

1 Cell Types of Origin in the Cell Free Transcriptome in Human Health and Disease

2
3 Sevahn K. Vorperian^{1,2}, Mira N. Moufarrej³, Tabula Sapiens Consortium⁴, Stephen R. Quake^{3,5,6*}

4 ¹Department of Chemical Engineering, Stanford University, Stanford CA, USA

5 ²ChEM-H, Stanford University, Stanford CA, USA

6 ³Department of Bioengineering, Stanford University, Stanford CA, USA

7 ⁴Members listed at the end

8 ⁵Department of Applied Physics, Stanford University, Stanford CA, USA

9 ⁶Chan-Zuckerberg Biohub, Stanford CA, USA

10 *steve@quake-lab.org

11 Abstract

12 Liquid biopsies using cell-free RNA (cfRNA) can noninvasively measure dynamic physiological changes
13 throughout the body. While there is much effort in the liquid biopsy field to determine disease tissue-of-
14 origin, pathophysiology occurs at the cellular level. Here, we show that it is possible to determine cell
15 type-of-origin from cfRNA by leveraging single cell transcriptomic atlases to perform computational
16 deconvolution. We derived cell type gene signatures by combining the whole-body single cell atlas
17 *Tabula Sapiens*, individual tissue single cell atlases, and bulk tissue atlases. Using deconvolution, we
18 identified cell types-of-origin in the healthy human cell-free transcriptome, including contributions from
19 multiple cell types in the brain, liver, lung, intestine, kidney, and pancreas in addition to hematopoietic
20 cell types. We further showed that it is possible not only to detect cell types implicated in the pathology of
21 chronic kidney disease (CKD) and Alzheimer's disease (AD), but also to measure changes in these cell
22 types as a function of disease state. Altogether, our results show that cfRNA measurements reflect cellular
23 contributions in health and disease from diverse tissue-specific cell types. These findings underscore the
24 resolution at which one can monitor pathophysiological changes and the broad potential prognostic utility
25 afforded by non-invasive transcriptomic measurement.
26

27 Introduction

28 Cell-free RNA (cfRNA) in blood plasma enables dynamic and longitudinal phenotypic insight
29 into diverse physiological conditions, spanning oncology and bone marrow transplantation¹, obstetrics^{2,3},
30 neurodegeneration⁴, and liver disease⁵. Liquid biopsies that measure cfRNA afford broad clinical utility
31 since cfRNA represents a mixture of transcripts that reflects the health status of multiple tissues.
32 However, several aspects about the physiologic origins of cfRNA including the contributing cell types-of-
33 origin remain unknown, and most current assays focus on tissue level contributions^{2-4,6,7}. Although
34 information about tissue-of-origin can provide insight into transcriptional changes at a disease site, it
35 would be even more powerful to incorporate knowledge from cellular pathophysiology which often forms
36 the basis of disease⁸. This would also more closely match the resolution afforded by invasive biopsy.

37 Single cell transcriptomics (scRNA-seq) enable insight into the heterogeneous cellular
38 transcriptional landscapes of tissues in health and disease⁹. Numerous scRNA-seq tissue atlases provide
39 powerful reference data for defining cell type specific gene profiles in the context of an individual tissue.
40 However, the starting set of cell types influences a differential expression analysis, which guides the
41 assignment of a gene as cell type specific. cfRNA originates from cell types across the human body.
42 Therefore, interpreting a measured gene in cfRNA as cell type specific relies on the completeness of
43 relevant atlases. The *Tabula Sapiens* (TSP) cell atlas¹⁰ from 14 tissues enables the most comprehensive
44

45 derivation of cell type specific gene profiles in the context of a single individual to date, all determined
46 with uniform methods and sequencing, and we used this resource for our deconvolution process. For cell
47 types originating from tissues absent from the draft TSP atlas, we derived specific gene profiles by
48 combining a given single tissue cell atlas with comprehensive bulk transcriptomic datasets, including the
49 Genotype-Tissue Expression (GTEx) project¹¹ and the Human Protein Atlas (HPA)¹².

50 In this work, we defined cell type specific gene profiles in the context of the whole body to
51 identify the cell types comprising the cf-transcriptome. First, we computationally deconvolved the cell
52 types-of-origin in the healthy human cf-transcriptome using the TSP cell atlas and individual scRNA-seq
53 tissue atlases. Next, we measured striking cfRNA changes associated with cell types implicated in chronic
54 kidney disease (CKD) and Alzheimer's disease (AD) that are consistent with observed clinical pathology.
55 Altogether, we demonstrate that it is possible to decompose the cf-transcriptome into distinct cell type
56 contributions even in the absence of a complete whole body single cell reference, and demonstrate that
57 cell type specific changes in disease can be measured noninvasively using cfRNA.

58

59 **Results**

60 **Deconvolution of cell type specific signals in the healthy cell free transcriptome**

61 We used published exome-enriched cf-transcriptome data¹ to characterize the landscape of cell
62 type specific signal in the plasma of healthy individuals (Fig. 1A). After eliminating low-quality cfRNA
63 samples (Fig. S1, Methods), we intersected the set of genes detected in healthy individuals ($n = 5$) with a
64 database of cell-type specific markers defined in context of the whole body¹³ with stringent expression
65 requirements (Fig. 1B, Methods). Marker genes for cell types originating from the blood, brain, and liver
66 were readily detected, as previously observed at tissue level^{1,3-5}. Kidney, GI track, and pancreas cell type
67 markers were additionally detected (Fig. 1B).

68 Given the robust detection of several cell types contributing to the cf-transcriptome, we then
69 deconvolved the fractions of cell-type specific RNA using TSP. We defined the cf-transcriptome as a
70 linear combination of cell type specific RNA contributions using a deconvolution method, nu-SVR,
71 originally developed to decompose bulk tissue transcriptomes into fractional cell type components^{14,15}
72 (fig. S2). This required specifying a basis matrix with a representative gene set (rows) that could
73 accurately and simultaneously resolve the distinct cell types (columns). To reduce multicollinearity, we
74 grouped transcriptionally similar cell types (Methods). We observed that the basis matrix appropriately
75 described cell types as most similar to others from the same organ compartment, where cell types
76 originating from the same compartment cluster together and correspond to the highest off-diagonal
77 similarity (Fig. 1C). We also confirmed that the defined basis matrix can correctly deconvolve cell type
78 specific RNA fractional contributions from several GTEx bulk tissue samples (fig. S3, S4, Supplementary
79 Note).

80 We then deconvolved the cell types-of-origin contributing to the plasma cf-transcriptome (Fig.
81 1D). We observed a large signal from hematopoietic cell types, as well as smaller, distinct transcriptional
82 contributions from tissue-specific cell types from the large and small intestine, lungs, and pancreas (Fig.
83 1D, fig. S5A, B). The highest cell type contributors were monocytes ($18.6 \pm 2.3\%$), platelets (13.6 ± 3.5
84 $\%$), erythrocytes and erythroid progenitors ($15.8 \pm 9.1\%$), and lymphocytes ($15.7 \pm 2.7\%$). There was
85 good pairwise similarity amongst all biological replicates ($r \geq 0.66$, fig. S5C). The predominant cell types
86 and their respective proportions we observe are generally consistent with recently published estimates for
87 serum cfRNA¹ and plasma cfDNA¹⁶. We also observed small fractional contributions from endothelial
88 cells, pancreatic cells, intestinal enterocytes, kidney epithelia, club cells, goblet cells, pancreatic acinar

89 cells, and other pancreatic cells (Fig. 1D), underscoring the contributions of non-hematopoietic cell types
90 to the cf-transcriptome.

91 Some cell types likely present in the plasma cf-transcriptome were not found in this
92 decomposition because the source tissues were absent from the TSP version 0.9. Deconvolution
93 performance yielded an elevated root mean square error (RMSE) (86 ± 5.2 CPM) and reduced Pearson
94 correlation (0.43 ± 0.08) compared to deconvolved GTEx tissues whose cell types were completely in the
95 basis matrix (fig. S3). To understand which cell type contributions might be absent from this present
96 analysis, we intersected the genes measured in cfRNA but absent from the basis matrix with tissue
97 specific genes annotated by the HPA transcriptomic atlas¹² (Methods). This identified brain, liver, testis,
98 skeletal muscle, and cardiac muscle as tissues with several reliably measured genes in the plasma cfRNA
99 (Fig. 1E, Methods) whose cell types were not found during systems-level deconvolution.

100 We then defined cell type specific gene profiles for these tissues in context of the whole body. To
101 do so, we leveraged individual tissue cell atlases¹⁷⁻¹⁹ but only considered cell types unique to a given
102 tissue (fig. S6, Methods). This formulation allowed us to apply bulk GTEx and HPA transcriptomic data
103 to ensure whole body specificity using stringent expression specificity constraints. First, we required a
104 given gene to be differentially expressed in a given cell type against all others within an individual tissue
105 cell atlas (Fig 2A, fig. S7A) (Methods). Second, we required high expression inequality across tissues
106 measured by the Gini coefficient²⁰ (Fig 2B, fig. S7B & fig. S8) (Methods). We validated the specificity of
107 a given gene profile to its corresponding cell type by comparing the aggregate expression of a given cell
108 type signature in its native tissue compared to that of the average across remaining GTEx tissues (Fig 2C,
109 fig. S7C). We uniformly observed a median fold change greater than one in the signature score of a cell
110 type gene profile in its native tissue relative to the mean expression in other tissues, confirming high
111 specificity.

112 Next, we estimated a signature score for each cell type in the cf-transcriptome using its specific
113 gene profile by summing the measured level for all included genes (Fig. 2D), and observed contributions
114 from multiple brain cell types, hepatocytes, and cardiomyocytes. Specifically, we measured a strong
115 signature score from excitatory neurons and a reduced signature score from inhibitory neurons. We also
116 observed strong signals from astrocytes, oligodendrocytes, and oligodendrocyte precursor cells. These
117 glial cells facilitate brain homeostasis, form myelin, and provide neuronal structure and support⁸.
118 Evidence of RNA transport across the blood brain barrier (BBB)²¹, BBB permeability²², and brain regions
119 in direct contact with the blood²³ help rationalize brain cell type signature detection in the cf-
120 transcriptome.

121 We additionally observed a strong hepatocyte signature score, which is consistent with their high
122 turnover rate and cellular mass²⁴, a small signal for atrial cardiomyocytes, and negligible signal from
123 ventricular cardiomyocytes, consistent with the low level of cardiomyocyte death in healthy adults²⁵ (Fig.
124 2D). These observations augment the resolution of previously observed brain-^{3,4}, liver-⁵, and heart²⁶-
125 specific genes reported to date in cfRNA.

126

127 **Plasma cfRNA measurement reflects cellular pathophysiology**

128 Cell type specific changes drive disease etiology⁸, and we asked whether cfRNA reflected
129 changes in pathological cell types. We considered trophoblasts in preeclampsia^{27,28}, proximal tubules in
130 CKD^{29,30}, and multiple brain cell types in AD^{17,31}. We utilized published cell atlases for the placenta^{32,33},
131 kidney³⁴, and brain¹⁷ to define cell specific gene profiles in context of the whole body following the
132 approach outlined above (fig. S7 and S9).

133 In pregnancy, extravillous trophoblast invasion is a stage in uteroplacental arterial
134 remodeling^{27,33}. Arterial remodeling occurs to ensure adequate maternal blood flow to the growing
135 fetus^{27,33} and is sometimes reduced in preeclampsia²⁷. Previously, the extravillous trophoblast was
136 reported by Tsang et al to be noninvasively resolvable and elevated in early onset preeclampsia
137 (gestational age < 34 weeks) compared to healthy pregnancy²⁸. The syncytiotrophoblast, involved in
138 nutrient exchange³³, showed no difference²⁸. We measured the respective signature scores in two
139 previously published preeclampsia cohorts³⁵. In contrast to Tsang et al, we observed no significant
140 difference in either trophoblast signature score in cfRNA samples collected at diagnosis for mothers with
141 early-onset preeclampsia (Fig. 3A) ($p = 0.703, 0.794$ respectively, two-sided Mann Whitney U) and for
142 mothers with either early- or late-onset preeclampsia ($p = 0.24, 0.54$ respectively, Kruskal Wallace, fig.
143 S10A) as compared to samples from mothers with no complications at a matched gestational age. Though
144 we generally recapitulate the observed signature score directionality for both cell types²⁸, examination of
145 the cell type gene profiles used by Tsang et al. in two independent placental atlases^{32,33} revealed genes
146 that were not cell type specific (fig. S10B,C). The presence of non-cell type specific genes in a cell type
147 gene profile likely impaired Tsang et al's signature score interpretation. The role of extravillous
148 trophoblast invasion and the ubiquity of its cellular pathophysiology in preeclampsia thus remains an
149 open question.

150 As a second example, the proximal tubule is a highly metabolic, predominant kidney cell type
151 and is a major source for injury and disease progression in CKD^{29,30}. Tubular atrophy is a hallmark of
152 CKD nearly independent of CKD etiology³⁶. Atubular formation results from proximal tubule damage
153 and death in renal pathologies^{29,30}. We discovered a striking decrease in proximal tubule intensity of CKD
154 ($n = 9$) patients (ages 67-91, CKD stage 3-5 or peritoneal dialysis) compared to healthy controls ($n = 7$) (p
155 $= 7.45 * 10^{-4}$, one-sided Mann Whitney U) (Fig. 3B, Methods). Random sampling of a gene set of
156 equivalent length as the proximal tubule gene profile in 10,000 trials yielded no discriminatory power
157 between CKD and healthy 96.6% of the time (Benjamini-Hochberg correction, FDR = 0.05, fig. S11)
158 and an adjusted proximal tubule signature score of 0.038, validating the specificity of our proximal tubule
159 gene profile. These findings demonstrate the ability to noninvasively resolve proximal tubule
160 deterioration observed in CKD histology³⁷ and is consistent with findings from routine invasive biopsy.

