

## **Isolation of nucleic acids from low biomass samples: detection and removal of sRNA contaminants**

Anna Heintz-Buschart<sup>1a\*</sup>, Dilmurat Yusuf<sup>1b</sup>, Anne Kaysen<sup>1c</sup>, Alton Etheridge<sup>2</sup>, Joëlle V. Fritz<sup>1c</sup>, Patrick May<sup>1</sup>, Carine de Beaufort<sup>1,3</sup>, Bimal B. Upadhyaya<sup>1</sup>, Anubrata Ghosal<sup>1d</sup>, David J. Galas<sup>2</sup>, and Paul Wilmes<sup>1\*</sup>

<sup>1</sup> Luxembourg Centre for Systems Biomedicine, University of Luxembourg, 4362 Esch-sur-Alzette, Luxembourg

<sup>2</sup> Pacific Northwest Research Institute, Seattle, WA, 98122, USA

<sup>3</sup> Centre Hospitalier de Luxembourg, 1210 Luxembourg, Luxembourg

\* To whom correspondence should be addressed. Paul Wilmes; Tel: 00352-466644-6188; Fax: 00352-466644-6949; Email: paul.wilmes@uni.lu; Anna Heintz-Buschart; Tel: 0049-345-558-5225; Email: anna.heintz-buschart@idiv.de

Present Addresses:

<sup>a</sup> Anna Heintz-Buschart, German Centre for Integrative Biodiversity Research (iDiv) Leipzig-Halle-Jena, 04103 Leipzig, Germany, and Department of Soil Ecology, Helmholtz-Centre for Environmental Research GmbH (UFZ), 06120 Halle (Saale), Germany

<sup>b</sup> Dilmurat Yusuf, Bioinformatics Group, Department of Computer Science, University of Freiburg, Freiburg, 79110, Germany

<sup>c</sup> Anne Kaysen and Joëlle V. Fritz, Centre Hospitalier de Luxembourg, 1210 Luxembourg, Luxembourg

<sup>d</sup> Anubrata Ghosal, Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139, United States

**Running title:** sRNA contaminants in microRNA extraction kits

**Keywords:** RNA sequencing, artefact removal, exogenous RNA in human blood plasma, contaminant RNA, spin columns

## ABSTRACT

Sequencing-based analyses of low-biomass samples are known to be prone to misinterpretation due to the potential presence of contaminating molecules derived from laboratory reagents and environments. Due to its inherent instability, contamination with RNA is usually considered to be unlikely. Here we report the presence of small RNA (sRNA) contaminants in widely used microRNA extraction kits and means for their depletion. Sequencing of sRNAs extracted from human plasma samples was performed and significant levels of non-human (exogenous) sequences were detected. The source of the most abundant of these sequences could be traced to the microRNA extraction columns by qPCR-based analysis of laboratory reagents. The presence of artefactual sequences originating from the confirmed contaminants were furthermore replicated in a range of published datasets. To avoid artefacts in future experiments, several protocols for the removal of the contaminants were elaborated, minimal amounts of starting material for artefact-free analyses were defined, and the reduction of contaminant levels for identification of *bona fide* sequences using ‘ultra-clean’ extraction kits was confirmed. In conclusion, this is the first report of the presence of RNA molecules as contaminants in laboratory reagents. The described protocols should be applied in the future to avoid confounding sRNA studies.

## INTRODUCTION

The characterization of different classes of small RNAs (sRNAs) in tissues and bodily fluids holds great promise in understanding human physiology as well as in health-related applications. In blood plasma, microRNAs and other sRNAs are relatively stable, and microRNAs in particular are thought to reflect a system-wide state, making them potential biomarkers for a multitude of human diseases (Mitchell *et al.* 2008). Different mechanisms of sRNA delivery as a means of long-distance intercellular communication have been recognized in several eukaryotes (Valadi *et al.* 2007; Zerneck *et al.* 2009; Pegtel *et al.* 2010; Molnar *et al.* 2010; Vickers *et al.* 2011; Arroyo *et al.* 2011). In addition, inter-individual, inter-species and even inter-kingdom communications via sRNAs have been proposed (Tomilov *et al.* 2008; Kosaka *et al.* 2010; Knip *et al.* 2014; Fritz *et al.* 2016; Koeppen *et al.* 2016), and some cases of microRNA-based control by the host (LaMonte *et al.* 2012; Liu *et al.* 2016) or pathogens (Weiberg *et al.* 2013; Buck *et al.* 2014) have been demonstrated.

As exogenous RNAs have been detected in the blood plasma of humans and mice (Wang *et al.* 2012; Zhang *et al.* 2012), the potential for exogenous RNA-based signalling in mammals is the subject of significant current debate (Zhang *et al.* 2011; Zhou *et al.* 2015). Diet-derived exogenous microRNAs have been proposed to exert an influence on human physiology (Liang *et al.* 2014; Baier *et al.* 2014), as have bacterial RNAs, which can be secreted in the protective environment of outer membrane vesicles (Ghosal *et al.* 2015; Celluzzi and Masotti 2016; Blenkiron *et al.* 2016). However, a heated discussion has at the same time been triggered around the genuineness of the observations of these exogenous sRNAs in human blood (Witwer 2015; Kang *et al.* 2017; Witwer and Zhang 2017) and the possibility of dietary uptake of sRNAs (Dickinson *et al.* 2013; Witwer and Hirschi 2014; Title *et al.* 2015). This discussion happens at a time where DNA sequencing-based analyses of low-biomass samples have been recognized to be prone to confounding by contaminants (Lusk 2014). From initial sample handling (Salzberg *et al.* 2016), to extraction kits (Naccache *et al.* 2013), to sequencing reagents (Salter *et al.* 2014), multiple sources of DNA contamination and artefactual sequencing data have been described.

Here, we report the contamination of widely used silica-based columns for the isolation of micro- and other small RNAs with RNA, which was apparent from sRNA sequencing data and was subsequently validated by qPCR. These artefactual sRNA sequences were also apparent in numerous published datasets. Furthermore, approaches for the depletion of the contaminants from the columns as well as an evaluation of a newer ultra-clean kit are presented, along with the determination of a minimum safe input volume to suppress the signal of the contaminant sequences in RNA sequencing data of human blood plasma samples. The potential presence of *bona fide* exogenous sRNA species in human plasma is examined. Finally, recommendations for the control and interpretation of sRNA sequencing data from low-biomass samples are provided.

## RESULTS

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

sRNA was extracted from 100  $\mu$ l blood plasma samples of ten healthy individuals and sequenced using regular RNeasy columns (workflow in **Figure 1**). The read profiles were mined for putative exogenous (non-human) sequences (Material and Methods). Among the potential exogenous sequences were 19 sequences that occurred with more than 1,000 counts per million (cpm) in all samples. To rule out sequencing errors or contamination during sequencing library preparation, a qPCR approach was developed to assess the presence of non-human sequences in the sRNA preparations from plasma. Six of the 19 highly abundant sRNA sequences from plasma that could not be mapped to the human genome were chosen for validation by qPCR (**Table 1**).

