1 **DevKidCC allows for robust classification and direct comparisons of**

2 **kidney organoid datasets**

3 Sean B. Wilson[1], Sara E. Howden[1,2], Jessica M. Vanslambrouck[1], Aude Dorison[1], Jose

4 Alquicira-Hernandez[3], Joseph E. Powell[3,4], Melissa H. Little[1,2,5*]

5

6

7 1. Murdoch Children's Research Institute, Flemington Rd, Parkville, VIC, Australia

8 2. Department of Paediatrics, The University of Melbourne, VIC, Australia.

9 3. Garvan-Weizmann Centre for Cellular Genomics, The Kinghorn Cancer Centre, NSW,

10 Australia

11 4. UNSW Cellular Genomics Futures Institute, University of New South Wales, NSW,

12 Australia

13 5. Department of Anatomy and Neuroscience, The University of Melbourne, VIC, Australia.

14

15 * Author for correspondence:

16 M.H.L.: +61 3 9936 6206; melissa.little@mcri.edu.au

17 Key words: cell identity prediction, human developing kidney, kidney organoid

18

19

20

21

22

23

**Abstract**

Kidney organoids provide a valuable resource to understand kidney development and disease. Clustering algorithms and marker genes fail to accurately and robustly classify cellular identity between human pluripotent stem cell (hPSC)-derived organoid datasets. Here we present a new method able to accurately classify kidney cell subtypes, a hierarchical machine learning model trained using comprehensive reference data from single cell RNA-sequencing of human fetal kidney (HFK). We demonstrate the tool's (*DevKidCC*) performance by application to all published kidney organoid datasets and a novel dataset. *DevKidCC* is available on Github and can be used on any kidney single cell RNA-sequence data.

**Background**

36  Single cell RNA sequencing has reformed how we approach biological questions at the

37  transcriptional level, facilitating accurate evaluation of cellular heterogeneity within complex

38  samples, including entire tissues. When coupled with approaches for molecular lineage

39  tagging[1] and computational approaches to analyse pseudotime[2–4] and RNA velocity[5,6], gene

40  expression in complex tissues such as kidney can be studied at an unprecedented resolution.

41  Despite these advantages, classification of cellular identity remains challenging and variable

42  between data, even when analysing similar cellular systems. Currently, a common approach

43  for identifying cell populations within single cell data is to first cluster cells, compute

44  differentially expressed genes between clusters, and label clusters of cells based on

45  expression of known marker genes[4,7,8]. The choice of clusters can be arbitrary, with users

46  defining the number of clusters, thereby raising the potential for biases in the reproducibility

47  of cell-type labels[9]. Placement of cells into a cluster relies on transcriptional similarity[10],

48  hence there needs to be a large enough population with a distinct gene signature for this to

49  occur. Cell clusters are also commonly defined based upon one or a few known differentially

50  expressed genes rather than their global transcriptional signature. Finally, technical

51  challenges such as batch variation can impact definitive cellular identification.

52  The application of single cell profiling to developmental biology presents unique challenges

53  due to the presence of intermediate cell types undergoing differentiation during

54  morphogenesis. The mammalian kidney contains more than 25 cell types in the mature

55  postnatal tissue, arising from a smaller number of progenitor cell types including nephron,

56  stromal, endothelial and collecting duct progenitors. Organogenesis is driven via reciprocal

57  signalling and self-organisation with many intermediate transcriptional states that are less

58  well defined, making the classification of cell types at the single cell level both extremely

59  useful but particularly difficult (reviewed in Little and Combes, 2019[11]). This is further

60    complicated with hPSC-derived kidney organoid datasets. While protocols for differentiating

61    kidney organoids from hPSC attempt to replicate *in vivo* kidney differentiation, they are

62    likely to be limited and contain emerging non-specific, off-target or synthetic cell types[12–15].

63    Here, unbiased classification of cellular identity is a computational challenge. Indeed, recent

64    single cell profiling of cell human fetal kidney (HFK) datasets have shown that the classical

65    canonical markers for many cell identities within the kidney are not unique to these cell types

66    but are also expressed at lower levels within other populations[15–18]. This makes cell

67    classification in organoids more challenging when analysing gene expression of these

68    markers in the single cell clusters. The ability to robustly identify and classify cells in hPSC-

69    derived organoid data is crucial to facilitate useful comparisons between datasets, particularly

70    data generated using different differentiation protocols and cell lines but also in response to

71    mutation or perturbation. These analyses will also help to improve and refine protocols

72    towards a more accurate endpoint tissue.

73    One approach to cellular identification is to apply a small set of 'known' genes to identify

74    clusters within a dataset based upon an existing reference dataset that has been accurately

75    classified. Of the 12 kidney organoid single cell RNA-seq datasets published to date (Table

76    1), seven used a HFK reference to find congruence with their clustered organoid populations

77    either through integration or training a unique random forest classifier. However, there have

78    been many different references used across these publications. Cell classifications may be

79    inconsistent when using various references containing different proportions of cells, possibly

80    captured at different ages or regions of the tissue. Indeed, the most commonly used HFK

81    reference only contained cells from the cortex of a 16-week kidney and hence was reported to

82    contain few nephron cells and no ureteric epithelium[19]. There have been many tools

83    developed to utilise reference data to classify a related query dataset, with scrna-tools.org[4]

84    listing 85 tools in the "Classification" category. These tools extract cell type information

85    from an annotated reference and apply that to a query dataset. Most rely upon the user to

86    supply the reference data and for those that supply a reference, none are directly relevant to

87    hPSC-derived kidney organoids. The *R* packages *scTyper*[20] and *scClassify*[21] are trained on

88    existing datasets. These are not ideal for human developing kidney classification as

89    *scClassify* is trained on mouse cell data, while *scTyper* contains gene sets of limited cell types

90    of the adult kidney and is thus not a developing kidney cell population. As such there is a no

91    tool that can be used to directly and accurately classify the cell types present within the

92    developing human kidney.

93    Here we have taken reference HFK datasets from three publications that span multiple ages

94    and kidney regions, performed individual annotations of the cells present based on prior

95    information, then used all confidently classified cells to train classification models using the

96    *R* package *scPred*[22], a generalizable method which has showed high accuracy in different

97    experiments and datasets from multiple tissues, and considered a top performer in

98    benchmarking studies[9]. The resulting model, referred to as *DevKidCC*, provides a robust and

99    accurate classification of cells in novel single cell datasets generated from developing human

100   kidney or stem cell-derived kidney organoids. *DevKidCC* defines a model of cellular identity

101   organised in a hierarchical manner to represent the key developmental trajectories of lineages

102   within the developing kidney. The classification method is complemented with custom

103   visualization tools in the *DevKidCC* package. This classifier was then used to investigate

104   published kidney organoid datasets to compare organoid patterning and gene expression

105   profiles across these datasets. We present a variety of applications of *DevKidCC* to the

106   reanalysis of existing data. This analysis revealed differences in nephron progenitor

107   proportion and nephron patterning and maturation between kidney organoid protocols. We

108   also apply *DevKidCC* to investigate approaches for directed differentiation to ureteric

109   epithelium and dissect the effect of all-trans retinoic acid on nephron patterning and podocyte

5

110    maturation. While *DevKidCC* is specifically trained on HFK for application to kidney

111    organoid models, the framework presented here could be applied for any tissue system to

112    generate a cell classification model.

113

**Results**

115    **Generation of the model hierarchy for complete cell classification**

116    We first build a comprehensive reference dataset on which to train the probabilistic

117    classification models. We used single cell RNA-sequence datasets from three publications

118    currently generated on HFK (Table 1). Samples range from 9 to 19 weeks' gestation across

119    which time the developing human kidney undergoes both growth and maturation, with week

120    16 being most frequently represented. Cells in all were originally annotated using clustering

121    and cluster labelling using marker gene expression. One dataset was a recently published high

122    quality HFK dataset[23] (8,987 cells) that included both medulla and cortex regions and

123    including a 96-day male and 108-day female sample. Of note, this dataset contained ureteric

124    epithelium, which had not been thoroughly analysed to this point[24]. This data was combined

125    with data from 17,759 HFKs cells ranging from week 11 to 18 of gestation[25] to increase the

126    developmental range of the training set. A further 8,317 cells from gestational week 17 which

127    had been microdissected into cortex, inner and outer medullary zones[26] were combined to

128    complete the comprehensive reference single cell RNA-sequencing HFK dataset. Cells from

129    all datasets were integrated using *Harmony*[27] (Figure 1A) before performing a supervised

130    clustering and annotation, using the original annotations of each dataset as a guide. This led

131    to a reference dataset containing three ureteric epithelial subpopulations (Tip, OuterStalk,

132    InnerStalk), four stromal subpopulations (Stromal Progenitor Cells (SPC), Cortex, Medullary,

133    Mesangial), endothelium, the nephron progenitor cells (NPC) and the nephron including

Table 1: Summary of existing kidney related single cell datasets

### Human Fetal Kidney

| Reference | Age (post coitum) | Sample Details |
|---|---|---|
| Lindstrom et al.[19] | 16 weeks | MARIS dissociation used to isolate cortical regions |
| Menon et al.[28] | 87-132 days | Cells from cortical nephgenic zone to inner medullary region |
| Young et al.[29] | 1.9 and 2.1 months | Biopsies of 1.9 and 2.1 month fetal kidneys |
| Hochane et al.[25] | 9, 11, 13, 15, 18 weeks | Week 9, 11, 13, 16 and 18 kidney pieces |
| Tran et al.[26] | 15 and 17 weeks | Regions dissected from both inner and outer cortex |
| Holloway et al.[23] | 16 weeks | Wedge biopsy including both medulla and cortex, one day 96 male and one day 108 female sample |

### Kidney Organoids

| Reference | Age (days) | Sample information | ID | Classification |
|---|---|---|---|---|
| Wu et al.[12] | 26 | 4 batches of iPS and 2 batches of ES derived organoids using Takasato[30] protocol | Wu_T | Clustering & DE genes, Integration with self-generated adult snRNA dataset, Lindstrom[19] trained random forest classifier |
| | 26 | 3 batches of iPS and 1 batch of ES derived organoids using Morizane[31] protocol | Wu_M | |
| | 34 | Older iPS derived organoid using Takasato protocol | Wu_TO | |
| | 7, 12, 19, 26 | Time course of iPS derived organoids using Takasato protocol | Wu_TC | |
| | 26 | 2 batches of iPS derived Takasato organoids with BDNF inhibition | Wu_TB | |
| Czerniecki et al.[32] | 25 | Freedman iPS and ES derived organoids, modified protocol for High Throughput Sequencing, comparing with/without VEGF | Cz_F | Clustering & DE genes, Menon |
| Howden et al.[13] | 18, 25 | Takasato iPS derived organoids using E6 base media | How_T | Clustering & DE genes |
| Phipson et al.[33], Combes et al.[15] | 25 | Takasato iPS derived organoids generated in two batches. Same dataset in both publications | PC_T | Clustering & DE genes; Integration with Lindstrom[19] |
| Harder et al.[34] | 19 | Freedman ES derived organoids, 6 datasets generated from the all organoids in a well, 3 separate batches | Har_F | Clustering & DE genes, integration and trajectory analysis with Menon[28] |
| | 20 | A single Freedman ES derived organoid isolated from a full well | Har_F_SO | |
| Subramanian et al.[14] | 7, 15, 29, 32 | Takasato iPS derived organoids with 3 pooled replicates per time using iPS cell line designated "ThF" | Sub_T_L1 | Clustering & DE genes, organoid trained random forest classifier, integration with Young[29], Lindstrom[19] and self-generated kidney tissue |
| | 7, 15, 29 | Takasato iPS derived organoids with 3 pooled replicates per time using iPS cell line designated "AS" | Sub_T_L2 | |
| Kumar et al.[35] | 25 | Modified Takasato iPS derived micro-organoid | Ku_TMO | Integration with organoid[15,33] with clustering & DE genes |
| Low et al.[36] | 10, 12, 14 | Modified Takasato ES derived organoids | Low_TMod | Clustering & DE genes |
| Tran et al.[26] | 16, 28 | Morizane ES derived organoids | Tran_M | Clustering & DE genes individually & after integrating with self-generated kidney tissue |
| Lawlor, Vanslambrouck, Higgins et al.[37] | 25 | Takasato iPS derived organoids generated by bioprinting. Organoids were compared with three different biophysical properties. | LVH_T | Clustering & DE genes compared to Hochane[25] trained machine learning model using scPred |
| Howden, Wilson et al.[24] | NA | Takasato iPS derived organoids dissociated and GATA3+EPCAM+ cells isolated. These cells cultured in ureteric epithelium promoting conditions. | HW_iUB | Seurat Label Transfer using reanalysed Holloway |
| Mae et al.[38] | NA | Induced Ureteric Bud cultures | Mae_iUB | Clustering & DE genes |