161 As a third example, AD pathogenesis results in neuronal death and synaptic loss³¹. We derived
162 brain cell type gene profiles for cell types in the AD brain in the same way as the NCI cell type profiles
163 (Methods). We then intersected a given cell type gene profile in AD with the equivalent NCI profile for
164 comparative analysis. Microglia, though often implicated in AD pathogenesis were excluded given their
165 high overlapping transcriptional profile with non-central nervous system macrophages³⁸. Inhibitory
166 neurons were also excluded given the low number of cell type specific genes intersecting between AD
167 and NCI phenotypes. Intersection of reportedly differentially downregulated AD genes (DEG) in plasma⁴
168 with the derived cell type specific gene profiles, identified several genes as cell type specific (Fig. 3C).
169 Astrocyte specific genes included filament protein (*GFAP*³⁹) and ion channels (*GRIN2C*¹⁷). Excitatory
170 neuron specific genes solute carrier proteins (*SLC17A7*¹⁷ & *SLC8A2*⁴⁰), cadherin proteins (*CDH8*⁴¹ &
171 *CDH22*⁴²), and a glutamate receptor stimulated a major excitatory neuron neurotransmitter (*GRM1*^{31,43}).
172 Neuronal death in AD phenotypes³¹, is likely the biochemical basis for the observed downregulation of
173 these cell type specific markers. Oligodendrocyte-specific genes encode proteins for myelin sheath
174 stabilization (*MOBP*³¹) and a synaptic/axonal membrane protein (*CNTN2*³¹). Oligodendrocyte precursor
175 cell-specific genes included transcription factors (*OLIG2*⁴⁴ & *MYT1*⁴⁵), neural growth and differentiation
176 factor (*CSPG5*⁴⁶), and a protein putatively involved in brain extracellular matrix formation (*BCAN*⁴⁷). A

177 permutation test using the Gini coefficients computed on the average single cell expression of the brain
178 specific and cell type specific DEGs corroborated that the DEGs assigned as cell type specific were more
179 specific to a given brain cell type than a brain specific DEG ($p < 1e-4$, 10000 trials, fig. S12) (Methods).
180 Taken together, these findings underscore the consistent detection cell-type specific changes in pathology
181 using noninvasive transcriptomic measurement of blood plasma and the resolution at which we can assert
182 the origins of cfRNA.

183

184 **Discussion**

185 We have shown that the cfRNA transcriptome can be reframed from a sum of tissue
186 transcriptomes to a sum of cell type transcriptomes. Using nu-SVR, we determined the fractional
187 contributions of cell type specific RNA relative to other cell types considered in the cell type column
188 space of the basis matrix⁴⁸. Ideally reference gene profiles for all possible cell types in the human body
189 would be simultaneously considered in nu-SVR deconvolution. However, a complete reference dataset of
190 all cell types in the human body does not yet exist. Despite an incomplete cell atlas, we demonstrate the
191 ability to decompose the cf-transcriptome into distinct cell type contributions by further leveraging
192 individual cell atlases and bulk transcriptomic data from GTEx and HPA to define specificity in context
193 of the whole body.

194 The decomposition of the cf-transcriptome reveals platelets, lymphocytes, and monocytes as the
195 predominant cell type specific RNA contributors. This is consistent with what is known about the cf-
196 transcriptome, and may also reflect a bias in nu-SVR deconvolution, which uses highly expressed genes
197 as support vectors, and consequently assigns a reduced fractional contribution to cell types expressing
198 genes at lower levels or that are smaller in size, such as neutrophils. Furthermore, our finding that
199 platelets are a majority cell type, rather than megakaryocytes¹, likely reflects annotation differences in
200 reference data. Megakaryocytes are absent from the TSP v0.9 annotations; however, they are responsible
201 for platelet production⁸. Comparison of nu-SVR to quadratic programming³ and non-negative linear least
202 squares⁴⁹ yielded similar deconvolution RMSE and slightly increased Pearson correlation. However, the
203 determined fractions with these methods excluded contributions from cell types with markers detected
204 using PanglaoDB, and so we chose nu-SVR for the comprehensive deconvolution in this work. Taken
205 together, these findings are consistent with prior work considering specific cell types¹ and the
206 hematopoietic tissues^{1,3}. Shared features among the cell types contributing to the to the cf-transcriptome
207 are large volume and/or increased turnover²⁴, suggesting cell death as the possible predominant entry
208 mechanism of cfRNA to the bloodstream.

209 To identify contributions from cell types absent from the basis matrix, we derived individual cell
210 type gene profiles from individual tissue scRNA-seq atlases. By considering cell types unique to a given
211 tissue, we could leverage bulk RNA transcriptomic datasets from GTEx and HPA to ensure specificity in
212 context of the whole body. This directed approach could enable the application of many single tissue cell
213 atlases, whose meaningful integration into approaches like nu-SVR is otherwise limited by batch effects⁵⁰.
214 Defining cell type specific genes from single tissue atlases in the context of the whole body surmounts the
215 problem of missing cell types from TSP and more generally enables a means to address missing cell types
216 from the basis matrix column space required for deconvolution approaches like nu-SVR^{14,15}.

217 Previous work has demonstrated cell types-of-origin identification from cell free nucleosomes;¹⁶
218 in particular of enriched cardiomyocyte signal in patients with acute myocardial infarction, and liver cell
219 type specific signatures in patients with various hepatic disorders. However, reference ChIP-seq data
220 from pure cell types is limited, thereby reducing the scope of resolvable cell types with this approach⁵¹

221 and limiting interpretation of some aspects of the data to tissue resolution, and sensitivity appears to be
222 limited. For example, in healthy (non-pathological) cases it was only possible to deconvolve blood cell
223 types, and other than a generalized signal from liver it was not possible to detect cell types from solid
224 tissues. Another interesting distinction which will potentially affect the overall applicability is that
225 cfDNA measurement requires cell death, whereas cfRNA additionally incorporates information from
226 living cells which secrete RNA by various mechanisms¹.

227 The results here reinforce the importance of reliable reference data annotation at both bulk tissue
228 and single cell level; differences in either impact the ability to meaningfully integrate in cfRNA analysis.
229 Cell type annotation differences across distinct cell atlases for the same tissue may impact the assignment
230 of a gene as cell-type specific when considering a single dataset. Specifically, we observed that several
231 genes reported as specific to a single trophoblast cell type²⁸ were not validated in two independent
232 placental cell atlases^{32,33}. Annotation discrepancies between atlases impacts the assignment of genes as
233 cell type specific in context of the whole body, and consequently impact the interpretation of a cell type
234 signature score in cfRNA.

235 This method for detecting cellular pathophysiology in the cf-transcriptome is most robust for
236 diseases with cell type changes independent of disease etiology. CKD cfRNA samples reveal striking
237 differences in proximal tubule signature score compared to healthy controls, in contrast with the minimal
238 effect size of extravillous trophoblast signature score in preeclampsia. Multiple differentially
239 downregulated genes in AD phenotypes are cell type specific, again reinforcing the ability to
240 noninvasively resolve pathologically implicated cell types. Cellular pathology in CKD and AD are
241 proximal tubule atrophy and neuronal death respectively, which occur irrespective of disease etiology,
242 whereas preeclampsia etiology may be multifactorial and the extent of underlying cellular pathogenesis
243 remains to be explored⁵².

244 CKD impacts 9.1% of the global population⁵³ and invasive biopsy (often multiple) is the clinical
245 gold-standard for diagnosis. The sensitivity and specificity of some glomerular filtration estimates for
246 various patient populations and kidney function may result in clinical confusion⁵⁴. Standard lab tests for
247 serum creatinine and urine albumin levels merely indicate renal dysfunction and provide limited clinically
248 actionable insight into the source(s) of renal pathophysiology. Tubular atrophy has been repeatedly shown
249 to be superior to glomerular pathology as a predictor of CKD progression³⁷. The ability to noninvasively
250 measure proximal tubule pathophysiology and the general detection of multiple renal cell types could
251 enable noninvasive classification of various renal disorders in future work and augment patient treatment
252 plans. Given the small sample size used here (n = 9, CKD; n = 7, control), we emphasize that our
253 findings, although striking, must still be validated in a larger follow-up study.

254 This work shows that one can apply cell atlases to measure disparate cell types that are disease-
255 implicated in the blood, relevant to a myriad of questions impacting human health. Unlike model
256 organisms which lack full translatability to human health, cf-transcriptomic measurement provides direct,
257 immediate insights into patient health. Readily measurable cell types in cfRNA, including those specific
258 to the brain, lung, intestine, liver, and kidney, have vast prognostic and clinical importance given the
259 multitude of diseases in these tissues. Single cell RNA-seq reveals numerous cell type specific changes in
260 pathologies within these tissues for investigation with cfRNA ranging from cancer to Crohn's disease,
261 drug or vaccine response, and aging. Consistent detection of cell types responsible for drug metabolism
262 (e.g. liver and renal cell types) as well as cell types that are drug targets, such as neurons or
263 oligodendrocytes for Alzheimer's-protective drugs, could provide powerful clinical trial end-point data in
264 evaluating drug toxicity. Chemotherapy regimens are known to have severe systemic side-effects. We

265 anticipate that the ability to noninvasively resolve cell type signatures in plasma cfRNA will both enhance
266 existing clinical knowledge in addition to enabling increased resolution in monitoring disease progression
267 and drug response.

268

269 **Acknowledgements:** We thank M. Chen for single cell analysis input, feedback and helpful discussions.
270 We thank E. Sattely for helpful discussions. We thank G. Loeb for kidney discussions.

271

272 **Funding:** This work is supported by the Chan Zuckerberg Biohub. S.K.V. is supported by a NSF
273 Graduate Research Fellowship (Grant # DGE 1656518), the Benchmark Stanford Graduate Fellowship,
274 and the Stanford ChEM-H Chemistry Biology Interface (CBI) training program.

275

276 **Author contributions:** S.K.V. and S.R.Q. conceptualized the study. S.K.V. and S.R.Q. designed the
277 study in collaboration with M.N.M. S.K.V. performed all analyses; M.N.M. wrote the bioinformatic
278 preprocessing pipeline to map reads to the human genome and cell free sample QC. S.K.V, M.N.M,
279 S.R.Q wrote the manuscript. All authors revised the manuscript and approved it for publication.

280

281 **Competing interests:** S.R.Q is a founder and shareholder of Molecular Stethoscope and Mirvie. S.K.V,
282 M.N.M, and S.R.Q are inventors on a patent application covering the methods and compositions to detect
283 specific cell types using cfRNA submitted by the Chan Zuckerberg Biohub and Stanford University.

284

285 **Data and materials availability:** Code for the work in this manuscript will be made available on Github.
286 All datasets used for this work were publicly available, downloaded with permission, or directly requested
287 from authors.

288

289 **Materials & Methods**

290

291 ***Data Processing***

292 **Data acquisition**

293 Cell free RNA: For samples from Ibarra et al, raw sequencing data was obtained from the SRA
294 (PRJNA517339). For samples from Munchel et al, processed counts tables were directly downloaded.

295

296 For all individual tissue single cell atlases, Seurat objects or AnnData objects were downloaded or
297 directly received from authors. Data from Mathys et al. were downloaded with appropriate approvals
298 from Synapse.

299

300 HPA v19 transcriptomic data, GTEx v8 raw counts, and Tabula Sapiens v0.9 were downloaded directly.

301

302 **Bioinformatic processing**

303 For each sample for which raw sequencing data were downloaded, we trimmed reads using trimmomatic
304 (v 0.36) and then mapped them to the human reference genome (hg38) with STAR (v 2.7.3a). Duplicate
305 reads were then marked and removed by GATK's (v 4.1.1) MarkDuplicates tool. Finally, mapped reads
306 were quantified using htseq-count (v 0.11.1), and read statistics were estimated using FastQC (v 0.11.8).

307

308 The bioinformatic pipeline was managed using snakemake (v 5.8.1). Read and tool performance statistics
309 were aggregated using MultiQC (v 1.7).

310

311

312 **Sample quality filtering**

313 For every sample for which raw sequencing data was available, we estimated three quality parameters as
314 previously described^{55,56}. To estimate RNA degradation, we calculated a 3' bias ratio. Specifically, we
315 first counted the number of reads per exon and then annotated each exon with its corresponding gene ID
316 and exon number using htseq-count. Using these annotations, we measured the frequency of genes for
317 which all reads mapped exclusively to the 3' most exon as compared to the total number of genes
318 detected. We approximate RNA degradation for a given sample as the fraction of genes where all reads
319 mapped to the 3' most exon.

320 To estimate ribosomal read fraction, we compared the number of reads that mapped to the ribosome
321 (Region GL00220.1:105424-118780, hg38) relative to the total number of reads (Samtools view). To
322 estimate DNA contamination, we used an intron to exon ratio and quantified the number of reads that
323 mapped to intronic as compared to exonic regions of the genome.

324 We then identified outlier samples using the 95th percentile bound within a given cfRNA dataset. We
325 considered any given sample a low quality sample if its value for any metric was greater than or equal to
326 the 95th percentile bound.

327

328 **Data Normalization**

329 All gene counts were adjusted to counts per million reads and per milliliter of plasma used. For a given
330 sample (i denotes gene index and j denotes sample index):

331

$$332 \quad \eta_{ij} = \frac{Gene_{ij}}{(Library\ Size_j) * (mL\ Plasma_j)} \quad \text{where } Library\ Size_j = \sum_i G_{ij} \quad (1)$$

333

334 For subjects who had samples with multiple technical replicates, these plasma volume CPM counts were
335 averaged prior to nu-SVR deconvolution.

