



Heriot-Watt University  
Research Gateway

## A Roadmap for the Human Gut Cell Atlas

### Citation for published version:

Zilbauer, M, James, KR, Kaur, M, Pott, S, Li, Z, Burger, A, Thiagarajah, JR, Burclaff, J, Jahnsen, FL, Perrone, F, Ross, AD, Matteoli, G, Stakenborg, N, Sujino, T, Moor, A, Bartolome-Casado, R, Bækkevold, ES, Zhou, R, Xie, B, Lau, KS, Din, S, Magness, ST, Yao, Q, Beyaz, S, Arends, M, Denadai-Souza, A, Coburn, LA, Gaublomme, JT, Baldock, R, Papatheodorou, I, Ordovas-Montanes, J, Boeckxstaens, G, Hupalowska, A, Teichmann, SA, Regev, A, Xavier, RJ, Simmons, A, Snyder, MP, Wilson, KT, Gut Cell Atlas Consortium & Human Cell Atlas Gut Biological Network Consortium 2023, 'A Roadmap for the Human Gut Cell Atlas', *Nature Reviews Gastroenterology and Hepatology*, vol. 20, no. 9, pp. 597-614.  
<https://doi.org/10.1038/s41575-023-00784-1>

### Digital Object Identifier (DOI):

[10.1038/s41575-023-00784-1](https://doi.org/10.1038/s41575-023-00784-1)

### Link:

[Link to publication record in Heriot-Watt Research Portal](#)

### Document Version:

Peer reviewed version

### Published In:

Nature Reviews Gastroenterology and Hepatology

### Publisher Rights Statement:

This version of the article has been accepted for publication, after peer review and is subject to Springer Nature's AM terms of use, but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <http://dx.doi.org/10.1038/s41575-023-00784-1>

### General rights

Copyright for the publications made accessible via Heriot-Watt Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

Heriot-Watt University has made every reasonable effort to ensure that the content in Heriot-Watt Research Portal complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [open.access@hw.ac.uk](mailto:open.access@hw.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

# Title: A Roadmap for the Human Gut Cell Atlas

**Author(s):** Matthias Zilbauer<sup>\*‡1,2,3</sup>, Kylie R James<sup>#‡4,5</sup>, Mandeep Kaur<sup>#6</sup>, Sebastian Pott<sup>‡#7</sup>, Zhixin Li<sup>#8</sup>, Albert Burger<sup>‡#9</sup>, Jay R Thiagarajah<sup>10</sup>, Joseph Burclaff<sup>11,12</sup>, Frode L Jahnsen<sup>13</sup>, Francesca Perrone<sup>1,2</sup>, Alexander D. Ross<sup>1,2,14</sup>, Gianluca Matteoli<sup>15</sup>, Nathalie Stakenborg<sup>‡15</sup>, Tomohisa Sujino<sup>16</sup>, Andreas Moor<sup>‡17</sup>, Raquel Bartolome-Casado<sup>13,18</sup>, Espen S Bækkevold<sup>13</sup>, Ran Zhou<sup>‡19</sup>, Bingqing Xie<sup>‡20</sup>, Ken S Lau<sup>‡21</sup>, Shahida Din<sup>‡22</sup>, Scott T Magness<sup>12,23</sup>, Qiuming Yao<sup>24</sup>, Semir Beyaz<sup>25</sup>, Mark Arends<sup>‡26</sup>, Alexandre Denadai-Souza<sup>‡27</sup>, Lori A. Coburn<sup>‡28,29</sup>, Jellert T Gaublomme<sup>‡30</sup>, Richard Baldock<sup>‡31</sup>, Irene Papatheodorou<sup>‡32</sup>, Jose Ordovas-Montanes<sup>10</sup>, Guy Boeckstaens<sup>‡15</sup>, Anna Hupalowska<sup>33</sup>, Sarah A Teichmann<sup>‡18,34</sup>, Aviv Regev<sup>33,35</sup>, Ramnik J Xavier<sup>‡36</sup>, Alison Simmons<sup>37</sup>, Michael P Snyder<sup>38</sup>, Keith T. Wilson<sup>‡28,29</sup>, Gut Cell Atlas<sup>‡</sup> & Human Cell Atlas Gut Biological Network

\*Corresponding author

#Authors that acted as section leads.