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

Synthetic sRNAs with the putative exogenous sequences found in plasma were poly-adenylated and reverse transcribed to yield cDNA, used for optimisation of PCR primers and conditions (**Table 1**). All primer sets yielded amplicons with single peaks in melting temperature analysis and efficiency values above 80 %. The optimised qPCR assays were then employed to test for the presence of the highly abundant sRNAs potentially representing exogenous sequences (workflow in **Figure 1**) in the

human plasma samples used for the initial sequencing experiment. The qPCR assays confirmed the presence of these sRNAs in the sRNA preparations used for sequencing (**Figure 2A**), yielding amplicons with melting temperatures expected from the synthetic sRNAs. To rule out contamination of the water used in the sRNA preparations, a water control was also examined. No amplification was observed in all but one assay, where amplification of a product with a different melting temperature occurred (**Figure 2A**). Thus, for the assays, contamination of the water could be ruled out.

#### *Non-human sequences derived from column contaminants*

To analyse whether the validated non-human sequences occurring in the sRNA extracts of plasma were present in any lab wear, a series of control experiments were carried out (**Supplementary Figure 1**). When nucleic acid- and RNase-free water (QIAGEN) was used as input to the miRNeasy Serum/Plasma kit (QIAGEN) instead of plasma (“mock-extraction”), all tested non-human sequences could be amplified from the mock-extract (**Figure 2B**). This indicates that one of the components of the extraction kit or lab-ware was contaminated with the non-human sequences. To locate the source of contamination, mock-extractions were performed by omitting single steps of the RNA-isolation protocol except for the elution step. Amplification from the resulting mock-extracts was tested for the most abundant non-human sequence (sRNA 1). In all cases, the sRNA 1 could be amplified (data not shown). We therefore carried out a simple experiment, in which nucleic acid- and RNase-free water was passed through an otherwise untreated spin column. From this column eluate, all target sequences could be amplified, in contrast to the nucleic acid- and RNase-free water (**Figure 2B**). The most abundant non-human sequences in the plasma sequencing experiments were therefore most likely contaminants originating from the untreated RNeasy columns.

#### *Detection of contaminant sequences in public datasets*

To assess whether our observation of contaminant sRNAs was also pertinent in other sequencing datasets of low-input samples, the levels of confirmed contaminant sRNA sequences in published datasets (Huang et al. 2013; Wang et al. 2012; Spornraft et al. 2014; Beatty et al. 2014; Dickinson et

al. 2013; Santa-Maria *et al.* 2015; Taft *et al.* 2010; Chen *et al.* 2011; Liu *et al.* 2012; Lebedeva *et al.* 2011; Kuchen *et al.* 2010; Wei *et al.* 2009; Mayr and Bartel 2009; Su *et al.* 2010; Chen *et al.* 2010; Legeai *et al.* 2010; Vaz *et al.* 2010; Liu *et al.* 2010; Lian *et al.* 2012; Nolte-'t Hoen *et al.* 2012; Zhang *et al.* 2012) were assessed. Irrespective of the RNA isolation procedure applied, non-target sequences were detected (making up between 5 and over 99 % of the sequencing libraries for the human samples; **Supplementary Table 1**). As shown in **Figure 3**, the six contaminant sequences which had been confirmed by qPCR were found in all analysed samples of low biomass samples which were extracted with regular miRNeasy kits, but the sequences were found at lower levels in studies with more biomass input (Spornraft *et al.* 2014; Dickinson *et al.* 2013; Santa-Maria *et al.* 2015) and hardly ever (Taft *et al.* 2010) in studies where samples were extracted using other methods (**Supplementary Table 1**). Within each study where the confirmed contaminant sequences were detected, the relative levels of the contaminant sequences were remarkably stable (**Supplementary Figure 2**).

#### *Depletion of contaminants from isolation columns*

In order to eliminate contamination from the columns to allow their use in studies of environmental samples or potential exogenous sRNAs from human samples, we were interested in the nature of these contaminants. The fact that they can be poly-adenylated by RNA-poly-A-polymerase points to them being RNA. Treatment of the eluate with RNase prior to cDNA preparation also abolished amplification (data not shown), but on-column DNase digest did not reduce their levels (**Figure 2C**). These findings suggest that the contaminants were RNAs.

Contaminating sequences could potentially be removed from the RNeasy columns using RNase, but as RNases are notoriously difficult to inactivate and RNases remaining on the column would be detrimental to sRNA recovery, an alternative means of removing RNA was deemed desirable. Loading and incubation of RNeasy columns with the oxidant sodium hypochlorite and subsequent washing with RNase-free water to remove traces of the oxidant reduced amplifiability of unwanted sRNA by at least 100 times (**Figure 2D**), while retaining the columns' efficiency to isolate sRNAs from samples applied afterwards. Elimination of contaminant sRNAs from the RNeasy columns by

washing with RNase-free water (**Figure 2D**; average +/- standard deviation of the contaminant reduction by 80 +/- 10 %) or treatment with sodium hydroxide (average +/- standard deviation of the contaminant reduction by 70 +/- 15 %) was not sufficient to remove the contaminants completely.

#### *Ultra-clean extraction kits*

Recently, RNeasy columns from an ultra-clean production have become available from QIAGEN within the miRNeasy Serum/Plasma Advanced Kit. We compared the levels of the previously analysed contaminant sequences in the flow-through of mock-extractions using 4 batches of ultra-clean RNeasy columns to 2 batches of the regular columns by qPCR. In all cases, marked reductions in the contaminant levels were observed in the clean columns (**Figure 4A**; 4 to 4,000 fold; median 60). To obtain an overview over potential other contaminants, sRNA sequencing of the mock-extracts from these six batches of spin columns was performed. With regards to the six previously analysed contaminant sequences, the results were similar to those of the qPCR assays (**Supplementary Figure 3**). Additionally, for the ultra-clean RNeasy columns, a smaller spectrum of other potential contaminant sequences was observed (**Figure 4B&C**) and those sequences made up a smaller proportion of the eluate sequences (**Figure 4D**).

As our initial analyses of plasma samples extracted using regular RNeasy spin columns had revealed contaminant levels of up to 7000 cpm, we were interested to define a safe input amount for human plasma for both column types that would be sufficient to suppress the contaminant signals to below 100 cpm. For this, we performed a titration experiment (**Supplementary Figure 3B**), isolating sRNA from a series of different input volumes of the same human plasma sample on four batches of RNeasy columns (2 batches of regular columns, 2 batches of ultra-clean columns) with subsequent sequencing. As expected from reagent contaminants, the observed levels of the contaminant sequences were generally inversely dependent on the plasma input volume (**Figure 5A**). In addition and in accordance with the earlier mock-extraction results, the levels of contaminant sequences were lower or they were completely absent in the ultra-clean columns (see levels for 100  $\mu$ l input in **Figure 5B**). An input

volume of 100  $\mu$ l plasma was sufficient to reduce all contaminant sequences to below 100 cpm when using the ultra-clean spin columns.

