

1 **Charting extracellular transcriptomes in The Human Biofluid RNA Atlas**

2

3 Eva Hulstaert (1,2,3), Annelien Morlion (1,2), Francisco Avila Cobos (1,2), Kimberly Verniers
4 (1,2), Justine Nuytens (1,2), Eveline Vanden Eynde (1,2), Nurten Yigit (1,2), Jasper Anckaert
5 (1,2), Anja Geerts (4), Pieter Hindryckx (4), Peggy Jacques (5,6), Guy Brusselle (7), Ken R. Bracke
6 (7), Tania Maes (7), Thomas Malfait (7), Thierry Derveaux (8), Virginie Ninclaus (8), Caroline
7 Van Cauwenbergh (8), Kristien Roelens (9), Ellen Roets (9), Dimitri Hemelsoet (10), Kelly
8 Tilleman (11), Lieve Brochez (2,3), Scott Kuersten (12), Lukas Simon (13), Sebastian Karg (14),
9 Alexandra Kautzky-Willers (15), Michael Leutner (15), Christa Nöhammer (16), Ondrej Slaby
10 (17,18,19), Roméo Willinge Prins (20), Jan Koster (20), Steve Lefever (1,2), Gary P. Schroth
11 (12), Jo Vandesompele (1,2)*, Pieter Mestdagh (1,2)*

12

13 1 Center for Medical Genetics, Department of Biomolecular Medicine, Ghent University, C. Heymanslaan 10,
14 9000, Ghent, Belgium

15 2 Cancer Research Institute Ghent (CRIG), Ghent University, C. Heymanslaan 10, 9000, Ghent, Belgium

16 3 Department of Dermatology, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

17 4 Department of Gastroenterology, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

18 5 Department of Rheumatology, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

19 6 VIB Inflammation Research Center, Ghent University, Technologiepark 71, 9052, Ghent, Belgium

20 7 Department of Respiratory Medicine, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

21 8 Department of Ophthalmology, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

22 9 Department of Obstetrics, Women's Clinic, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent,
23 Belgium

24 10 Department of Neurology, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

25 11 Department of Reproductive Medicine, Ghent University Hospital, C. Heymanslaan 10, 9000, Ghent, Belgium

26 12 Illumina, San Diego, CA, USA

27 13 University of Texas Health Science Center, School of Biomedical Informatics, Center for Precision Health,
28 Houston, TX, USA

29 14 Helmholtz Zentrum München, German Research Center for Environmental Health, Institute of Computational
30 Biology, Neuherberg, Germany

31 15 Department of Internal Medicine III, Clinical Division of Endocrinology and Metabolism, Unit of Gender
32 Medicine, Medical University of Vienna, Vienna, Austria

33 16 Austrian Institute of Technology, Center for Health and Bioresources, Molecular Diagnostics, Vienna, Austria

34 17 Masaryk Memorial Cancer Institute, Department of Comprehensive Cancer Care, Brno, Czech Republic

35 18 Department of Pathology, University Hospital Brno, Brno, Czech Republic

36 19 Central European Institute of Technology, Masaryk University, Brno, Czech Republic

37 20 Department of Oncogenomics, Amsterdam University Medical Centers (AUMC), University of Amsterdam, the
38 Netherlands

39 * equally contributing authors

40

41

42 **Abstract**

43 Extracellular RNAs present in biofluids have emerged as potential biomarkers for disease.
44 Where most studies focus on plasma or serum, other biofluids may contain more informative
45 RNA molecules, depending on the type of disease. Here, we present an unprecedented atlas
46 of messenger, circular and small RNA transcriptomes of a comprehensive collection of 20
47 different human biofluids. By means of synthetic spike-in controls, we compared RNA content
48 across biofluids, revealing a more than 10 000-fold difference in RNA concentration. The
49 circular RNA fraction is increased in nearly all biofluids compared to tissues. Each biofluid
50 transcriptome is enriched for RNA molecules derived from specific tissues and cell types. In
51 addition, a subset of biofluids, including stool, sweat, saliva and sputum, contains high levels
52 of bacterial RNAs. Our atlas enables a more informed selection of the most relevant biofluid
53 to monitor particular diseases. To verify the biomarker potential in these biofluids, four
54 validation cohorts representing a broad spectrum of diseases were profiled, revealing
55 numerous differential RNAs between case and control subjects. Taken together, our results
56 reveal novel insights in the RNA content of human biofluids and may serve as a valuable
57 resource for future biomarker studies. All spike-normalized data is publicly available in the R2
58 web portal and serve as a basis to further explore the RNA content in biofluids.

59

60 **Keywords**

61 RNA-sequencing, biofluids, circular RNA, messenger RNA, small RNA, biomarker, extracellular
62 RNA, cell-free RNA

63

64

65 Introduction

66 Extracellular RNAs (exRNAs) in blood and other biofluids are emerging as potential
67 biomarkers for a wide range of diseases¹⁻⁶. These so-called liquid biopsies may offer a non-
68 invasive alternative to tissue biopsies for both diagnosis and treatment response monitoring.
69 Previous studies have extensively profiled the small RNA content of several biofluids and
70 identified large differences in the small RNA content amongst different biofluids.¹⁻¹² These
71 efforts were gathered by the NIH Extracellular RNA Communication Consortium in the
72 exRNA Atlas Resource (<https://exrna-atlas.org>).⁸ Besides microRNAs (miRNAs), the most
73 studied small RNA biotype in biofluids, other small RNAs, such as piwi-interacting RNAs
74 (piRNAs), small nuclear RNAs (snRNAs), small nucleolar RNAs (snoRNAs), ribosomal RNAs
75 (rRNAs), transfer RNA fragments (tRNAs) and Y-RNAs have also been identified^{5-7,9,12,13}.
76 Weber et al.¹³ was the first to compare the miRNA content in 12 different human biofluids
77 (pooled samples of plasma, saliva, tears, urine, amniotic fluid, colostrum, breast milk,
78 bronchial lavage fluid, cerebrospinal fluid, peritoneal fluid, pleural fluid and seminal plasma)
79 using reverse transcription quantitative polymerase chain reaction (RT-qPCR) of selected
80 miRNAs. Large variations in RNA concentration were observed among the different biofluids,
81 with the highest small RNA concentrations measured in breast milk and seminal fluid. Since
82 the advent of small RNA sequencing, other small RNA biotypes were characterized in various
83 biofluids, such as plasma, serum, stool, urine, amniotic fluid, bronchial lavage fluid, bile,
84 cerebrospinal fluid (CSF), saliva, seminal plasma and ovarian follicle fluid^{5,7,9,9,12}. The
85 distribution of small RNA biotypes clearly varies across these biofluids, with a high
86 abundance of piRNAs and tRNAs reported in urine and a high abundance of Y-RNAs in
87 plasma^{6,7,12}. Also non-human RNA sequences, mapping to bacterial genomes, were reported
88 in plasma, urine and saliva⁶.
89 A systematic RNA-sequencing analysis of biofluids to explore the messenger RNAs (mRNA) and
90 circular RNA (circRNA) transcriptome is challenging due to low RNA concentration and RNA
91 fragmentation in biofluids. As such, most studies have explored the abundance of individual
92 mRNAs in one specific biofluid by RT-qPCR¹⁴⁻²⁰. CircRNAs have been reported in saliva²¹,
93 semen²², blood²³ and urine^{24,25}. Recently, the mRNA content of plasma and serum has been
94 investigated using dedicated sequencing approaches like Phospho-RNA-Seq, SILVER-seq and
95 SMARTer Stranded Total RNA-Seq method²⁶⁻²⁹. Studies comparing the small RNA, mRNA and

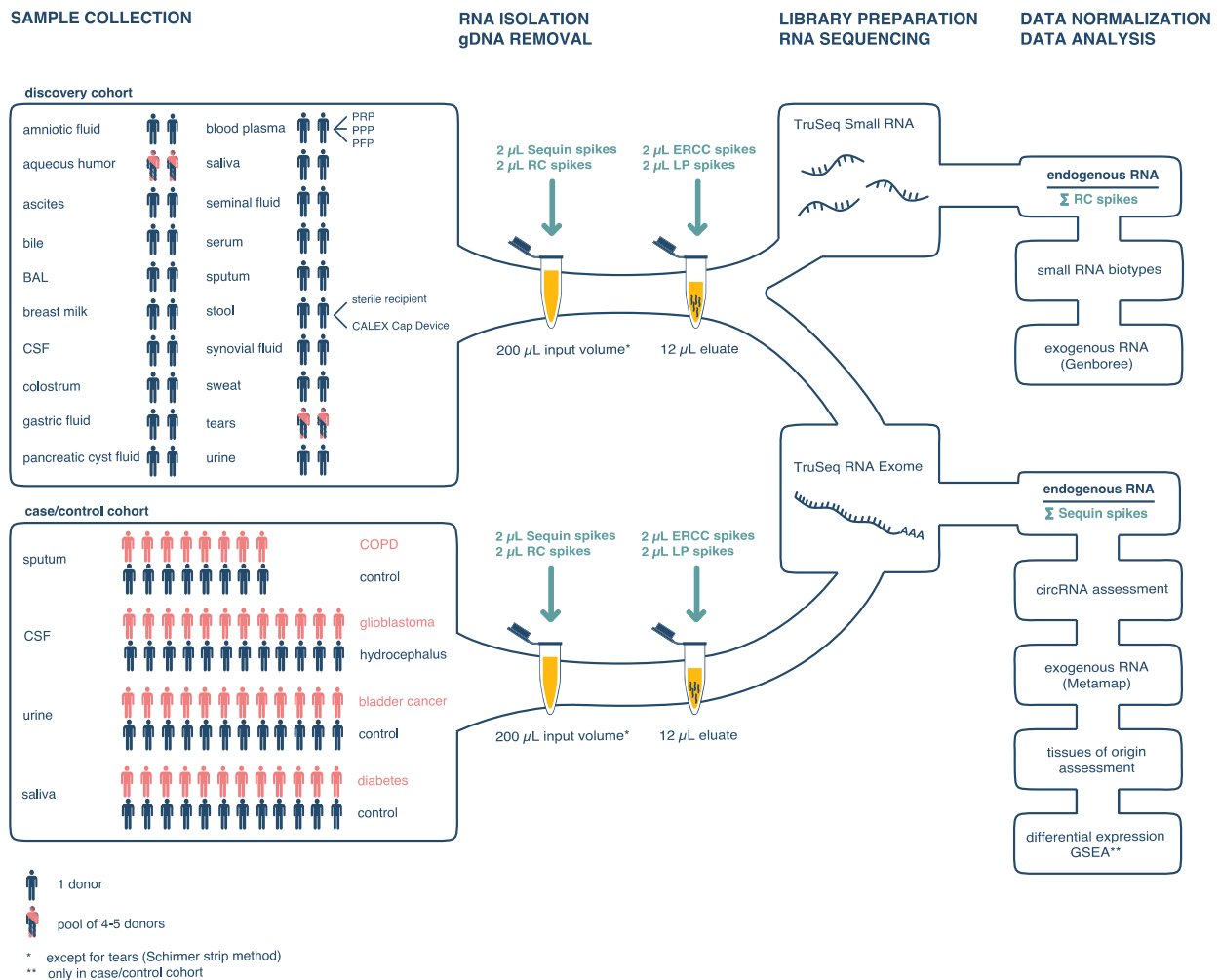
96 circRNA content in a wide range of human biofluids are currently lacking and are essential to
97 explore the biomarker potential of exRNAs.

98 The goal of the Human Biofluid RNA Atlas is to define the extracellular transcriptome across a
99 wide range of human biofluids (amniotic fluid, aqueous humor, ascites, bile, bronchial lavage
100 fluid, breast milk, cerebrospinal fluid, colostrum, gastric fluid, pancreatic cyst fluid, plasma,
101 saliva, seminal fluid, serum, sputum, stool, synovial fluid, sweat, tear fluid and urine) and to
102 assess biomarker potential in selected case-control cohorts. We used small RNA-sequencing
103 to quantify different small RNA species and present a dedicated mRNA-capture sequencing
104 workflow to simultaneously quantify mRNAs and circRNAs.

105 In the first phase of our study, small RNA sequencing and mRNA capture sequencing was
106 performed in a discovery cohort of 20 different biofluids (Fig. 1). The goal of this phase was to
107 assess the technical feasibility of the methodology and to generate a comprehensive set of
108 mRNAs, circRNAs and small RNAs in which the contributing tissues and cell types per biofluid
109 were assessed.

110 In the second phase of our study, we aimed to investigate the biological relevance of exRNAs
111 in various biofluids. Therefore, mRNA capture sequencing was applied to four different
112 case/control cohorts, each consisting of 16-24 samples (Fig. 1). These samples included
113 sputum samples from 8 patients with chronic obstructive pulmonary disease (COPD) versus 8
114 controls, urine samples from 12 bladder cancer patients versus 12 controls, CSF samples from
115 12 glioblastoma patients versus 12 hydrocephalus patients and saliva samples from 12
116 diabetes mellitus patients versus 12 controls.