134     subpopulations of CellCycle (CC), EarlyNephron (EN), early distal and medial tubule

135     (EDT_EMT), distal tubule (DT), Loop of Henle (LOH), early proximal tubule (EPT),

136     proximal tubule (PT), parietal epithelial cells (PEC), early podocytes (EPod) and podocytes

137     (Pod) (Supplementary 1). These populations have been further classified in the original

138     publications, such as the DT being split into distal straight, distal convoluted and connecting

139     segment or classifying populations in relation to morphological features, such renal vesicle,

140     comma shaped body and S-shaped body segmentation[24–26]. While morphologically there is a

141     consistency in segment identification, this is less clear in single cell data and has led to

142     inconsistency in classification terminology. As such, here we have classified cell populations

143     based on expression of known differentiation markers as cells take on a more distinct identity

144     (Figure 1B).

145     The complex and dynamic nature of the developing kidney, with multiple cell lineages and

146     waves of nephrogenesis, means that cells of many stages of differentiation can be present at

147     all developing timepoints within the same single cell data. This is one of the main challenges

148     in classifying cells in the HFK single cell data, as the cells are in transitional flux. The

149     multiple lineages within the kidney also make classifying cell types difficult, as the

150     differences between lineages mask the subtle differences in gene expression between cell

151     types within a lineage, such as those of the epithelial sub-types. To minimise the impact of

152     this transcriptional variance on classification, we took a hierarchal approach by training three

153     tiers of models (Figure 1C). The first tier classified cells based on their lineage; nephron

154     progenitor cells (NPC), nephron, ureteric epithelium (UrEp), stroma and endothelial. The

155     second tier for the UrEp lineage classified cells into the highly proliferative Tip cells, *AQP2*-

156     expressing outer stalk and uroplakin-expressing inner stalk. The second tier for the stroma

157     lineage classified cells into the *FOXD1*-expressing stromal progenitors (SPC), the cortical

158     and medullary stroma clearly identifiable in the outer and inner zones[26] (*DCN* low/high

159    respectively) and the mesangial cells which express *GATA3*. The nephron segmentation

160    required an extra tier due to the complexity of cell types present and their transcriptional

161    similarity. Here, the second tier classified the early nephron (EN) that could not be clearly

162    identified as polarised, the proximal (PN) and distal (DN) nephron epithelium and the renal

163    corpuscle (RC) lineage. These were then further classified in third tier models (Figure 1C).

164    The models were trained with the package *scPred*[22] using a support vector machine with a

165    radial basis kernel and 100 principal components. The *scPred* package utilises a machine

166    learning approach to train predictive models on a reference single cell dataset. This model

167    can estimate the similarity of a cell within a query dataset to the identities classified within

168    the model. This has been shown to be a robust method to classify cells of a novel dataset

169    based on a known reference[9,37]. We created wrapper functions of all the models into a single

170    use function (*DevKidCC*), which takes an input of a *Seurat* object. To determine cells in the

171    first tier, we use a probability threshold of 0.7, while at all other tiers the threshold is

172    removed. This enables all cells that are classified at the top tier to be given an assigned

173    identity regardless of the highest level of similarity predicted by the lower tier models.

174    Further investigation of the calculated similarity value can be interrogated as every cell has a

175    record in the metadata of the scores from each classification. No pre-processing is required as

176    data is normalised during the function call. The recommended pipeline is to read in raw

177    counts data using the *Seurat* pipeline, filter out poor quality cells and then run *DevKidCC*.

178    The classifications for each tier and the final identities can be accessed within the metadata

179    slot for further investigation. The package contains custom in-build functions *ComparePlot*,

180    *DotPlotCompare* and *SankeyPlot* to investigate the cell populations within the classified

181    sample.

182

183    **DevKidCC classification rapidly and accurately reproduces published annotations**

9

184    While this tool was designed to classify cells within kidney organoids, we first confirm the

185    capacity to accurately classify developing kidney cell types by applying it to other HFK

186    datasets. We applied *DevKidCC* to the dataset of Lindstrom[19] of which the original cell

187    classification identities were equivalent to those of the first classification tier within

188    *DevKidCC* (Figure 2A). *DevKidCC* classified 90% of the 2945 cells that passed quality

189    control, while the remaining cells expressed markers for immune cells (*HLA-DRA*, *CCL3*,

190    *SRGN*) which are not represented in the model and so were not assigned an identity. 14 cells

191    were classified as UrEp, positioned at the tips of one end of the nephron cluster, which

192    *DevKidCC* further classified as DN epithelium. While these two cell populations arise from

193    distinct precursors, they share a very similar transcriptional profile, making them very

194    difficult to distinguish at single cell level[15–18,24]. The ability to identify and classify these two

195    populations separately, even with a small contribution of one population within a dataset,

196    demonstrates the power of *DevKidCC* as a classification tool, particularly in comparison to

197    clustering algorithms. The expression of marker genes used by Lindstrom[19] to annotate cell

198    identities were shown as enriched in the same populations classified using *DevKidCC* (Figure

199    2B), affirming the accuracy and relevance of our classification tool.

200    The arbitrary nature of classifying cells using clustering algorithms is challenged when

201    identifying cells transitioning between populations, often represented as the "borders" of

202    clusters. The cluster-based classification of such cells will change with different approaches

203    to analysis. The application of a cell-centered identification approach circumvents this

204    challenge. To investigate this, *DevKidCC* classification of two published kidney organoid

205    single cell datasets was compared to their original cluster-based annotations. Howden[13]

206    contained samples from two differentiation timepoints; intermediate (18 day) and late (25

207    day) stage organoids (Figure 2C) while Wu[12] contained datasets from two distinct protocols

208    for deriving kidney organoids, labelled as Takasato[30] and Morizane[31] after the original

10

209    authors. *scPred*[22] allows for the setting of a threshold of minimum similarity for a cell to be

210    assigned a given identity. The distribution of the maximum scores for cells in the HFK and

211    organoid datasets showed very similar patterns, however in the HFK there are more distinct

212    peaks at the higher end of similarity (Supplementary Figure 2). In organoids we see a more

213    gradual decrease in scores, meaning there is no set point at which the threshold should

214    obviously be set (Supplementary Figure 2). Organoid datasets from Howden[13] and Wu[12]

215    displayed a similar distribution of similarity scores for Stroma and NPC (Supplementary

216    Figure 2). A sensitivity analysis was performed by comparing threshold points of 0.7 and 0.9

217    with the Howden[13] dataset where we had access to the original annotation for each cell. When

218    mapping the *DevKidCC* classification at both 0.7 and 0.9 thresholds onto the UMAP plot and

219    comparing this to the original classification, *DevKidCC* accurately replicated the original

220    annotation in both settings (Figure 2D). Only a small number of cells did not get classified at

221    the top tier model, defining them as "unassigned" cell types. Such cell types may represent

222    non-renal off target cell types not normally present in HFK or cells in which identity is not

223    sufficiently strong to definitively classify. While all original clusters contained cells that were

224    reclassified as unassigned, the largest contribution was from clusters previously annotated as

225    neuron and muscle, illustrating the specificity with which the model classifies renal cell types

226    (Figure 2E).

227    Both stroma and NPC are mesenchymal cell types. The mesenchymal cells present within

228    kidney organoids have been difficult to accurately classify due to their gene expression

229    profiles being different to those of characterized developing kidney stroma[15]. The previous

230    analysis of the Howden[13] dataset identified seven clusters as stromal (Figure 2C), of which

231    almost all of those assigned an identity using *DevKidCC* remained classified as a stromal sub-

232    type (Figure 2E). However, of the two largest stromal clusters initially identified (see Figure

233    2A[13]), the *TCF21*+ cluster showed a higher number of classified stromal cells while the

234 second was more 'unassigned'. When looking at the similarity scores for stroma and NPC

235 identity in tier 1, the *TCF21* expressing population showed a stromal similarity > 0.9 but very

236 low NPC scores, while the other population had lower stromal scores, albeit still

237 predominantly >0.7, and higher NPC scores, indicating these represent a less defined

238 mesenchymal population (Figure 2D, Supplementary Figure 2). Within the nephron, cells

239 previously identified as "Committed and Early Nephron" were reclassified by *DevKidCC* to

240 comprise a smaller population of NPC together with a larger population of cells identified as

241 Early Nephron (Figure 2E). To examine this further, we analysed organoids generated from

242 either embryonic (ES) or induced pluripotent (iPS) stem cells using two different protocols[12].

243 Using *DevKidCC* we were able to rapidly reproduce the initial classification of these

244 organoids, accounting for the differences in the nomenclature (Figure 2F). Using *DevKidCC*

245 classification we identified cells which do not match the reference (termed "unassigned")

246 enabling further investigation. Here, *DevKidCC* could again distinguish kidney stroma from

247 likely off target cell types like muscle and neural that may represent artefacts of *in vitro*

248 culture[12,13]. Together this reanalysis demonstrates the accuracy with which *DevKidCC* can

249 classify renal cell types within organoid datasets.

250

251 ***DevKidCC* provides a method for direct comparison between protocols**

252 A major challenge for the field has been to compare between datasets generated from

253 different labs, lines, batches or from different protocols due to differences in the analyses that

254 were used. This is particularly pertinent given the use of several distinct protocols for

255 generating kidney tissue from hPSCs (Takasato[30], Morizane[31] and Freedman[39], see Table 1).

256 Direct comparisons between studies and protocols requires an integration of all existing

257 samples to allow re-clustering and differential gene expression analysis on the combined

12

258    dataset. This is challenging due to the noise between samples, the majority of which relates to

259    technical or batch effects[33] which can confound biological variations of interest during data

260    integration[40]. To avoid these challenges, *DevKidCC* was used to directly identify all cell

261    types present within multiple datasets enabling direct comparisons without the need for

262    integration. As *DevKidCC* will compare all cells to the same comprehensive reference, the

263    biological information for each sample can be directly compared without prior dimensional

264    reduction and clustering. To demonstrate this, we applied *DevKidCC* to all available single

265    cell kidney organoid datasets (summarised in Table 1) irrespective of the cell line, organoid

266    age, differentiation protocol or laboratory. The resulting comprehensive analysis (Figure 3)

267    allows a direct comparison of cell proportions across all samples at each tier of classification,

268    grouped into the three main differentiation protocols represented in the literature (Figure 3).