#### *Potential plasma-derived exogenous RNAs*

Finally, to detect potential exogenous sRNAs, we mined the plasma datasets used in the well-controlled titration experiment for sequences that do not originate from the human genome and were not detected in any of the mock-extracts. On average, 5 % of the sequencing reads of sRNA isolated from plasma did not map to the human genome. 127 sequences which did not map to the human genome assembly hg38 were detected in the majority of the plasma samples and were not represented in the control samples (empty libraries, column eluates or water). Out of these, 3 sequences had low complexity and 81 could be matched to sequences in the NCBI-nr that are not part of the current version of the human genome assembly (hg38) but annotated as human sequences or to sequences from other vertebrates. Of the 43 remaining sequences which matched to bacterial, fungal or plant sequences, 22 matched best to genera which have previously been identified as a source of contaminations of sequencing kits (Salter et al. 2014). The remaining 21 sequences displayed very low (up to 47 cpm), yet consistent relative abundances in the 28 replicates of a plasma sample from the one healthy individual. Their potential origins were heterogeneous, including fungi and bacteria, with a notable enrichment in *Lactobacillus* sequences (**Supplementary Table 2**).

## **DISCUSSION**

Several instances of contamination of laboratory reagents with DNA, which can confound the analysis of sequencing data, have been reported in recent years (Salter et al. 2014; Lusk 2014; Lauder et al. 2016; Glassing et al. 2016). In contrast, the contamination of reagents with RNA has not yet been reported. Contamination with RNA is usually considered very unlikely, due to the ubiquitous presence of RNases in the environment and RNA's lower chemical stability due to being prone to hydrolysis, especially at higher pH. However, our results suggest that the detected contaminants were not DNA, but RNA, because treatment with RNase and not DNase could decrease the contaminant load. In



addition, the contaminating molecules could not be amplified without poly-adenylation and reverse-transcription. The stability of the contaminants is likely due to the extraction columns being RNase-free and their silica protecting loaded sRNAs from degradation. While the results presented here focused on one manufacturer's spin column-based extraction kit, for which contaminants were validated, other RNA-stabilizing or extraction reagents may carry RNA contaminations. This is suggested by previously observed significant batch effects of sequencing data derived from samples extracted with a number of different extraction kits (Kang et al. 2017). Based on the analysis of the published data sets, where significant numbers of sequences that did not map to the source organism's genome were found independent of the RNA extraction kit used, the potential contaminants in other extraction kit would have different sequences than the ones confirmed by qPCR here.

The reported contaminant sequences can confound studies of organisms whose transcriptomes contain sequences similar to the contaminants. They can also give rise to misinterpretation in studies without *a priori* knowledge of the present organisms as well as lead to the overestimation of miRNA yields in low-biomass samples. Therefore, based on the present study, care has to be taken when analysing low-input samples, in particular for surveys of environmental or otherwise undefined sources of RNAs. A number of recommendations can be conceived based on the presented data (**Figure 6**): Extraction columns should be obtained as clean as possible. Simple clean-up procedures can also reduce contaminants. The input mass of sRNA should be as high as possible, e.g. for human plasma volumes above 100  $\mu$ l are preferable. Extraction controls should always be sequenced with the study samples. To facilitate library preparation for the extraction controls, spike-in RNAs with defined sequences can be used. They should be applied at concentrations similar to the levels of RNA found in the study samples. As the spike-in signal can drown out the contaminants, it is necessary to avoid too high concentrations for the spike-ins. Sequences found in the extraction controls should be treated as artefacts and removed from the sequencing data. Independent techniques that are more robust to low input material, such as qPCR or ddPCR, should be applied to both study samples and controls in case of doubt.

The results presented here should also help to assess the question whether exogenous sRNA species derived from oral intake (Zhang et al. 2012) or the human microbiome (Wang et al. 2012; Beatty et al. 2014; Yeri et al. 2017) really occur frequently in human plasma or are merely artefacts (Witwer 2015). While the limited data from this study (one healthy person) points to very low levels and a small spectrum of potential foreign sRNAs, properly controlled studies using laboratory materials without contaminants on individuals or animals with conditions that limit gastrointestinal barrier function will shed more light on this important research question in the future.

## **MATERIAL AND METHODS**

### *Blood plasma sampling*

Written informed consent was obtained from all blood donors. The sample collection and analysis was approved by the Comité d'Ethique de Recherche (CNER; Reference: 201110/05) and the National Commission for Data Protection in Luxembourg. Blood was collected by venepuncture into EDTA-treated tubes. Plasma was prepared immediately after blood collection by centrifugation (10 min at 1,000 x g) and platelets were depleted by a second centrifugation step (5 min at 10,000 x g). The blood plasma was flash-frozen in liquid nitrogen and stored at -80 °C until extraction.

### *Use of sRNA isolation columns*

Unless stated otherwise, 100 µl blood plasma was lysed using the QIAzol (QIAGEN) lysis reagent prior to binding to the column, as recommended by the manufacturer. RNeasy MinElute spin columns from the miRNeasy Serum/Plasma Kit (QIAGEN) were then loaded, washed and dried, and RNA was eluted as recommended by the manufacturer's manual. We further tested four batches of ultra-clean RNeasy MinElute columns, which underwent an ultra-clean production process (UCP) to remove potential nucleic acid contaminations, including environmental sRNAs. These columns were treated as recommended in the manual of the miRNeasy Serum/Plasma Advanced Kit (QIAGEN). All eluates were stored at -80 °C until analysis.

For the mock-extractions, ultra-clean or regular RNeasy columns were loaded with the aqueous phase from a QIAzol extraction of nucleic acid- and RNase-free water (QIAGEN) instead of plasma. For mock-extractions with a defined spike-in, the aqueous phase was spiked with synthetic *hsa*-miR-486-3p RNA (Eurogentec) to yield 40,000 copies per  $\mu$ l eluate. To obtain column eluates, spin columns were not loaded, washed or dried. Instead, 14  $\mu$ l of RNase-free water (QIAGEN) was applied directly to a new column and centrifuged for 1 min.

To eliminate environmental sRNAs from the regular RNeasy columns, the columns were incubated with 500  $\mu$ l of a sodium hypochlorite solution (Sigma; diluted in nuclease free water (Invitrogen) to approx. 0.5 %) for 10 min at room temperature. Columns were subsequently washed 10 times with 500  $\mu$ l nuclease free water (Invitrogen), before use. Similarly, in the attempt to remove sRNAs by application of sodium hydroxide, 500  $\mu$ l 50 mM NaOH were incubated on the spin columns for 5 min, followed by incubation with 50 mM HCl for 5 min, prior to washing the columns 10 times with 500  $\mu$ l nuclease-free water (Invitrogen) before use.