117 The resulting catalog of extracellular transcriptomes of 185 human samples can guide
118 researchers in the biomarker field to investigate other biofluids besides the well-studied
119 blood-derived ones and is a first step to more dedicated mRNA and circRNA profiling of
120 biofluids in larger cohorts.



121

122 **Fig.1 Study flow chart**

123 *In the discovery cohort, 20 different biofluids were collected in two donors or in a pool of 4-5*
 124 *donors. In the case/control cohorts, selected biofluids (sputum, CSF, urine and saliva) were*
 125 *collected in 8-12 patients and an equal number of healthy controls. Both small RNA sequencing*
 126 *and mRNA capture sequencing were performed in the discovery cohort. In the case/control*
 127 *cohorts, mRNA capture sequencing was performed. To compare the RNA content across the*
 128 *different biofluids, the RC spikes and the Sequin spikes are used for normalization of small RNA*
 129 *and mRNA data, respectively.*

130 *BAL, bronchoalveolar lavage fluid; CSF, cerebrospinal fluid; PRP, platelet-rich plasma; PPP,*
 131 *platelet-poor plasma; PFP, platelet-free plasma*

132

133 **Results**

134 ***RNA spike-in controls enable process control of the RNA sequencing workflow***

135 Synthetic spike-in RNA sequences are crucial to control the process from RNA isolation to RNA
136 sequencing, especially when working with challenging and low input material. We applied 4
137 different mixes of synthetic RNA spike-in controls (in total 189 RNAs) as workflow processing
138 and normalization controls that enable direct comparison of the RNA profiles across the
139 different biofluids. Sequin and Small RNA extraction Control (RC) spikes were added prior to
140 RNA isolation whereas External RNA Control Consortium (ERCC) spikes and small RNA Library
141 Prep (LP) spikes were added to the RNA eluate prior to genomic DNA (gDNA) removal (Fig. 1).
142 Of note, every spike mix consists of multiple RNA molecules of different lengths over a wide
143 concentration range. Detailed information is provided in Supplementary Note 1. Besides
144 normalization, the spike-in controls enabled quality control of the RNA extraction and library
145 preparation steps in the workflow and relative quantification of the RNA yield and
146 concentration across the different biofluids.

147 First, the correlation between the expected and the observed relative quantities for all four
148 spike mixes can be used to assess quantitative linearity. In the discovery cohort, the expected
149 and the observed relative quantities for all four spike mixes were well correlated (Pearson
150 correlation coefficients range from 0.50 to 1.00 for Sequin spikes, 0.92 to 1.00 for ERCC spikes,
151 0.44 to 0.98 for RC spikes and 0.40 to 0.96 for LP spikes). In some biofluids (e.g. seminal plasma
152 and tears), the sequencing coverage of spikes was low due to a high concentration of
153 endogenous RNA. Detailed information per sample is provided in Supplementary Fig. 1.

154 The spike-in controls can also be used to assess the RNA isolation efficiency. The Sequin/ERCC
155 ratio and the RC/LP ratio reflect the relative mRNA and microRNA isolation efficiency,
156 respectively. A 170-fold and 104-fold difference in RNA isolation efficiency across the samples
157 was observed when assessing long and small RNAs, respectively (Supplementary Fig. 2). These
158 differences underline the challenges of working with heterogenous samples and the
159 importance of spike-in controls for proper data normalization and cross-sample comparison
160 of results.

161 Finally, the spikes can be utilized to normalize the endogenous RNA abundance data. In this
162 study, we applied a biofluid volume-based normalization by dividing the RNA reads consumed
163 by the endogenous transcripts by the sum of the Sequin spikes for mRNA data and by the sum

164 of the RC spikes for small RNA data. The spike-normalized data represent relative abundance
165 values of RNA molecules proportional to the input volume. Of note, there is an inverse
166 relationship between the number of spike-in RNA reads and the number of endogenous RNA
167 reads. As such, the ratio between the sum of the reads consumed by the endogenous
168 transcripts and the total number of spike-in reads is a relative measure for the RNA
169 concentration of the various samples.

170 ***Highly variable mRNA and small RNA content among biofluids in the discovery cohort***

171 Both small RNAs and mRNAs were quantified in each of the 20 biofluids in the discovery
172 cohort. Mapping rates varied substantially across the different biofluids (Fig. 2A). In general,
173 the proportion of mapped reads was higher for the mRNA capture sequencing data (further
174 referred to as mRNA data) than for the small RNA sequencing data, in line with the fact that
175 human mRNAs were enriched using biotinylated capture probes during the library
176 preparation. The fraction of mapped reads in the mRNA data ranged from 16% in stool to 97%
177 in seminal plasma. Low mapping rates were observed in stool, in one of the bile samples and
178 in saliva. Mapping rates for samples in the case/control cohorts are in line with these of the
179 discovery cohort (Supplementary Fig. 3A). In the small RNA sequencing data, the proportion
180 of mapped reads ranged from ~7% in stool, saliva and CSF to 95% in platelet-rich plasma (PRP).
181 A 10 000-fold difference in mRNA and small RNA concentration was observed between the
182 lowest concentrated fluids, i.e. platelet-free plasma, urine and CSF, and the highest
183 concentrated biofluids, i.e. tears, seminal plasma and bile (Fig. 2B). The generalizability of the
184 difference in mRNA concentration between highly concentrated biofluids (seminal plasma)
185 and lowly concentrated biofluids (CSF) was confirmed in additional samples (Supplementary
186 Fig. 3B). In the discovery cohort a 5547-fold difference in mRNA concentration is observed
187 between seminal plasma and CSF; in independent validation samples, a similarly large 19 851-
188 fold difference in mRNA concentration is observed between both biofluids. In the discovery
189 cohort, the mRNA and miRNA concentrations were significantly correlated across biofluids
190 (Pearson correlation coefficient 0.76, p-value = 8.5e-10, Fig. 2D). Normalized abundance levels
191 of exRNAs were significantly correlated between biological replicates within each biofluid
192 (Supplementary Fig. 4). The median Pearson correlation coefficient of the mRNA and the small
193 RNA data was 0.84 and 0.92, respectively. While the mRNA and miRNA data was well
194 correlated in most biofluids (e.g. tears, colostrum, saliva), correlation in other biofluids (e.g.
195 bile, pancreatic cyst fluid) was poor. These biofluids are obtained with a more challenging

196 collection method involving echo-endoscopy, impacting the reproducibility of collection and
 197 the correlation of the RNA content between biological replicates.

198 The likelihood of identifying RNA biomarkers in a given biofluid will not only depend on its
 199 relative RNA concentration, but also on its RNA diversity, here approximated by the fraction
 200 of read counts consumed by the top 10 most abundant mRNAs/miRNAs (Fig. 2C). In aqueous
 201 humor, the top 10 mRNAs represent up to 70% of all reads, indicating that this fluid does not
 202 contain a rich mRNA repertoire. In both PRP and PPP, about 50% of all reads go to the top 10
 203 mRNAs. While amniotic fluid has a median RNA concentration, this fluid seems to contain a
 204 diverse mRNA profile, with only 7% of all reads going to the top 10 mRNAs. When looking into
 205 the miRNA data, the top 10 miRNAs represent more than 90% of all reads in PFP, urine and
 206 serum. BAL contains the most diverse miRNA repertoire, with 57% of all reads going to the
 207 top 10 miRNAs. Similar conclusions with respect to biofluid exRNA diversity can be drawn
 208 based on the number of miRNAs/mRNAs representing 50% of the counts (Supplementary Fig.
 209 5). RNA diversity is also reflected by the number of detected exRNAs. The total number of
 210 mRNAs and miRNAs detected with at least 4 counts in both samples of the same biofluid
 211 ranged from 13 722 mRNAs in pancreatic cyst fluid to 107 mRNAs in aqueous humor and from
 212 231 miRNAs in tears to 18 miRNAs in stool (Table 1).

213

214 **Table 1 Number of mRNAs and miRNAs per biofluid.**

215 *The number of mRNAs and miRNAs with at least 4 unique read counts in both replicates is*
 216 *shown per biofluid.*

217

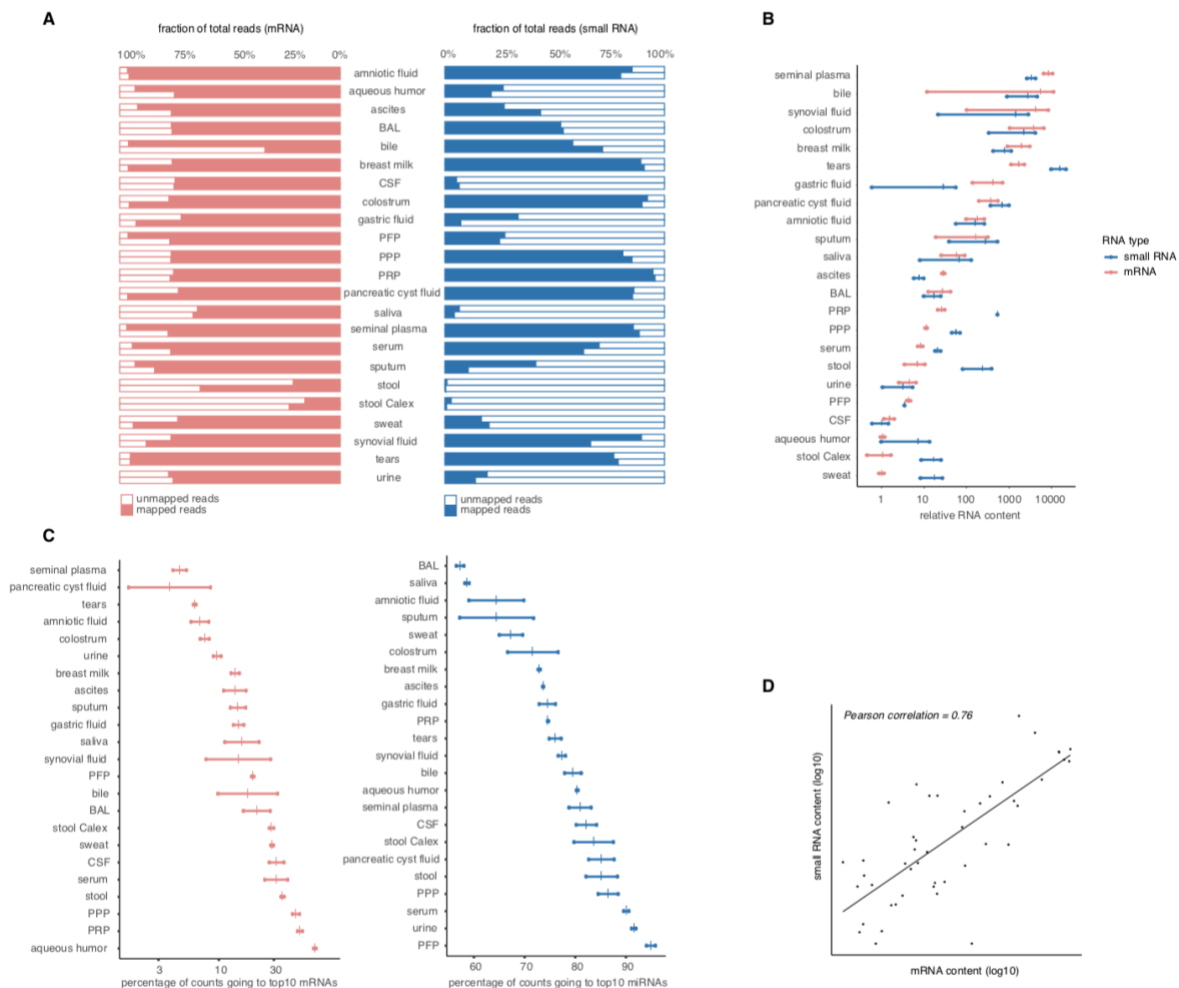
biofluid	number of mRNAs
amniotic fluid	10 531
aqueous humor	107
ascites	5578
BAL	3565
bile	2279
breastmilk	11 607
colostrum	11 914
CSF	438
gastric fluid	9288
pancreatic cyst fluid	13 722
PFP	2699
PPP	4548
PRP	5440
saliva	6353
seminal plasma	11 868
serum	4152
sputum	7738

biofluid	number of miRNAs
amniotic fluid	119
aqueous humor	20
ascites	75
BAL	126
bile	45
breastmilk	213
colostrum	229
CSF	32
gastric fluid	21
pancreatic cyst fluid	129
PFP	95
PPP	113
PRP	192
saliva	110
seminal plasma	211
serum	122
sputum	91

stool	134
stool Calex	135
sweat	410
synovial fluid	1614
tears	13 366
urine	2094

stool	19
stool Calex	18
sweat	45
synovial fluid	122
tears	231
urine	41

218
219



220

221 **Fig. 2 mRNA and small RNA content varies across the 20 biofluids**

222 (A) Percentage of the total read count mapping to the human transcriptome.