269    What is immediately evident is both the variation in the proportions of "unassigned" cells

270    across all datasets and the lack of nephron maturation even in the oldest organoids regardless

271    of protocol. The maturation of nephron cell types was limited in all protocols and samples,

272    although the Morizane[31] protocol produced organoids with the highest number of cells

273    reflective of a more mature podocyte stage. While there are a small number of mature

274    podocytes, there are almost no mature proximal tubule cells generated with any organoid

275    protocol, but rather being classified as less mature EPT. These have expression of proximal

276    markers such as *CUBN*, *LRP2* and *HNF4A* but lack the specific solute channels such as

277    *SLC47A1*, *SLC22A2* and *SLC22A8* (Supplementary Figure 3). In clustering-based analyses,

278    these cell populations are often split into two or more groups which are interpreted to have

279    varying degrees of maturation, whereas the *DevKidCC* classification indicates that these are

280    mostly immature.

281    *DevKidCC* analysis revealed differences in cell proportion and nephron patterning between

282    organoids generated with different protocols. Organoids generated using the Freedman[39]

283    protocol show a small stromal population in comparison to other protocols. The Morizane[31]

284    organoids show little early nephron cell identity while the Freedman[39] organoids tend to have

285    more early-stage nephron cells. In the Morizane[31] organoids we identify limited distal tubule

286    regions, having less than 25% of the nephrons classified as distal whereas in the Takasato[30]

287    and Freedman[39] protocols is more evenly segmented across nephron components. The

288    Takasato[30] protocol generates the most distal tubule, including some cells classed as a more

289    mature DT segment as well as a Loop of Henle population (Figure 3). The DT expressed

290    *GATA3* and *TMEM52B* but lacked the distal convoluted tubule (DCT)-specific marker

291    *SLC12A3*. However, in some cases the connecting segment (CS)-specific marker *CALB1* is

292    expressed. This would indicate that the connecting segment, which represents the most distal

293    region of the nephron and which invades and fuses into the ureteric tip to form a contiguous

294    tube, is being generated in some organoids. This is promising as it would indicate that there is

295    the potential to promote fusion of these nephrons to any separately induced collecting duct

296    structure, potentially through engineering methods.  In summary, while nephrons are forming

297    and showing evidence of patterning and identifiable segmentation in all protocols, one should

298    keep in mind their relative proximo-distal patterning and evident immaturity prior to their

299    application in disease modelling and drug screen studies.

300

301    **Identifying nephron progenitor cell variation between protocols using *ComponentPlot***

302    **and *DotPlotCompare***

303    To further investigate relative gene expression between datasets, we extracted gene

304    expression profiles and proportions of cells in each classified population, in all available

305    organoid datasets (see Table 1) and the comprehensive reference. A modified version of the

306    *DotPlot* function from the *Seurat*[7,8] package was included to compare gene expression

14

307  between datasets. The direct comparison between kidney organoids (Figure 3) revealed

308  substantial variation in the proportion of NPC, which we further investigated applying the

309  modified function named *DotPlotCompare* to visualization relative gene expression in NPCs

310  across all protocols.

311  The nephron develops from NPCs which are a heterogeneous population of mesenchyme that

312  undergo a mesenchyme to epithelial transition (MET) in response to signals from the ureteric

313  epithelium, giving rise to the entire nephron epithelium[41,42]. *In vivo* analysis has shown

314  markers like *SIX1*, *SIX2*, *CITED1*, *DAPL1* and *LYPD1* are expressed in this population and

315  can be used to reliably identify these cells from the surrounding stromal mesenchyme *in*

316  *situ*[17,19]. These markers have also been used to identify the NPC populations of cells in both

317  HFK and organoids in single cell datasets. NPCs express a posterior HOX code, particularly

318  the HOX10 and HOX11 paralogues[43,44]. Visualising the NPC populations from within the

319  reference HFK dataset using *DevKidCC*, we can see that 44.9% of cells express *SIX2*, 56.3%

320  express *SIX1*, 53.3% *CITED1* while over 70% express *DAPL1* and *LYPD1* (Figure 4A). The

321  posterior HOX genes are also expressed, with *HOXA10* most abundant and *HOXC10*,

322  *HOXD10*, *HOXA11* and *HOXD11* at lower levels and in less cells (Figure 4A). The surprising

323  heterogeneity of gene expression within this population could be explained by technical

324  challenges, including data sparseness, dropout levels and capture bias. It may also be

325  explained by transcriptional bursting[45], where genes are not constantly being transcribed and

326  so the sample harvesting may occur during a transcriptional lull. However, this does provide

327  a true reference for comparison to the expression profiles expected within these cell

328  populations in organoids.

329  When we compare organoid NPCs to the HFK reference, we again note variance between

330  publications and protocols. While NPCs constitute 5-10% of the total cells for the Freedman

331  protocol (Figure 3, 4B), these populations have almost no expression of *SIX2*, *CITED1* and

15

332    *DAPL1*. Similarly, they do not express the posterior HOX code, and only express low levels

333    of *LYPD1* and *SIX2* (Figure 4A). Organoids generated from the Morizane protocol have a

334    more similar profile to the reference NPCs, including some posterior HOX gene expression

335    but little *SIX1* and *LYPD1* expression and almost no *SIX2*, *CITED1* and *DAPL1* present.

336    Takasato-derived organoid NPCs have the most similar profiles to the reference, with NPCs

337    in some samples coexpressing *SIX1*, *SIX2*, *CITED1*, *DAPL1* and *LYPD1*. There is some

338    variance between publications generating organoids from the same protocol (Figure 4B),

339    concurring with earlier studies showing that batch differences are a notable source of

340    variation[12,33]. In the "unassigned" populations generated in organoids, expression of the

341    muscle markers, including *MYOG* and *MYOD1* was sometimes evident. A subset of

342    individual cells within such a published 'muscle' cluster[13] were re-classified by *DevKidCC* as

343    NPC but do show expression of these muscle genes (Figure 4A). Indeed, muscle gene

344    expression is detectable in kidney organoid clusters previously labelled as NPC from multiple

345    protocols and publications[12–15,32]. However, there is no evidence the expression of these genes

346    in the HFK reference, suggesting that their consistent expression in organoid populations is

347    an artifact of the *in vitro* culture conditions. This demonstrates how using *DevKidCC* to

348    classify and directly compare all published organoids datasets can improve our understanding

349    of NPC population generated across multiple kidney organoid protocols. We have identified

350    an *in vitro* culture artefact muscle gene signature within the NPC population present across

351    multiple protocols, giving a target to modulate for improving NPC identity within organoids.

352

353    ***DevKidCC* classification highlights distinct expression profiles of organoid podocytes**

354    **compared to human fetal kidney**

355    Estimating the maturity of cells within a single cell dataset is commonly performed by

356    combining an analysis of cell specific maturation markers within clusters and placing cells or

357    clusters along pseudotime trajectories. However without a time-stamped reference to align

358    the transcriptional profiles these results can be open to interpretation. *DevKidCC* classifies

359    cells based on a reference dataset with a range of maturation states, enabling us to directly

360    compare maturation levels across samples.

361    The glomerular epithelial cells or podocytes are a non-dividing architecturally constrained

362    cell type surrounding the fenestrated capillaries within the glomerulus of each nephron.

363    Podocytes arise from the proximal nephron with trajectory analysis suggesting a distinct

364    transition from the NPCs to that of the remaining nephron epithelium[12,25,28]. Forming a

365    component cell type of the renal corpuscle / glomerulus, the podocytes are anatomically

366    surrounded by a Bowman's capsule comprised of parietal epithelial cells (PECs) which show

367    transcriptional overlap with both podocytes and proximal tubule[46]. Hochane[25] defined a

368    pattern of differential expression across podocytes during maturation. Here, *OLFM3* was

369    expressed in the proximal end of the early nephron (S-Shaped body stage) preceding

370    podocyte patterning with expression of this gene decreasing during podocyte maturation and

371    upregulation of *NPHS1* and *NPHS2*. While expressing markers of podocyte at a lower level,

372    including *MAFB*, *TCF21*, *NPHS1* and *NPHS2*, PECs showed specific expression of *CLDN1*

373    and enriched expression of *PAX8*[24,25,47]. *DevKidCC* analysis of organoid protocols classified

374    most renal corpuscle components as immature podocytes (EPod), with most protocols

375    containing cells classified as PEC and EPod. Organoid EPod and Pod populations had

376    varying levels of *CLDN1*, while *OLFM3* and *PAX8* were co-expressed with more mature

377    podocyte markers like *NPHS1* and *NPHS2* in the PEC and Pod populations. (Figure 5). This

378    may indicate that *in vitro* podocyte differentiation does not progress in the same manner as *in*

379    *vivo* or that these cells are undergoing maturation. The key collagen genes expressed by the

380     podocytes to form a mature glomerular basement membrane are *COL4A3*, *COL4A4* and

381     *COL4A5*[48]. Organoid podocytes again show low expression of these genes compared to

382     podocytes in the HFK reference data. The exception to this observation was seen in organoids

383     subjected to a longer period of time in culture[14,33] suggesting a capacity to mature with time.

384     A critical switch in podocyte maturation is suppression of proliferation, with this post-mitotic

385     state maintained via the expression of key cell cycle regulators including *CDKN1A* (p21) and

386     *CDKN1C* (p57). This is seen in the reference with an increase in expression in the EPod and

387     Pod populations, paired with a decrease in mitotic markers such as *TOP2A,* however in the

388     organoid podocytes there is little decrease in mitosis markers, but expression of *CDKN1A* and

389     *CDKN1C* do increase (Supplementary Figure 4). As such, *DevKidCC* can also be employed

390     as a tool to gain biologically relevant insights into kidney organoids generated from different

391     protocols and users. This is promising for the application of such a tool to compare between

392     wildtype and mutant organoid datasets.

393

394     **Application of *DevKidCC* to investigate the impact of retinoic acid on kidney organoid**

395     **maturation**

396     Accurately identifying the cell types present within an organoid is crucial for the analysis of

397     disease states or the optimization of the differentiation protocols. To evaluate the application

398     of *DevKidCC* in analyzing functional differences between methods, we analysed unpublished

399     data in which kidney organoids from the same starting cell line generated in the same batch

400     were treated with 5µM retinoic acid after removal of all other growth factors at day 12 (7+5)

401     of the protocol to promote maturation. Mammalian nephrogenesis *in vivo* occurs in waves

402     with new nephrons constantly forming up to 36 weeks gestation[49,50] in humans and into the

403     first week of life in mice[51]. This is facilitated by the presence of a peripheral nephrogenic

18

404    niche within which the NPC balance self-renewal versus nephron commitment. Once

405    differentiated, NPCs exist throughout the duration of organoid culture and deplete with time,

406    although a population does remain in mature organoids able to undergo nephrogenesis when

407    induced with a canonical Wnt agonist[13] (Figure 4A). Retinoic acid signaling plays many roles

408    in kidney development depending on spatiotemporal expression[52–54], and is also known to

409    promote the differentiation of progenitor cell populations[55]. We investigated adding all-trans

410    retinoic acid (RA) to organoids at multiple time points to see what effect this would have on

411    organoids. The addition of 1 -5 µM RA before day five of 3D organoid culture, substantially

412    impaired nephron formation, whereas addition at day five onwards led to organoids with fully

413    segmented nephrons similar to organoids without RA (data not shown). The *DevKidCC*

414    classification identified an increase in the percentage of classified stromal cells, seemingly at

415    the expense of the 'unassigned' population. In contrast to control organoids at day 25, the

416    addition of RA resulted in a complete depletion of NPC cells (Figure 5A). While the

417    percentage of nephron cells did not change, there was a shift towards more proximal tubule

418    (EPT) than early distal tubule (EDT_EMT) (Figure 5B). These comparisons show evidence

419    that RA caused the depletion of NPCs and proximalisation of nephrons within forming

420    organoids. NPC depletion can be seen 6 days after addition of RA, when organoids generated

421    using a SIX2[EGFP] reporter line[13,56] were analysed by flow cytometry. The control organoids

422    had 31.44% EGFP+ cells while the organoids with RA had less than 0.5% (Figure 5C).  This

423    confirms that RA acts directly or indirectly on the NPC population, forcing them to either

424    undergo commitment to form nephrons or differentiate away from NPC identity down a

425    stromal pathway.