336

337 For all analyses except nu-SVR (e.g. all work except Fig. 1b), we next applied trimmed mean of M values
338 (TMM) normalization as previously described⁵⁷:

339

$$340 \quad \frac{\eta_{ij}}{TMM_j} \quad (2)$$

341

342 CPM TMM normalized gene counts across technical replicates for a given biological replicate were
343 averaged for the count tables used in the analyses performed in Figures 2 & 3.

344

345 Sequencing batches and plasma volumes were obtained from the authors in Toden et al for per-sample
346 normalization. For samples from Ibarra et al., plasma volume was assumed to be constant at 1 mL as we
347 were unable to attain this information from the authors. Sequencing batches were inferred based on the
348 figure and confirmed with authors that sequencing strategy was figure-dependent (personal
349 communication).

350 All samples from Munchel et al were used to compute TMM scaling factors and 4.5 mL plasma was used
351 to normalize all samples within a given dataset (both PEARL-PEC and iPEC)

352

353 For the work in Figure 3B, longitudinal samples for a given CKD patient were averaged, given that the
354 timescale over which renal cell type changes would occur were longer than the patient samples (~30
355 days). These samples were collected alongside two healthy patient biological replicates passing sample
356 QC in this sequencing batch. We additionally used the five plasma biological replicates from above. This
357 ultimately yielded $n = 7$ healthy biological replicates and $n = 9$ CKD patients.

358

359

360

361 **Zero-Centered Batch Normalization**

362 To account for center-specific effects that could impact meaningful comparison of data across centers in
363 Figure 3B, we subtracted the mean normalized value across all samples measurements for given gene
364 within a given batch from the measured normalized value for a given sample⁵⁸:

$$365 \quad \quad \quad 366 \quad \quad \quad \overline{G_{ij}} = G_{ij} - \mu_{ik} \quad (3)$$

367
368 Where the gene index is i , the sample is j , and k is the batch. The mean expression of the i^{th} gene in the k^{th}
369 batch is denoted by μ_{ik} .

370
371 We defined a ‘batch’ of samples to reflect the experimental workflow for each of the corresponding
372 analyses:

- 373 • For samples from Ibarra et al, given that only two control biological replicates were
374 sequenced with the CKD samples and the other healthy controls came from another batch, we
375 did not directly compare CKD vs. healthy samples. A difference in the raw median value of
376 the proximal tubule signature score between the two sick and the nine healthy samples
377 sequenced in the same batch was observed prior to grouping healthy plasma data from a
378 different batch, consistent with the observed difference post-zero centered normalization. All
379 CKD and healthy samples were treated as a single batch from which normalization was
380 performed.
- 381 • For the datasets from Munchel et al., zero-centered batch normalization was not performed
382 given that the data were compared within the same sequencing studies (e.g. iPEC and
383 PEARL-PEC)

384 **Cell Type Marker Identification using PanglaoDB**

385 The PanglaoDB cell type marker database was downloaded on March 27, 2020. Markers were filtered for
386 human (“Hs”) only. Specificity (how often marker was not expressed in a given cell type) and sensitivity
387 (how frequently marker is expressed in cells of this type) thresholds determined the total gene space for
388 intersection across the cfRNA samples. Gene synonyms from Panglao were determined using MyGene
389 version 3.1.0 to ensure full gene space.

390
391 The intersection of this space with each cfRNA sample were then determined, where the error bars reflect
392 the differences in number of markers detected across the samples for a given cell type across samples. A
393 given cell type marker was counted in a given healthy cfRNA sample its gene expression was greater than
394 zero in log + 1 transformed CPM-TMM gene count space.

395
396 Cell types with markers filtered by sensitivity = 0.9 and specificity = 0.2 and samples with ≥ 5 cell type
397 markers are shown in Fig 1B.

398
399 The samples used for nu-SVR deconvolution were the five healthy donor plasma samples as in Figure 1D
400 of Ibarra et al.

401 **Basis Matrix Formation**

402 Only cells from droplet sequencing (“10X”) were used in analysis. Disassociation genes as reported¹⁰ and
403 cell types too granular (i.e. fast muscle cell, smooth muscle cell, LYVE1 aortic macrophage,
404 differentiated basal cell of epithelium of trachea, etc) or too broad (i.e. granulocyte, lymphocyte,
405 monocyte subtypes, innate lymphoid cells, etc) in annotation were excluded from subsequent analysis.

406
407
408 Of the remaining cell types, either 30 observations were randomly sampled or the maximum number of
409 available observations if less than 30 were subsampled, whichever was greater.

410

411 Cells were assigned broader labels to enable linear independence of the matrix column space, as several
412 cell types that are very transcriptionally similar with few distinct gene would be challenging to resolve
413 noninvasively and hence would impact nu-SVR deconvolution. In any sort of regression, multicollinearity
414 between features will impact the learned coefficients.

415

416 Scanpy (version 1.6.0) was used to analyze all single cell data. Hierarchical clustering was performed on
417 PCA-transformed (scanpy pca) CPM cell counts (scanpy normalize_total, target_sum = 1e6) log
418 transformed (scanpy log1p) counts using the scanpy dendrogram function. Cell types that were close in
419 clustering were grouped together, including:

- 420 • 'pancreatic A/B/PP cell' comprised 'pancreatic A cell', 'pancreatic PP cell', 'type B pancreatic cell'
- 421 • Lung pneumocyte comprised 'type I pneumocyte' and 'type II pneumocyte'
- 422 • "intestinal crypt stem cell + transient amplifying cell" comprised 'intestinal crypt stem cell',
423 'intestinal transient amplifying cell', 'paneth cell'
- 424 • "vascular smooth muscle cell" comprised 'aortic smooth muscle cell' and 'vascular associated
425 smooth muscle cell'

426 All cell types annotated as some type of 'B cell', 'T cell', 'NK cell', 'dendritic cell', 'thymocyte', etc
427 were labelled respectively with the broader category.

428

429 All groupings are available on Github in the script entitled 'coarsegraincells_forBMGeneration.py'.

430

431 This subsampled counts matrix was then passed to the 'Create Signature Matrix' analysis module on
432 available at cibersortx.stanford.edu, with the following parameters:

- 433 • Disable quantile normalization = False
- 434 • Min. expression = 0.25
- 435 • Replicates = 5
- 436 • Sampling = 0.5
- 437 • Kappa = 999
- 438 • q-value = 0.01
- 439 • No. barcode genes = 2450 - 5000
- 440 • Filter non-hematopoietic genes = False

441 The resulting basis matrix was then saved as a .txt file and used in nu-SVR deconvolution

442

443 *Nu-SVR deconvolution*

444 We formulated the cell free transcriptome as a linear summation of the cell types from which it
445 originates^{3,59}. With this formulation, we adapted existing deconvolution methods developed with the
446 objective of decomposing a bulk tissue sample into its single cell constituents^{14,15}, where the
447 deconvolution problem is formulated as:

448

$$449 \quad A\theta = b \quad (4)$$

450

451 Here, A is the representative basis matrix ($g \times c$) of g genes for c cell types, which represent the gene
452 expression profiles of the c cell types. θ is a vector ($c \times 1$) of the contributions of each of the cell types
453 and b is the measured expression of the genes observed in from blood plasma ($g \times 1$). The goal here is to
454 learn θ such that the matrix product $A\theta$ predicts the measured signal b . The derivation of the basis matrix
455 A is described in the section 'Basis Matrix Formation'.

456

457 We performed nu-SVR using a linear kernel to learn θ from a subset of genes from the signature matrix
458 to best recapitulate the observed signal b , where nu denotes the lower bound on the fraction of support
459 vectors and the upper bound on the fraction of errors at the margin⁶⁰. Here, the support vectors are the
460 genes used from the basis matrix from which to learn θ ; θ reflects the weights of the cell types in the

461 basis matrix column space. For each sample, we learned coefficients for six values of ν , $\nu \in$
462 $\{0.1, 0.15, 0.25, 0.5, 0.75, 0.9\}$ and estimated the resulting deconvolution error using the root mean square
463 error (RMSE). We determined the product of the basis matrix with the learned coefficients ($A\theta$), which
464 reflected some predicted expression value for each of the genes in a given cfRNA mixture. The RMSE
465 was then computed using the predicted expression values and the measured values across all the non-zero
466 CPM genes in a cfRNA mixture.

467
468 Only CPM counts > 0 were considered in the mixture. The values in the basis matrix were also in CPM
469 space. Prior to deconvolution, the mixture and basis matrix were scaled to zero mean and unit variance for
470 improved runtime performance. We emphasize that we did not log-transform counts in b or in A , as this
471 would destroy the requisite linearity assumption in equation 4. Specifically, the concavity of the log
472 function would result in the consistent underestimation of θ during deconvolution⁶¹.

473
474 Using the θ resulting from the value of ν whose coefficients yielded the smallest RMSE was transformed
475 to. Specifically, the relative fractional contributions of cell type specific RNA from θ , we repeat what was
476 previously described^{14,15}:

$$477 \quad \forall \theta_j < 0 \in \{\theta_1, \dots, \theta_c\} \rightarrow 0 \quad (5)$$

478 All non-zero coefficients were then normalized by their sum to result in the relative fractions to determine
479 the relative fractional contributions of cell type specific RNA.

480
481 We used the function nuSVR from scikitlearn version 0.23.2.
482 The samples used for nu-SVR deconvolution were the five healthy donor plasma samples as in Figure 1D
483 of Ibarra et al.

484
485 ***Evaluating Basis Matrix on GTEx samples***
486 Bulk RNA-seq samples from GTEx v8 were deconvolved with the derived basis matrix from tissues that
487 were present (kidney cortex, whole blood, small intestine – terminal ileum, lung, and spleen) or absent
488 (kidney medulla and liver) from the basis matrix derived using Tabula Sapiens version 0.9. For each
489 tissue type, the maximum number available samples or ten samples, whichever was smaller, was
490 deconvolved. Please see Supplementary Note 1 for additional discussion.

491
492 ***Identifying tissue specific genes in cfRNA absent from basis matrix***

493 To identify cell type specific genes in cfRNA that were distinct to a given tissue, we considered the set
494 difference of the non-zero genes measured in a given cfRNA sample with the row space of the basis
495 matrix and intersected this with HPA tissue specific genes:

$$496 \quad (G_j - R) \cap HPA \quad (6)$$

497 Where G_j is the gene set in the j^{th} deconvolved sample, where a given gene in the set's expression was \geq
498 5 TMM-CPM. R is the set of genes in the row space of the basis matrix used for nu-SVR deconvolution.
499 HPA denotes the total set of tissue specific genes from HPA.

500
501 The HPA tissue specific gene set (HPA) were genes across all tissues with Tissue Specificity assignments
502 'Group Enriched', 'Tissue Enhanced', 'Tissue Enriched' and NX expression ≥ 10 .

503
504 This approach yielded tissues with several distinct genes present in cfRNA which could then be
505 subsequently interrogated using single cell data.

506
507 ***Derivation of cell type specific gene profiles in context of the whole body using single cell data***

508 For this analysis, only cell types unique to a given tissue (i.e. hepatocytes unique to the liver, or excitatory
509 neurons unique to the brain) were considered so that bulk transcriptomic data could be used to ensure
510 specificity in context of the whole body. A gene was asserted to be cell type specific if it was (i)

511 differentially expressed within a given tissue cell type atlas (ii) had a Gini coefficient ≥ 0.6 , indicating
512 comprehensive tissue specificity in context of the whole body.

513

514 **(1) Single cell differential expression**

515 For data received as a Seurat object, conversion to AnnData was performed by saving as an intermediate
516 loom objects (Seurat version 3.1.5) and converting to AnnData (loompy version 3.0.6). Scanpy (version
517 1.6.0) was used for all other single cell analysis. Reads per cell were normalized for library size (scanpy
518 `normalize_total`, `target_sum = 1e4`), then logged (scanpy `log1p`). Differential expression was performed
519 using the Wilcoxon rank sum test in Scanpy's `filter_rank_genes_groups` with the following arguments:
520 `min_fold_change = 1.5`, `min_in_group_fraction = 0.2`, `max_out_group_fraction = 0.5`, `corr_method =`
521 `"benjamini-hochberg"`. For differentially expressed genes (DEG) with Benjamini Hochberg adjusted p-
522 values < 0.01 , the ratio of the highest out_group percent expressed to in_group percent expressed < 0.5 to
523 ensure high specific expression in the cell type of interest within a given cell type atlas.

524

525 **(2) Quantifying comprehensive whole body tissue specificity using the Gini coefficient**

526 The distribution of all the Gini coefficients and Tau values across all genes belonging to cell type gene
527 profiles for cell types native to a given tissue were compared using the HPA gene expression Tissue
528 Specificity and Tissue Distribution assignments¹² (fig S8). The Gini coefficient better reflected the
529 underlying distribution of gene expression tissue-specificity than Tau (fig. S8) and was hence used for
530 subsequent analysis. As the Gini coefficient approaches unity, this indicates extreme gene expression
531 inequality, or equivalently high specificity. A single threshold (Gini coefficient ≥ 0.6) was applied across
532 all atlases to facilitate a generalizable framework from which to define tissue specific cell type gene
533 profiles in context of the whole body in a principled fashion for signature scoring in cfRNA.

534

535 For the following definitions, n denotes the total number of tissues and x_i is the expression of a given
536 gene in the i^{th} tissue.