#### *Real-time PCR*

5  $\mu$ l of eluted RNA was polyadenylated and reverse-transcribed to cDNA using the qScript microRNA cDNA Synthesis Kit (Quanta BIOSCIENCES). 1  $\mu$ l of cDNA (except for the initial plasma experiment, where 0.2  $\mu$ l cDNA were used) was amplified by use of sequence-specific forward primers (see **Table 1**, obtained from Eurogentec) or the miR486-5p specific assay from Quanta BIOSCIENCES, PerfeCTa Universal PCR Primer and PerfeCTa SYBR Green SuperMix (Quanta BIOSCIENCES) in a total reaction volume of 10  $\mu$ l. Primers were added at a final concentration of 0.2  $\mu$ M. Primer design and amplification settings were optimised with respect to reaction efficiency and specificity. Efficiency was calculated using a dilution series covering seven orders of magnitude of template cDNA reverse transcribed from synthetic sRNA. Real-time PCR was performed on a LightCycler® 480 Real-Time PCR System (Roche) including denaturation at 95 °C for 2 min and 40 cycles of 95 °C for 5 sec, 54-60 °C for 15 sec (for annealing temperatures see **Table 1**), and 72 °C for 15 sec. All reactions were carried out in duplicates. No-template-controls were

performed analogously with water as input. Cp values were obtained using the second derivative procedure provided by the LightCycler® 480 Software, Version 1.5. Cp data were analysed using the comparative  $C_T$  method ( $\Delta\Delta C_T$ ).

#### *sRNA seq: library preparation and sequencing*

sRNA libraries were made using the TruSeq small RNA library preparation kit (Illumina) according to the manufacturer's instructions, except that the 3' and 5' adapters were diluted 1:3 before use. PCR-amplified libraries were size selected using a PippinHT instrument (Sage Science), collecting the range of 121-163 bp. Completed, size-selected libraries were run on a High Sensitivity DNA chip on a 2100 Bioanalyzer (Agilent) to assess library quality. Concentration was determined by qPCR using the NEBNext Library Quant kit (NEB). Libraries were pooled, diluted and sequenced with 75 cycle single-end reads on a NextSeq 500 (Illumina) according the manufacturer's instructions. The sequencing reads can be accessed at NCBI's short read archive via PRJNA419919 (for sample identifiers and accessions see **Supplementary Table 1**).

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

For the initial analysis of plasma-derived sRNA sequencing data, FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc>) was used to determine over-represented primer and adapter sequences, which were subsequently removed using cutadapt (<http://dx.doi.org/10.14806/ej.17.1.200>). This step was repeated recursively until no over-represented primer or adapter sequences were detected. 5'-Ns were removed using fastx\_clipper of the FASTX-toolkit. Trimmed reads were quality-filtered using fastq\_quality\_filter of the FASTX-toolkit (with -q 30 -p 90; [http://hannonlab.cshl.edu/fastx\\_toolkit](http://hannonlab.cshl.edu/fastx_toolkit)). Finally, identical reads were collapsed, retaining the read abundance information using fastx\_collapser of the FASTX-toolkit. The collapsed reads were mapped against the human genome (GRCh37), including RefSeq exon junction sequences, as well as prokaryotic, viral, fungal, plant and animal genomes from Genbank (Benson et al. 2012) and the Human Microbiome Project (The NIH HMP Working Group et al. 2009) using Novoalign V2.08.02

(<http://www.novocraft.com>; **Supplementary Table 3**). These organisms were selected based on the presence in the human microbiome, human nutrition and the public availability of the genomes. As reads were commonly mapping to genomic sequences of multiple organisms, and random alignment can easily occur between short sequences and reference genomes, the following approach was taken to refine their taxonomic classification: First, reads were attributed to the human genome if they mapped to it. Secondly, reads mapping to each reference genome was compared to mapping of a shuffled decoy read set. Based on this, the list of reference genomes was limited to the genomes recruiting at least one read with a minimum length of 25 nt. Loci on non-human genomes were established by the position of the mapping reads. The number of mapping reads per locus was adjusted using a previously established cross-mapping correction (de Hoon et al. 2010). Finally, the sequences of the loci, the number of mapping reads and their potential taxonomy were extracted.

#### *sRNA sequence analysis of controls*

For the subsequent analysis of the mock-extractions, column eluates and nucleic acid- and RNase-free water, and no-template controls as well as human plasma samples, extracted using either regular or ultra-clean RNeasy columns, the trimming and quality check of the reads was done analogously to the description above. Collapsed reads were mapped against the most recent version of the human genome (hg38) either to remove operator-derived sequences or to distinguish the reads mapping to the human genome in the different datasets. Sequencing was performed in two batches, with one batch filling an entire flow cell, and one mixed with other samples. The latter batch of samples was sequenced on the same flow cell as sRNAs extracted from *Salmonella typhimurium* LT2. To avoid misinterpretations due to multiplexing errors, reads mapping to *Salmonella typhimurium* LT2 (McClelland et al. 2001) (Genbank accession AE006468) were additionally removed in this batch. To limit the analysis to only frequently occurring sequences and therefore avoid over-interpretation of erroneous sequences, only read sequences that were found at least 30 times in all analysed samples together were retained for further analysis. Public sRNA datasets of low-input samples (see **Supplementary Table 1**) were analysed in a fashion analogous to the study's control and plasma

samples. As the published studies consisted of different numbers of samples, no overall threshold was imposed, but to limit the analysis to frequently occurring sequences, singleton reads were removed.

To compare the sequencing results to the qPCR-based results and to detect the same sequences in public datasets, reads matching the sequences assayed by qPCR were determined by clustering the trimmed, filtered and collapsed sRNA reads with 100 % sequence identity and 14 nt alignment length with the primer sequences, while allowing the sRNA reads to be longer than the primer sequences, using CD-HIT-EST-2D (parameters -c 1 -n 8 -G 0 -A 14 -S2 40 -g 1 -r 0) (Fu et al. 2012).