‡Members of the Gut Cell Atlas, an Initiative Supported by the Helmsley Charitable Trust.

## Author affiliations:

<sup>1</sup>Wellcome-MRC Cambridge Stem Cell Institute, University of Cambridge, Cambridge, UK,

<sup>2</sup>University Department of Paediatrics, University of Cambridge, Cambridge, UK

<sup>3</sup>Department of Paediatric Gastroenterology, Hepatology and Nutrition, Cambridge University Hospitals, Cambridge, UK

<sup>4</sup>Garvan Institute of Medical Research, NSW, Australia,

<sup>5</sup>School of Biomedical Sciences, University of New South Wales, NSW, Australia

<sup>6</sup>School of Molecular and Cell Biology, University of the Witwatersrand, Johannesburg, South Africa

<sup>7</sup>Section of Genetic Medicine, Department of Medicine, University of Chicago, Chicago, USA

<sup>8</sup>Dana-Farber Cancer Institute, Boston, USA

<sup>9</sup>Department of Computer Science, Heriot-Watt University, Edinburgh, UK

<sup>10</sup>Division of Gastroenterology, Hepatology and Nutrition, Boston Children's Hospital, Harvard Medical School, Boston, USA.

34 <sup>11</sup>Joint Department of Biomedical Engineering, University of North Carolina at Chapel Hill,  
35 North Carolina State University, Chapel Hill, North Carolina, USA  
36 <sup>12</sup>Center for Gastrointestinal Biology and Disease, University of North Carolina at Chapel  
37 Hill, Chapel Hill, North Carolina, USA  
38 <sup>13</sup>Department of Pathology, Oslo University Hospital and University of Oslo, Oslo, Norway  
39 <sup>14</sup>University Department of Medical Genetics, University of Cambridge, Cambridge, UK  
40 <sup>15</sup>Translational Research Center for Gastrointestinal Disorders (TARGID), Department of  
41 Chronic Diseases, Metabolism and Ageing, KU Leuven, Leuven, Belgium  
42 <sup>16</sup>Center for the Diagnostic and Therapeutic Endoscopy, School of Medicine, Keio  
43 University, Tokyo, Japan  
44 <sup>17</sup>Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland  
45 <sup>18</sup>Wellcome Sanger Institute, Hinxton, Cambridge, UK  
46 <sup>19</sup>Section of Genetic Medicine, University of Chicago, Chicago, USA  
47 <sup>20</sup>Department of medicine, University of Chicago, Chicago, USA  
48 <sup>21</sup>Epithelial Biology Center and Department of Cell and Developmental Biology, Vanderbilt  
49 University School of Medicine, Nashville TN, USA  
50 <sup>22</sup>Edinburgh IBD Unit, Western General Hospital, NHS Lothian  
51 <sup>23</sup>Joint Department of Biomedical Engineering, University of North Carolina at Chapel  
52 Hill/North Carolina State University, Chapel Hill, North Carolina, USA  
53 <sup>24</sup>Department of Computer Science and Engineering, University of Nebraska Lincoln,  
54 Lincoln, Nebraska, USA  
55 <sup>25</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, 11724 USA  
56 <sup>26</sup>Edinburgh Pathology, University of Edinburgh, Institute of Genetics and Cancer,  
57 Edinburgh, UK  
58 <sup>27</sup>Laboratory of Mucosal Biology, Department of Chronic Diseases, Metabolism and Ageing,  
59 KU Leuven, Belgium  
60 <sup>28</sup>Vanderbilt University Medical Center, Nashville, USA  
61 <sup>29</sup>Veterans Affairs Tennessee Valley Healthcare System, Nashville, USA  
62 <sup>30</sup>Department of Biological Sciences, Columbia University, New York, USA  
63 <sup>31</sup>University of Edinburgh, UK  
64 <sup>32</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, EMBL-EBI,  
65 Wellcome Genome Campus, Hinxton, UK.  
66 <sup>33</sup>Genentech, San Francisco, CA, USA  
67

68 <sup>34</sup>Theory of Condensed Matter Group, Cavendish Laboratory/Department of Physics,  
69 University of Cambridge, Cambridge, UK

70 <sup>35</sup>Klarman Cell Observatory, Broad Institute of Harvard and MIT, Cambridge, MA, USA.