426    To investigate the maturation of the nephrons we visualized maturation markers for all

427    segments using the *DotPlotCompare* function within the package. Only the podocytes

428    showed evidence of maturation, with an increase in the expression of genes such as *WT1,*

19

429     *MAFB, TCF21* and *NPHS2* with RA addition while also showing a decrease in *OLFM3*

430     expression, a marker of the immature podocytes[25]. Interestingly, the PEC marker *CLDN1*

431     remained expressed in the podocytes, although immunofluorescence showed more specific

432     localization to the epithelial cells surrounding the podocytes, which is the normal location of

433     PECs (Figure 5D, 5E). These results may indicate that the podocytes and potential PEC cells

434     increase in maturity when RA is added. The expression of both PEC and podocyte markers in

435     cells assigned to all three renal corpuscle identities is consistent with the previous analysis of

436     these populations and may indicate that while maturation is occurring, the delineation of

437     specific gene signatures within these cells is not.

438     **Analysis of existing protocols for the development of ureteric epithelium**

439     *DevKidCC* is able to predict stromal, endothelial, ureteric and nephron cell identity based upon the

440     reference data from HFK. Our analysis of existing standard organoid protocols confirms the absence

441     of populations classified as ureteric epithelium. The ureteric epithelium in the mammalian kidney

442     arises as a side branch of the mesonephric duct that grows into the presumptive kidney mesenchyme.

443     Hence it has been suggested that it is not possible to generate ureteric epithelium using the same

444     differentiation protocol able to generate the nephron lineages[57]. To date, a number of protocols have

445     been published that report the generation of ureteric epithelium[24,38,57,58] with all of all these methods

446     involving the isolation of cellular fractions that are then cultured separately to form ureteric

447     epithelium. Single cell analyses have recently revealed the significant transcriptional congruence

448     between the distal nephron and the ureteric epithelium in both human and mouse[12,16]. It has also been

449     established that distal nephron from standard organoids remains plastic and can be induced to adopt a

450     ureteric epithelial fate[18]. To investigate how accurately *DevKidCC* can identify this kidney cell

451     identity, we applied this analysis to a single cell dataset available from a specific UE protocol[38] and

452     the single cell dataset recently generated from UE that had been derived from DN[24]. *DevKidCC*

453     classified 20% and 28% of cells as UE respectively (Figure 6A), with most of these classified as Tip

454     (Figure 6B). The DN-derived sample[24] contains a population of cells classed as nephron while the UE

455     sample directly differentiated from hPSC[38] contains cells classed as NPC (Figure 6A, 6B). The

456     different cell types present between these two samples may be explained by the different protocols

457     used to generate UE and kidney developmental biology. Cultures differentiated towards an

458     anteriorised intermediate mesoderm population directly from hPSCs are likely to generate a

459     proportion of NPC-like cells as a *bona fide* posterior intermediate mesoderm of a more anterior

460     nephrogenic cord. In contrast, the DN-derived cultures contain nephron-like cells that have not

461     become UE. In both samples the overwhelming majority of cells were "unassigned" (Figure 6A).

462     However, when visualizing the distribution of scores there is an even distribution of UrEp scores in

463     both samples between 0.1 and 1 (Figure 6C) with the majority of cells being most similar to the

464     ureteric population over any other lineage. This indicates a spectrum of similarity to the true ureteric

465     epithelium. The implications of this when attempting to classify these cells are that a clustering

466     analysis will break them up without appreciating the overall transcriptional similarity while

467     *DevKidCC* will classify each cell based on its own merit, giving a more accurate overall picture of

468     cell identity compared to the true HFK profile. **Discussion**

469     The question of cell identity is one that is difficult to answer. Histologically we can try to

470     define a cell type based on its morphology, gene expression or protein expression, the latter

471     typically being read by immunohistochemistry and immunofluorescence assays. In many

472     cellular states, particularly those present during organogenesis, evaluation of cellular identity

473     by functional assays is challenging and marker expression is rarely unique. This challenge is

474     significant when evaluating cell identity using single cell RNA sequencing data. Such data is

475     sparse, providing an incomplete snapshot rather than a comprehensive picture. As capture

476     technology and bioinformatics tools have improved, increased levels of information can be

477     extracted from this data, providing an overall synergy of expression profile for groups of cells

478     within a sample. This can be combined with the pseudotime trajectory or even molecular

479     lineage tagging to relate cells within a sample by history, assisting in likely classification of

480     cell type.  Such inferences are much more difficult in a synthetic *in vitro* system such as

481    hPSC-derived organoids. Such protocols direct cells to undergo a series of changes that

482    attempt to replicate the *in vivo* process. However, in reality, hPSC-derived lineages often do

483    not completely recapitulate their *in vivo* counterparts, at least at the level of the transcriptome.

484    We can often identify a gene, or a number of genes, that are expressed in a cell that give us

485    some information of what it can be classified as, but in many cases there is ambiguity. This is

486    compounded by our knowledge that hPSC-derived organoid models replicate early

487    developmental cell states that are frequently in flux, not present in adult tissue and are less

488    well defined.

489    The classification of cells within all single cell data has been inconsistent as clustering and

490    classification decisions vary between individual researchers and the limitations within each

491    dataset. *DevKidCC* represents a method of specifically classifying individual cellular identity

492    within hPSC-derived kidney organoids based upon models trained on a comprehensive

493    reference dataset. Our tool facilitates direct comparisons between kidney organoid datasets by

494    classifying cells based on the reference data. While the base package, *scPred*[22], includes a

495    way to integrate the data within the models using *Harmony*[27], this can introduce false

496    correlations between similar cell populations such as the mesenchymal cells that have

497    intermediate to high scores for both stroma and NPC. Hence, *DevKidCC* provides an option

498    to run the harmonization step, but this is not required or recommended for kidney organoid

499    datasets. The classification for all datasets has been integrated into functions allowing for

500    plotting any novel dataset in direct comparison using the classification from *DevKidCC*.

501    Gene expression can be visualised using the *DotPlotCompare* function, while sample

502    annotation can be visualised using *ComponentPlot* or *SankeyPlot.* These tools included in

503    *DevKidCC* provide a classification and visualization toolset to investigate cell identity and

504    gene expression within novel and existing kidney organoids.

505    *DevKidCC* was developed so that it could be applied to novel datasets facilitating direct

506    comparisons to those previously generated. This will make comparative studies much easier,

507    facilitating the analysis of genetic variants, disease states or methodological variation in new

508    protocols. While this system has developed a model with three tiers of subclassification, the

509    complexity of the human nephron, even in the fetal kidney, is such that there is scope to

510    interrogate individual cellular identity even further within this and other subcomponents. As

511    these models were trained using developing HFK, the ability of the tool to accurately classify

512    cell identity during earlier stages of mesoderm patterning or mature kidney is limited. The

513    adult kidney shows significant specification of functional cell types within all segments of the

514    final nephron, many of which have distinct functional roles in renal filtration and fluid

515    homeostasis but are not present in the fetal organ. Indeed, the ratio of epithelium to stroma is

516    dramatically shifted in the adult. While the fetal kidney begins to form some more mature

517    cellular states, such as the intercalated and principal cells of the distal nephron / collecting

518    duct, it is likely that a distinct cellular identity tool will be required for the accurate

519    identification of cellular identity in postnatal kidney tissue. Conversely, the use of HFK from

520    Trimester 1 and 2 as the reference dataset limits the ability to identify earlier stages of

521    morphogenesis. This may explain the large percentage of unassigned cell calls in datasets in

522    early stages of kidney organoid differentiation protocols (Figure 3, Figure 4A). However,

523    *DevKidCC* applied to early-stage differentiations (day 7, intermediate mesoderm) split cell

524    identity between NPC and unassigned, suggesting that the tool is able to identify those cells

525    beginning to commit to the mesenchymal precursors of the kidney. Indeed, in a dataset that

526    includes day 7, 15 and 29 organoids between two cell lines[14], there is a direct relationship

527    between the proportion of cells classified as NPC at day 7 to the proportion of nephron cells

528    at day 15 and 29 (Figure 4A). We conclude that at this early stage the cells identified as NPC

529    at this early stage could be the percentage of the differentiation correctly patterned to

530    intermediate mesoderm and are still the cells that will go on to form the nephron population.

531    **Conclusions**

532    DevKidCC provides a robust, reproducible and computationally efficient tool for the

533    classification of kidney single cell data, in both human and organoid-derived tissue. Using

534    DevKidCC we can now directly compare between kidney samples regardless of batch and

535    have done so for all available published datasets. This important advance has provided

536    insights into differences in organoids derived using different protocols and allows for any

537    novel dataset to be directly compared to all previous datasets. The included custom functions

538    simplify visualisation of cell identity proportion and gene expression within samples and

539    between multiple samples. Any novel dataset can be classified using the framework provided

540    in this package, allowing for direct comparison to all previous datasets, all of which are

541    included within the package. For visualisation of gene expression profiles and organoid cell

542    identities, the gene expression profiles of all datasets have been built into an *R* Shiny app

543    available at https://sbwilson91.shinyapps.io/devkidcc_interactive/[59] that does not require the

544    use of *R* directly, allowing for easy access to this information. Finally, while this package has

545    been built using HFK data to classify kidney cells, the framework can be transferred to any

546    tissue type where adequate single cell data is available.

547    **Methods**

548    *DevKidCC* **algorithm**

549    *DevKidCC* (Developing Kidney Cell Classifier) is a function written in *R* designed to provide

550    an accurate, robust and reproducible method to classify single cell RNA-sequencing datasets

551    containing human developing kidney-like cells. The algorithm has two steps: data pre-

24

552    processing and cell classification. Below we describe the development and utilisation of these

553    steps.

**Data pre-processing**

555    The required input is a scRNA-seq dataset as a *Seurat*[7,8] object. This object is first normalised

556    by dividing the total expression of each gene by the total gene expression per cell then

557    multiplied by a scale factor of 10,000 and natural log-transformed with pseudocount of 1.

**Cell classification**

559    We generated a comprehensive developing kidney reference single cell dataset by

560    harmonising the raw data from multiple high quality human fetal kidney datasets. The

561    annotation of the reference included three tiers with increasing specificity, with a clear

562    hierarchical structure between the tiers. This dataset was then used to train machine learning

563    models using the *R* package *scPred*[22]. One model was trained for each node of identities

564    within the classification hierarchy.

565    Utilising *scPred*[22] the models were trained using the same parameters, with the relevant cells

566    inputted for each. The feature space used was the top 100 principal components. The models

567    were trained using a support vector machine with a radial kernel. The models are stored as a

568    *scPred*[22]object and can be used to classify cells within a *Seurat*[7,8] object using the *scPred*[22]

569    package. For classification, these models will calculate the similarity of a cell to each of the

570    trained identities within that model, giving a probability score between 0 and 1 for each

571    identity. It will then assign an identity of the highest similarity score above the set threshold,

572    or call the cell unassigned if no identity scores above the threshold.