537

538 The Gini coefficient was computed as defined in²⁰:

$$539 \text{Gini} = \frac{n+1}{n} - \frac{2 \sum_{i=1}^n (n+1-i)x_i}{n \sum_{i=1}^n x_i}; x_i \text{ is ordered from least to greatest. (7)}$$

540

541 Tau, as defined in²⁰:

$$542 \tau = \frac{\sum_{i=1}^n 1 - \bar{x}}{n-1} \text{ where } \bar{x} = \frac{x_i}{\max(x_i) \forall i \in \{1 \dots n\}} \text{ (8)}$$

543

544 HPA NX Counts from the HPA object entitled 'rna_tissue_consensus.tsv' accessed on July 1, 2019 were
545 used for computing Gini coefficients and Tau.

546

547 Note for brain cell type gene profiles: given that there are multiple sub brain-regions in the HPA data, the
548 determined Gini coefficients are lower (e.g. not as close to unity compared to other cell type gene
549 profiles) since there are multiple regions of the brain with high expression, which would result in reduced
550 count inequality.

551

552 **Gene Expression in GTEx**

553 We used the raw GTEx data v8 (accessed August 26 2019) and converted to $\log(\text{CPM} + 1)$ counts. The
554 signature score was determined by summing the expression of the genes in a given bulk RNA sample for
555 a given cell type gene profile. Since only gene profiles were derived for cell types that correspond to a
556 given tissue, the mean signature score of a cell type profile across the non-native tissues was then
557 computed and used to determine the log fold change

558

559

560 ***Estimating signature scores for each cell type***

561 The signature score is defined as the sum of genes asserted to be cell type specific, where i denotes the
562 index of the gene in a cell type signature gene profile in the j^{th} patient sample.

563

$$564 \text{Signature Score}_j = \sum_i G_{ij} \quad (9)$$

565

566 For signature scoring of syncytiotrophoblast and extravillous trophoblast gene profiles in PEARL-PEC
567 and iPEC³⁵. The genes in a respective profile used for signature scoring were derived as described in
568 ‘Derivation of cell type specific gene profiles in context of the whole body using single cell data’
569 independently using two different placental single cell datasets^{32,33}. Only the intersection of the cell type
570 specific gene profiles was considered for signature scoring.

571

572 ***Comparison of proximal tubule signature score to random for discriminatory between CKD 3+ and***
573 ***Healthy***

574 To assess the discriminatory power of a given cell type signature score with a statistically significant
575 difference in Fig 3B, we randomly sampled an equivalent gene length as the proximal tubule gene profile
576 in 10,000 trials and performed a one-sided Mann Whitney U with the alternative hypothesis that healthy
577 would be greater than CKD 3+. For a given trial, the signature score of the random gene list was
578 computed across all samples and tested. Benjamini-Hochberg correction at FDR = 0.05 was performed
579 using ‘multitest’ function in statsmodels version 0.10.1 with the following arguments: alpha=0.05
580 method=‘fdr_bh’.

581

582 ***Cell Type specific differentially expressed neuronal and glial cell type specific genes in Alzheimer’s***
583 ***plasma***

584 To assess whether DEGs in AD/NCI plasma⁴ that intersected with a brain cell type gene profile were
585 more specific to a given brain cell type than DEGs in AD/NCI plasma that was generally brain tissue
586 specific, we performed a permutation test. Specifically, we compared the Gini coefficient for genes in
587 these two groups, computed using the mean expression of a given gene across brain cell types from
588 healthy brain single cell data¹⁷.

589

590 The starting set of brain specific genes were defined using in the HPA brain transcriptional data annotated
591 as either ‘Tissue enriched’, ‘Group enriched’, or ‘Tissue enhanced’ (accessed January 13, 2021). These
592 requirements ensured the specificity of a given brain gene in context of the whole body. This formed the
593 initial set of brain specific genes B .

594

595 The union of all brain cell type specific genes is the set C .

596

597 Genes in B that that did not intersect with C (e.g. any brain cell type gene profile (‘brain cell type
598 specific’)) and intersected with DEG-up (U) or DEG-down genes (D)⁴ were then defined as ‘brain tissue
599 specific’.

600

$$601 T = (B \cap U) + (B \cap D) \quad (10)$$

602

603 All genes belonging to brain cell type gene profiles (‘brain cell type specific’) were a subset of the initial
604 set of brain specific genes.

605

$$C - B = 0 \quad (11)$$

606 Genes defined as ‘brain cell type specific’ for signature scoring in Fig. 3C were intersected with
607 differentially upregulated (DEG-up) and differentially downregulated genes (DEG-down) reported⁴. No
608 DEG-up genes intersected with any of the brain signatures used in Fig 3C. Only DEG-down were
609 considered in the subsequent analysis as ‘brain cell type specific’.

610

611 The Gini coefficients reflecting the gene expression inequality across brain cell types were computed for
612 the gene sets labelled as ‘brain cell type specific’ and ‘brain tissue specific’. Brain reference data to
613 compute Gini coefficients was the single cell brain atlas with diagnosis as ‘Normal’¹⁷. All Gini
614 coefficients were computed using the mean log transformed CPTT (counts per ten thousand) gene
615 expression per cell type.

616

617 A permutation test was then performed on the union of the Gini coefficients for the genes labeled as
618 ‘brain cell type specific’ and ‘brain tissue specific’. The purpose of this test was to assess probability that
619 the observed mean difference in Gini coefficient for these two groups yielded no difference in specificity
620 (e.g. $H_0: \mu_{cell\ type\ Gini\ Coefficient} = \mu_{brain\ tissue\ Gini\ coefficient}$).

621

622 Gini coefficients were permuted and reassigned to the list of ‘brain tissue’ or ‘brain cell type’ genes, then
623 the difference in mean of the two groups was computed. This procedure was repeated 10,000 times. The
624 p-value was determined as follows:

625

$$626 \quad p = \frac{\# \text{ trials with permuted } (\mu_{cell\ type} - \mu_{tissue}) \geq \mu_{observed}}{10,000+1} \quad (12)$$

627 Where $\mu_{observed} := (\mu_{cell\ type\ Gini\ Coefficient} - \mu_{brain\ tissue\ Gini\ coefficient})$.

628

629 The additional 1 in the denominator reflects the original test between the true difference in means (e.g. the
630 true comparison yielding $\mu_{observed}$)

631

632 **Supplementary Note 1: Deconvolution of bulk GTEx tissues using the *Tabula Sapiens*-derived basis** 633 **matrix**

634 To assess the ability of the basis matrix to deconvolve tissues whose cell types were wholly present in the
635 cell type column space, we deconvolved a subset of bulk RNA-seq GTEx samples. The determined
636 fractions of cell type specific RNA generally recapitulated the predominant cell types within a given
637 tissue (fig. S4). Kidney cortex majority fractions were from kidney epithelia and vascular endothelia (fig.
638 4A,B); small intestine, smooth muscle cells and intestinal enterocytes (fig. S4E); whole blood,
639 erythrocytes (fig. S4G). Cells with larger volume yielded larger deconvolved fractions for all tissues (fig.
640 S3). Variance in the relative cell type fractional contributions across the deconvolved bulk samples within
641 a given tissue reflects the underlying cell type heterogeneity. GTEx kidney medulla samples recorded to
642 be contaminated with renal cortex reflect the presence of the kidney epithelia, the majority cell type in the
643 renal cortex (fig. S4A). Tissues absent from the cell type column space, such as liver, yielded cell types
644 that are transcriptionally similar (kidney epithelia) (fig. S4C) and a higher deconvolution error (fig. S3) as
645 expected.

646

647 **References**

- 648 1. Ibarra, A. *et al.* Non-invasive characterization of human bone marrow stimulation and
649 reconstitution by cell-free messenger RNA sequencing. *Nat. Commun.* **11**, 400 (2020).
- 650 2. Ngo, T. T. M. *et al.* Noninvasive blood tests for fetal development predict gestational age and
651 preterm delivery. *Science* **360**, 1133–1136 (2018).
- 652 3. Koh, W. *et al.* Noninvasive in vivo monitoring of tissue-specific global gene expression in
653 humans. *Proc Natl Acad Sci USA* **111**, 7361–7366 (2014).

- 654 4. Toden, S. *et al.* Noninvasive characterization of Alzheimer's disease by circulating, cell-free
655 messenger RNA next-generation sequencing. *Sci. Adv.* **6**, (2020).
- 656 5. Chalasani, N. *et al.* Non-invasive Stratification of Non-Alcoholic Fatty Liver Disease by Whole-
657 transcriptome Cell-free mRNA Characterization. *Am. J. Physiol. Gastrointest. Liver Physiol.*
658 (2021) doi:10.1152/ajpgi.00397.2020.
- 659 6. Larson, M. H. *et al.* A comprehensive characterization of the cell-free transcriptome reveals
660 tissue- and subtype-specific biomarkers for cancer detection. *Nat. Commun.* **12**, 2357 (2021).
- 661 7. Basu, M., Wang, K., Ruppin, E. & Hannehalli, S. Predicting tissue-specific gene expression from
662 whole blood transcriptome. *Sci. Adv.* **7**, (2021).
- 663 8. Klatt, E. C. *Robbins & Cotran Atlas of Pathology.* (Elsevier, 2021).
- 664 9. Gawad, C., Koh, W. & Quake, S. R. Single-cell genome sequencing: current state of the science.
665 *Nat. Rev. Genet.* **17**, 175–188 (2016).
- 666 10. Tabula Sapiens: human transcriptome reference at single cell resolution. <https://tabula-sapiens-portal.ds.czbiohub.org>.
- 668 11. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**,
669 204–213 (2017).
- 670 12. Uhlen, M. *et al.* A genome-wide transcriptomic analysis of protein-coding genes in human blood
671 cells. *Science* **366**, (2019).
- 672 13. Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of
673 mouse and human single-cell RNA sequencing data. *Database (Oxford)* **2019**, (2019).
- 674 14. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat.*
675 *Methods* **12**, 453–457 (2015).
- 676 15. Newman, A. M. *et al.* Determining cell type abundance and expression from bulk tissues with
677 digital cytometry. *Nat. Biotechnol.* **37**, 773–782 (2019).
- 678 16. Sadeh, R. *et al.* ChIP-seq of plasma cell-free nucleosomes identifies gene expression programs of
679 the cells of origin. *Nat. Biotechnol.* (2021) doi:10.1038/s41587-020-00775-6.
- 680 17. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* **570**, 332–337
681 (2019).
- 682 18. Aizarani, N. *et al.* A human liver cell atlas reveals heterogeneity and epithelial progenitors. *Nature*
683 **572**, 199–204 (2019).
- 684 19. Litviňuková, M. *et al.* Cells of the adult human heart. *Nature* **588**, 466–472 (2020).
- 685 20. Kryuchkova-Mostacci, N. & Robinson-Rechavi, M. A benchmark of gene expression tissue-
686 specificity metrics. *Brief. Bioinformatics* **18**, 205–214 (2017).

- 687 21. András, I. E. & Toborek, M. Extracellular vesicles of the blood-brain barrier. *Tissue Barriers* **4**,
688 e1131804 (2016).
- 689 22. Abbott, N. J. Inflammatory mediators and modulation of blood-brain barrier permeability. *Cell.*
690 *Mol. Neurobiol.* **20**, 131–147 (2000).
- 691 23. Ganong, W. F. Circumventricular organs: definition and role in the regulation of endocrine and
692 autonomic function. *Clin. Exp. Pharmacol. Physiol.* **27**, 422–427 (2000).
- 693 24. Sender, R. & Milo, R. The distribution of cellular turnover in the human body. *Nat. Med.* **27**, 45–
694 48.
- 695 25. Zemmour, H. *et al.* Non-invasive detection of human cardiomyocyte death using methylation
696 patterns of circulating DNA. *Nat. Commun.* **9**, 1443 (2018).
- 697 26. Danielson, K. M. *et al.* Plasma Circulating Extracellular RNAs in Left Ventricular Remodeling
698 Post-Myocardial Infarction. *EBioMedicine* **32**, 172–181 (2018).
- 699 27. Kaufmann, P., Black, S. & Huppertz, B. Endovascular trophoblast invasion: implications for the
700 pathogenesis of intrauterine growth retardation and preeclampsia. *Biol. Reprod.* **69**, 1–7 (2003).
- 701 28. Tsang, J. C. H. *et al.* Integrative single-cell and cell-free plasma RNA transcriptomics elucidates
702 placental cellular dynamics. *Proc Natl Acad Sci USA* **114**, E7786–E7795 (2017).
- 703 29. Nakhoul, N. & Batuman, V. Role of proximal tubules in the pathogenesis of kidney disease.
704 *Contrib. Nephrol.* **169**, 37–50 (2011).
- 705 30. Chevalier, R. L. The proximal tubule is the primary target of injury and progression of kidney
706 disease: role of the glomerulotubular junction. *Am. J. Physiol. Renal Physiol.* **311**, F145–61
707 (2016).
- 708 31. Grubman, A. *et al.* A single-cell atlas of entorhinal cortex from individuals with Alzheimer’s
709 disease reveals cell-type-specific gene expression regulation. *Nat. Neurosci.* **22**, 2087–2097
710 (2019).
- 711 32. Suryawanshi, H. *et al.* A single-cell survey of the human first-trimester placenta and decidua. *Sci.*
712 *Adv.* **4**, eaau4788 (2018).
- 713 33. Vento-Tormo, R. *et al.* Single-cell reconstruction of the early maternal-fetal interface in humans.
714 *Nature* **563**, 347–353 (2018).
- 715 34. Stewart, B. J. *et al.* Spatiotemporal immune zonation of the human kidney. *Science* **365**, 1461–
716 1466 (2019).
- 717 35. Munchel, S. *et al.* Circulating transcripts in maternal blood reflect a molecular signature of early-
718 onset preeclampsia. *Sci. Transl. Med.* **12**, (2020).
- 719 36. Dhillon, P. *et al.* The Nuclear Receptor ESRRA Protects from Kidney Disease by Coupling
720 Metabolism and Differentiation. *Cell Metab.* **33**, 379–394.e8 (2021).