To compare the diversity and levels of putative contaminant sequences in the different samples, identical reads derived from all study samples (that did not map to the human genome) were clustered using CD-HIT-EST (Fu et al. 2012), and a table with the number of reads sequenced for each sample per sequence was created using R v.3.0.2. This table was also used to extract candidate sequences from the study plasma samples that are likely exogenous plasma sRNAs, based on the following criteria: for a sequence to be considered a potential exogenous plasma sRNA, it had to be non-identical to any of the sequences assigned to the confirmed contaminant sequences (**Table 1**), and it had to be absent from at least 90 % of the controls (no-library controls, water and spike-in controls, eluates and mock-extracts) and never detected in any of these controls with at least 10 copy numbers, and it had to be detected by more than 3 reads in more than 7 of the 28 libraries generated from the plasma titration experiment. These thresholds were chosen in order to make the analysis robust against multiplexing errors (e.g. which would result in false-negative identifications if a sequence that is very dominant in a plasma sample is falsely assigned to the control-samples), while at the same time making it sensitive to low-abundant sequences (which would not be detected in every library). To confirm the non-human origin and find potential microbial taxa of origin for these sequences, they were subsequently searched within the NCBI nr database using megablast and blastn web tools, with parameters auto-set for short inputs (Altschul et al. 1990; Morgulis et al. 2008; 2016). All sequences with best hits to human sequences or other vertebrates were removed, because they were potentially human. The remaining sequences were matched against a set of genera previously reported (Salter et

al. 2014) to be common sequencing kit contaminants. Sequences with better hits to non-contaminant taxa than contaminant taxa were kept as potential exogenous sequences.

#### *Data and code availability*

Scripts for the analysis of the data from sRNA sequencing of column eluates and the plasma titration experiment is available at <https://git.ufz.de/metaOmics/contaminomics>.

RNA sequencing data has been deposited at NCBI Bioproject PRJNA419919.

#### *Supplementary Information*

The following supplementary information are available online: **Supplementary Figures 1-3**; **Supplementary Table 1**: list of the generated datasets and analysed published datasets; **Supplementary Table 2**: potential exogenous sRNA sequences detected in human plasma after removal of contaminants; **Supplementary Table 3**: list of the species whose reference genomes and cDNA collections were used in the initial analysis.

#### **ACKNOWLEDGEMENT**

*In silico* analyses presented in this paper were carried out using the HPC facilities of the University of Luxembourg (Varrette et al. 2014) whose administrators are acknowledged for excellent support.

#### **FUNDING**

This work was supported by the Luxembourg National Research Fund (FNR) through an ATTRACT programme grant (ATTRACT/A09/03), CORE programme grant (CORE/15/BM/10404093) and Proof-of-Concept Programme Grant (PoC/13/02) to P.W., an Aide à la Formation Recherche grant (Ref. no. 1180851) to D.Y., an Aide à la Formation Recherche grant (Ref. no. 5821107) and a CORE grant (CORE14/BM/8066232) to J.V.F., a National Institutes of Health Extracellular RNA

Communication Consortium award (1U01HL126496) to D.J.G., and by the University of Luxembourg (ImMicroDyn1).

## CONFLICT OF INTEREST

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

## REFERENCES

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* **215**: 403–410.
- Arroyo JD, Chevillet JR, Kroh EM, Ruf IK, Pritchard CC, Gibson DF, Mitchell PS, Bennett CF, Pogossova-Agadjanyan EL, Stirewalt DL, et al. 2011. Argonaute2 complexes carry a population of circulating microRNAs independent of vesicles in human plasma. *Proc Natl Acad Sci USA* **108**: 5003–5008.
- Baier SR, Nguyen C, Xie F, Wood JR, Zempleni J. 2014. MicroRNAs Are Absorbed in Biologically Meaningful Amounts from Nutritionally Relevant Doses of Cow Milk and Affect Gene Expression in Peripheral Blood Mononuclear Cells, HEK-293 Kidney Cell Cultures, and Mouse Livers. *Journal of Nutrition* **144**: 1495–1500.
- Beatty M, Guduric-Fuchs J, Brown E, Bridgett S, Chakravarthy U, Hogg RE, Simpson DA. 2014. Small RNAs from plants, bacteria and fungi within the order Hypocreales are ubiquitous in human plasma. *BMC Genomics* **15**: 933.
- Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. 2012. GenBank. *Nucleic Acids Res* **41**: D36–D42.
- Blenkiron C, Simonov D, Muthukaruppan A, Tsai P, Dauros P, Green S, Hong J, Print CG, Swift S,



- Phillips AR. 2016. Uropathogenic *Escherichia coli* Releases Extracellular Vesicles That Are Associated with RNA ed. E. Cascales. *PLoS ONE* **11**: e0160440–16.
- Buck AH, Coakley G, Simbari F, McSorley HJ, Quintana JF, Le Bihan T, Kumar S, Abreu-Goodger C, Lear M, Harcus Y, et al. 2014. Exosomes secreted by nematode parasites transfer small RNAs to mammalian cells and modulate innate immunity. *Nature Communications* **5**: 5488.
- Celluzzi A, Masotti A. 2016. How Our Other Genome Controls Our Epi-Genome. *Trends Microbiol* **24**: 777–787.
- Chen C, Ai H, Ren J, Li W, Li P, Qiao R, Ouyang J, Yang M, Ma J, Huang L. 2011. A global view of porcine transcriptome in three tissues from a full-sib pair with extreme phenotypes in growth and fat deposition by paired-end RNA sequencing. *BMC Genomics* **12**: 448.
- Chen X, Yu X, Cai Y, Zheng H, Yu D, Liu G, Zhou Q, Hu S, Hu F. 2010. Next-generation small RNA sequencing for microRNAs profiling in the honey bee *Apis mellifera*. *Insect Mol Biol* **19**: 799–805.
- de Hoon MJL, Taft RJ, Hashimoto T, Kanamori-Katayama M, Kawaji H, Kawano M, Kishima M, Lassmann T, Faulkner GJ, Mattick JS, et al. 2010. Cross-mapping and the identification of editing sites in mature microRNAs in high-throughput sequencing libraries. *Genome Research* **20**: 257–264.
- Dickinson B, Zhang Y, Petrick JS, Heck G, Ivashuta S, Marshall WS. 2013. Lack of detectable oral bioavailability of plant microRNAs after feeding in mice. *Nat Biotechnol* **31**: 965–967.
- Fritz JV, Heintz-Buschart A, Ghosal A, Wampach L, Etheridge A, Galas D, Wilmes P. 2016. Sources and Functions of Extracellular Small RNAs in Human Circulation. *Annu Rev Nutr* **36**: 301–336.
- Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**: 3150–3152.

Ghosal A, Upadhyaya BB, Fritz JV, Heintz-Buschart A, Desai MS, Yusuf D, Huang D, Baumuratov

A, Wang K, Galas D, et al. 2015. The extracellular RNA complement of *Escherichia coli*.

*MicrobiologyOpen* **4**: 252–266.

Glassing A, Dowd SE, Galandiuk S, Davis B, Chiodini RJ. 2016. Inherent bacterial DNA

contamination of extraction and sequencing reagents may affect interpretation of microbiota in

low bacterial biomass samples. *Gut Pathog* **8**: 24.

Huang X, Yuan T, Tschannen M, Sun Z, Jacob H, Du M, Liang M, Dittmar RL, Liu Y, Liang M, et

al. 2013. Characterization of human plasma-derived exosomal RNAs by deep sequencing. *BMC*

*Genomics* **14**: 319.