573    Cells are classified using these models, organised in a biologically relevant hierarchy so as to

574    optimally and accurately identify the cellular identity of all analysed cells. All cells are first

575    classified using the first-tier model, contains generalised lineage identities of stroma, nephron

25

576   progenitors, nephron, ureteric epithelium and endothelium. After similarity calculation using

577   the first-tier model, cells that do not pass the threshold are classified as unassigned. The

578   threshold is set to 0.7 by default but can be adjusted by the user, which can be useful if the

579   user wants to classify cells with at decreasing levels of similarity. Cells assigned to stroma,

580   nephron and ureteric epithelium are passed into a second tier of classification specific to these

581   identities. It is important to note that at the second and third classification tiers, there is no

582   thresholding, i.e., all cells are assigned an identity with no cells classed as unassigned. The

583   second-tier ureteric epithelium model is trained on the tip, cortical, outer and inner medullary

584   cell identities. The second-tier stroma model is trained on the stromal progenitors, cortex,

585   medullary and mesangial cell identities. The second-tier nephron model is trained on the early

586   nephron, distal nephron, proximal nephron, renal corpuscle and nephron cell cycle

587   population. The distal nephron, proximal nephron and renal corpuscle are then further

588   classified into more specific identities in a third tier of models. The third-tier distal nephron

589   model is trained on early distal/medial cells, distal tubule and loop of Henle cells. The third-

590   tier proximal nephron model is trained on early proximal tubule and proximal tubule cells.

591   The third-tier renal corpuscle model is trained on parietal epithelial cells, early podocytes and

592   podocytes. Each stage of the classification step is recorded as a metadata column, as is the

593   final classification for each cell. All the similarity scores and tier classifications are readily

594   accessible within the *Seurat*[7,8] object for further analysis.

595   **Comprehensive reference generation**

596   Raw data was downloaded from GEO database from repositories GSE114530[60] and

597   GSE124472[61], or provided to us directly by the authors, since made available at EMBL-EBI

598   ArrayExpress under accession number E-MTAB-9083[62]. The data as *CellRanger* output was

599   read into *R* and processed using *Seurat*[7,8] (v3.2.2), using *SCTransform*[63] for pre-processing.

600   Clustering and manual annotation was performed on each dataset individually, referring back

26

601    to the original papers and using established markers enriched in clusters to classify each

602    cluster. Once annotated, datasets were integrated using *Harmony*[27] with 100 PCAs and 10000

603    variable features.

**Organoid gene expression database**

605    A reference database of all available kidney organoid datasets (Table 1) was generated by

606    running *DevKidCC*, extracting summaries of the gene expression information at each

607    classification tier, and combining these into a database. This database can be used to directly

608    compare gene expression between existing datasets, also novel datasets classified using

609    *DevKidCC*. The link to download this database is available at the package Github

610    repository[64].

**Downstream visualisation functions**

612    To facilitate data visualisation and analysis of *DevKidCC* classified datasets, three

613    customised functions were included in the package. *DotPlotCompare* is a modified version of

614    the *DotPlot* function from the Seurat package. A gene expression profile of the reference is

615    present within the function and can be used for direct comparisons to an existing or novel

616    dataset. There is an option to visualise the organoid database within this function as well, the

617    downloading instructions for this are available at the package Github repository[64]. The

618    proportions of cells classified using *DevKidCC* can be visualised as a bar chart using the

619    *ComparePlot* function. This can also take as input a gene and show the expression of that

620    gene in each segment. The *SankeyPlot* function utilises the *networkD3* package to generate an

621    interactive Sankey chart showing the flow of cell classification.

**DevKidCC Kidney Organoid Gene Explorer shiny app**

623    To make visualisation of the organoid database possible outside of using *R,* a shiny app was

624    developed[59]. This allows for an interactive way to visualise and analyse gene expression

625    within published organoid datasets.

626    **iPSC-derived organoid differentiation**

627    The day prior to differentiation, cells were dissociated with TrypLE (Thermo Fisher

628    Scientific), counted using a hemocytometer, and seeded onto Laminin 521-coated 6-well

629    plates at a density of 50 x $10^3$ cells per well in Essential 8 () medium. Intermediate mesoderm

630    induction was performed by culturing iPSCs in TeSR-E6 medium (Stem Cell Technologies)

631    containing 4-8 µM CHIR99021 (R&D Systems) for 4 days. On day 4, cells were switched to

632    TeSR-E6 medium supplemented with 200ng/ml FGF9 (R&D Systems) and 1 µg/ml Heparin

633    (Sigma-Aldrich). On day 7, cells were dissociated with TrypLE, diluted fivefold with TeSR-

634    E6 medium, transferred to a 15-ml conical tube, and centrifuged for 5 min at 300 x g to pellet

635    cells. The supernatant was discarded, and cells were resuspended in residual medium and

636    transferred directly into a syringe for bioprinting. Syringes containing the cell paste were

637    loaded onto a NovoGen MMX Bioprinter, primed to ensure cell material was flowing, with

638    100,000 cells deposited per organoid onto a 0.4-µm Transwell polyester membranes in 6-well

639    plates (Corning). Following bioprinting, organoids were cultured for 1h in presence of 6µM

640    CHIR99021 in TeSR-E6 medium in the basolateral compartment and subsequently cultured

641    until day 12 in TeSR-E6 medium supplemented with 200 ng/ml FGF9 and 1 µg/ml Heparin.

642    From day 12 to day 25, organoids were grown in TeSR-E6 medium either without additional

643    supplement, or with additional 5uM all-trans retinoic acid (). Unless otherwise stated, kidney

644    organoids were cultured until harvest at day 25.

645    **Flow cytometry**

646    Prior to analysis, single kidney organoids were dissociated with 0.2 ml of a 1:1

647    TrypLE/Accutase solution in 1.5-ml tubes at 37°C for 15–25 min, with occasional mixing

648    (flicking) until large clumps were no longer clearly visible. 1 ml of HBBS supplemented with

649    2% FBS was added to the cells before passing through a 40-lM FACS tube cell strainer

650    (Falcon). Flow cytometry was performed using a LSRFortessa Cell Analyzer (BD

651    Biosciences). Data acquisition and analysis were performed using FACSDiva (BD) and

652    FlowLogic software (Inivai). Gating was performed on live cells based on forward and

653    side-scatter analysis.

654    **Whole mount immunostaining**

655    Fixed kidney organoids were incubated in blocking buffer (PBS 1X donkey serum 10% triton

656    X100 0.3%) at 4°C for 3h before adding primary antibodies against HNF4α (Life

657    Technologies 1:300, cat# MA1-199), Nephrin (NPHS1 1:300, Bioscientific, cat# AF4269)

658    and Claudin-1 (CLDN1 1:100, Thermo Fisher Scientific, cat# 71-7800) at 4°C for 2 days.

659    After washing in PBS 1X triton X-100 0.1%, organoids were incubated in secondary

660    antibodies 1:400 at 4°C for 2 days: Alexa fluor 405 donkey anti-mouse (Abcam, cat#

661    ab175659), Alexa fluor 488 donkey anti-goat (Molecular Probes, cat# A11055), Alexa fluor

662    568 donkey anti-rabbit (Life Technologies, cat# A10042). Samples were then washed before

663    blocking at 4°C for 3h with PBS 1X mouse serum 10µg/ml triton X-100 0.3%, and adding an

664    APC-conjugated CD31 antibody (1:50, Biolegend, cat# 303115) at 4°C for 2 days. Finally,

665    samples were washed and imaged in 50:50 glycerol:PBS 1X using a Dragonfly Spinning Disc

666    Confocal Microscope (Andor Technology).

667    **Single-cell transcriptional profiling and data analysis**

668    Organoids were dissociated as described above (for flow cytometry) and passed through a 40-

669    µM FACS tube cell strainer. Following centrifugation at 300 g for 3 min, the supernatant was

670  discarded and cells resuspended in 50 µl TeSR-E6 medium. Viability and cell number were

671  assessed, and samples were run across separate runs on a Chromium Chip Kit (10×

672  Genomics). Libraries were prepared using Chromium Single Cell Li sequenced on an

673  Illumina HiSeq with 100-bp paired-end reads. Cell Ranger (v1.3.1) was used to process and

674  aggregate raw data from each of the samples returning a count matrix. Quality control and

675  analysis was performed in *R* using the *Seurat* package (v3.2.2). Classification was performed

676  using *DevKidCC* (v0.1.6) as described in this manuscript.

677

678  **Declarations**

679  **Availability of data and materials**

680  *DevKidCC* is available from Github at https://github.com/KidneyRegeneration/DevKidCC[64]

681  under the MIT licence. *DevKidCC Kidney Organoid Gene Expression* interactive shiny

682  dashboard is available at https://sbwilson91.shinyapps.io/devkidcc_interactive/[59] and from

683  Github at https://github.com/KidneyRegeneration/DevKidCC_Interactive[59]

684  Single cell RNA-sequencing human fetal kidney datasets can be found in GEO (GSE102596,

685  GSE114530) and EMBL-EBI ArrayExpress (E-MTAB-9083) [60,62,65]. Single cell RNA-

686  sequencing organoid datasets can be found in GEO (GSE118184, GSE109718, GSE119561,

687  GSE114802, GSE115986, GSE132026, GSE124472, GSE152014, GSE161255,

688  GSE152685)[61,66–73]. The single cell RNA-sequencing organoid dataset generated in this study

689  will be available from GEO upon manuscript publication.

690  **Competing interests**

691  The authors declare that they have no competing interests.

692  **Funding**

698     **Authors' contributions**

699     SBW, MHL and JEP conceived the study. SBW, JAH and JEP contributed to method

700     development. SBW performed bioinformatics analysis. SBW, SEH, JMV and AD performed

701     kidney differentiation experiments, immunofluorescence and FLOW analysis. SBW and

702     MHL wrote the manuscript while all authors assisted in manuscript preparation.

707     **References**

708     1.      Wagner DE, Klein AM. Lineage tracing meets single-cell omics: opportunities and

709             challenges. *Nat Rev Genet*. 2020;21(7):410-427. doi:10.1038/s41576-020-0223-2

710     2.      Gitter A. Single-cell RNA-seq pseudotime estimation algorithms. Zenodo.

711             doi:10.5281/zenodo.1297422

712     3.      Saelens W, Cannoodt R, Todorov H, Saeys Y. A comparison of single-cell trajectory

713             inference methods. *Nat Biotechnol*. 2019;37(5):547-554. doi:10.1038/s41587-019-

714             0071-9

715    4.    Zappia L, Phipson B, Oshlack A. Exploring the single-cell RNA-seq analysis

716          landscape with the scRNA-tools database. *PLOS Comput Biol*. 2018;14(6):e1006245.

717          https://doi.org/10.1371/journal.pcbi.1006245.

718    5.    La Manno G, Soldatov R, Zeisel A, et al. RNA velocity of single cells. *Nature*.

719          2018;560(7719):494-498. doi:10.1038/s41586-018-0414-6

720    6.    Bergen V, Lange M, Peidli S, Wolf FA, Theis FJ. Generalizing RNA velocity to

721          transient cell states through dynamical modeling. *Nat Biotechnol*. 2020.

722          doi:10.1038/s41587-020-0591-3

723    7.    Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell

724          transcriptomic data across different conditions, technologies, and species. *Nat*

725          *Biotechnol*. 2018;36(5):411-420. doi:10.1038/nbt.4096

726    8.    Stuart T, Butler A, Hoffman P, et al. Comprehensive Integration of Single-Cell Data.

727          *Cell*. 2019;177(7):1888-1902.e21. doi:10.1016/j.cell.2019.05.031

728    9.    Abdelaal T, Michielsen L, Cats D, et al. A comparison of automatic cell identification

729          methods for single-cell RNA sequencing data. *Genome Biol*. 2019;20(1):194.

730          doi:10.1186/s13059-019-1795-z

731    10.   Kiselev VY, Yiu A, Hemberg M. scmap: projection of single-cell RNA-seq data across

732          data sets. *Nat Methods*. 2018;15:359. https://doi.org/10.1038/nmeth.4644.

733    11.   Little MH, Combes AN. Kidney organoids: accurate models or fortunate accidents.