- 721 37. Schelling, J. R. Tubular atrophy in the pathogenesis of chronic kidney disease progression.
722 *Pediatr. Nephrol.* **31**, 693–706 (2016).
- 723 38. van Rossum, D. & Hanisch, U.-K. Microglia. *Metab. Brain Dis.* **19**, 393–411 (2004).
- 724 39. McCall, M. A. *et al.* Targeted deletion in astrocyte intermediate filament (Gfap) alters neuronal
725 physiology. *Proc Natl Acad Sci USA* **93**, 6361–6366 (1996).
- 726 40. Lytton, J. Na⁺/Ca²⁺ exchangers: three mammalian gene families control Ca²⁺ transport.
727 *Biochem. J.* **406**, 365–382 (2007).
- 728 41. Friedman, L. G. *et al.* Cadherin-8 expression, synaptic localization, and molecular control of
729 neuronal form in prefrontal corticostriatal circuits. *J. Comp. Neurol.* **523**, 75–92 (2015).
- 730 42. Arlotta, P. *et al.* Neuronal subtype-specific genes that control corticospinal motor neuron
731 development in vivo. *Neuron* **45**, 207–221 (2005).
- 732 43. Shigemoto, R., Nakanishi, S. & Mizuno, N. Distribution of the mRNA for a metabotropic
733 glutamate receptor (mGluR1) in the central nervous system: an in situ hybridization study in adult
734 and developing rat. *J. Comp. Neurol.* **322**, 121–135 (1992).
- 735 44. Zhou, Q., Choi, G. & Anderson, D. J. The bHLH transcription factor Olig2 promotes
736 oligodendrocyte differentiation in collaboration with Nkx2.2. *Neuron* **31**, 791–807 (2001).
- 737 45. Nielsen, J. A., Berndt, J. A., Hudson, L. D. & Armstrong, R. C. Myelin transcription factor 1
738 (Myt1) modulates the proliferation and differentiation of oligodendrocyte lineage cells. *Mol. Cell.*
739 *Neurosci.* **25**, 111–123 (2004).
- 740 46. Ichihara-Tanaka, K., Oohira, A., Rumsby, M. & Muramatsu, T. Neuroglycan C is a novel midkine
741 receptor involved in process elongation of oligodendroglial precursor-like cells. *J. Biol. Chem.*
742 **281**, 30857–30864 (2006).
- 743 47. Levine, J. M., Reynolds, R. & Fawcett, J. W. The oligodendrocyte precursor cell in health and
744 disease. *Trends Neurosci.* **24**, 39–47 (2001).
- 745 48. Zaitsev, K., Bambouskova, M., Swain, A. & Artyomov, M. N. Complete deconvolution of cellular
746 mixtures based on linearity of transcriptional signatures. *Nat. Commun.* **10**, 2209 (2019).
- 747 49. Qiao, W. *et al.* PERT: a method for expression deconvolution of human blood samples from
748 varied microenvironmental and developmental conditions. *PLoS Comput. Biol.* **8**, e1002838
749 (2012).
- 750 50. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-
751 sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–
752 427 (2018).
- 753 51. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human
754 epigenomes. *Nature* **518**, 317–330 (2015).

- 755 52. Huppertz, B. Placental origins of preeclampsia: challenging the current hypothesis. *Hypertension*
756 **51**, 970–975 (2008).
- 757 53. GBD Chronic Kidney Disease Collaboration. Global, regional, and national burden of chronic
758 kidney disease, 1990-2017: a systematic analysis for the Global Burden of Disease Study 2017.
759 *Lancet* **395**, 709–733 (2020).
- 760 54. Levin, A. & Stevens, P. E. Early detection of CKD: the benefits, limitations and effects on
761 prognosis. *Nat. Rev. Nephrol.* **7**, 446–457 (2011).
- 762 55. Moufarrej, M. N., Wong, R. J., Shaw, G. M., Stevenson, D. K. & Quake, S. R. Investigating
763 Pregnancy and Its Complications Using Circulating Cell-Free RNA in Women’s Blood During
764 Gestation. *Front. Pediatr.* **8**, 605219 (2020).
- 765 56. Pan, W. Development of diagnostic methods using cell-free nucleic acids . (Stanford University,
766 2016).
- 767 57. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression
768 analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
- 769 58. Nygaard, V., Rødland, E. A. & Hovig, E. Methods that remove batch effects while retaining group
770 differences may lead to exaggerated confidence in downstream analyses. *Biostatistics* **17**, 29–39
771 (2016).
- 772 59. Shen-Orr, S. S., Tibshirani, R. & Butte, A. J. Gene expression deconvolution in linear space. *Nat.*
773 *Methods* **9**, 9–9 (2011).
- 774 60. Chang, C.-C. & Lin, C.-J. Training nu-support vector regression: theory and algorithms. *Neural*
775 *Comput.* **14**, 1959–1977 (2002).
- 776 61. Zhong, Y. & Liu, Z. Gene expression deconvolution in linear space. *Nat. Methods* **9**, 8–9; author
777 reply 9 (2012).

778

779

780

781

782 **Tabula Sapiens Consortium**

783 **Overall Project Direction and Coordination**

784 Robert C Jones, Jim Karkanas, Mark Krasnow, Angela Oliveira Pisco, Stephen Quake, Julia Salzman,
785 Nir Yosef

786

787 **Donor Recruitment**

788 Bryan Bulthaupt, Phillip Brown, Will Harper, Marisa Hemenez, Ravikumar Ponnusamy, Ahmad Salehi,

789 Bhavani Sanagavarapu, Eileen Spallino

790

791 **Surgeons**

792 Ksenia A. Aaron, Waldo Concepcion, James Gardner, Burnett Kelly, Nikole Neidlinger, Zifa Wang

793

794 **Logistical coordination**

795 Sheela Crasta, Saroja Kolluru, Maurizio Morri, Angela Oliveira Pisco, Serena Y. Tan, Kyle J. Travaglini,

796 Chenling Xu

797

798 **Organ Processing**

799 Marcela Alcántara-Hernández, Nicole Almanzar, Jane Antony, Benjamin Beyersdorf, Deviana Burhan,

800 Lauren Byrnes, Kruti Calcuttawala, Mathew Carter, Charles K. F. Chan, Charles A. Chang, Alex Colville,

801 Sheela Crasta, Rebecca Culver, Ivana Cvijović, Jessica D'Addabbo, Gaetano D'Amato, Camille Ezran,

802 Francisco Galdos, Astrid Gillich, William R. Goodyer, Yan Hang, Alyssa Hayashi, Sahar Houshdaran,

803 Xianxi Huang, Juan Irwin, SoRi Jang, Julia Vallve Juanico, Aaron M. Kershner, Soochi Kim, Bernhard

804 Kiss, Saroja Kolluru, William Kong, Maya Kumar, Rebecca Leylek, Baoxiang Li, Shixuan Liu, Gabriel

805 Loeb, Wan-Jin Lu, Shruti Mantri, Maxim Markovic, Patrick L. McAlpine, Ross Metzger, Antoine de

806 Morree, Maurizio Morri, Karim Mrouj, Shravani Mukherjee, Tyler Muser, Patrick Neuhöfer, Thi Nguyen,

807 Kimberly Perez, Ragini Phansalkar, Angela Oliveira Pisco, Nazan Puluca, Zhen Qi, Poorvi Rao, Hayley

808 Raquer, Koki Sasagawa, Nicholas Schaum, Bronwyn Lane Scott, Bobak Seddighzadeh, Joe Segal,

809 Sushmita Sen, Sean Spencer, Lea Steffes, Varun R. Subramaniam, Aditi Swarup, Michael Swift, Kyle J

810 Travaglini, Will Van Treuren, Emily Trimm, Maggie Tsui, Sivakamasundari Vijayakumar, Kim Chi Vo,

811 Sevahn K. Vorperian, Hannah Weinstein, Juliane Winkler, Timothy T.H. Wu, Jamie Xie, Andrea

812 R.Yung, Yue Zhang

813

814 **Sequencing**

815 Angela M. Detweiler, Honey Mekonen, Norma Neff, Rene V. Sit, Michelle Tan, Jia Yan

816

817 **Histology**

818 Gregory R. Bean, Gerald J. Berry, Vivek Charu, Erna Forgó, Brock A. Martin, Michael G. Ozawa, Oscar

819 Silva, Serena Y. Tan, Pranathi Vemuri

820

821 **Computational Data Analysis**

822 Shaked Afik, Rob Bierman, Olga Botvinnik, Ashley Byrne, Michelle Chen, Roozbeh Dehghannasiri,

823 Angela Detweiler, Adam Gayoso, Qiqing Li, Gita Mahmoudabadi, Aaron McGeever, Antoine de Morree,

824 Julia Olivieri, Madeline Park, Angela Oliveira Pisco, Neha Ravikumar, Julia Salzman, Geoff Stanley,

825 Michael Swift, Michelle Tan, Weilun Tan, Sevahn K. Vorperian, Sheng Wang, Galen Xing, Chenling Xu,

826 Nir Yosef

827

828 **Expert Cell Type Annotation**

829 Marcela Alcántara-Hernández, Jane Antony, Charles A. Chang, Alex Colville, Sheela Crasta, Rebecca

830 Culver, Camille Ezran, Astrid Gillich, Yan Hang, Juan Irwin, SoRi Jang, Aaron M. Kershner, William

831 Kong, Rebecca Leylek, Gabriel Loeb, Ross Metzger, Antoine de Morree, Patrick Neuhöfer, Kimberly

832 Perez, Ragini Phansalkar, Zhen Qi, Hayley Raquer, Bronwyn Lane Scott, Rahul Sinha, Hanbing Song,

833 Sean Spencer, Aditi Swarup, Michael Swift, Kyle J. Travaglini, Jamie Xie

834

835 **Tissue Expert Principal Investigators**

836 Steven E. Artandi, Philip Beachy, Michael F. Clarke, Linda Giudice, Franklin Huang, KC Huang, Juliana

837 Idoyaga, Seung K Kim, Mark Krasnow, Christin Kuo, Patricia Nguyen, Stephen Quake, Thomas A.

838 Rando, Kristy Red-Horse, Jeremy Reiter, Justin Sonnenburg, Bruce Wang, Albert Wu, Sean Wu, Tony
839 Wyss-Coray

840

841

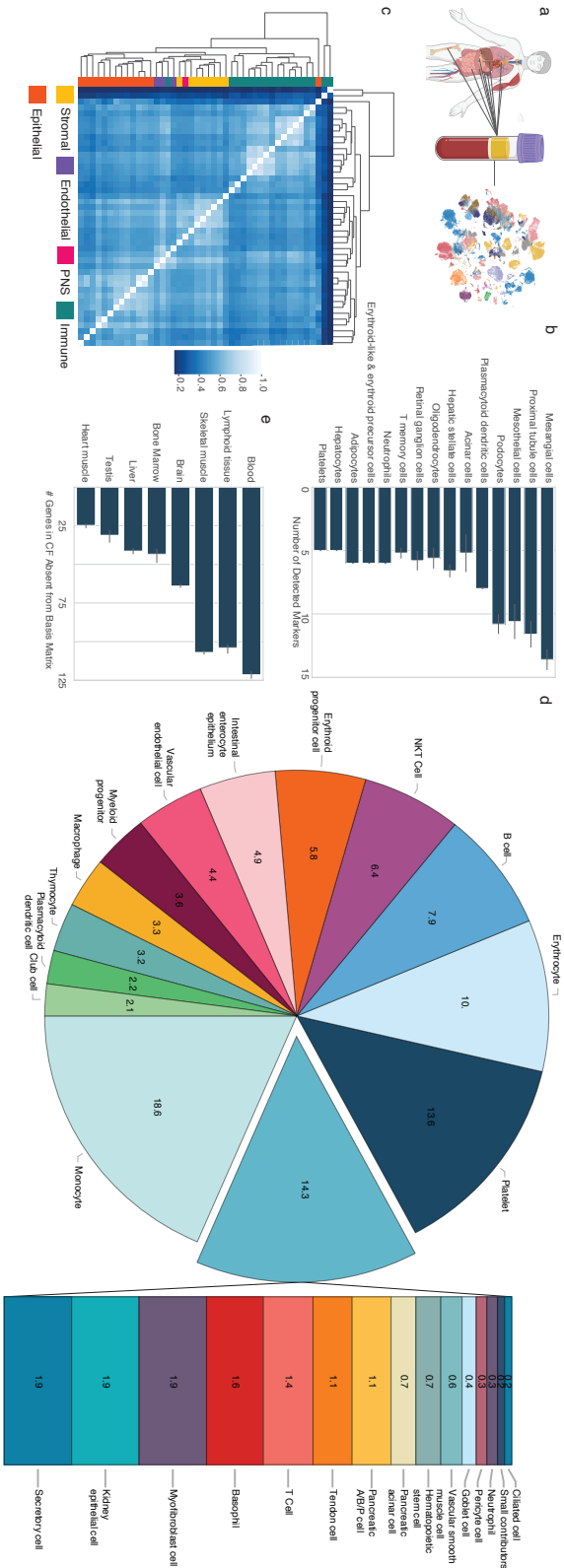


Fig. 1 : Cell type decomposition of the plasma cell free transcriptome using Tabula Sapiens reveals the comprehensive cell type-of-origin landscape

- Integration of tissue-of-origin and single cell transcriptomics to identify cell types-of-origin in plasma cfRNA.
- Cell type specific markers defined in context of human body identified in cfRNA (n = 5). Error bar denotes variance among intersection size across biological replicates.
- Cluster heatmap of Spearman correlations of cell types in basis matrix derived from Tabula Sapiens. PNS denotes 'peripheral nervous system'.
- Mean fractional contributions of cell type specific RNA in plasma cf-transcriptome (n = 5). Full distributions of learned coefficients across biological samples are available in fig. S5.
- Top tissues in cfRNA not captured by basis matrix (e.g. the set difference of all genes detected in a given cfRNA sample and the row space of the basis matrix intersected with HPA tissue specific genes).