Kang W, Bang-Berthelsen CH, Holm A, Houben AJS, Müller AH, Thymann T, Pociot F, Estivill X,

Friedländer MR. 2017. Survey of 800+ data sets from human tissue and body fluid reveals

xenomiRs are likely artifacts. *RNA* **23**: 433–445.

Knip M, Constantin ME, Thordal-Christensen H. 2014. Trans-kingdom cross-talk: small RNAs on the

move. *PLoS Genet* **10**: e1004602.

Koeppen K, Hampton TH, Jarek M, Scharfe M, Gerber SA, Mielcarz DW, Demers EG, Dolben EL,

Hammond JH, Hogan DA, et al. 2016. A Novel Mechanism of Host-Pathogen Interaction through

sRNA in Bacterial Outer Membrane Vesicles. *PLoS Pathog* **12**: e1005672.

Kosaka N, Izumi H, Sekine K, Ochiya T. 2010. microRNA as a new immune-regulatory agent in

breast milk. *Silence* **1**: 7.

Kuchen S, Resch W, Yamane A, Kuo N, Li Z, Chakraborty T, Wei L, Laurence A, Yasuda T, Peng S,

et al. 2010. Regulation of MicroRNA Expression and Abundance during Lymphopoiesis.

*Immunity* **32**: 828–839.

LaMonte G, Philip N, Reardon J, Lacsina JR, Majoros W, Chapman L, Thornburg CD, Telen MJ,

Ohler U, Nicchitta CV, et al. 2012. Translocation of sickle cell erythrocyte microRNAs into

Plasmodium falciparum inhibits parasite translation and contributes to malaria resistance. *Cell Host Microbe* **12**: 187–199.

Lauder AP, Roche AM, Sherrill-Mix S, Bailey A, Laughlin AL, Bittinger K, Leite R, Elovitz MA, Parry S, Bushman FD. 2016. Comparison of placenta samples with contamination controls does not provide evidence for a distinct placenta microbiota. *Microbiome* **4**: 29.

Lebedeva S, Jens M, Theil K, Schwanhäusser B, Selbach M, Landthaler M, Rajewsky N. 2011. Transcriptome-wide Analysis of Regulatory Interactions of the RNA-Binding Protein HuR. *Mol Cell* **43**: 340–352.

Legeai F, Rizk G, Walsh T, Edwards O, Gordon K, Lavenier D, Leterme N, Méreau A, Nicolas J, Tagu D, et al. 2010. Bioinformatic prediction, deep sequencing of microRNAs and expression analysis during phenotypic plasticity in the pea aphid, *Acyrtosiphon pisum*. *BMC Genomics* **11**: 281.

Lian L, Qu L, Chen Y, Lamont SJ, Yang N. 2012. A Systematic Analysis of miRNA Transcriptome in Marek's Disease Virus-Induced Lymphoma Reveals Novel and Differentially Expressed miRNAs ed. M. Watson. *PLoS ONE* **7**: e51003–13.

Liang G, Zhu Y, Sun B, Shao Y, Jing A, Wang J, Xiao Z. 2014. Assessing the survival of exogenous plant microRNA in mice. *Food Sci Nutr* **2**: 380–388.

Liu J-L, Liang X-H, Su R-W, Lei W, Jia B, Feng X-H, Li Z-X, Yang Z-M. 2012. Combined Analysis of MicroRNome and 3'-UTRome Reveals a Species-specific Regulation of Progesterone Receptor Expression in the Endometrium of Rhesus Monkey. *J Biol Chem* **287**: 13899–13910.

Liu S, da Cunha AP, Rezende RM, Cialic R, Wei Z, Bry L, Comstock LE, Gandhi R, Weiner HL. 2016. The Host Shapes the Gut Microbiota via Fecal MicroRNA. *Cell Host Microbe* **19**: 32–43.

Liu S, Li D, Li Q, Zhao P, Xiang Z, Xia Q. 2010. MicroRNAs of *Bombyx mori* identified by Solexa sequencing. *BMC Genomics* **11**: 148.

- Lusk RW. 2014. Diverse and widespread contamination evident in the unmapped depths of high throughput sequencing data. *PLoS ONE* **9**: e110808.
- Mayr C, Bartel DP. 2009. Widespread Shortening of 3'UTRs by Alternative Cleavage and Polyadenylation Activates Oncogenes in Cancer Cells. *Cell* **138**: 673–684.
- McClelland M, Sanderson KE, Spieth J, Clifton SW, Latreille P, Courtney L, Porwollik S, Ali J, Dante M, Du F, et al. 2001. Complete genome sequence of *Salmonella enterica* serovar Typhimurium LT2. *Nature* **413**: 852–856.
- Mitchell PS, Parkin RK, Kroh EM, Fritz BR, Wyman SK, Pogosova-Agadjanian EL, Peterson A, Noteboom J, O'Briant KC, Allen A, et al. 2008. Circulating microRNAs as stable blood-based markers for cancer detection. *Proc Natl Acad Sci USA* **105**: 10513–10518.
- Molnar A, Melnyk CW, Bassett A, Hardcastle TJ, Dunn R, Baulcombe DC. 2010. Small silencing RNAs in plants are mobile and direct epigenetic modification in recipient cells. *Science* **328**: 872–875.
- Morgulis A, Coulouris G, Raytselis Y, Madden TL, Agarwala R, Schäffer AA. 2008. Database indexing for production MegaBLAST searches. *Bioinformatics* **24**: 1757–1764.
- Naccache SN, Greninger AL, Lee D, Coffey LL, Phan T, Rein-Weston A, Aronsohn A, Hackett J, Delwart EL, Chiu CY. 2013. The Perils of Pathogen Discovery: Origin of a Novel Parvovirus-Like Hybrid Genome Traced to Nucleic Acid Extraction Spin Columns. *J Virol* **87**: 11966–11977.
- Nolte-'t Hoen ENM, Buermans HPJ, Waasdorp M, Stoorvogel W, Wauben MHM, 't Hoen PAC. 2012. Deep sequencing of RNA from immune cell-derived vesicles uncovers the selective incorporation of small non-coding RNA biotypes with potential regulatory functions. *Nucleic Acids Research* **40**: 9272–9285.
- Pegtel DM, Cosmopoulos K, Thorley-Lawson DA, van Eijndhoven MAJ, Hopmans ES, Lindenberg