223 (B) Relative RNA concentration per biofluid; every dot represents the relative RNA
224 concentration in one sample, every vertical mark indicates the mean per biofluid.

225 (C) The diversity of the RNA content expressed as fraction of read counts consumed by the
226 top 10 most abundant mRNAs/miRNAs. Only genes with at least 4 unique reads are
227 taken into account. Every dot represents the fraction in one sample, every vertical mark
228 indicates the mean percentage per biofluid.

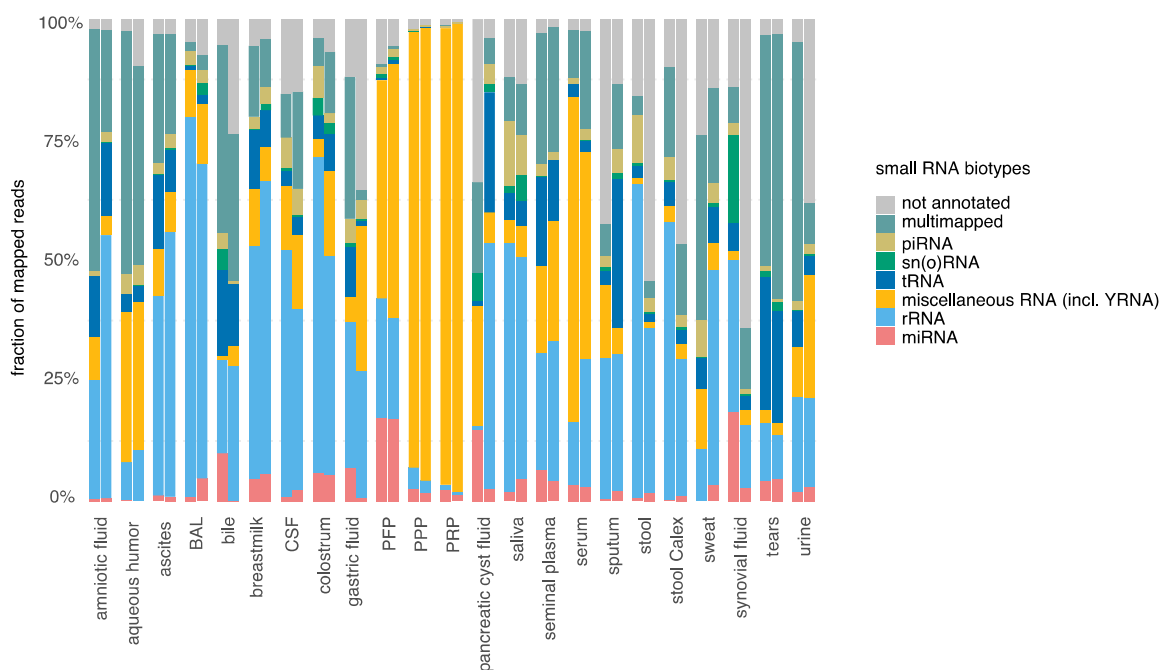
229 (D) Correlation between the small RNA and the mRNA relative concentration. The Pearson
 230 correlation coefficient is 0.76 (p -value = 8.58×10^{-10}). The correlation coefficients is
 231 calculated on log10 transformed data.

232 BAL, bronchoalveolar lavage fluid; CSF, cerebrospinal fluid; PRP, platelet-rich plasma; PPP,
 233 platelet-poor plasma; PFP, platelet-free plasma

234

235 **The distribution of small RNA biotypes varies across the different biofluids**

236 The distribution of small RNA biotypes shows distinct patterns among the 20 different
 237 biofluids (Fig. 3). The exceptionally high percentage of miscellaneous RNAs (mainly Y-RNAs)
 238 observed in blood-derived fluids is in line with a previous study¹² and with the Y-RNA function
 239 in platelets. The fraction of reads mapping to miRNAs is lower than 15% in all samples but
 240 platelet-free plasma and one synovial fluid sample. Tears, bile and amniotic fluid have the
 241 highest fraction of tRNA fragments while saliva has the highest fraction of piRNAs. The rRNA
 242 fraction is higher than 15% in all samples but tears, aqueous fluid and the three plasma
 243 fractions. The majority of these reads map to the 45S ribosomal RNA transcript. The not
 244 annotated read fraction contains uniquely mapped reads that could not be classified in one of
 245 the small RNA biotypes. These reads most likely originate from degraded longer RNAs, such
 246 as mRNAs and long non-coding RNAs.



247

248 **Fig. 3 Distinct small RNA biotype patterns are present across the different biofluids**

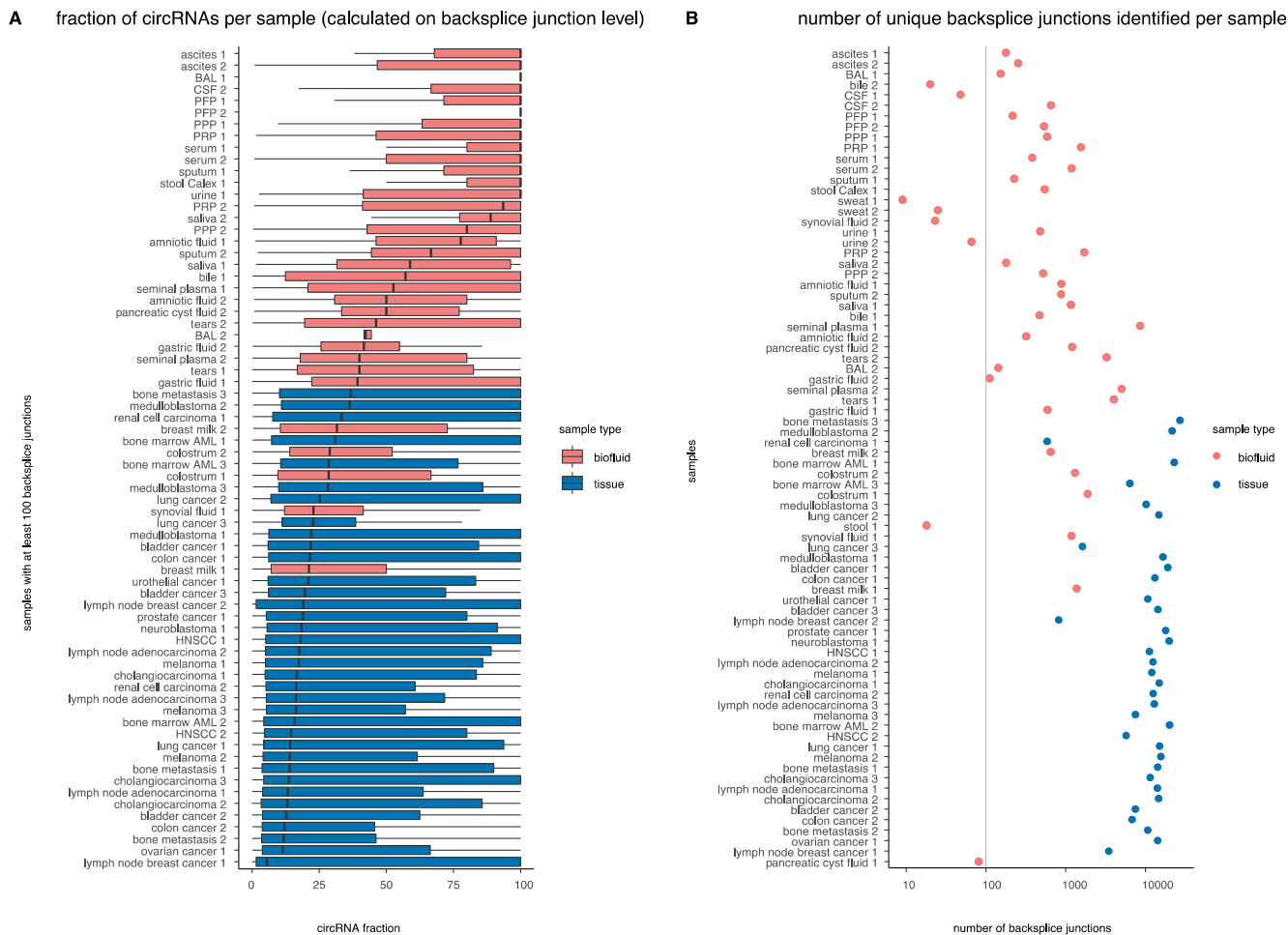
249 *The fraction of reads that align to small RNA biotypes are shown per biofluid. Only mapped*
250 *reads of the small RNA sequencing data are taken into account. BAL, bronchoalveolar lavage*
251 *fluid; CSF, cerebrospinal fluid; miRNA: microRNA; PFP, platelet-free plasma; PPP, platelet-poor*
252 *plasma; PRP, platelet-rich plasma; piRNAs: piwi-interacting RNA; sn(o)RNAs: small nuclear and*
253 *nucleolar RNAs; tRNAs: transfer RNA.*

254

255 ***Circular RNAs are enriched in biofluids compared to tissues***

256 CircRNAs are produced from unspliced RNA through a process called back-splicing where a
257 downstream 5' donor binds to an upstream 3' acceptor. CircRNAs are resistant to endogenous
258 exonucleases that target free 5' or 3' terminal ends. As a result, circRNAs are highly stable and
259 have extended half-lives compared to linear mRNAs.³⁰ CircRNAs have been reported to be
260 present in numerous human tissues²⁴ and in a few biofluids such as saliva²¹, blood³¹, semen²²
261 and urine^{24,25}. A direct comparison of the circRNA read fraction between biofluids and tissues
262 is currently lacking in literature. We compared the circRNA fraction, for genes that produce
263 both linear and circular transcripts, identified through mRNA capture sequencing of the 20
264 biofluids in this study with the circRNA fraction identified in mRNA capture sequencing of 36
265 cancerous tissue types obtained from the MiOncoCirc Database²⁴. While more unique
266 backsplice junctions were identified in tissues compared to biofluids, in line with the higher
267 RNA concentration in tissues (Fig. 4B), the circRNA read fraction is clearly higher in biofluid
268 exRNA compared to cellular RNA (Fig. 4A). The median circRNA read fraction in biofluids is
269 84.4%, which is significantly higher than the median circRNA read fraction in tissues of 17.5%
270 (Mann-Whitney-U test, two-sided, p-value = 5.36×10^{-12}). For genes that produce both linear
271 and circular transcripts, the stable circRNAs are more abundant than the linear mRNAs in
272 biofluids, while it is the other way around in tissues.

273 We used two different methods to define the circRNA read fraction (see "*Circular RNA*
274 *detection*" in methods; Supplementary Fig. 6): one based on individual backsplice junctions
275 (shown in Fig. 4) and another method based on backsplice junctions aggregated at gene-level
276 (Supplementary Fig. 7). Both methods clearly point towards a substantial enrichment of
277 circRNAs in biofluids.



278

279 **Fig. 4 CircRNAs are enriched in biofluids compared to tissues**

280 (A) The circRNA fraction, calculated at the backsplice junction level, is plotted per sample
 281 and is higher in cell-free biofluid RNA than in tissue RNA. Only samples with at least
 282 100 backsplice junctions are plotted.

283 (B) The number of unique backsplice junctions per sample is higher in tissues compared to
 284 biofluids, in line with the higher input concentration of RNA into the library prep.