734          *Genes Dev*. 2019;33(19-20):1319-1345. doi:10.1101/gad.329573.119

735    12.   Wu H, Uchimura K, Donnelly EL, Kirita Y, Morris SA, Humphreys BD. Comparative

736          Analysis and Refinement of Human PSC-Derived Kidney Organoid Differentiation

737          with Single-Cell Transcriptomics. *Cell Stem Cell*. 2018;23(6):869-881.e8.

738     doi:10.1016/j.stem.2018.10.010

739     13.     Howden SE, Vanslambrouck JM, Wilson SB, Tan KS, Little MH. Reporter-based fate

740             mapping in human kidney organoids confirms nephron lineage relationships and

741             reveals synchronous nephron formation. *EMBO Rep*. 2019;0(0):e47483.

742             doi:10.15252/embr.201847483

743     14.     Subramanian A, Sidhom E-H, Emani M, et al. Single cell census of human kidney

744             organoids shows reproducibility and diminished off-target cells after transplantation.

745             *Nat Commun*. 2019;10(1). doi:10.1038/s41467-019-13382-0

746     15.     Combes AN, Zappia L, Er PX, Oshlack A, Little MH. Single-cell analysis reveals

747             congruence between kidney organoids and human fetal kidney. *Genome Med*.

748             2019;11(1). doi:10.1186/s13073-019-0615-0

749     16.     Combes AN, Phipson B, Lawlor KT, et al. Single cell analysis of the developing

750             mouse kidney provides deeper insight into  marker gene expression and ligand-

751             receptor crosstalk. *Development*. 2019;146(12). doi:10.1242/dev.178673

752     17.     Lindström NO, Tran T, Guo J, et al. Conserved and Divergent Molecular and

753             Anatomic Features of Human and Mouse Nephron Patterning. *J Am Soc Nephrol*.

754             2018:ASN.2017091036. doi:10.1681/asn.2017091036

755     18.     Ransick A, Lindstrom NO, Liu J, et al. Single-Cell Profiling Reveals Sex, Lineage,

756             and Regional Diversity in the Mouse Kidney. *Dev Cell*. 2019;51(3):399-413.e7.

757             doi:10.1016/j.devcel.2019.10.005

758     19.     Lindström NO, Guo J, Kim AD, et al. Conserved and Divergent Features of

759             Mesenchymal Progenitor Cell Types within the Cortical Nephrogenic Niche of the

760             Human and Mouse Kidney. *J Am Soc Nephrol*. 2018:ASN.2017080890.

761       doi:10.1681/asn.2017080890

762   20.  Choi J-H, In Kim H, Woo HG. scTyper: a comprehensive pipeline for the cell typing

763       analysis of single-cell RNA-seq data. *BMC Bioinformatics*. 2020;21(1):342.

764       doi:10.1186/s12859-020-03700-5

765   21.  Lin Y, Cao Y, Kim HJ, et al. scClassify: hierarchical classification of cells. *bioRxiv*.

766       January 2019:776948. doi:10.1101/776948

767   22.  Alquicira-Hernandez J, Sathe A, Ji HP, Nguyen Q, Powell JE. scPred: accurate

768       supervised method for cell-type classification from single-cell  RNA-seq data. *Genome*

769       *Biol*. 2019;20(1):264. doi:10.1186/s13059-019-1862-5

770   23.  Holloway EM, Wu JH, Czerwinski M, et al. Differentiation of Human Intestinal

771       Organoids with Endogenous Vascular Endothelial Cells. *Dev Cell*. 2020;54(4):516-

772       528.e7. doi:https://doi.org/10.1016/j.devcel.2020.07.023

773   24.  Howden SE, Wilson SB, Groenewegen E, et al. Plasticity of distal nephron epithelia

774       from human kidney organoids enables the induction of ureteric tip and stalk. *Cell Stem*

775       *Cell*. December 2020. doi:10.1016/j.stem.2020.12.001

776   25.  Hochane M, van den Berg PR, Fan X, et al. Single-cell transcriptomics reveals gene

777       expression dynamics of human fetal kidney development. *PLOS Biol*.

778       2019;17(2):e3000152. https://doi.org/10.1371/journal.pbio.3000152.

779   26.  Tran T, Lindstrom NO, Ransick A, et al. In Vivo Developmental Trajectories of

780       Human Podocyte Inform In Vitro Differentiation of Pluripotent Stem Cell-Derived

781       Podocytes. *Dev Cell*. 2019;50(1):102-116.e6. doi:10.1016/j.devcel.2019.06.001

782   27.  Korsunsky I, Millard N, Fan J, et al. Fast, sensitive and accurate integration of single-

783       cell data with Harmony. *Nat Methods*. 2019;16(12):1289-1296. doi:10.1038/s41592-

784        019-0619-0

785  28.   Menon R, Otto EA, Kokoruda A, et al. Single-cell analysis of progenitor cell dynamics

786        and lineage specification in the human fetal kidney. *Development*.

787        2018;145(16):dev164038. doi:10.1242/dev.164038

788  29.   Young MD, Mitchell TJ, Vieira Braga FA, et al. Single-cell transcriptomes from

789        human kidneys reveal the cellular identity of renal tumors. *Science (80- )*.

790        2018;361(6402):594-599. doi:10.1126/science.aat1699

791  30.   Takasato M, Er PX, Chiu HS, et al. Kidney organoids from human iPS cells contain

792        multiple lineages and model human nephrogenesis. *Nature*. 2015;526(7574):564-568.

793        doi:10.1038/nature15695

794  31.   Morizane R, Lam AQ, Freedman BS, Kishi S, Valerius MT, Bonventre J V. Nephron

795        organoids derived from human pluripotent stem cells model kidney development and

796        injury. *Nat Biotechnol*. 2015;33:1193. http://dx.doi.org/10.1038/nbt.3392.

797  32.   Czerniecki SM, Cruz NM, Harder JL, et al. High-Throughput Screening Enhances

798        Kidney Organoid Differentiation from Human Pluripotent Stem Cells and Enables

799        Automated Multidimensional Phenotyping. *Cell Stem Cell*. 2018;22(6):929-940.e4.

800        doi:10.1016/j.stem.2018.04.022

801  33.   Phipson B, Er PX, Combes AN, et al. Evaluation of variability in human kidney

802        organoids. *Nat Methods*. 2019;16(1):79-87. doi:10.1038/s41592-018-0253-2

803  34.   Harder JL, Menon R, Otto EA, et al. Organoid single cell profiling identifies a

804        transcriptional signature of glomerular disease. *JCI insight*. 2019;4(1):e122697.

805        doi:10.1172/jci.insight.122697

806  35.   Kumar S V, Er PX, Lawlor KT, et al. Kidney micro-organoids in suspension culture as

807     a scalable source of human pluripotent stem cell-derived kidney cells. *Development*.

808     2019;146(5):dev172361. doi:10.1242/dev.172361

809  36.  Low JH, Li P, Chew EGY, et al. Generation of Human PSC-Derived Kidney

810     Organoids with Patterned Nephron Segments and a De Novo Vascular Network. *Cell*

811     *Stem Cell*. 2019;25(3):373-387.e9. doi:https://doi.org/10.1016/j.stem.2019.06.009

812  37.  Lawlor KT, Vanslambrouck JM, Higgins JW, et al. Cellular extrusion bioprinting

813     improves kidney organoid reproducibility and conformation. *Nat Mater*. 2020.

814     doi:10.1038/s41563-020-00853-9

815  38.  Mae S-I, Ryosaka M, Sakamoto S, et al. Expansion of Human iPSC-Derived Ureteric

816     Bud Organoids with Repeated Branching Potential. *Cell Rep*. 2020;32(4):107963.

817     doi:https://doi.org/10.1016/j.celrep.2020.107963

818  39.  Freedman BS, Brooks CR, Lam AQ, et al. Modelling kidney disease with CRISPR-

819     mutant kidney organoids derived from human pluripotent epiblast spheroids. *Nat*

820     *Commun*. 2015;6:8715. http://dx.doi.org/10.1038/ncomms9715.

821  40.  Tran HTN, Ang KS, Chevrier M, et al. A benchmark of batch-effect correction

822     methods for single-cell RNA sequencing data. *Genome Biol*. 2020;21(1):12.

823     doi:10.1186/s13059-019-1850-9

824  41.  Self M, Lagutin O V, Bowling B, et al. Six2 is required for suppression of

825     nephrogenesis and progenitor renewal in the developing kidney. *EMBO J*.

826     2006;25(21):5214-5228. doi:10.1038/sj.emboj.7601381

827  42.  Kobayashi A, Valerius MT, Mugford JW, et al. Six2 Defines and Regulates a

828     Multipotent Self-Renewing Nephron Progenitor Population throughout Mammalian

829     Kidney Development. *Cell Stem Cell*. 2008;3(2):169-181.

830            doi:10.1016/J.STEM.2008.05.020

831    43.    Wellik DM, Hawkes PJ, Capecchi MR. Hox11 paralogous genes are essential for

832            metanephric kidney induction. *Genes Dev*. 2002;16(11):1423-1432.

833            doi:10.1101/gad.993302

834    44.    Yallowitz AR, Hrycaj SM, Short KM, Smyth IM, Wellik DM. Hox10 genes function

835            in kidney development in the differentiation and integration of the cortical stroma.

836            *PLoS One*. 2011;6(8):e23410-e23410. doi:10.1371/journal.pone.0023410

837    45.    Chubb JR, Trcek T, Shenoy SM, Singer RH. Transcriptional Pulsing of a

838            Developmental Gene. *Curr Biol*. 2006;16(10):1018-1025.

839            doi:10.1016/j.cub.2006.03.092

840    46.    Shankland SJ, Smeets B, Pippin JW, Moeller MJ. The emergence of the glomerular

841            parietal epithelial cell. *Nat Rev Nephrol*. 2014;10(3):158-173.

842            doi:10.1038/nrneph.2014.1

843    47.    Ohse T, Chang AM, Pippin JW, et al. A new function for parietal epithelial cells: a

844            second glomerular barrier. *Am J Physiol Physiol*. 2009;297(6):F1566-F1574.

845            doi:10.1152/ajprenal.00214.2009

846    48.    Abrahamson DR, Hudson BG, Stroganova L, Borza D-B, St. John PL. Cellular Origins

847            of Type IV Collagen Networks in Developing Glomeruli. *J Am Soc Nephrol*.

848            2009;20(7):1471 LP - 1479. doi:10.1681/ASN.2008101086

849    49.    Hartman HA, Lai HL, Patterson LT. Cessation of renal morphogenesis in mice. *Dev

850            Biol*. 2007;310(2):379-387. doi:10.1016/j.ydbio.2007.08.021

851    50.    Rumballe BA, Georgas KM, Combes AN, Ju AL, Gilbert T, Little MH. Nephron

852            formation adopts a novel spatial topology at cessation of nephrogenesis. *Dev Biol*.

853        2011;360(1):110-122. doi:10.1016/j.ydbio.2011.09.011

854    51.    Short KM, Combes AN, Lefevre J, et al. Global Quantification of Tissue Dynamics in

855        the Developing Mouse Kidney. *Dev Cell*. 2014;29(2):188-202.

856        doi:https://doi.org/10.1016/j.devcel.2014.02.017

857    52.    Burrow CR. Retinoids and Renal Development. *Nephron Exp Nephrol*. 2000;8(4-

858        5):219-225. doi:10.1159/000020672

859    53.    Janesick A, Tang W, Shioda T, Blumberg B. RARγ is required for mesodermal gene

860        expression prior to gastrulation. *Development*. 2018:dev.147769.

861        doi:10.1242/dev.147769

862    54.    Janesick A, Nguyen TTL, Aisaki K-I, et al. Active repression by RAR  signaling is

863        required for vertebrate axial elongation. 2014;141(11):2260-2270.

864        doi:10.1242/dev.103705

865    55.    Gudas LJ, Wagner JA. Retinoids regulate stem cell differentiation. *J Cell Physiol*.