- JL, de Grujil TD, Wurdinger T, Middeldorp JM. 2010. Functional delivery of viral miRNAs via exosomes. *Proc Natl Acad Sci USA* **107**: 6328–6333.
- Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, Moffatt MF, Turner P, Parkhill J, Loman NJ, Walker AW. 2014. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol* **12**: 87.
- Salzberg SL, Breitwieser FP, Kumar A, Hao H, Burger P, Rodriguez FJ, Lim M, Quiñones-Hinojosa A, Gallia GL, Tornheim JA, et al. 2016. Next-generation sequencing in neuropathologic diagnosis of infections of the nervous system. *Neurol Neuroimmunol Neuroinflamm* **3**: e251.
- Santa-Maria I, Alaniz ME, Renwick N, Cela C, Fulga TA, Van Vactor D, Tuschl T, Clark LN, Shelanski ML, McCabe BD, et al. 2015. Dysregulation of microRNA-219 promotes neurodegeneration through post-transcriptional regulation of tau. *J Clin Invest* **125**: 681–686.
- Spornraft M, Kirchner B, Haase B, Benes V, Pfaffl MW, Riedmaier I. 2014. Optimization of Extraction of Circulating RNAs from Plasma – Enabling Small RNA Sequencing ed. C. Antoniewski. *PLoS ONE* **9**: e107259.
- Su R-W, Lei W, Liu J-L, Zhang Z-R, Jia B, Feng X-H, Ren G, Hu S-J, Yang Z-M. 2010. The Integrative Analysis of microRNA and mRNA Expression in Mouse Uterus under Delayed Implantation and Activation ed. H. Wang. *PLoS ONE* **5**: e15513–18.
- Taft RJ, Simons C, Nahkuri S, Oey H, Korbie DJ, Mercer TR, Holst J, Ritchie W, Wong JJ-L, Rasko JE, et al. 2010. Nuclear-localized tiny RNAs are associated with transcription initiation and splice sites in metazoans. *Nat Struct Mol Biol* **17**: 1030–1034.
- The NIH HMP Working Group, Peterson J, Garges S, Giovanni M, McInnes P, Wang L, Schloss JA, Bonazzi V, McEwen JE, Wetterstrand KA, et al. 2009. The NIH Human Microbiome Project. *Genome Research* **19**: 2317–2323.
- Title AC, Denzler R, Stoffel M. 2015. Uptake and Function Studies of Maternal Milk-derived

- MicroRNAs. *Journal of Biological Chemistry* **290**: 23680–23691.
- Tomilov AA, Tomilova NB, Wroblewski T, Michelmore R, Yoder JI. 2008. Trans-specific gene silencing between host and parasitic plants. *Plant J* **56**: 389–397.
- Valadi H, Ekström K, Bossios A, Sjöstrand M, Lee JJ, Lötvall JO. 2007. Exosome-mediated transfer of mRNAs and microRNAs is a novel mechanism of genetic exchange between cells. *Nat Cell Biol* **9**: 654–659.
- Varrette S, Bouvry P, Cartiaux H, Georgatos F. 2014. Management of an academic HPC cluster: The UL experience. 959–967.
- Vaz C, Ahmad HM, Sharma P, Gupta R, Kumar L, Kulshreshtha R, Bhattacharya A. 2010. Analysis of microRNA transcriptome by deep sequencing of small RNA libraries of peripheral blood. *BMC Genomics* **11**: 288.
- Vickers KC, Palmisano BT, Shoucri BM, Shamburek RD, Remaley AT. 2011. MicroRNAs are transported in plasma and delivered to recipient cells by high-density lipoproteins. *Nat Cell Biol* **13**: 423–433.
- Wang K, Li H, Yuan Y, Etheridge A, Zhou Y, Huang D, Wilmes P, Galas D. 2012. The complex exogenous RNA spectra in human plasma: an interface with human gut biota? eds. K. Wang, H. Li, Y. Yuan, A. Etheridge, Y. Zhou, D. Huang, P. Wilmes, and D. Galas. *PLoS ONE* **7**: e51009.
- Wei Y, Chen S, Yang P, Ma Z, Kang L. 2009. Characterization and comparative profiling of the small RNA transcriptomes in two phases of locust. *Genome Biology* **10**: R6.
- Weiberg A, Wang M, Lin F-M, Zhao H, Zhang Z, Kaloshian I, Huang H-D, Jin H. 2013. Fungal small RNAs suppress plant immunity by hijacking host RNA interference pathways. *Science* **342**: 118–123.
- Witwer KW. 2015. Contamination or artifacts may explain reports of plant miRNAs in humans. *The*

*Journal of Nutritional Biochemistry* **26**: 1685.

Witwer KW, Hirschi KD. 2014. Transfer and functional consequences of dietary microRNAs in vertebrates: Concepts in search of corroboration. *Bioessays* **36**: 394–406.

Witwer KW, Zhang C-Y. 2017. Diet-derived microRNAs: unicorn or silver bullet? *Genes & Nutrition* **12**: 15.

Yeri A, Courtright A, Reiman R, Carlson E, Beecroft T, Janss A, Siniard A, Richholt R, Balak C, Rozowsky J, et al. 2017. Total Extracellular Small RNA Profiles from Plasma, Saliva, and Urine of Healthy Subjects. *Sci Rep* **7**: 44061.

Zernecke A, Bidzhekov K, Noels H, Shagdarsuren E, Gan L, Denecke B, Hristov M, Koppel T, Jahantigh MN, Lutgens E, et al. 2009. Delivery of microRNA-126 by apoptotic bodies induces CXCL12-dependent vascular protection. *Sci Signal* **2**: ra81.

Zhang L, Hou D, Chen X, Li D, Zhu L, Zhang Y, Li J, Bian Z, Liang X, Cai X, et al. 2011. Exogenous plant MIR168a specifically targets mammalian LDLRAP1: evidence of cross-kingdom regulation by microRNA. *Nature Publishing Group* **22**: 107–126.

Zhang Y, Wiggins BE, Lawrence C, Petrick J, Ivashuta S, Heck G. 2012. Analysis of plant-derived miRNAs in animal small RNA datasets. *BMC Genomics* **13**: 381.

Zhou Z, Li X, Liu J, Dong L, Chen Q, Liu J, Kong H, Zhang Q, Qi X, Hou D, et al. 2015. Honeysuckle-encoded atypical microRNA2911 directly targets influenza A viruses. *Cell Research* **25**: 39–49.

2016. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res* **44**: D7–D19.