285 AML, acute myeloid leukemia; BAL, bronchoalveolar lavage fluid; CSF, cerebrospinal fluid;
 286 HNSCC: head and neck squamous-cell carcinoma; PFP, platelet-free plasma; PPP, platelet-
 287 poor plasma; PRP, platelet-rich plasma

288

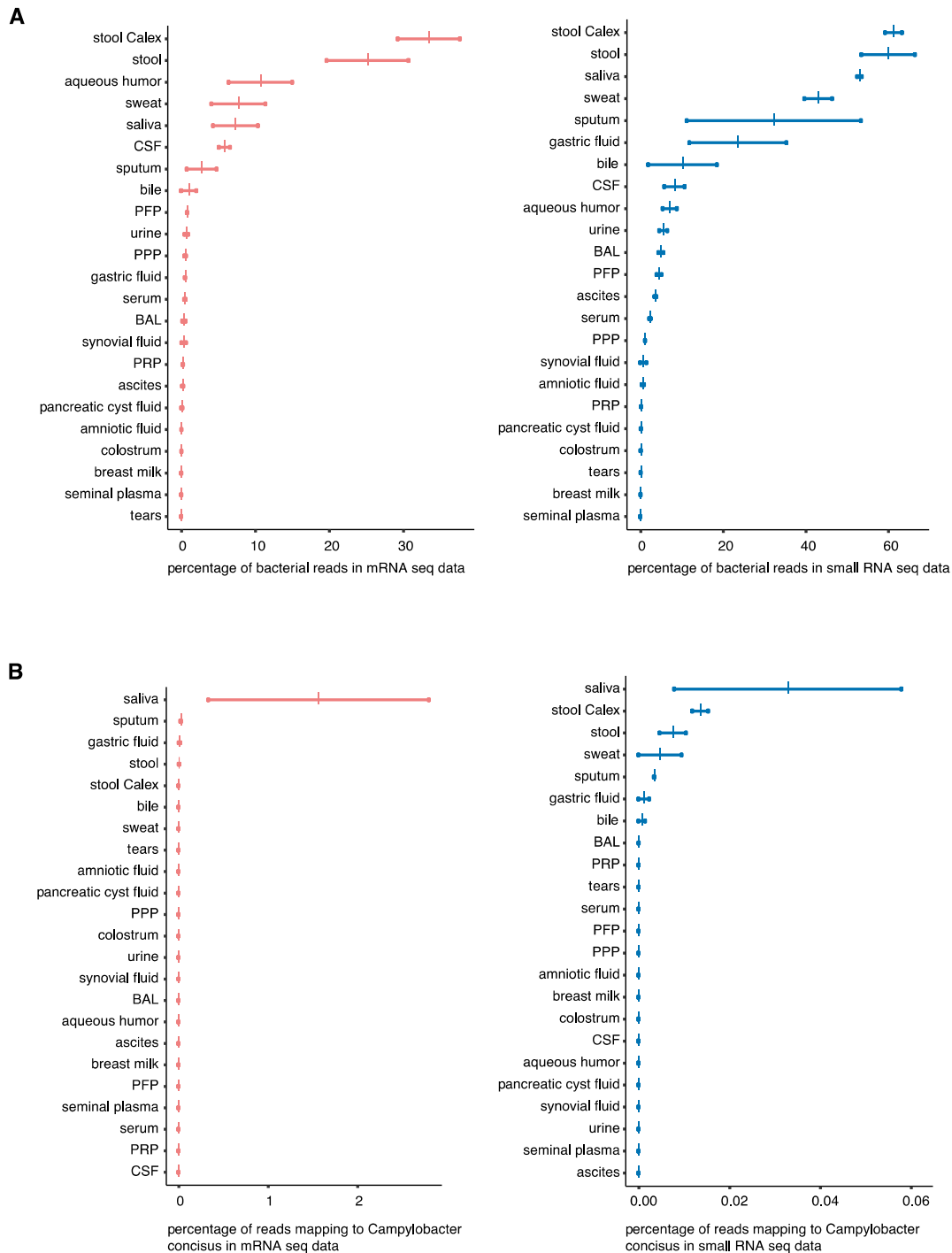
289 **Assessment of exogenous RNA in human biofluids**

290 Two dedicated pipelines were used for the non-trivial assessment of the presence of microbial
 291 or viral RNA in human biofluid extracellular RNA. Overall, the fraction of bacterial reads is

292 higher in small RNA sequencing data than in the mRNA data, in line with the unbiased nature
293 of small RNA sequencing and the targeted hybrid capture enrichment using probes against
294 human RNA during the mRNA capture library preparation. Stool (both collection methods),
295 sweat, saliva and sputum are among the biofluids with the highest fraction of bacterial RNA
296 in both the small RNA sequencing data and the mRNA data. The percentage of bacterial reads
297 in mRNA data and in small RNA data are significantly correlated across biofluids (Pearson
298 correlation coefficient 0.78, p-value = 1.94e-10).

299 Bacterial reads in aqueous humor and CSF, two fluids with very low endogenous RNA content
300 that were collected in a sterile setting (and thus presumed to be sterile), most likely reflect
301 background contamination during the workflow³². To illustrate the biological relevance of the
302 bacterial signal, we looked into reads mapping to *Campylobacter concisus*, a gram-negative
303 bacterium that is known to primarily colonize the human oral cavity, with some strains
304 translocated to the intestinal tract³³. We confirm the selective presence of reads mapping to
305 *Campylobacter concisus* in saliva in both the small RNA and the mRNA data (Fig. 5B). In all
306 samples and for both the small RNA and the mRNA data, the percentage of the total reads
307 that maps to viral transcriptomes is less than 1%.

308



309

310 **Fig. 5. Reads mapping to bacterial genomes**

311 (A) Percentage of reads mapping to bacteria in mRNA data (pink) and in small RNA
312 sequencing data (blue).

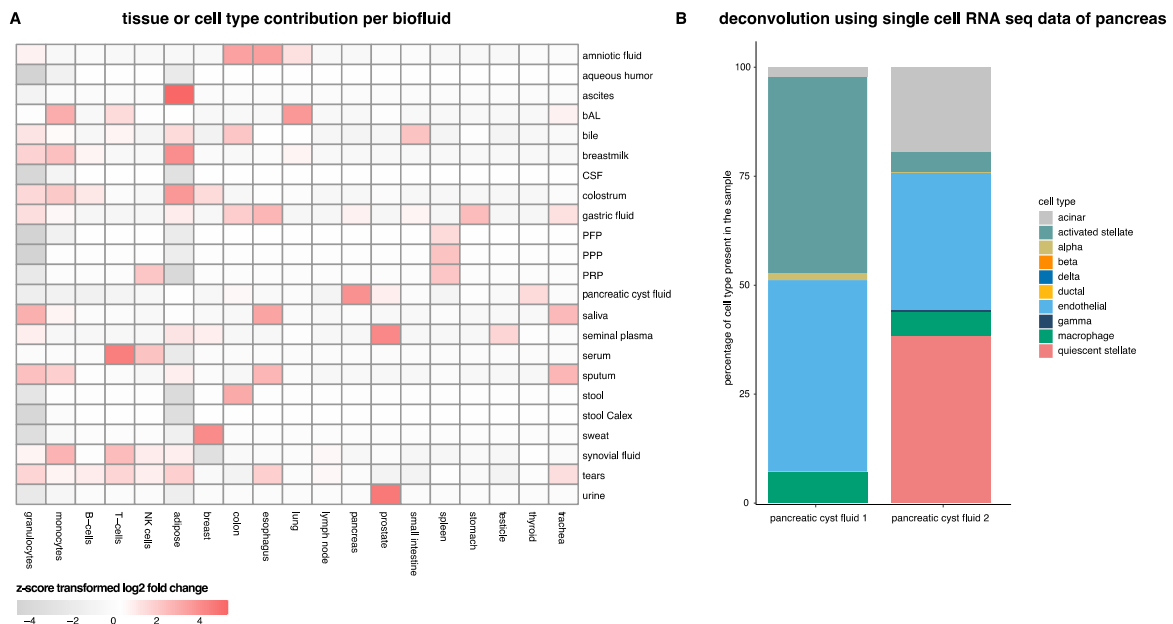
313 (B) Percentage of reads mapping to *Campylobacter concisus* in mRNA data (pink) and in
314 small RNA sequencing data (blue). *Campylobacter concisus* is known to be present in
315 saliva.

316 ***Assessment of the tissues of origin and deconvolution of pancreatic cyst fluid***

317 Gaining insights in tissue contribution to biofluid RNA profiles may guide the selection of the
318 most appropriate biofluid to investigate a given disease. To define tissues that specifically
319 contribute RNA molecules to individual biofluids, we explored the relationship between
320 extracellular mRNA levels and tissue or cell type specific mRNA signatures. The heatmap in
321 Fig. 6A highlights the relative contribution of tissues and cell types to a specific biofluid
322 compared to the other biofluids. More detailed results per biofluid are shown in
323 Supplementary Fig. 8. The results of this analysis were validated in an independent sample
324 cohort for CSF, saliva, sputum, seminal plasma and urine (Supplementary Fig. 3C). As
325 expected, prostate tissue RNA markers are more abundant in urine and in seminal plasma
326 than in any other biofluid. Both sputum and saliva contain mRNAs specific for trachea and
327 esophagus. In amniotic fluid, markers for esophagus, small intestine, colon and lung are more
328 abundant than the other tissues and cell types, probably reflecting organs that actively shed
329 RNA (at the gestational age of sampling) into the amniotic cavity. These data strongly suggest
330 that biofluid mRNA levels, at least to some degree, reflect intracellular mRNA levels from cells
331 that produce or transport the fluid. To further investigate the origin of biofluid RNA at the
332 cellular level, we applied computational deconvolution of the pancreatic cyst fluid RNA
333 profiles using single cell RNA sequencing data from 10 pancreatic cell types³⁴. Fig. 6B reveals
334 that pancreatic cyst fluid 1 consists of 45% of activated stellate cells and 43% of endothelial
335 cells, while pancreatic cyst fluid 2 mainly consists of quiescent stellate cells (38%), endothelial
336 cells (31%) and acinar cells (19%).

337

338



339

340 **Fig. 6 Identification of the tissues of origin per biofluid and deconvolution of pancreatic cyst**
 341 **fluid**

342 (A) Assessment of the tissues of origin in the biofluids of the discovery cohort.

343 Heatmap showing tissues and cell types that contribute more specifically to a certain biofluid
 344 compared to the other biofluids. Rows depict the biofluids of the discovery cohort and the
 345 columns are the tissues or cell types for which markers were selected based on the RNA Atlas³⁵.
 346 For visualization purposes, only tissues and cell types with a z-score transformed log₂ fold
 347 change $\geq |1|$ in at least one biofluid are shown.

348 (B) Composition of pancreatic cyst fluid samples based on deconvolution using sequencing
 349 data from 10 pancreatic cell types.

350

351 **Biomarker potential of mRNA in sputum, urine, CSF and saliva in selected case/control**
 352 **cohorts**

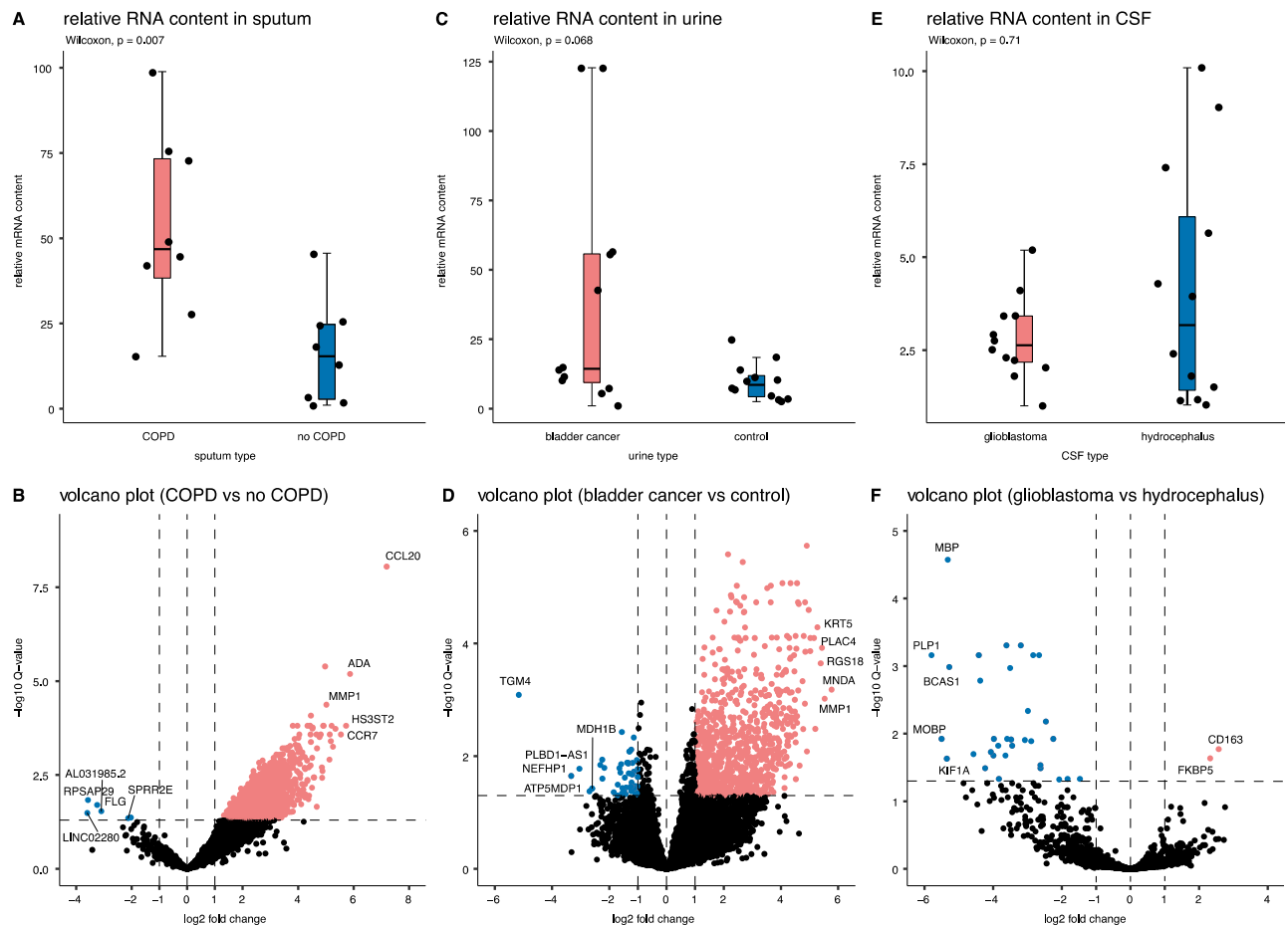
353 Additional biofluid samples were collected in patients with a specific disease or in healthy
 354 controls to investigate potential biologically relevant differences in mRNA content between
 355 both groups. Sequin RNA spikes were used for biofluid volume-based data normalization.
 356 Strikingly, the relative RNA concentration in sputum of COPD patients was higher than in non-
 357 COPD patients, probably reflecting the high turnover of immune cells during the state of
 358 chronic inflammation (Fig. 7A). Differential expression analysis revealed 5513 and 6 mRNAs
 359 that were significantly up- and downregulated, respectively, in sputum from COPD patients