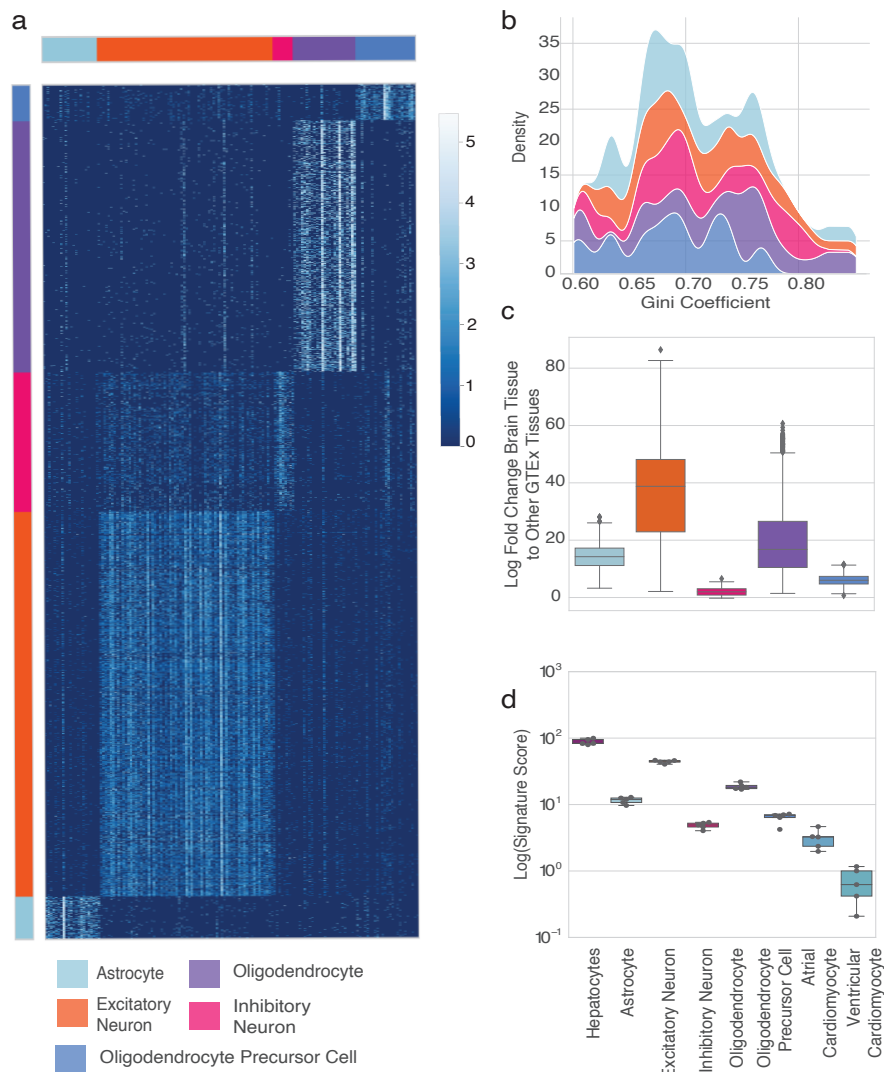


Fig. 2: Individual single cell atlases identify additional cell types-of-origin in the plasma cell free transcriptome

a. Heatmap of gene expression in NCI prefrontal cortex single cell data indicating cell type specificity across brain cell types. Rows are individual cells; columns are a given marker gene. Values are log-transformed counts per ten thousand.

b. Density plot of the Gini coefficients for the genes within a given cell type gene profile. The Gini coefficient for a given gene was computed using normalized counts across all HPA tissues. Area under curve for a given cell type sums to one.

c. Validation of specificity of cell type specific genes across the tissues of the human body. Log-fold change in total log-transformed counts-per-ten-thousand expression of genes for a given cell type signature across all GTEx tissues over the mean non-brain aggregate cell type signature score assess overall abundance in expression.

d. Hepatocyte, neuronal, glial, and cardiomyocyte cell type signature scores in healthy cfRNA plasma (n = 5) on a logarithmic scale.

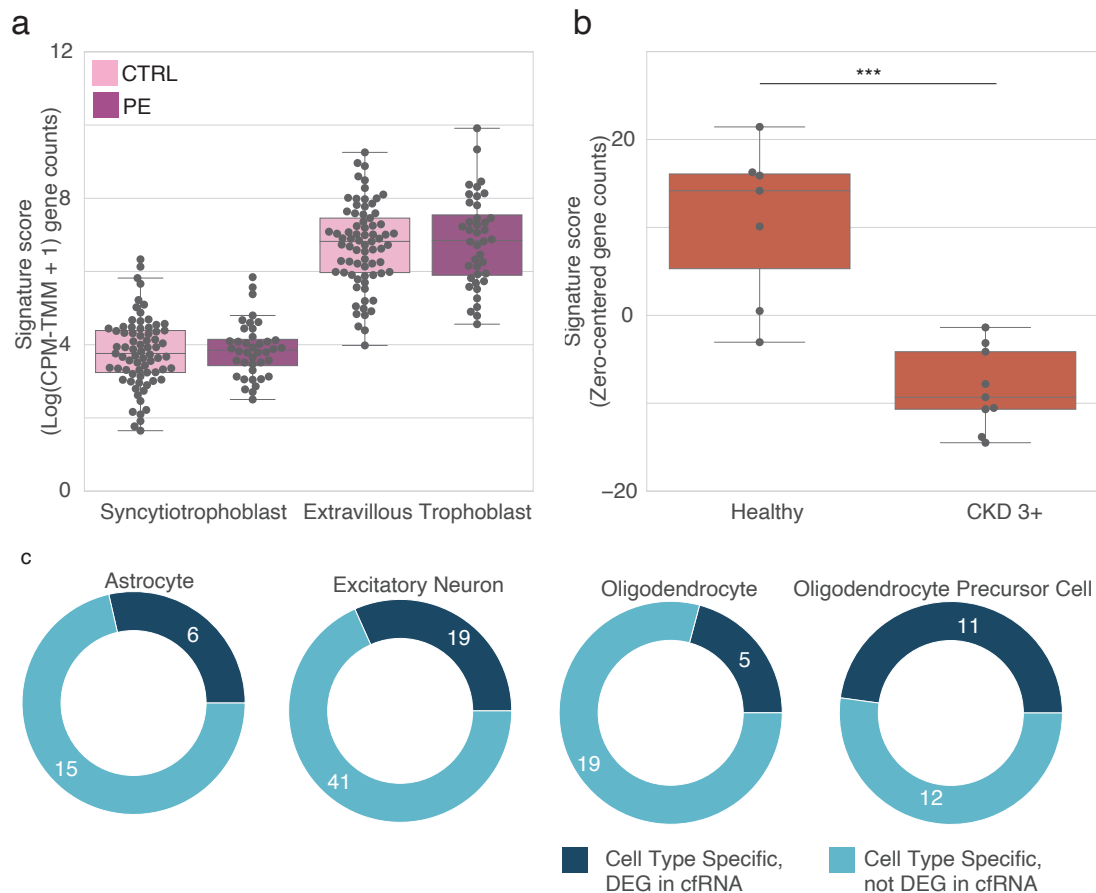
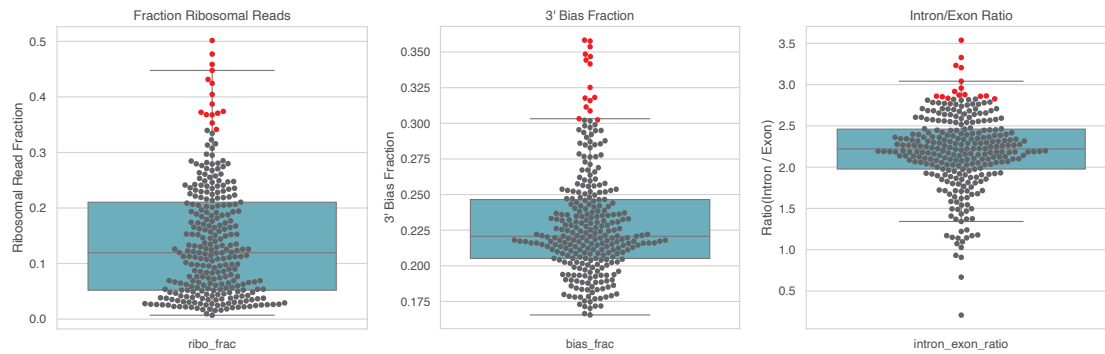


Fig. 3: Cell type signatures reveal implicated cell types in pathologic cfRNA samples

- Comparison of signature scores for extravillous trophoblasts and syncytiotrophoblasts in cfRNA samples (Munchel et al. 2020, iPEC cohort) from mothers with no complications ($n = 73$) as compared to mothers with preeclampsia ($n = 40$). Gene counts were CPM-TMM normalized and then natural log transformed and summed to determine the signature score.
- Comparison of signature scores for proximal tubule cells in cfRNA samples from CKD stages 3+ ($n = 9$) and healthy patients ($n = 7$) using data from Ibarra et al (***) denotes $p < 10^{-3}$). Gene counts were CPM-TMM normalized and then zero centered across conditions and summed to determine the signature score.
- Proportion of cell type specific genes from brain cell type single cell data intersecting with reported DEGs downregulated in AD-derived plasma⁴.

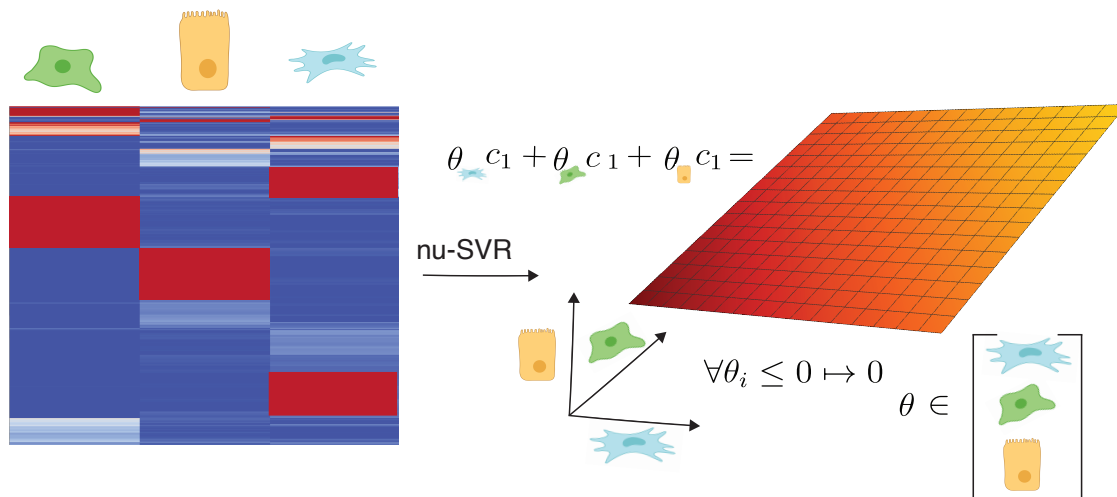
844

845



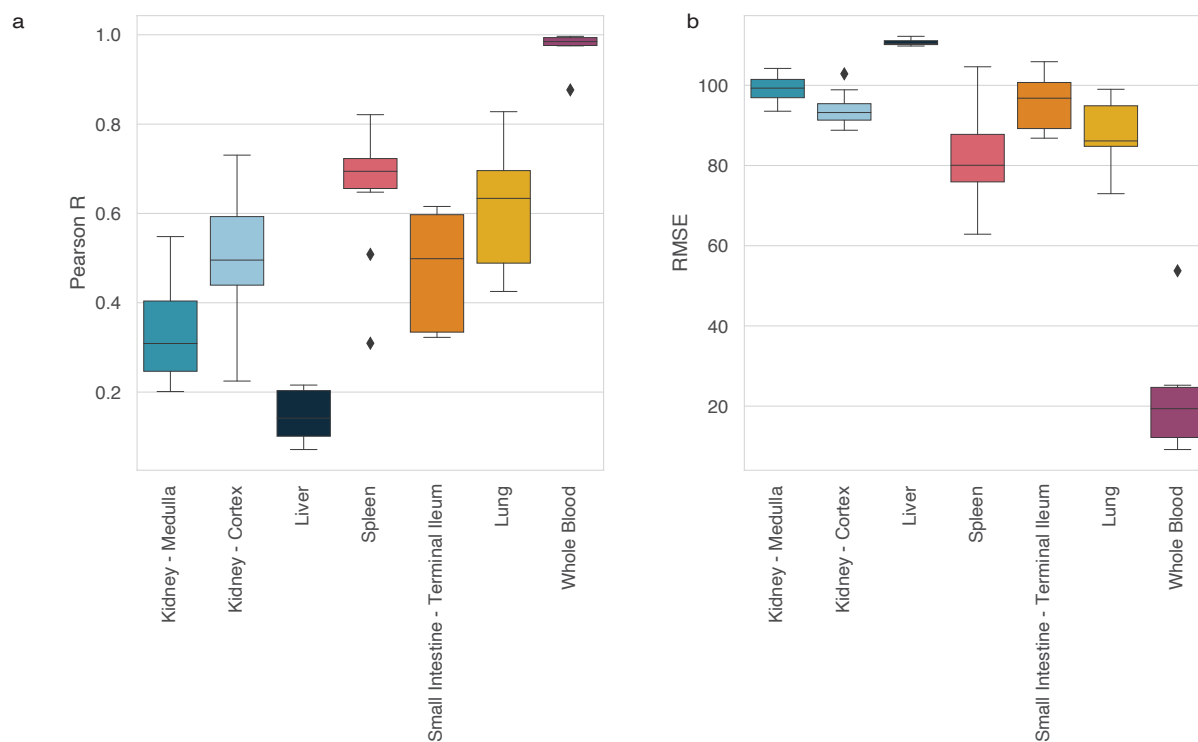
846
847
848
849
850
851

Fig S1 Identification of samples with outlier values for at least one quality control metric including a measure of RNA degradation, ribosomal fraction, and DNA contamination from Ibarra et al. Samples with outlier values are highlighted in red. (See Methods section ‘Data Processing’ for details)