866        2011;226(2):322-330. doi:10.1002/jcp.22417

867    56.    Vanslambrouck JM, Wilson SB, Tan KS, et al. A Toolbox to Characterize Human

868        Induced Pluripotent Stem Cell-Derived Kidney Cell Types and Organoids. *J Am Soc*

869        *Nephrol*. 2019;30(10):1811-1823. doi:10.1681/ASN.2019030303

870    57.    Taguchi A, Nishinakamura R. Higher-Order Kidney Organogenesis from Pluripotent

871        Stem Cells. *Cell Stem Cell*. 2017;21(6):730-746.e6. doi:10.1016/j.stem.2017.10.011

872    58.    Xia Y, Nivet E, Sancho-Martinez I, et al. Directed differentiation of human pluripotent

873        cells to ureteric bud kidney progenitor-like cells. *Nat Cell Biol*. 2013;15(12):1507-

874        1515. doi:10.1038/ncb2872

875    59.    Wilson SB, Little MH. DevKidCC Kidney Organoid Gene Expression Shiny

876    application. Shiny App. https://sbwilson91.shinyapps.io/devkidcc_interactive/.

877    Published 2021.

878  60.  Hochane M, van den Berg PR, Fan X, et al. Single cell RNA-sequencing of human

879    fetal kidneys. Gene Expression Omnibus.

880    https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE114530. Published 2019.

881  61.  Ransick A, Tran T, Lindstrom NO, De Sena Brandine G, McMahon AP. Single Cell

882    RNA-Seq profiling of human embryonic kidney outer and inner cortical cells and

883    kidney organoid cells. Gene Expression Omnibus.

884    https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE124472. Published 2019.

885  62.  Holloway EM, Spence JR, Wu JH. scRNA-seq of human fetal kidney tissue. EMBL-

886    EBI ArrayExress. https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-9083/.

887    Published 2020.

888  63.  Hafemeister C, Satija R. Normalization and variance stabilization of single-cell RNA-

889    seq data using regularized negative binomial regression. *Genome Biol*. 2019;20(1):296.

890    doi:10.1186/s13059-019-1874-1

891  64.  Wilson SB, Little MH. DevKidCC: Developing Kidney Cell Classifier. Github.

892    https://github.com/KidneyRegeneration/DevKidCC. Published 2021.

893  65.  Ransick A, Kim AD, De Sena Brandine G, Lindstrom NO, McMahon AP. Single Cell

894    RNA-Seq profiling human embryonic kidney cortex cells. Gene Expression Omnibus.

895    https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE102596. Published 2018.

896  66.  Menon R, Harder JL, Kretzler M, Otto EA, Freedman BS. Enhancing human kidney

897    organoid differentiation from pluripotent stem cells with high-throughput automation.

898    Gene Expression Omnibus.

899          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE109718. Published 2018.

900    67.   Howden SE, Vanslambrouck JM, Little MH, Lonsdale A, Wilson SB. Fate-mapping

901          within human iPSC-derived kidney organoids reveals conserved mammalian nephron

902          progenitor lineage relationships. Gene Expression Omnibus.

903          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE119561. Published 2019.

904    68.   Phipson B, Zappia L, Combes AN. Single cell RNA-Seq of four human kidney

905          organoids. Gene Expression Omnibus.

906          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE114802. Published 2018.

907    69.   Menon R, Harder JL, Otto EA, Kretzler M. Single-cell analysis of human kidney

908          organoids. Gene Expression Omnibus.

909          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE115986. Published 2019.

910    70.   Low JH, Li P, Chew EGY, Zhou B. Generating Patterned Kidney Organoids for

911          Studying Development and Diseases. Gene Expression Omnibus.

912          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE132026. Published 2019.

913    71.   Lawlor KT, Vanslambrouck JM, Little MH. Comparison manual and two types of

914          bioprinted kidney organoids by single cell RNA-seq. Gene Expression Omnibus.

915          https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152014. Published 2020.

916    72.   Wilson SB, Howden SE, Little MH. Distal nephron plasticity allows the induction of

917          ureteric tip and stalk for the modelling of collecting duct disease. Gene Expression

918          Omnibus. https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE161255.

919          Published 2020.

920    73.   Mae S-I, Ryosaka M, Sakamoto S, et al. Expansion of human iPSC derived ureteric

921          bud organoids with repeated branching potential. Gene Expression Omnibus.

922        https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152685. Published 2020.
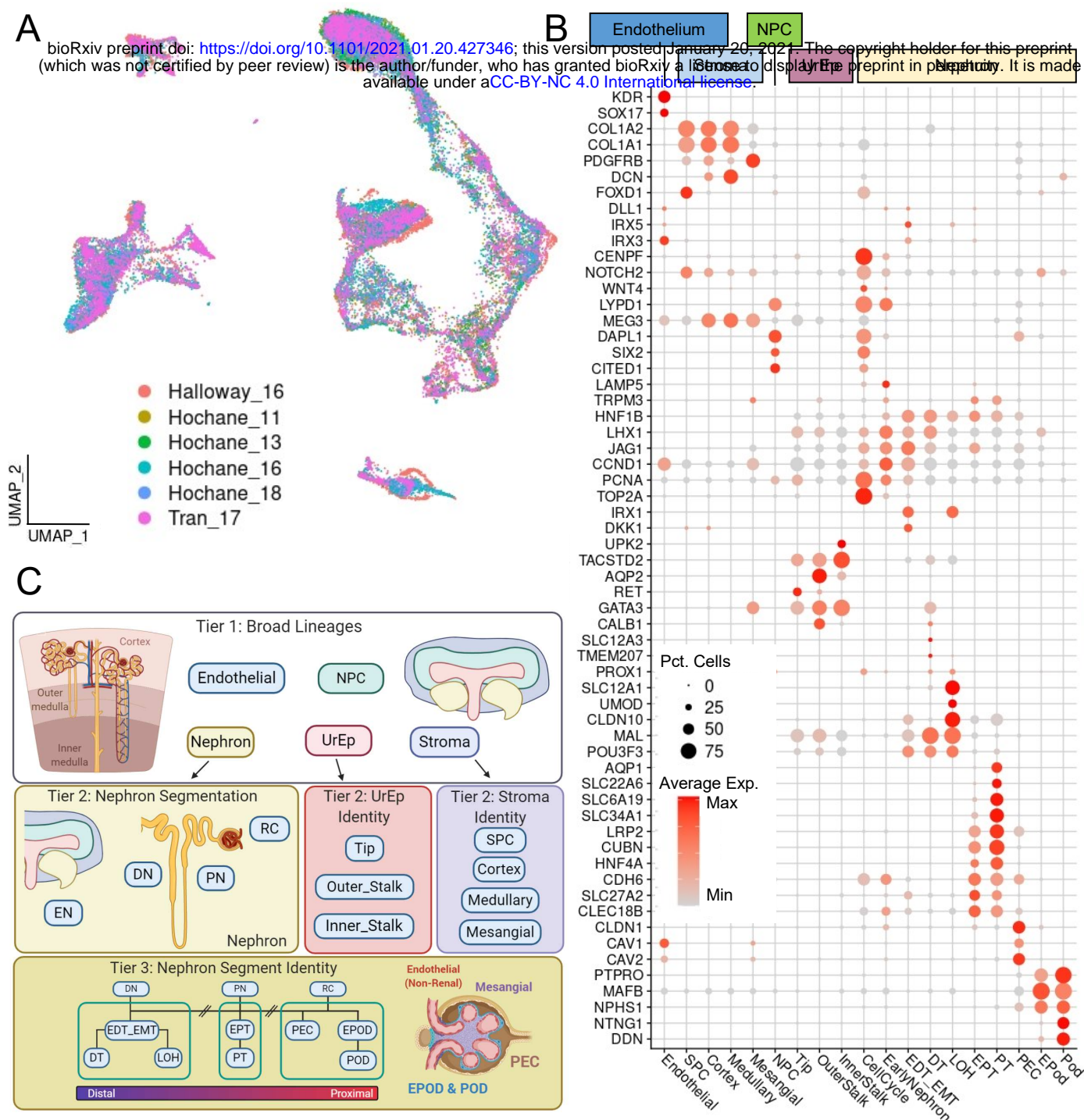
923

**Figure 1: Generating a comprehensive reference to train the classification model hierarchy**

A) The comprehensive human fetal kidney reference displayed in the first two UMAP dimensions grouped by their dataset of origin. The samples were integrated using *Harmony.* B) A DotPlot showing the expression of known marker and important genes present in each kidney segment. C) Graphical outline of the model hierarchy employed by *DevKidCC* to classify cells from single cell RNA sequencing kidney organoid datasets.

**Figure 2:** *DevKidCC* **classification in human fetal kidney and organoid datasets**

A) UMAP representation of the Lindstrom 2018 human fetal kidney dataset, grouped by the *DevKidCC* classification. B) A *DotPlot* showing the expression of key marker genes from the original Lindstrom 2018 analysis and their expression in the *DevKidCC* tier 1 annotated cell types. C) UMAP representation of the original annotation from Howden[13] kidney organoids. D) UMAP representation of cell classification using *DevKidCC* using thresholds of 0.7 and 0.9 similarity scores. E) A *ComparePlot* showing the reclassification of cells from the original Howden[13] annotation using *DevKidCC*. Cell classification is well conserved when considering differences in nomenclature. F) Directly comparing the original annotation of four organoid samples from Wu[12] to that of *DevKidCC* shows the congruence of classification with increased accuracy in determining kidney-like mesenchymal cells from non kidney-like cells.
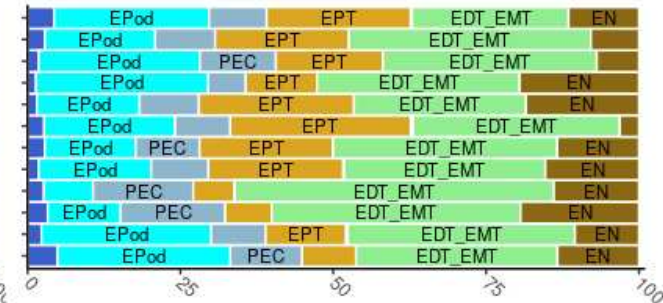
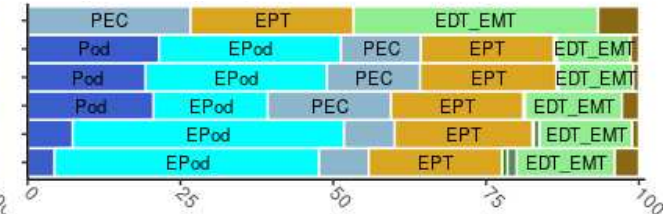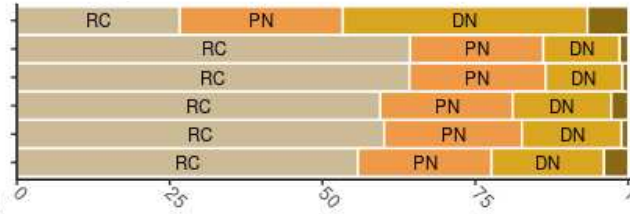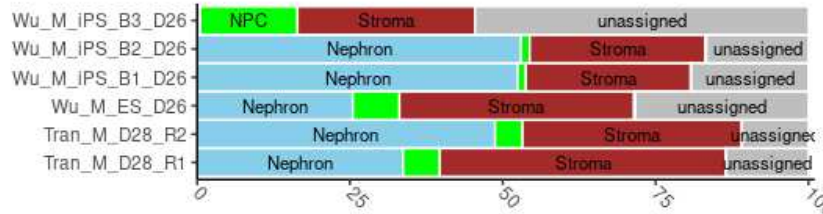# Freedman protocol