360 compared to healthy controls (Fig. 7B). CCL20, the most differential mRNA, showed a 146-fold
361 upregulation in COPD patients compared to healthy donors. This potent chemokine attracting
362 dendritic cells has previously been linked to the pathogenesis of COPD^{36,37}. ADA and MMP1,
363 also among the most differential mRNAs, have also been associated with the pathogenesis of
364 COPD³⁸⁻⁴⁰. To verify the RNA-seq findings, 8/8 of the most differentially abundant mRNAs
365 were validated by RT-qPCR (Supplementary Fig. 9A-B).

366 In contrast to COPD, the relative RNA content is comparable in urine from bladder cancer
367 patients and healthy volunteers, in CSF from glioblastoma patients and hydrocephalus patient,
368 and in saliva from diabetes patients and healthy volunteers (Fig. 7C/E, Supplementary Fig. 10).
369 A higher RNA yield in CSF from glioblastoma patients compared to CSF from healthy controls
370 has been reported by Saugstad et al.⁴¹, however the collection method of CSF differed
371 between both groups and it is therefore not possible to assess whether the reported
372 difference in RNA yield between both groups is due to the different CSF collection site (lumbar
373 puncture versus craniotomy) or due to the neurological disease. In urine from patients with a
374 muscle invaded bladder cancer, 529 mRNAs and 9 mRNAs were significantly upregulated and
375 downregulated, respectively, compared to urine from healthy volunteers (Fig. 7D). Some of
376 the upregulated mRNAs, such as MDK, SLC2A1, GPRC5A, KRT17 and KRT5, have been reported
377 in urine and were suggested as biomarker for the accurate detection and classification of
378 bladder cancer⁴²⁻⁴⁵. In CSF from glioblastoma patients, only 2 mRNAs are significantly
379 upregulated compared to CSF from hydrocephalus patients. CD163, one of the upregulated
380 genes in glioblastoma, has been linked with glioblastoma pathogenesis⁴⁶. In saliva from
381 diabetes patients and saliva from healthy volunteers, no differentially expressed genes could
382 be identified. A list with differentially expressed genes in all case/control cohorts can be found
383 in Supplementary Data 5.

384 Differential abundance analysis was performed for circular RNAs as well, but in none of the
385 case/control cohorts differentially abundant circRNAs could be detected (data not shown). As
386 circular RNAs can only be identified based on their backsplice junction, the read coverage is
387 generally (too) low for biomarker discovery based on mRNA capture sequencing data. When
388 applying a similar strategy for mRNAs by looking at the reads of only one "linear only" junction
389 per gene (outside every detected back-splice junction) a significantly lower number of
390 differentially abundant mRNAs were detected (sputum: 13 out of 5519 mRNAs; urine: 0 out

391 of 538 mRNAs; CSF: 0 out of 35 mRNAs). These results strongly suggest that a dedicated
 392 circRNA enrichment strategies may be needed to assess circRNA biomarker potential.
 393 To validate the identification of the 10 most abundant circRNAs detected by mRNA capture
 394 sequencing in sputum, an orthogonal validation by RT-qPCR of the backsplice sequence region
 395 was performed. For 9 of the 10 circRNAs, the RNA-sequencing results could be validated.
 396 (Supplementary Fig. 9C)



397
 398 **Fig. 7 Relative RNA concentration and volcano plot in case/control cohorts**

399 *Top: Boxplots of relative mRNA content, bottom: Volcano plots of differentially expressed*
 400 *mRNAs ($q < 0.05$; pink up; blue down in patient vs. control) with labeling of up to 5 most*
 401 *differential genes. (A) Sputum from COPD patients ($n = 8$) compared to sputum from healthy*
 402 *donors ($n = 8$; Wilcoxon rank test, two-sided, $p = 0.007$); (B) 5513 and 6 mRNAs up and down,*
 403 *respectively in COPD samples. (C) Urine from bladder cancer patients ($n = 12$) compared to*
 404 *urine from healthy donors ($n = 12$; Wilcoxon signed-rank test, two-sided, $p = 0.068$). (D) 529*
 405 *and 9 mRNAs up and down, respectively in bladder cancer samples. (E) CSF from glioblastoma*
 406 *cancer patients ($n = 12$) compared to CSF from hydrocephalus patients ($n = 12$); Wilcoxon*

407 *signed-rank test, two-sided, $p = 0.71$); (D) 2 and 33 mRNAs up and down, respectively in*
408 *glioblastoma samples*

409 **Discussion**

410 By applying two complementary RNA-sequencing technologies on 20 different biofluids, we
411 assembled the most comprehensive human biofluid transcriptome, covering small RNAs,
412 mRNAs and circRNAs. Until now, most efforts to investigate and compare the RNA content
413 within biofluids focused on small RNA sequencing, most likely because of technical limitations
414 and unawareness of the abundance of extracellular mRNA (fragments)^{5-7,9,12,13}.

415 The availability of both small RNA sequencing data and mRNA data allows a more in-depth
416 characterization of the human transcriptome in biofluids. To our knowledge, this is the first
417 study reporting on the mRNA content, generated through a dedicated mRNA enrichment
418 sequencing method, in tear fluid, amniotic fluid, aqueous humor, bile, bronchial lavage fluid,
419 gastric fluid, saliva, seminal plasma, synovial fluid, sweat and urine. Selected mRNAs were
420 previously studied by means of RT-qPCR in amniotic fluid¹⁴, pancreatic cyst fluid^{15,18}, seminal
421 plasma¹⁶, sputum¹⁷, stool¹⁹ and in extracellular vesicles isolated from cell-free urine²⁰. In
422 saliva, selected mRNAs were detected using microarrays⁴⁷. We have demonstrated that it is
423 technically feasible to generate mRNA data from low input biofluid samples. This is expected
424 to accelerate biomarker research in these fluids. Further efforts to profile and share the mRNA
425 and circRNA content in larger sample cohorts of biofluids, comparable to the exRNA Atlas
426 Resource for small RNAs, are necessary to move this scientific field forward.⁸

427 Our small RNA results confirm previous studies observing high miRNA concentration in tears¹³,
428 low mapping rates in CSF^{5,48} and low miRNA concentration in cell-free urine¹². A direct
429 comparison of the absolute number of detected miRNAs, mRNAs and circRNAs detected per
430 sample in our study with the numbers in published literature is hampered by the fact that the
431 absolute read count is dependent on the input volume of the biofluids, the RNA isolation kit
432 and library preparation method used, the sequencing depth and data-analysis settings (e.g.
433 mapping without mismatches, filtering of the data). In addition, different pre-analytical
434 variables when preparing the biofluid samples may also affect the sequencing results.
435 However, on a higher level, we can look into the most abundant miRNAs detected in specific
436 biofluids. The majority of the 10 most abundant miRNAs detected in 9 specific biofluids

437 reported by Godoy et al. are also detected amongst the most abundant miRNAs in the samples
438 from the discovery cohort (Supplementary Data 10)⁵.

439 We compared the mRNA results of the discovery cohort with these of the case/control
440 cohorts. Mapping rates for samples in the discovery cohort are in the same range for saliva,
441 sputum and seminal plasma. The mapping rates for CSF and urine are about 15% higher in the
442 case/control cohorts compared to the discovery cohort. These differences may be due to
443 different pre-analytical variables between both cohorts (collection tube, centrifugation speed
444 and the portion of urine collected (Supplementary Fig. 3A; Supplemental material and
445 methods).

446 In the discovery cohort on average 53% of all small RNA reads in saliva can be traced to
447 bacteria, perfectly in line with the average of 45.5% reads mapping to bacteria reported by
448 Yeri et al.⁶ Aqueous humor and CSF, although collected in a sterile setting and presumed to
449 be sterile, contain up to 11% of reads mapping to bacteria, in line with previous studies^{5,48}.
450 However, bacterial cultures of our two CSF samples were negative. As both CSF and aqueous
451 humor display a very low relative RNA content, the exogenous sequences may represent
452 bacterial contaminants introduced during the sample processing workflow. Contaminants can
453 derive from contaminated spin columns used during RNA purification³², enzymes produced in
454 microorganisms⁴⁹, or various environmental sources⁵⁰. Such contaminant signals are likely
455 underrepresented in samples with high concentration of endogenous exRNAs.

456 Although we collected a broad range of biofluids, only two samples per biofluid were studied,
457 limiting our ability to assess donor variability. The input volume for the RNA isolations in all
458 biofluids was set to 200 μ L and a volume-based comparison of the RNA content was made
459 among the biofluids. We did not explore if higher input volumes would result in higher RNA
460 yields in biofluids where this could have been possible (e.g. urine). We also note that the
461 results in Table 1 are impacted by biofluid input volume in the RNA purification, RNA input in
462 the sequencing library prep, and the sequencing depth.

463 Biofluid data normalization with synthetic spike-in controls is a unique and powerful approach
464 and reflects more accurately the biological situation compared to classic normalization
465 approaches where global differences on overall abundance are neutralized. For instance, the
466 relative mRNA concentration in sputum from COPD patients is higher than in sputum from
467 healthy donors. Typically, RNA sequencing data is subsampled or normalized based on the
468 library size before performing a differential expression analysis, resulting in an artificially more

469 balanced volcano plot, an overcorrection of the biological situation and a loss of information,
470 which is not the case when the data is normalized based on spike-in controls.

471 Our results highlighting tissues and cell types that contribute more specifically to a certain
472 biofluid compared to the other biofluids (Fig. 6A) can be used as a roadmap to formulate
473 hypotheses when initiating biomarker research. Not surprisingly, the RNA signal from prostate
474 is reflected in urine and seminal plasma. Both fluids can be collected in a non-invasive way
475 and may be of value to investigate further in prostate cancer patients. Of interest, the mRNA
476 concentration in seminal plasma is 1000-fold higher than in urine and seminal plasma contains
477 more unique mRNAs compared to urine, suggesting that the biomarker potential of seminal
478 plasma is higher. However, one should also be cautious in interpreting the tissue enrichment
479 results: while the RNA signal of breast seems relatively enriched in sweat, this biofluid has the
480 lowest RNA concentration. The limited number of detected mRNAs in sweat show overlap
481 with mRNAs related to secretion (MCL1 gene, SCGB2A2 gene, SCGB1D2 gene) that also appear
482 as markers in breast tissue.

483 The pancreatic tissue RNA signal appears to be enriched in pancreatic cyst fluid and a different
484 cell type composition is observed when both samples are deconvoluted using single cell RNA
485 sequencing data of pancreatic cell types (Fig. 6B). Pancreatic cyst fluid was collected in these
486 donors to investigate a cystic lesion in the pancreas. The routine cytological analysis of these
487 fluid samples was inconclusive at the moment of sample collection. By following up both
488 patients, we discovered that the first patient developed a walled off necrosis collection after
489 necrotizing pancreatitis. The incipient high fraction of activated stellate cells in the first cyst
490 fluid sample may have been an indication pointing towards the inflammation and necrosis
491 that finally occurred. The second patient was diagnosed with a side-branch intra papillary
492 mucinous neoplasia, probably reflected by the relative high fraction of acinar cells. Pancreatic
493 cysts are often detected on abdominal imaging, resulting in a diagnostic and treatment
494 dilemma. Furthermore, pancreatic cysts represent a broad group of lesions, ranging from
495 benign to malignant entities. The main challenge in their management is to accurately predict
496 the malignant potential and to determine the risk to benefit of a surgical resection⁵¹. Our
497 results show that the cellular contribution to the RNA content of pancreatic cyst fluids can be
498 estimated through deconvolution and that these results may be associated with clinical
499 phenotypes. Larger cohorts are necessary to investigate the clinical potential of this approach

500 and pancreatic tumor cells may also need to be added to the reference set with single cell
501 RNA sequencing data to improve the accuracy of the prediction.