852
853
854
855
856

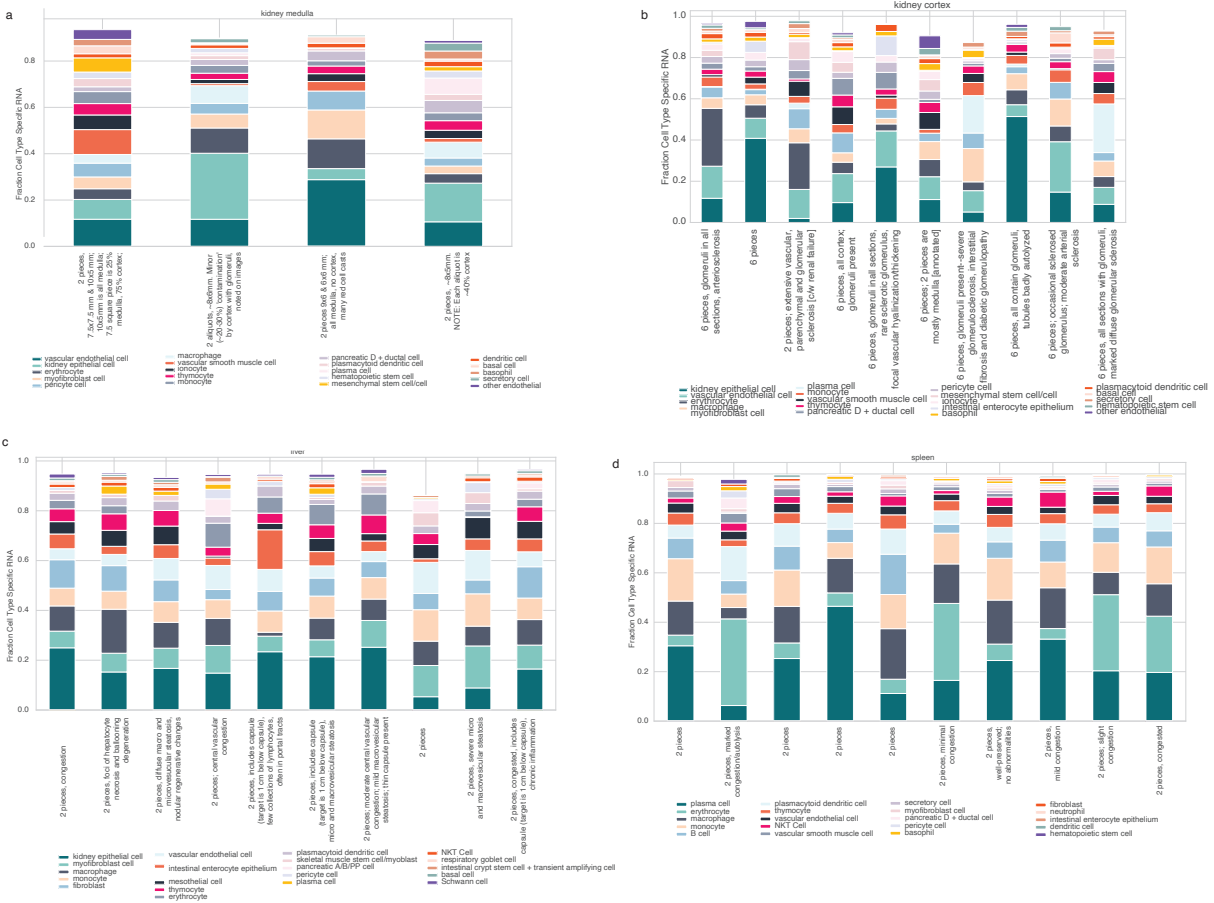
Fig. S2. Schematic overview of nu-SVR deconvolution in trivial 3-dimensional cell type dimensional space denoting the learning of a hyperplane in cell type dimensional space and subsequent normalization to infer relative fractions of cell type specific RNA.

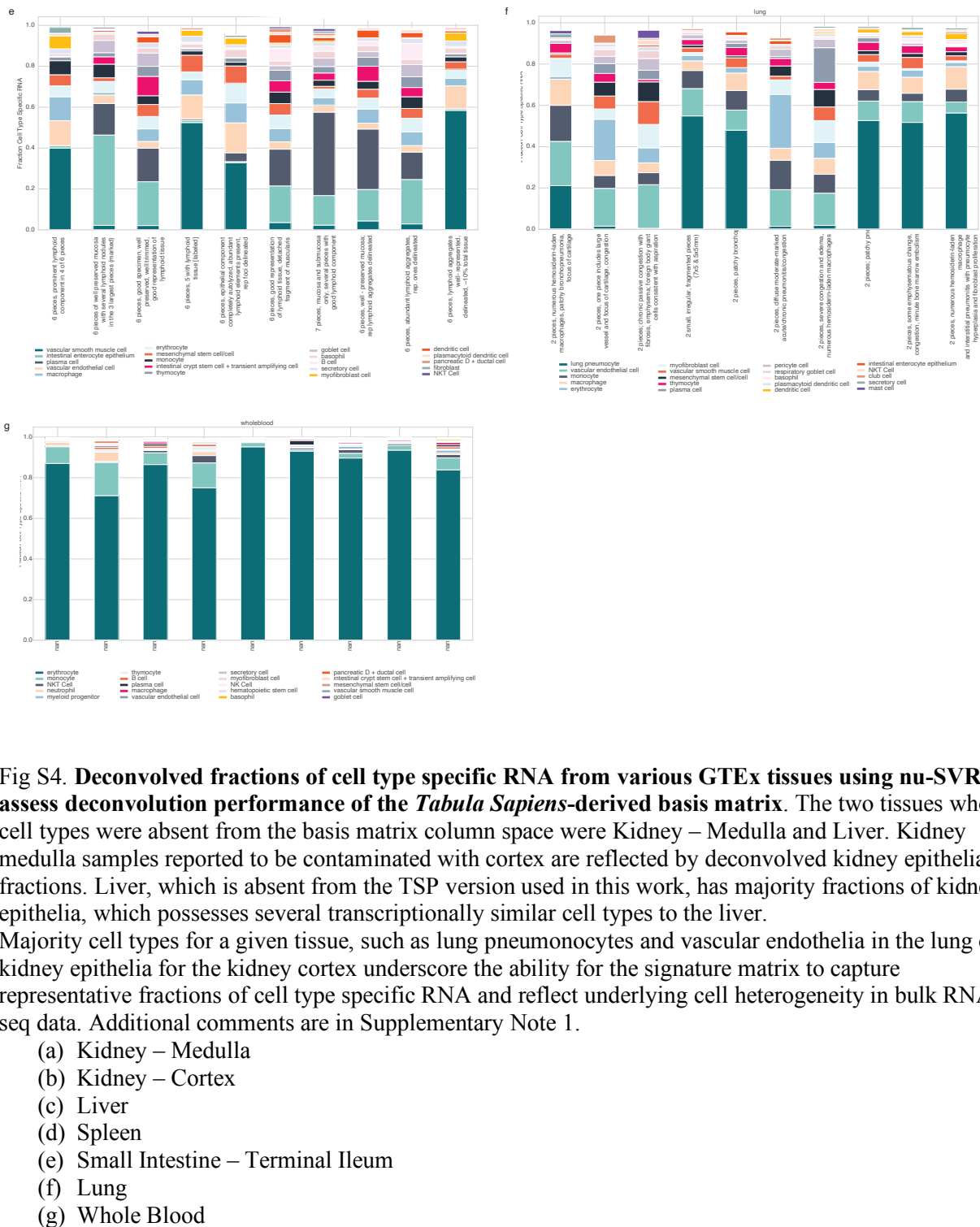


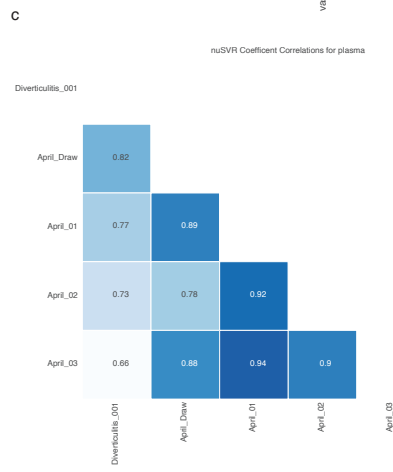
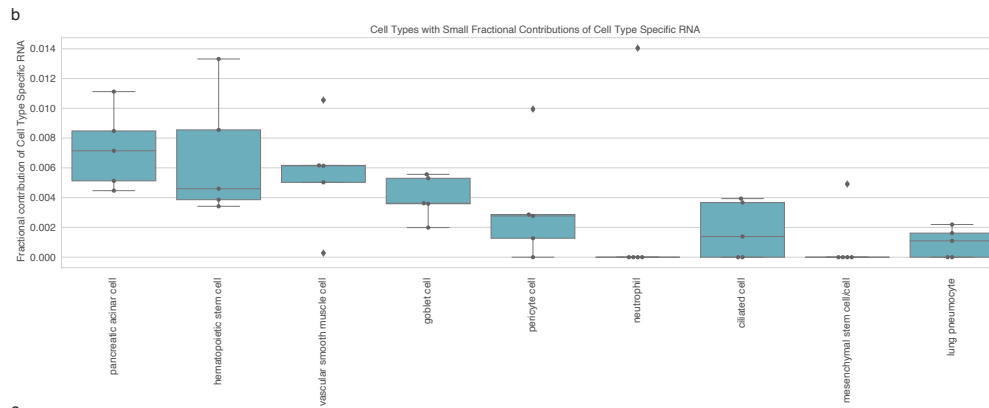
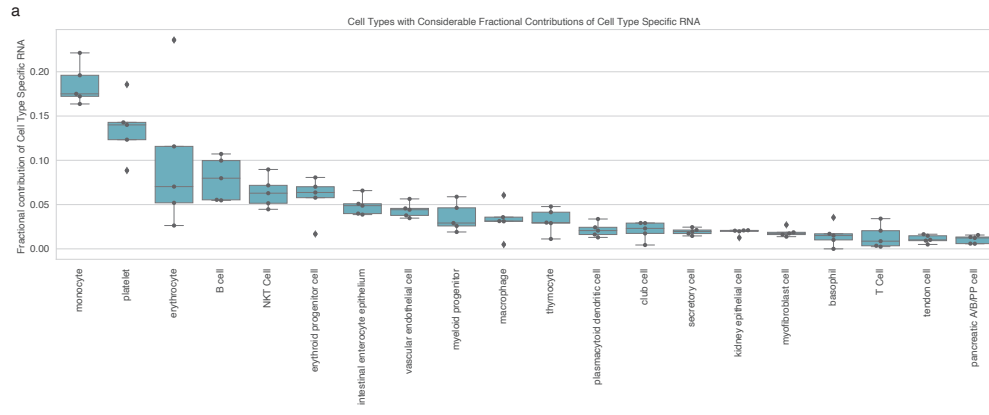
857
858
859
860
861
862
863
864

Fig. S3. Basis matrix performance on GTEx bulk RNA samples using nu-SVR. GTEx tissue samples possessing cell types wholly present (Kidney – Cortex, Spleen, Small Intestine – Terminal Ileum, Lung, Whole Blood) and absent from the basis matrix column space (Kidney – Medulla, Liver) were selected.