**Tier 1: All** — **Tier 2: Nephron** — **DevKidCC: Nephron**

| | Tier 1: All | Tier 2: Nephron | DevKidCC: Nephron |
|---|---|---|---|
| Har_F_SO_D20 | Nephron / unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |
| Har_F_D21 | Nephron, NPC, unassigned | RC, PN, DN | EPod, EPT, EDT_EMT |
| Har_F_B3R2_D19 | Nephron, unassigned | RC, PN, DN | EPod, PEC, EPT, EDT_EMT |
| Har_F_B3R1_D19 | Nephron, unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |
| Har_F_B2R2_D19 | Nephron, NPC, unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |
| Har_F_B2R1_D19 | Nephron, unassigned | RC, PN, DN | EPod, EPT, EDT_EMT |
| Har_F_B1R2_D19 | Nephron, unassigned | RC, PN, DN, EN | EPod, PEC, EPT, EDT_EMT, EN |
| Har_F_B1R1_D19 | Nephron, unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |
| Cz_F_VEGF_R2_D25 | Nephron, Stroma, unassigned | RC, DN, EN | PEC, EDT_EMT, EN |
| Cz_F_VEGF_R2_D25 | Nephron, Stroma, unassigned | RC, DN, EN | EPod, PEC, EDT_EMT, EN |
| Cz_F_Std_R2_D25 | Nephron, unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |
| Cz_F_Std_R1_D25 | Nephron, unassigned | RC, DN, EN | EPod, PEC, EDT_EMT, EN |

# Morizane protocol

| | Tier 1: All | Tier 2: Nephron | DevKidCC: Nephron |
|---|---|---|---|
| Wu_M_iPS_B3_D26 | NPC, Stroma, unassigned | RC, PN, DN | PEC, EPT, EDT_EMT |
| Wu_M_iPS_B2_D26 | Nephron, Stroma, unassigned | RC, PN, DN | Pod, EPod, PEC, EPT, EDT_EMT |
| Wu_M_iPS_B1_D26 | Nephron, Stroma, unassigned | RC, PN, DN | Pod, EPod, PEC, EPT, EDT_EMT |
| Wu_M_ES_D26 | Nephron, Stroma, unassigned | RC, PN, DN | Pod, EPod, PEC, EPT, EDT_EMT |
| Tran_M_D28_R2 | Nephron, Stroma, unassigned | RC, PN, DN | EPod, EPT, EDT_EMT |
| Tran_M_D28_R1 | Nephron, Stroma, unassigned | RC, PN, DN | EPod, EPT, EDT_EMT |

# Takasato protocol

| | Tier 1: All | Tier 2: Nephron | DevKidCC: Nephron |
|---|---|---|---|
| Wu_TO_iPS_D34 | Nephron, unassigned | RC, PN, DN, EN | EPT, EDT_EMT, EN |
| Wu_TC_iPS_D26 | Nephron, Stroma, unassigned | RC, PN, DN | EPT, EDT_EMT |
| Wu_TC_iPS_D19 | Nephron, unassigned | RC, PN, DN | EPT, EDT_EMT |
| Wu_T_iPS_B3_D26 | Nephron, Stroma, unassigned | RC, PN, DN | EPT, EDT_EMT |
| Wu_T_iPS_B2_D26 | Nephron, Stroma, unassigned | RC, PN, DN, EN | EPT, EDT_EMT, EN |
| Wu_T_iPS_B1_D26 | Nephron, Stroma, unassigned | RC, PN, DN, EN | EPT, EDT_EMT, EN |
| Wu_T_ES_B2_D26 | Nephron, Stroma, unassigned | PN, DN | EPT, EDT_EMT |
| Wu_T_ES_B1_D26 | Nephron, Stroma, unassigned | PN, DN | EPT, EDT_EMT |
| Sub_T_L2_D29 | Nephron, Stroma, unassigned | RC, PN, DN | EPT, EDT_EMT |
| Sub_T_L1_D29 | Nephron, Stroma, unassigned | PN, DN | EPT, EDT_EMT |
| Sub_T_L1_32 | Nephron, unassigned | PN, DN | EPT, EDT_EMT |
| PC_T_R4_D25 | NPC, Stroma, unassigned | RC, DN, EN | EPod, EDT_EMT, EN |
| PC_T_R3_D25 | Nephron, NPC, Stroma, unassigned | RC, PN | Pod, EPod, EPT |
| PC_T_R2_D25 | NPC, Stroma, unassigned | RC | Pod, EPod, PEC |
| PC_T_R1_D25 | Nephron, NPC, Stroma, unassigned | RC | Pod, EPod |
| LVH_T_PrintLine_D25 | Nephron, NPC, unassigned | RC, PN, DN | EPod, PEC, EPT, EDT_EMT |
| LVH_T_PrintBlob_D25 | Nephron, NPC, unassigned | RC, PN, DN | EPod, EPT, DT, EDT_EMT |
| LVH_T_ManualBlob_D25 | Nephron, NPC, unassigned | RC, PN, DN, EN | EPod, EPT, DT, EDT_EMT, EN |
| Ku_TMO_D25 | Nephron, NPC, unassigned | DN, EN | EDT_EMT, EN |
| How_T_D25 | Nephron, Stroma, unassigned | RC, PN, DN, EN | EPod, EPT, EDT_EMT, EN |

Legend Tier 1: UrEp, Nephron, NPC, Stroma, Endothelial, unassigned

Legend Tier 2: RC, PN, DN, EN, CC

Legend DevKidCC: Pod, EPod, PEC, PT, EPT, LOH, DT, EDT_EMT, EN, CC

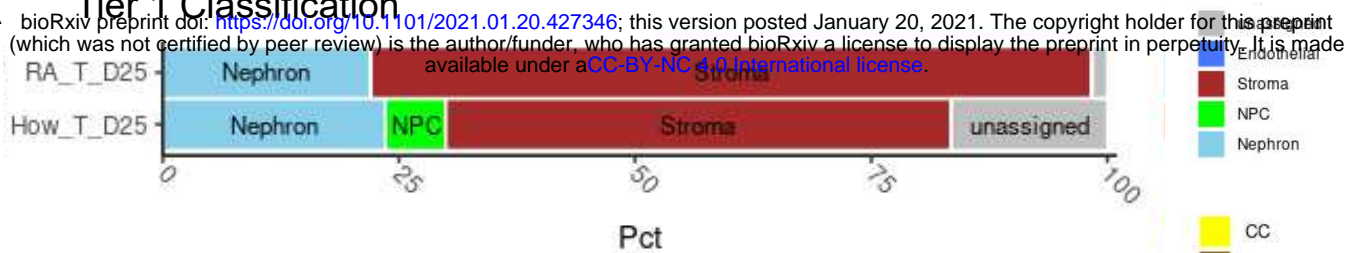**Figure 3: Direct comparison of organoids generated from different protocols**

All the mature age organoids were classified using DevKidCC. The proportions of identities are classified, the first column is all cells classified at the top tier, the second column is the nephron cells classed at the second tier, and the last column is the nephron cells classified at the third tier. Samples are grouped in rows by the protocol used to derive the organoids.
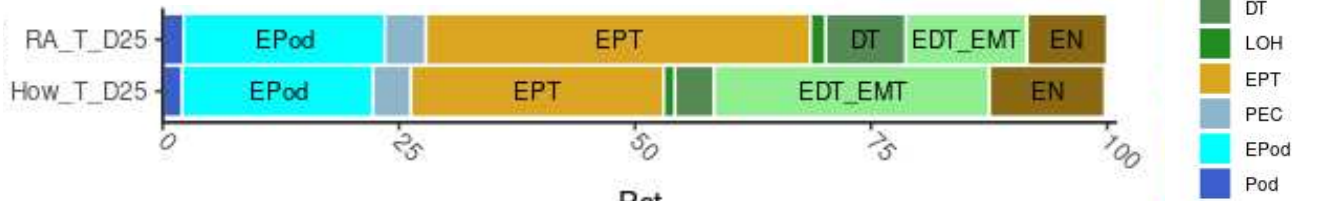
**Figure 4: NPCs deplete as organoids age and vary in transcriptional similarity to HFK NPCs**

A) Expression of NPC marker genes in the NPC cluster from each sample of each publication compared to the reference dataset. Takasato derived organoids show more congruence with the reference profile than Morizane or Freedman organoids. B) Proportion of tier one classification in i) organoids 17 days or more, ii) organoids 16 days or less, iii) monolayer differentiations showing the variation of NPC contribution across ages and datasets.

**Figure 5: Addition of Retinoic Acid during organoid development depletes NPCs and promotes glomerular maturation**

A) The proportions of cells classified between day 25 organoids of standard protocol or with Retinoic Acid (RA) added at D7+5 (D12). B) The proportion of nephron cells classified into their subpopulations shows an expansion of the proximal tubule at the expense of the distal tubule and early nephron when adding RA. C) FLOW plots showing the expression of SIX2$^{EGFP}$ cells from D7+11 (D18) organoids with and without RA addition at D7+5 (D12) of protocol. D) The expression and localisation of CLDN1 can clearly be seen to be improved in RA organoids by immunofluorescence. E) Comparative gene expression between the reference, standard organoid and RA organoid for informative genes. *HNF4A* and *SLC22A6* are expressed in immature and mature proximal tubule respectively. *OLFM3*, *MAFB* and *NPHS2* are expressed in precursor, immature/mature and mature podocytes respectively. *CDKN1C* is a post-mitotic marker. *PAX8* is expressed in the nephron epithelium but not mature podocytes. *CLDN1* is expressed in the parietal epithelial cells.

**Figure 6: Classification of Ureteric cell types in targeted cultures**

A) The *DevKidCC* classification for the Howden[24] and Mae[38] datasets show populations of ureteric cells (UrEp) classified. B) The complete classification of all cells not classed as "unassigned" shows interesting differences between the cell types present between the two datasets. Howden[24] has a mix roughly 50:50 split of Tip to Stalk cells of the ureteric cells classified, and also distal nephron cell types. Mae[38] has a higher proportion of Tip compared to stalk cells, and the nephron cell types are NPC not nephron epithelium. C) Density plots showing the spread of similarity scores for the UrEp, Nephron, NPC at the highest tier and the Tip and InnerStalk of the UrEp subsets.

**Supplementary Figure 1: Annotation of the comprehensive reference**

A) The reference annotation at Tier 1. B) The reference annotation at Tier 2.

**Supplementary Figure 2: Scoring outcomes**

A) UMAP plots showing the distribution of scores at the top tier for Stroma, NPC and Nephron, then the expression of *TCF21* which is kidney stromal marker. B) Density plots showing the distribution of cell scores for the Stroma, Nephron and NPC across the reference dataset, Lindstrom[19] human fetal kidney (HFK) and Howden[13] organoids datasets.

**Supplementary Figure 3: Epithelial maturation marker expression in end-stage organoids**

Gene expression profiles for immature and mature markers of A) Distal Tubule and B) Proximal Tubule present in the relevant epithelial segments. *CUBN, LRP2, HNF4A, LHX1* are immature markers, *SLC47A1, SLC22A2, SLC22A8, SLC12A1, SLC12A3, KCNJ1, SCNN1A* are mature markers.

**Supplementary Figure 4: Expression of glomerular maturation markers in end-stage organoids**

Gene expression profiles for markers of glomerular maturation. *WT1* is expressed from NPC through to mature podocytes. *OLFM3* is expressed in immature podocytes only. *NPHS2* is expressed in mature podocytes only. *CLDN1* is expressed in PECs. Immature podocytes express *LAMB1* and switch to *LAMB2* upon maturation. They also turn on *COL4A3*, *COL4A4* and *COL4A5* as they mature. Mature podocytes stop cycling and so are lowly expressing *TOP2A* and highly expressing *CDKN1A* (p21) and *CDKN1C* (p57).