502 In addition to linear mRNA transcripts, we also explored the circular RNA content in biofluids.
503 CircRNAs are a growing class of non-coding RNAs and a promising RNA biotype to investigate
504 in the liquid biopsy setting, as they are presumed to be less prone to degradation compared
505 to linear forms⁵². The circRNA fraction in tissues has previously been reported and is in line
506 with our findings⁵³. In our study, we demonstrated that for genes that produce both circRNAs
507 and linear mRNAs, the circRNAs are more abundant than the linear forms in biofluids. Further
508 assessment of the biomarker potential of circRNAs in biofluids require dedicated library
509 preparation methods with circRNA enrichment.

510 In conclusion, The Human Biofluid RNA Atlas provides a systematic and comprehensive
511 comparison of the extracellular RNA content in 20 different human biofluids. The results
512 presented here may serve as a valuable resource for future biomarker studies.

513

514

515

516 **Material and methods**

517 ***Donor material, collection and biofluid preparation procedure***

518 Sample collection for the discovery cohort and sputum collection for the case/control cohort
519 was approved by the ethics committee of Ghent University Hospital, Ghent, Belgium (no.
520 B670201734450) and written informed consent was obtained from all donors according to the
521 Helsinki declaration. Breast milk, colostrum, plasma, serum, sputum, seminal plasma, sweat,
522 stool, tears and urine were obtained in healthy volunteers. All other biofluids were collected
523 from non-oncological patients.

524 The collection of two case series of each 12 cases and 12 control samples was approved by
525 the Masaryk Memorial Cancer Institute, Brno, Czech Republic (no. 14-08-27-01 and no.
526 MOU190814). Urine was collected in healthy donors and muscle-invasive bladder cancer
527 patients; CSF was collected in hydrocephalus patients and glioblastoma patients.

528 Collection of saliva samples in 12 healthy donors and in patients with diabetes mellitus for the
529 case/control cohort was approved by the ethics committee of the Medical University of
530 Vienna, Vienna, Austria (no. 2197/2015). Written informed consent was obtained from all
531 donors. The demographic and clinical patient information is provided in Supplementary Table
532 1. Detailed information on the sample collection per biofluid is provided in Supplementary
533 Note 2. All samples, except tear fluid, plasma and serum, were centrifuged at 2000 g (rcf) for
534 10 minutes without brake at room temperature. All samples were processed within 2 hours
535 after collection. The cell-free supernatant was carefully pipetted into 2 mL LoBind tubes
536 (Eppendorf LoBind microcentrifuge tubes, Z666556-250EA) and stored at -80 °C.

537 ***RNA isolation and gDNA removal***

538 *RNA isolation from all biofluids, except tears*

539 In the discovery cohort, two RNA isolations per biofluid and per sample were simultaneously
540 performed by two researchers (E.V.E. and E.H.). In the end, RNA obtained from both RNA
541 isolations was pooled per biofluid and per sample and this pooled RNA was used as starting
542 material for both library preparations. Hence, small RNA and mRNA capture sequencing on
543 the discovery cohort were performed on the same batch of RNA. In the case/control cohorts,
544 one RNA isolation was performed per sample and the RNA was used as starting material for
545 mRNA capture sequencing.

546 RNA was isolated with the miRNeasy Serum/Plasma Kit (Qiagen, Hilden, Germany, 217184)
547 according to the manufacturer's instructions. An input volume of 200 μ L was used for all
548 samples, except for tear fluid, and total RNA was eluted in 12 μ L of RNase-free water. Tear
549 fluid was collected with Schirmer strips and RNA was isolated directly from the strips (see
550 further). Per 200 μ L biofluid input volume, 2 μ L Sequin spike-in controls (Garvan Institute of
551 Medical Research) and 2 μ L RNA extraction Control (RC) spike-ins (Integrated DNA
552 Technologies)⁵⁴ were added to the lysate for TruSeq RNA Exome Library Prep sequencing and
553 TruSeq Small RNA Library Prep sequencing, respectively. Details on the spike-in controls are
554 available in the Supplementary Note 1.

555 Briefly, 2 μ L External RNA Control Consortium (ERCC) spike-in controls (ThermoFisher
556 Scientific, Waltham, MA, USA, 4456740), 2 μ L Library Prep Control (LP) spike-ins (Integrated
557 DNA Technologies)⁵⁵, 1 μ L HL-dsDNase and 1.6 μ L reaction buffer were added to 12 μ L RNA
558 eluate, and incubated for 10 min at 37 °C, followed by 5 min at 55 °C. Per biofluid and per
559 donor the RNA after gDNA removal was pooled. RNA was stored at -80 °C and only thawed on
560 ice immediately before the start of the library prep. Multiple freeze/thaw cycles did not occur.

561 *RNA isolation from tear fluid*

562 Tear fluid was collected in 8 healthy donors with Schirmer strips (2 strips per eye per donor),
563 as previously described^{56,57}. RNA was isolated within two hours after tear collection with the
564 miRNeasy Serum/Plasma Kit (Qiagen, Hilden, Germany, 217184), starting from one 2 mL tube
565 containing each 4 Schirmer strips. The same reagent volumes as suggested by the
566 manufacturer for a 200 μ L input volume were used. Throughout the RNA isolation protocol,
567 the two final RNA samples each result from 4 tear fluid samples (each containing the 4 strips
568 of a single donor) that were pooled in a two-step method. First, the upper aqueous phase of
569 two tear fluid samples was put together (in step 8 of the RNA isolation protocol). Second, the
570 RNA eluate of these two samples was pooled into the final RNA that was used as input for the
571 library prep (in step 15 of the RNA isolation protocol).

572 ***TruSeq RNA Exome library prep sequencing***

573 Messenger RNA capture based libraries were prepared starting from 8.5 μ L DNase treated and
574 spike-in supplemented RNA eluate using the TruSeq RNA Exome Library Prep Kit (Illumina, San
575 Diego, CA, USA). Each sample underwent individual enrichment according to the
576 manufacturer's protocol. The quality and yield of the prepared libraries were assessed using
577 a high sensitivity Small DNA Fragment Analysis Kit (Agilent Technologies, Santa Clara, CA, USA)

578 according to manufacturer's instructions. The libraries were quantified using qPCR with the
579 KAPA Library Quantification Kit (Roche Diagnostics, Diegem, Belgium, KK4854) according to
580 manufacturer's instructions. Based on the qPCR results, equimolar library pools were
581 prepared.

582 Paired-end sequencing was performed on a NextSeq 500 instrument using a high output v2
583 kit (Illumina, San Diego, CA, USA) with a read length of 75 nucleotides to an average
584 sequencing depth of 11 million read pairs in the discovery cohort, 16.8 million read pairs in
585 the sputum case/control cohorts, 15.4 million read pairs in the urine case/control cohort, 15
586 million read pairs in the CSF case/control cohort and 18.8 million read pairs in the saliva
587 case/control cohort. Samples from the discovery cohort were randomly assigned over two
588 pools and sequenced with a loading concentration of 1.2 pM (5% PhiX) and 1.6 pM (5% PhiX),
589 respectively. Urine, CSF and saliva samples from the case/control cohorts were loaded in 3
590 separate runs at 2 pM (2% PhiX) and sputum samples from the case/control cohorts were
591 loaded at 1.6 pM (5% PhiX).

592 ***TruSeq Small RNA library prep sequencing***

593 Small RNA libraries were prepared starting from 5 μ L DNase treated and spike-in
594 supplemented RNA eluate using a TruSeq Small RNA Library Prep Kit (Illumina, San Diego, CA,
595 USA) according to the manufacturer's protocol with two minor modifications(1). The RNA 3'
596 adapter (RA3) and the RNA 5' adapter (RA5) were 4-fold diluted with RNase-free water(2) and
597 the number of PCR cycles was increased to 16.

598 First, a volume-based pool of all 46 samples of the discovery cohort was sequenced. After PCR
599 amplification, quality of libraries was assessed using a high sensitivity Small DNA Fragment
600 Analysis Kit (Agilent Technologies, Santa Clara, CA, USA) according to manufacturer's
601 instructions. Size selection of the pooled samples was performed using 3% agarose dye-free
602 marker H cassettes on a Pippin Prep (Sage Science, Beverly, MA, USA) following
603 manufacturer's instructions with a specified collection size range of 125–163 bp. Libraries
604 were further purified and concentrated by ethanol precipitation, resuspended in 10 μ L of
605 10 mM tris-HCl (pH = 8.5) and quantified using qPCR with the KAPA Library Quantification Kit
606 (Roche Diagnostics, Diegem, Belgium, KK4854) according to manufacturer's instructions. The
607 pooled library was quality controlled via sequencing at a concentration of 1.7 pM with 35%
608 PhiX on a NextSeq 500 using a mid-output v2 kit (single-end 75 nucleotides, Illumina, San
609 Diego, CA, USA), resulting in an average sequencing depth of 1 million reads, ranging from

610 3341 reads to 14 million reads. Twenty-three samples with less than 200 000 reads were
611 assigned to a low concentrated pool, 23 samples with more than 17 million reads were
612 assigned to a highly concentrated pool. Based on the read numbers from the mid output run,
613 two new equimolar pools were prepared, purified and quantified as described higher. Both
614 re-pooled libraries were then sequenced at a final concentration of 1.7 pM with 25% PhiX on
615 a NextSeq 500 using a high output v2 kit (single-end, 75 nucleotides, Illumina, San Diego, CA,
616 USA), resulting in an average sequencing depth of 9 million reads (range 817 469 – 41.7 million
617 reads).

618 ***RT-qPCR***

619 To validate findings observed in the RNA sequencing data, we performed a targeted mRNA
620 and circRNA expression profiling with RT-qPCR for 8 differentially expressed mRNAs in sputum
621 (COPD versus healthy control) and for the 10 most abundant circRNAs in sputum. As reference
622 RNAs for normalization purposes, we selected Sequin spikes stably detected in all samples
623 based on the available RNA sequencing data. The assays to measure mRNA, circRNA and
624 Sequin spike expression were custom designed using primerXL⁵⁸ (Supplementary Data 9) and
625 purchased from Integrated DNA Technologies, Inc. (Coralville, USA).

626 For cDNA synthesis, 5 µl of total RNA was reverse transcribed using the iScript Advanced cDNA
627 Synthesis Kit (BioRad, USA) in a 10 µL volume. 5 µL of cDNA was pre-amplified in a 12-cycle
628 PCR reaction using the Sso Advanced PreAmp Supermix (Bio-Rad, USA) in a 50 µL reaction.
629 Pre-amplified cDNA was diluted (1:8) and 2 µL was used as input for a 45-cycle qPCR reaction,
630 quantifying 8 mRNAs and 10 circRNAs of interest with the SsoAdvanced™ Universal SYBR
631 Green Supermix (BioRad, USA). All reactions were performed in 384-well plates on the
632 LightCycler480 instrument (Roche) in a 5 µL reaction volume using 250 nM primer
633 concentrations. Cq-values were determined with the LightCycler®480 Software (release 1.5.0,
634 Roche) with the “Abs Quant/2nd Derivative Max” method.

635 The geNorm analysis to select the optimal number of reference targets was performed using
636 Biogazelle’s qbase+ software (www.qbaseplus.com) using log2-transformed RNA count data.
637 We observed medium reference target stability (average geNorm M ≤ 1.0) with an optimal
638 number of reference targets in this experimental situation of two (geNorm V < 0.15 when
639 comparing a normalization factor based on the two or three most stable targets). As such, the
640 optimal normalization factor can be calculated as the geometric mean of reference targets
641 R2_150 and R2_65. These Sequin spike RNAs were considered as reference RNAs.

642 **Data analysis**

643 *Processing TruSeq RNA Exome sequencing data*

644 Read quality was assessed by running FastQC (v0.11.5) on the FASTQ files and reads shorter
645 than 35 nucleotides and with a quality (phred) score < 30 were removed. The reads were
646 mapped with STAR (v2.6.0). Mapped reads were annotated by matching genomic coordinates
647 of each read with genomic locations of mRNAs (obtained from UCSC GRCh38/hg38 and
648 Ensembl, v91) or by matching the spike-in sequences. Picard (v2.18.5) was used for duplicate
649 removal. HTSeq (v0.9.1) was used for quantification of PCR deduplicated reads. A cut-off for
650 filtering noisy genes was set based on historic data to remove noisy genes. Using a threshold
651 of 4 counts, at least 95% of the single positive replicate values are filtered out. A table with
652 the read count of mRNAs per sample is provided in Supplementary Data 6.