- (a) Pearson correlation between predicted expression and actual expression in cfRNA
(b) Root Mean Square Error between predicted expression and actual expression in cfRNA. Units are zero-mean unit variance scaled CPM counts; tissues present in TSP have reduced RMSE compared to those that are absent (e.g. Kidney – Medulla and Liver)

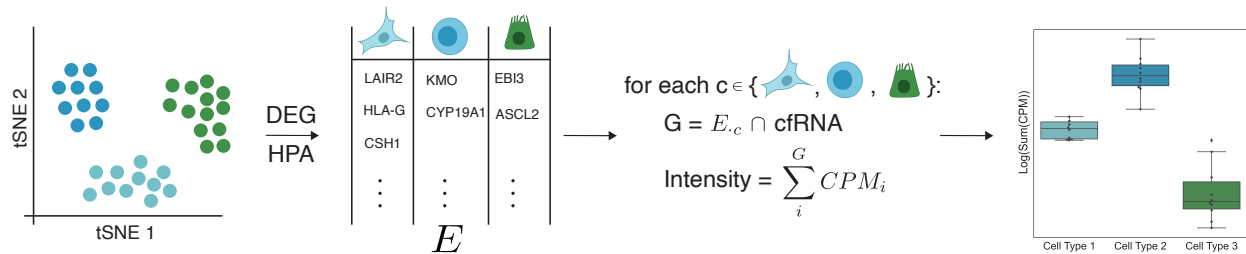






885 **Fig. S5. Deconvolving the plasma cell free transcriptome using nu-SVR**

- 886 (a) Distribution of learned relative fractional contributions with mean cell type proportions were >
 887 1%
- 888 (b) Distribution of learned coefficients of small contributions using nu-SVR across samples ($\leq 1\%$).
 889 ‘Small contributors’ slice in main text Fig 1D reflects cell types with cumulative fractions < 0.1%
 890 (mesenchymal stem cell/mesenchymal cell and lung pneumocyte).
- 891 (c) Pairwise Pearson correlations of nu-SVR learned coefficients between biological replicates
 892
 893
 894



895 **Fig S6. Cell type signature score derivation overview.** See ‘Signature Scoring’ section of methods for
 896 filtering criteria and thresholds.
 897
 898
 899

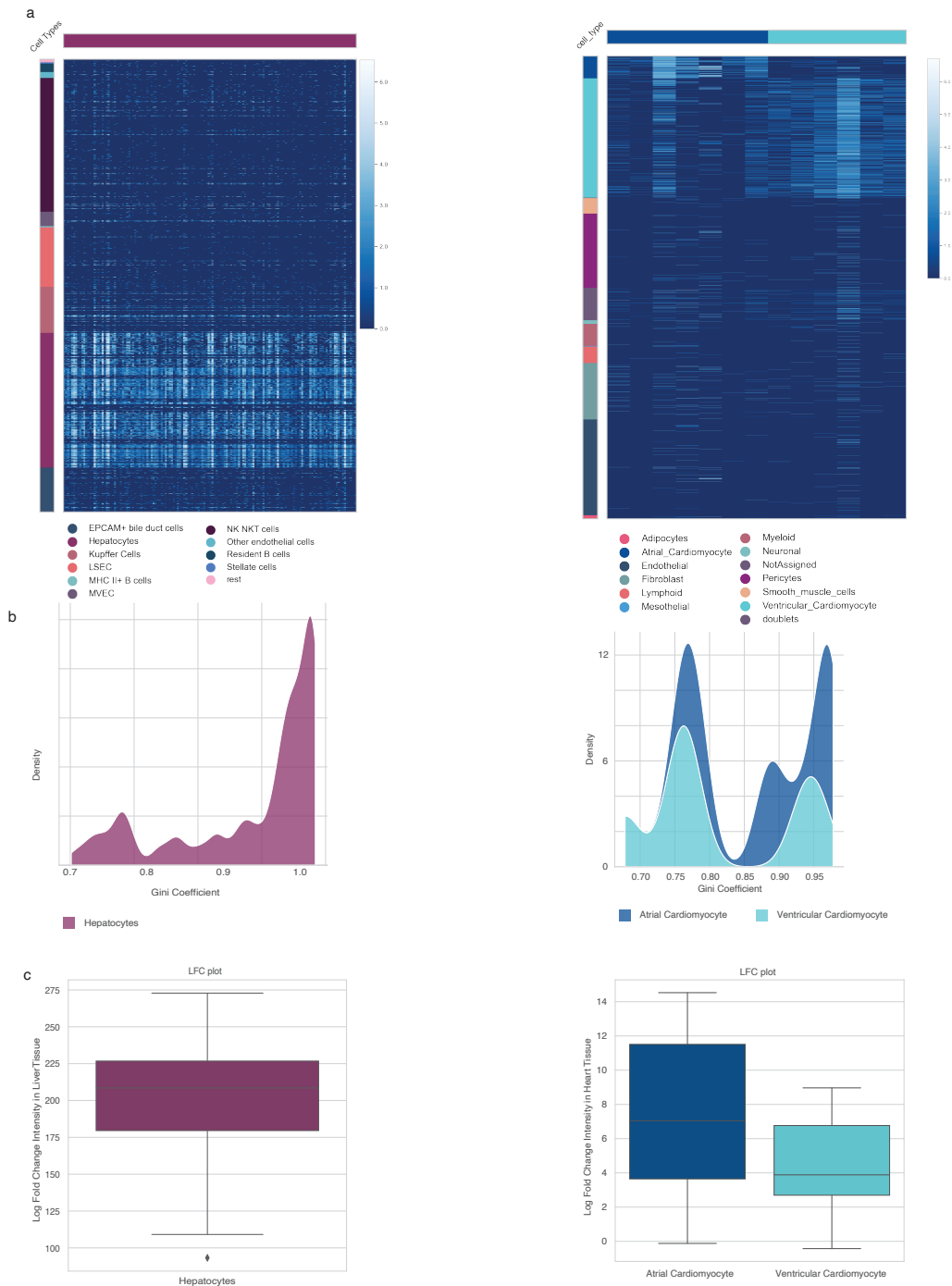


Fig S7. Establishing gene profile cell type specificity in context of the whole human body using single cell and bulk RNAseq data

(a) Single cell heatmaps for gene cell type profiles within the corresponding tissue cell atlas, demonstrating that predominant expression in bulk data is within the cell type of interest.

(b) Gini coefficient density plot for genes in cell type profiles derived from liver, testis, and heart single cell atlases using HPA NX counts.

(c) Log fold change in bulk RNA-seq data of the cell type profile

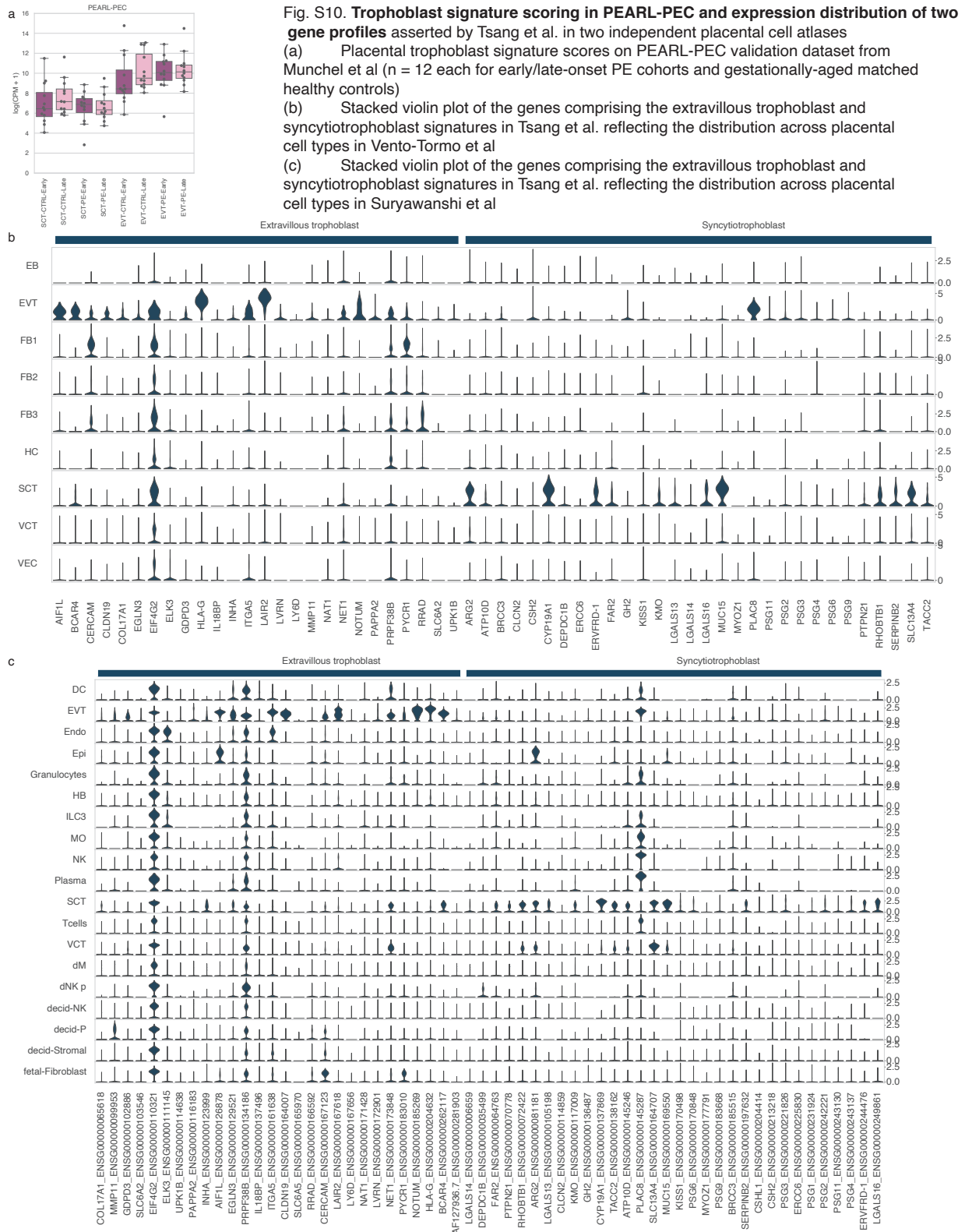


Fig S8. Distribution of Gini coefficient and Tau for all genes denoted by HPA as specific to the brain, liver, placenta, kidney, and heart



Fig. S9. **Comprehensive placental and renal cell type gene profile specificity** at single cell and whole body resolution
 (a) Violin plot of derived syncytiotrophoblast and extravillous trophoblast from Vento-Tormo et al.
 (b) Violin plot of derived syncytiotrophoblast and extravillous trophoblast from Suryawanshi et al.
 (c) Violin plot of derived proximal tubule and podocyte markers
 (d) Gini coefficient distribution for placental trophoblast cell types in (a) and (b)
 (e) Gini coefficient distribution for renal cell types in (c)
 (f) Distribution of placental trophoblast signature scores across all GTEx tissues, since the placenta is not in GTEx, so the values plotted are just the aggregate expression of genes in a given signature.
 (g) Log fold change of renal cell type intensity in GTEx Kidney Cortex/Medulla samples (sum of log-transformed counts-per-ten thousand) relative to the mean non-kidney signature score intensity.

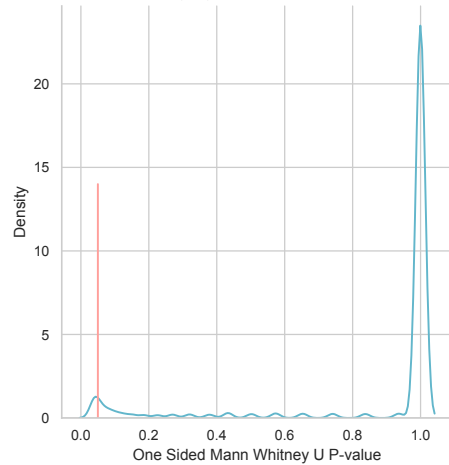
902
903
904
905



906

907

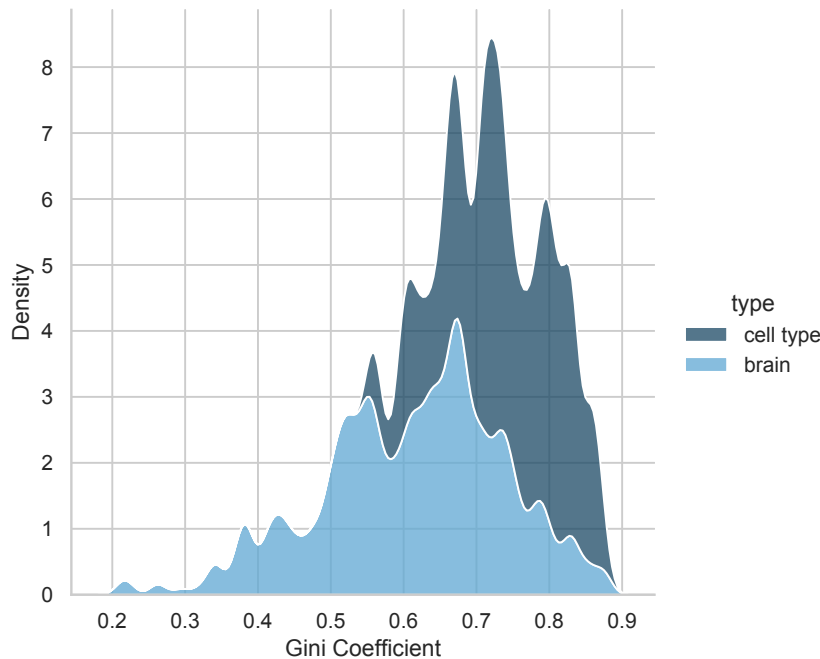
BENJAMINI HOCHBERG ADJUSTED ALPHA = 0.05, 10,000 trials of Difference in Random Signature between CKD and Healthy



908
909
910
911
912
913
914
915
916
917
918
919
920

Fig. S11. **Validation of the discriminatory power in the proximal tubule signature** in CKD stage 3+ vs healthy samples

Distribution of one-sided Mann Whitney U in 10,000 trials of comparing a random signature score of an equivalent number of genes as the proximal tubule signature with an alternative hypothesis that the signature in healthy is greater than CKD. Red vertical line denotes the $p = 0.05$ threshold. Multiple hypothesis testing correction using Benjamini Hochberg with $FDR = 0.05$ was performed and this yielded an adjusted p-value of 0.038 between the CKD 3+ and healthy groups using the actual proximal tubule signature score. In 96.6% of the trials, no significance was observed in the random signature discriminating between sick and healthy, indicating the specificity of our signature score in discriminating between CKD and healthy patients.



921
922
923
924

Fig. S12. **Comparison of brain specific DEG and cell type specific DEG.** Distribution in gini coefficients for AD downregulated DEG in Toden et al. that are brain-specific and cell type specific respectively. Area under curve for each group sums to 1.