**Table 1.** Sequences of non-human sRNAs found in plasma preparations, synthetic sRNA templates, 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)AACAGACCGAGGACU UGAA(U)	133,700	algae	AACAGACCGAGGACTTGAA	57 °C
sRNA 2	ACGGACAAGAAUAGGCUU CGGCU	8,000	fungi or plants	ACGGACAAGAATAGGCTTC	54 °C
sRNA 3	GCCUUGGUUGUAGGAUCU GU	8,200	plants	GCCTTGTTGTAGGATCTGT	57 °C
sRNA 4	GCCAGCAUCAGUUCGGUG UG	6,800	bacteria	CAGCATCAGTTCGGTGTG	57 °C
sRNA 5	GAGAGUAGGACGUUGCCA GGUU	3,900	bacteria	AGTAGGACGTTGCCAGGTT	57 °C
sRNA 6	UUGAAGGGUCGUUCGAGA CCAGGACGUUGAUAGGCU GGGUG	3,400	bacteria	GAAGGGTCGTTTCGAGACC	57 °C
<i>hsa</i> - miR486 -5p	UCCUGUACUGAGCUGCCC CGAG		human	-*	60 °C



## FIGURES LEGENDS

**Figure 1.** Workflow of the initial screen for and validation of exogenous sRNA sequences in human plasma samples.

**Figure 2.** Detection of non-human sRNA species in column eluates and their removal from columns. **A)** qPCR amplification of six non-human sRNA species in extracts from human plasma and qPCR control (water). **B)** Detection of the same sRNA species in mock-extracts without input to extract columns and water passed through extraction columns (“eluate”). **C)** Levels of the same sRNA species in mock-extracts without and with DNase treatment during the extraction. **D)** Relative levels of sRNA remaining after pre-treatment of extraction columns with bleach or washing ten times with water, detected after eluting columns with water. **All:** mean results of three experiments, measured in reaction duplicates; error bars represent one standard deviation. Experiments displayed in panels **B** and **D** were performed on the same batch of columns, **A** and **C** on independent batches.

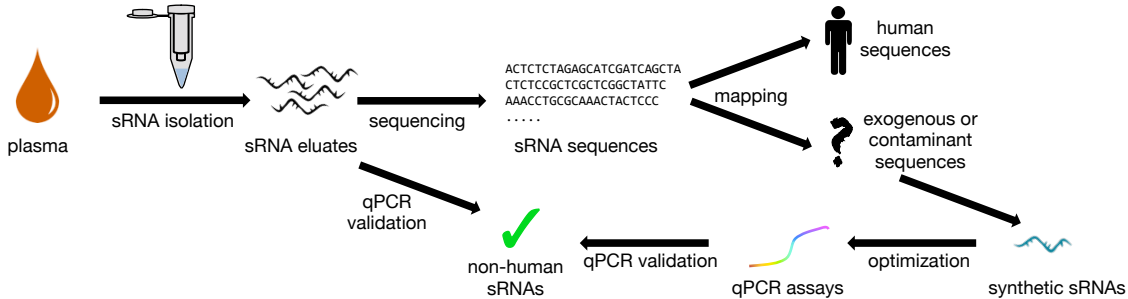
**Figure 3.** Detection of contaminant sequences in published sRNA sequencing datasets of low biomass samples. Datasets are referenced by NCBI bioproject accession or first author of the published manuscript. n: number of samples in the dataset. E: extraction kit used (if this information is available) – Q: regular miRNeasy (QIAGEN), T: TRIzol (Thermo Fisher), P: mirVana PARIS RNA extraction kit (Thermo Fisher), V: mirVana RNA extraction kit with phenol. rpm: reads per million. Error bars indicate one standard deviation.

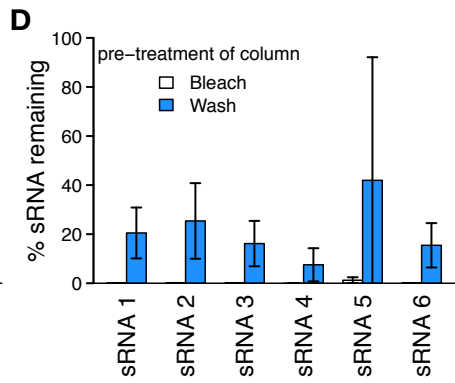
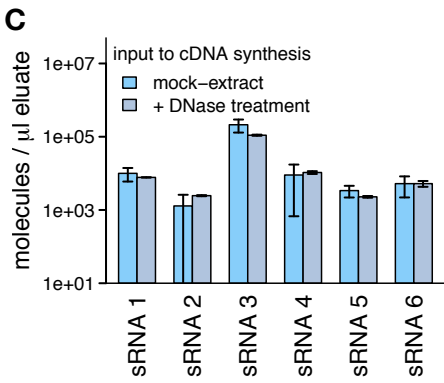
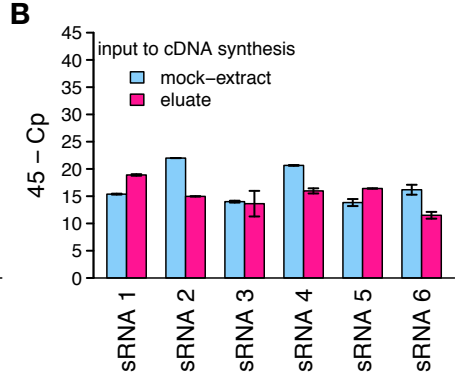
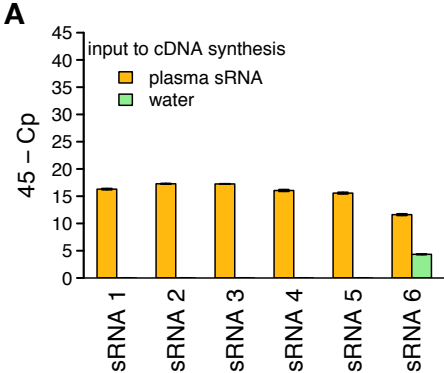
**Figure 4.** Levels of confirmed and potential contaminant sequences in eluates of regular and ultra-clean RNeasy spin columns. cpm: counts per million. **A)** Levels of contaminant sequences in eluates of two batches of regular and four batches of ultra-clean spin columns, based on qPCR; ultra-clean batches 1 and 2 are cleaned-up versions of regular batch 2 and ultra-clean batches 3 and 4 are cleaned-up versions of regular batch 3; error bars indicate one standard deviation. **B&C)** Numbers of different further potential contaminant sequences on the regular and ultra-clean spin columns from two

different batches. **D)** Total levels of further potential contaminant sequences, based on sRNA sequencing data normalized to spike-in levels.

**Figure 5.** Detection of contaminants in sRNA preparations of human plasma using different input volumes and extraction columns. **A)** Detected levels of the six contaminant sRNA sequences in sRNA sequencing data of preparations using 0 to 1115  $\mu$ l human plasma and regular or ultra-clean RNeasy spin columns. **B)** Detailed view of the data displayed in **A** for 100  $\mu$ l human plasma as input to regular and ultra-clean RNeasy spin columns. cpm: counts per million; error bars indicate one standard deviation.

**Figure 6.** Summary of recommendations for artefact-free analysis of sRNA by sequencing.





contaminant rpm

10000  
1000  
100  
10  
1

n = 9 33 3 9 16 1 3 6 3 2 45 2 2 3 2 2 1 4 2  
E : Q Q Q Q Q T Q T T T T T T T T T P V V

Wang et al.

PRJNA196121

PRJNA230622

PRJEB6683

PRJNA268092

PRJNA124945

PRJNA213914

PRJNA142489

PRJNA140779

PRJNA144887

PRJNA127103

PRJNA112823

PRJNA116139

PRJNA119935

PRJNA120351

PRJNA124019

PRJNA129787

PRJNA145935

PRJNA184379