653 *Processing TruSeq Small RNA sequencing data*

654 Adaptor trimming was performed using Cutadapt (v1.8.1) with a maximum error rate of 0.15.
655 Reads shorter than 15 nts and those in which no adaptor was found were discarded. For
656 quality control the FASTX-Toolkit (v0.0.14) was used, a minimum quality score of 20 in at least
657 80% of nucleotides was applied as a cutoff. The reads were mapped with Bowtie (v1.1.2)
658 without allowing mismatches. Mapped reads were annotated by matching genomic
659 coordinates of each read with genomic locations of miRNAs (obtained from miRBase, v22) and
660 other small RNAs (obtained from UCSC GRCh38/hg38 and Ensembl, v91) or by matching the
661 spike-in sequences. Reads assigned as “not annotated” represent uniquely mapped reads that
662 could not be classified in one of the small RNA biotype groups. As for the mRNA data, genes
663 with fewer than 4 counts were filtered out. A table with the read count of miRNAs per sample
664 is provided in Supplementary Data 7.

665 *Exogenous RNA characterization*

666 The exogenous RNA content in the mRNA data was assessed using the MetaMap pipeline⁵⁹.
667 Briefly, all reads were mapped to the human reference genome (hg38) using STAR (v2.5.2)⁶⁰.
668 Unmapped reads were subsequently subjected to metagenomic classification using CLARK-S
669 (v1.2.3)⁶¹. Reads were summed across all bacterial species.

670 The exogenous RNA content in the small RNA data was assessed using the exceRpt small RNA-
671 seq pipeline (v4.6.2) in the Genboree workbench with default settings⁶². Briefly, after adapter
672 trimming, read quality was assessed by FASTQC (v0.11.2). A minimum quality score of 20 in at
673 least 80% of nucleotides was applied as cutoff. The minimum read length after adapter

674 trimming was set to 18 nucleotides. Reads were first mapped to the custom spike-in
675 sequences using Bowtie2 (v2.2.6), followed by mapping the unmapped reads with STAR
676 (v2.4.2a) to UniVec contaminants and human ribosomal (rRNA) sequences to exclude them
677 before mapping (also with STAR) to the following databases: miRbase (v21), gtRNAdb,
678 piRNABank, GeneCode version 24 (hg38) and circBase (version last updated in July 2017). A
679 single mismatch was allowed during mapping to the human genome. Unmapped reads were
680 then mapped with STAR to exogenous miRNAs and rRNAs. In the end, the remaining
681 unmapped reads were mapped to the genomes of all sequenced species in Ensembl and NCBI.
682 No mismatches were allowed during exogenous alignment. Raw read counts obtained from
683 the Genboree workbench were further analyzed in R (v3.5.1) making use of tidyverse (v1.2.1).

684 *Circular RNA detection and circular/linear ratio determination*

685 Only TruSeq RNA Exome reads passing quality control (base calling accuracy of $\geq 99\%$ in at
686 least 80% of the nucleotides in both mates of a pair) were included in this analysis. Clumpify
687 dedupe (v38.26) was used to remove duplicates in paired-end mode (2 allowed substitutions,
688 kmer size of 31 and 20 passes). We used a two-step mapping strategy to identify forward
689 splice (further referred to as linear) junction reads and backsplice junction reads. First, reads
690 were aligned with TopHat2 (v2.1.0) to the GRCh38/hg38 reference genome (Ensembl, v91).
691 Micro-exons were included, a minimum anchor length of 6 nucleotides was required, and up
692 to two mismatches in the anchor region were allowed. The resulting output contains linear
693 junction information. Secondly, unmapped reads from the first mapping strategy were
694 realigned with TopHat2 (v2.1.0) to the same reference, but this time with the fusion search
695 option that can align reads to potential fusion transcripts. Processing the fusion search output
696 with CIRCexplorer2 parse (v2.3.3) results in backsplice junction information. Junction read
697 counts obtained with the mapping strategies described above were used as a measure for the
698 relative level of linear and circular RNA in each sample. Only genes with at least one detected
699 backsplice junction were considered. Junctions that could be part of both linear and circular
700 transcripts (ambiguous junctions) were filtered out. As there is currently no consensus on how
701 to calculate the circular to linear ratio (CIRC/LIN), we decided to calculate the ratio in two
702 different ways (Supplementary Fig. 8). The circRNA fraction is defined as $100 * \text{CIRC} / (\text{CIRC} + \text{LIN})$.
703 The first method (referred to as “backsplice junction-level method”) zooms in on each
704 particular backsplice junction. CIRC was defined as the backsplice junction read count of one
705 particular backsplice junction. LIN was defined as the average read count of all junctions

706 flanking the backsplice junction of interest. The second method (referred to as “gene-level
707 method”) considers all backsplice junctions within a given gene. CIRC was defined as the
708 average number of backsplice junction reads for a given gene. LIN was defined as the average
709 number of linear junction reads for a given gene. For both methods, CIRC > 3 was used as a
710 cut-off for filtering noisy backsplice junctions. To enable a comparison of the circular/linear
711 genic ratios in biofluids with those of tissues, the mRNA capture sequencing FASTQ files of 16
712 cancerous tissue types (34 samples in total) were downloaded from the MiOncoCirc database
713 (dbGaP Study Accession phs000673.v3.p1)²⁴. A list with the downloaded samples is attached
714 in Supplementary Table 2. A table with the read count of backsplice junctions per sample is
715 provided in Supplementary Data 8.

716 *Assessment of tissue and cell contribution to biofluid exRNA*

717 Using total RNA-sequencing data from 27 normal human tissue types and 5 immune cell types
718 from peripheral blood from the RNA Atlas³⁵, we created gene sets containing marker genes
719 for each individual entity (Supplementary Data 4). We removed redundant tissues and cell
720 types from the original RNA Atlas (e.g. granulocytes and monocytes were present twice; brain
721 was kept and specific brain sub-regions such as cerebellum, frontal cortex, occipital cortex and
722 parietal cortex were removed) and we used genes where at least one tissue or cell type had
723 expression values greater or equal to 1 TPM normalized counts. A gene was considered to be
724 a marker if its abundance was at least 5 times higher in the most abundant sample compared
725 to the others. For the final analysis, only tissues and cell types with at least 3 markers were
726 included, resulting in 26 tissues and 5 immune cell types.

727 Gene abundance read counts from the biofluids were normalized using Sequin spikes as size
728 factors in DESeq2 (v1.22.2). For all marker genes within each gene set, we computed the log₂
729 fold changes between the median read count of a biofluid sample pair versus the median read
730 count of all other biofluids. The median log₂ fold change of all markers in a gene set was
731 selected, followed by z-score transformation over all biofluids (Fig. 7). For visualization
732 purposes, only tissues and cell types with a z-score $\geq |1|$ in at least one biofluid were used.

733 *Cellular deconvolution of pancreatic cyst fluid samples*

734 To build the reference matrix for the computational deconvolution of pancreatic cyst fluid
735 samples, single cell RNA sequencing data of 10 pancreatic cell types³⁴ was processed with the
736 statistical programming language R (v3.6.0). For each gene, the mean count across all
737 individual cells from each cell type was computed. Next, this reference matrix was normalized

738 using the trimmed means of M values (TMM) with the edgeR package (v3.26.4)⁶³⁶⁴. Limma-
739 voom (v3.40.2)⁶⁵ was used for subsequent differential gene expression analysis and those
740 genes with an absolute fold change greater or equal to 2 and an adjusted p-value < 0.05
741 (Benjamini-Hochberg) were retained as markers⁶⁶. Finally, using these markers and both the
742 pancreatic cyst fluid samples and the reference matrix described above, the cell type
743 proportions were obtained through computational deconvolution using non-negative least
744 squares (nnls package; v1.4)⁶⁷⁶⁸.

745 *Differential expression analysis in case/control cohorts*

746 Further processing of the count tables was done with R (v3.5.1) making use of tidyverse
747 (v1.2.1). Gene expression read counts from the biofluids were normalized using Sequin spikes
748 as size factors in DESeq2 (v1.20.0)⁶⁹. To assess the biological signal in the case/control cohorts,
749 we performed differential expression analysis between the patients and control groups using
750 DESeq2 (v1.20.0). Genes were considered differentially expressed when the absolute log₂ fold
751 change > 1 and at q < 0.05.

752

753 **Data availability**

754 The raw RNA-sequencing data have been deposited at the European Genome-phenome
755 Archive (EGA) under accession number EGAS00001003917. All spike-normalized sequencing
756 data can be readily explored in the interactive web-based application R2: Genomics analysis
757 and visualization platform (<http://r2.amc.nl>), and via a dedicated accessible portal
758 (<http://r2platform.com/HumanBiofluidRNAAtlas>). This portal allows the analysis and
759 visualization of mRNA, circRNA and miRNA abundance, as illustrated in Supplementary Fig. 11.
760 All samples can be used for correlation, principle component, and gene set enrichment
761 analyses, and many more. All other data are available within the article and supplementary
762 information.

763

764 **Code availability**

765 The R scripts to reproduce the analyses and plots reported in this paper are available from the
766 corresponding authors upon request.

767

768 **Acknowledgments**

769 E.H. is a recipient of a grant of the Fund for Scientific Research Flanders (FWO). F.A.C. and
770 A.M. are supported by Special Research Fund (BOF) scholarship of Ghent University
771 (BOF.DOC.2017.0026.01, BOF.DOC.2019.0047.01). A.G. is senior clinical researcher of the
772 Fund for Scientific Research Flanders (FWO; 1805718N). This research is partly funded by
773 “RNA-MAGIC”, “LNCCA” Concerted Research Actions of Ghent University (BOF19/GOA/008,
774 BOF16/GOA/023), “Kom Op Tegen Kanker” and by the Czech Ministry of Health grants (NV18-
775 03-00398, NV18-03-00360). We thank Tim Mercer for providing the Sequin spikes.
776 The results published here are in part based upon data generated by the Clinical Sequencing
777 Exploratory Research (CSER) consortium established by the NHGRI. Funding support was
778 provided through cooperative agreements with the NHGRI and NCI through grant numbers
779 U01 HG006508 (Exploring Cancer Medicine for Sarcoma and Rare Cancers). Information about
780 CSER and the investigators and institutions who comprise the CSER consortium can be found
781 at <http://www.genome.gov/27546194>.

782

783 **Contributions**

784 J.V. and P.M. conceived and supervised the project; E.H., K.V., N.Y., E.V.E., J.N. designed and
785 performed the experiments; E.H., A.M., J.A. and F.A.C. analyzed the data; L.S. and S.K.
786 performed analysis using the MetaMap pipeline; G.S. and S.K. contributed technical support
787 and resources; E.H., A.G., P.H., P.J., G.B., K.B., T.M., T.D., V.N., C.V.C., K.R., E.R., D.H., K.T.,
788 O.S., C.N. collected samples; S.L. designed RT-qPCR primers; E.H., P.M and J.V. wrote the
789 paper; J.K. developed dedicated tools to analyze RNA atlas data and results and
790 implemented them in the online portal R2. All authors contributed to manuscript editing
791 and approved the final draft.

792

793

794 References

- 795 1. Weiland, M., Gao, X.-H., Zhou, L. & Mi, Q.-S. Small RNAs have a large impact: circulating
796 microRNAs as biomarkers for human diseases. *RNA Biol.* **9**, 850–859 (2012).
- 797 2. Max, K. E. A. *et al.* Human plasma and serum extracellular small RNA reference profiles and their
798 clinical utility. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E5334–E5343 (2018).
- 799 3. Freedman, J. E. *et al.* Diverse human extracellular RNAs are widely detected in human plasma.
800 *Nat. Commun.* **7**, 11106 (2016).
- 801 4. Yuan, T. *et al.* Plasma extracellular RNA profiles in healthy and cancer patients. *Sci. Rep.* **6**, 19413
802 (2016).
- 803 5. Godoy, P. M. *et al.* Large Differences in Small RNA Composition Between Human Biofluids. *Cell*
804 *Rep.* **25**, 1346–1358 (2018).
- 805 6. Yeri, A. *et al.* Total Extracellular Small RNA Profiles from Plasma, Saliva, and Urine of Healthy
806 Subjects. *Sci. Rep.* **7**, 44061 (2017).
- 807 7. Ferrero, G. *et al.* Small non-coding RNA profiling in human biofluids and surrogate tissues from
808 healthy individuals: description of the diverse and most represented species. *Oncotarget* **9**,
809 (2018).
- 810 8. Murillo, O. D. *et al.* exRNA Atlas Analysis Reveals Distinct Extracellular RNA Cargo Types and
811 Their Carriers Present across Human Biofluids. *Cell* **177**, 463–477.e15 (2019).
- 812 9. Srinivasan, S. *et al.* Small RNA Sequencing across Diverse Biofluids Identifies Optimal Methods for
813 exRNA Isolation. *Cell* **177**, 446–462.e16 (2019).
- 814 10. Fehlmann, T., Ludwig, N., Backes, C., Meese, E. & Keller, A. Distribution of microRNA biomarker
815 candidates in solid tissues and body fluids. *RNA Biol.* **13**, 1084–1088 (2016).
- 816 11. Umu, S. U. *et al.* A comprehensive profile of circulating RNAs in human serum. *RNA Biol.* **15**, 242–
817 250 (2018).
- 818 12. El-Mogy, M. *et al.* Diversity and signature of small RNA in different bodily fluids using next
819 generation sequencing. *BMC Genomics* **19**, 408 (2018).
- 820 13. Weber, J. A. *et al.* The microRNA spectrum in 12 body fluids. *Clin. Chem.* **56**, 1733–1741 (2010).
- 821 14. Welch, R. A., Shaw, M. K. & Welch, K. C. Amniotic fluid LPCAT1 mRNA correlates with the
822 lamellar body count. *J. Perinat. Med.* **44**, 531–532 (2016).
- 823 15. Marzioni, M. *et al.* PDX-1 mRNA expression in endoscopic ultrasound-guided fine needle
824 cytoaspirate: perspectives in the diagnosis of pancreatic cancer. *Dig. Liver Dis. Off. J. Ital. Soc.*
825 *Gastroenterol. Ital. Assoc. Study Liver* **47**, 138–143 (2015).
- 826 16. Tian, Y., Li, L., Zhang, F. & Xu, J. Seminal plasma HSPA2 mRNA content is associated with semen
827 quality. *J. Assist. Reprod. Genet.* **33**, 1079–1084 (2016).
- 828 17. Oreo, K. M. *et al.* Sputum ADAM8 expression is increased in severe asthma and COPD. *Clin. Exp.*
829 *Allergy J. Br. Soc. Allergy Clin. Immunol.* **44**, 342–352 (2014).
- 830 18. Maker, A. V. *et al.* Cyst Fluid Biosignature to Predict Intraductal Papillary Mucinous Neoplasms of
831 the Pancreas with High Malignant Potential. *J. Am. Coll. Surg.* **228**, 721–729 (2019).
- 832 19. Herring, E., Kanaoka, S., Tremblay, E. & Beaulieu, J.-F. A Stool Multitarget mRNA Assay for the
833 Detection of Colorectal Neoplasms. *Methods Mol. Biol. Clifton NJ* **1765**, 217–227 (2018).
- 834 20. Bazzell, B. G. *et al.* Human Urinary mRNA as a Biomarker of Cardiovascular Disease. *Circ.*
835 *Genomic Precis. Med.* **11**, e002213 (2018).
- 836 21. Bahn, J. H. *et al.* The landscape of microRNA, Piwi-interacting RNA, and circular RNA in human
837 saliva. *Clin. Chem.* **61**, 221–230 (2015).
- 838 22. Liu, B. *et al.* Characterization of tissue-specific biomarkers with the expression of circRNAs in
839 forensically relevant body fluids. *Int. J. Legal Med.* **133**, 1321–1331 (2019).
- 840 23. Li, S. *et al.* exoRBase: a database of circRNA, lncRNA and mRNA in human blood exosomes.
841 *Nucleic Acids Res.* **46**, D106–D112 (2018).
- 842 24. Vo, J. N. *et al.* The Landscape of Circular RNA in Cancer. *Cell* **176**, 869–881.e13 (2019).

- 843 25. Kölling, M. *et al.* Circular RNAs in Urine of Kidney Transplant Patients with Acute T Cell-Mediated
844 Allograft Rejection. *Clin. Chem.* (2019) doi:10.1373/clinchem.2019.305854.
- 845 26. Giraldez, M. D. *et al.* Phospho-RNA-seq: a modified small RNA-seq method that reveals
846 circulating mRNA and lncRNA fragments as potential biomarkers in human plasma. *EMBO J.* **38**,
847 (2019).
- 848 27. Zhou, Z. *et al.* Extracellular RNA in a single droplet of human serum reflects physiologic and
849 disease states. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 19200–19208 (2019).
- 850 28. Everaert, C. *et al.* Performance assessment of total RNA sequencing of human biofluids and
851 extracellular vesicles. *Sci. Rep.* **9**, (2019).
- 852 29. Metzenmacher *et al.* Plasma Next Generation Sequencing and Droplet Digital-qPCR-Based
853 Quantification of Circulating Cell-Free RNA for Noninvasive Early Detection of Cancer. *Cancers*
854 **12**, 353 (2020).
- 855 30. Li, X., Yang, L. & Chen, L.-L. The Biogenesis, Functions, and Challenges of Circular RNAs. *Mol. Cell*
856 **71**, 428–442 (2018).
- 857 31. Memczak, S., Papavasileiou, P., Peters, O. & Rajewsky, N. Identification and Characterization of
858 Circular RNAs As a New Class of Putative Biomarkers in Human Blood. *PLoS One* **10**, e0141214
859 (2015).
- 860 32. Heintz-Buschart, A. *et al.* Small RNA profiling of low biomass samples: identification and removal
861 of contaminants. *BMC Biol.* **16**, 52 (2018).
- 862 33. Liu, F., Ma, R., Wang, Y. & Zhang, L. The Clinical Importance of *Campylobacter concisus* and
863 Other Human Hosted *Campylobacter* Species. *Front. Cell. Infect. Microbiol.* **8**, 243 (2018).
- 864 34. Baron, M. *et al.* A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals
865 Inter- and Intra-cell Population Structure. *Cell Syst.* **3**, 346–360.e4 (2016).
- 866 35. Lorenzi, L. *et al.* The RNA Atlas, a single nucleotide resolution map of the human transcriptome.
867 *bioRxiv* (2019) doi:10.1101/807529.
- 868 36. Demedts, I. K. *et al.* Accumulation of dendritic cells and increased CCL20 levels in the airways of
869 patients with chronic obstructive pulmonary disease. *Am. J. Respir. Crit. Care Med.* **175**, 998–
870 1005 (2007).
- 871 37. Bracke, K. R. *et al.* Cigarette smoke-induced pulmonary inflammation and emphysema are
872 attenuated in CCR6-deficient mice. *J. Immunol. Baltim. Md 1950* **177**, 4350–4359 (2006).
- 873 38. Karmouty-Quintana, H. *et al.* Adenosine A2B receptor and hyaluronan modulate pulmonary
874 hypertension associated with chronic obstructive pulmonary disease. *Am. J. Respir. Cell Mol.*
875 *Biol.* **49**, 1038–1047 (2013).
- 876 39. Singh Patidar, B. *et al.* Adenosine Metabolism in COPD: A Study on Adenosine Levels, 5'-
877 Nucleotidase, Adenosine Deaminase and Its Isoenzymes Activity in Serum, Lymphocytes and
878 Erythrocytes. *COPD* **15**, 559–571 (2018).
- 879 40. Stankovic, M., Nikolic, A., Nagorni-Obradovic, L., Petrovic-Stanojevic, N. & Radojkovic, D. Gene-
880 Gene Interactions Between Glutathione S-Transferase M1 and Matrix Metalloproteinases 1, 9,
881 and 12 in Chronic Obstructive Pulmonary Disease in Serbians. *COPD* **14**, 581–589 (2017).
- 882 41. Saugstad, J. A. *et al.* Analysis of extracellular RNA in cerebrospinal fluid. *J. Extracell. Vesicles* **6**,
883 1317577 (2017).
- 884 42. Holyoake, A. *et al.* Development of a multiplex RNA urine test for the detection and stratification
885 of transitional cell carcinoma of the bladder. *Clin. Cancer Res. Off. J. Am. Assoc. Cancer Res.* **14**,
886 742–749 (2008).
- 887 43. Eckstein, M. *et al.* mRNA-Expression of KRT5 and KRT20 Defines Distinct Prognostic Subgroups of
888 Muscle-Invasive Urothelial Bladder Cancer Correlating with Histological Variants. *Int. J. Mol. Sci.*
889 **19**, (2018).
- 890 44. Murakami, T. *et al.* Bladder cancer detection by urinary extracellular vesicle mRNA analysis.
891 *Oncotarget* **9**, 32810–32821 (2018).
- 892 45. Lin, H., Zhou, Q., Wu, W. & Ma, Y. Midkine Is a Potential Urinary Biomarker for Non-Invasive
893 Detection of Bladder Cancer with Microscopic Hematuria. *OncoTargets Ther.* **Volume 12**, 11765–
894 11775 (2020).

- 895 46. Chen, T. *et al.* CD163, a novel therapeutic target, regulates the proliferation and stemness of
896 glioma cells via casein kinase 2. *Oncogene* **38**, 1183–1199 (2019).
- 897 47. Zhang, L. *et al.* Salivary transcriptomic biomarkers for detection of resectable pancreatic cancer.
898 *Gastroenterology* **138**, 949-957.e1–7 (2010).
- 899 48. Waller, R. *et al.* Small RNA Sequencing of Sporadic Amyotrophic Lateral Sclerosis Cerebrospinal
900 Fluid Reveals Differentially Expressed miRNAs Related to Neural and Glial Activity. *Front.*
901 *Neurosci.* **11**, (2018).
- 902 49. Salter, S. J. *et al.* Reagent and laboratory contamination can critically impact sequence-based
903 microbiome analyses. *BMC Biol.* **12**, (2014).
- 904 50. Strong, M. J. *et al.* Microbial contamination in next generation sequencing: implications for
905 sequence-based analysis of clinical samples. *PLoS Pathog.* **10**, e1004437 (2014).
- 906 51. Farrell, J. J. Pancreatic Cysts and Guidelines. *Dig. Dis. Sci.* **62**, 1827–1839 (2017).
- 907 52. Jeck, W. R. & Sharpless, N. E. Detecting and characterizing circular RNAs. *Nat. Biotechnol.* **32**,
908 453–461 (2014).
- 909 53. Guo, J. U., Agarwal, V., Guo, H. & Bartel, D. P. Expanded identification and characterization of
910 mammalian circular RNAs. **14** (2014).
- 911 54. Locati, M. D. *et al.* Improving small RNA-seq by using a synthetic spike-in set for size-range
912 quality control together with a set for data normalization. *Nucleic Acids Res.* **43**, e89 (2015).
- 913 55. Hafner, M. *et al.* RNA-ligase-dependent biases in miRNA representation in deep-sequenced small
914 RNA cDNA libraries. *RNA N. Y. N* **17**, 1697–1712 (2011).
- 915 56. Green-Church, K. B., Nichols, K. K., Kleinholz, N. M., Zhang, L. & Nichols, J. J. Investigation of the
916 human tear film proteome using multiple proteomic approaches. *Mol. Vis.* **14**, 456–470 (2008).
- 917 57. Pieragostino, D. *et al.* Tear Film Steroid Profiling in Dry Eye Disease by Liquid Chromatography
918 Tandem Mass Spectrometry. *Int. J. Mol. Sci.* **18**, 1349 (2017).
- 919 58. Lefever, S. *et al.* High-throughput PCR assay design for targeted resequencing using primerXL.
920 *BMC Bioinformatics* **18**, 400 (2017).
- 921 59. Simon, L. M. *et al.* MetaMap: an atlas of metatranscriptomic reads in human disease-related
922 RNA-seq data. *GigaScience* **7**, (2018).
- 923 60. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinforma. Oxf. Engl.* **29**, 15–21
924 (2013).
- 925 61. Ounit, R. & Lonardi, S. Higher classification sensitivity of short metagenomic reads with CLARK-S.
926 *Bioinforma. Oxf. Engl.* **32**, 3823–3825 (2016).
- 927 62. Rozowsky, J. *et al.* exceRpt: A Comprehensive Analytic Platform for Extracellular RNA Profiling.
928 *Cell Syst.* **8**, 352-357.e3 (2019).
- 929 63. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis
930 of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
- 931 64. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential
932 expression analysis of digital gene expression data. *Bioinforma. Oxf. Engl.* **26**, 139–140 (2010).
- 933 65. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
934 microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
- 935 66. Benjamini, Yoav, H., Yosef. Controlling the False Discovery Rate: A Practical and Powerful
936 Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)*
937 **57**, 289–300 (1995).
- 938 67. Mullen, K. M. & Van Stokkum, I. H. M. nns: The Lawson-Hanson algorithm for non-negative least
939 squares (NNLS). R package version 1.4. <https://CRAN.R-project.org/package=nns>. (2012).
- 940 68. Avila Cobos, F., Vandesompele, J., Mestdagh, P. & De Preter, K. Computational deconvolution of
941 transcriptomics data from mixed cell populations. *Bioinforma. Oxf. Engl.* **34**, 1969–1979 (2018).
- 942 69. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
943 RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 944