71 <sup>36</sup>Broad Institute and Department of Molecular Biology, Massachusetts General Hospital,  
72 Harvard Medical School, Boston, USA

73 <sup>37</sup>MRC Human Immunology Unit, MRC Weatherall Institute of Molecular Medicine,  
74 University of Oxford, Oxford, UK.

75 <sup>38</sup>Stanford University School of Medicine, Stanford, USA

76

77 **Corresponding author:**

78 Matthias Zilbauer

79 Email: mz304@cam.ac.uk

80

81

82 **Abstract:**

83 The number of studies investigating the human gastrointestinal tract using various single cell  
84 profiling methods have increased substantially in recent years. Although this provides a unique  
85 opportunity for the generation of a first comprehensive Human Gut Cell Atlas (HGCA), there  
86 are still a range of major challenges ahead. Above all, the ultimate success will largely depend  
87 on a structured and coordinated approach that aligns global efforts undertaken by a large  
88 number of research groups. Here we present a detailed roadmap for the generation of the HGCA  
89 on behalf of the Human Cell Atlas (HCA) Gut Biological Network. Based on the consensus  
90 opinion of experts from across the globe, we outline the main requirements for a first complete  
91 Gut Atlas by summarising existing datasets and highlighting anatomical regions and/or tissues  
92 with limited coverage. We provide recommendations for future studies, discuss key  
93 methodologies as well as the importance of integrating the healthy gut atlas with related  
94 diseases and gut organoids. Importantly, we critically discuss computational tools available  
95 and provide recommendations to overcome key challenges.

96

97 **Introduction:**

98 The intestinal tract is one of the most complex organs in the human body and serves a wide  
99 range of functions, including the digestion and absorption of nutrients as well as representing  
100 a major site of immune interactions. Different anatomical sections of the intestinal tract play

101 specific roles in the digestive process requiring the presence and complex interaction of various  
102 cell types. Additionally, the interaction with trillions of nearby microbes has been shown to be  
103 of critical importance adding further complexity to this finely tuned symbiosis. A detailed  
104 knowledge of gut physiology and cellular function in health is a prerequisite for investigating  
105 related diseases such as bowel cancer and inflammatory conditions.

106 The development of methodologies allowing genome wide molecular profiling on a single cell  
107 (sc) level has opened unprecedented opportunities to generate detailed anatomical maps of the  
108 human body. Although the number of studies investigating the human intestine using various  
109 sc profiling methods have increased substantially in recent years, generating a first  
110 comprehensive Human Gut Cell Atlas (HGCA) is associated with major challenges. Above all,  
111 the ultimate success will largely depend on a structured and coordinated approach that aligns  
112 global efforts. The Human Cell Atlas (HCA) Gut Biological Network connects leading experts  
113 in a wide range of related areas providing an ideal platform to lead the development of the first  
114 sc HGCA. Here we provide a detailed roadmap including recommendations on how to combine  
115 existing data and requirements for future studies. We discuss key methodologies as well as the  
116 importance of integrating the healthy gut atlas with related diseases and gut organoids.  
117 Importantly, we will critically reflect on computational tools available, highlight any existing  
118 limitations and provide recommendations to overcome key challenges. Importantly, generating  
119 a first comprehensive HGCA will provide researchers, scientist and clinicians with greater  
120 scope and power to enable novel discoveries into intestinal biology and pathophysiology with  
121 the ultimate goal of improving human health.

122

## 123 **Mapping the Human Gastrointestinal tract**

124 The human digestive tract reaches from the oral cavity to the rectum and includes the  
125 oesophagus, stomach, and the small and large intestine. Generating a complete map requires  
126 the inclusion of all organs associated with the digestive process and detailed sampling of each  
127 gut segment to capture distinct anatomical changes along the cephalad-caudal axis (Figure 1).

128 A unique feature of the human intestine is the presence of trillions of microbes that exist in a  
129 finely balanced symbiosis and are thought to be of critical importance to the host<sup>1</sup>. As such,  
130 profiling microbial communities will provide an opportunity to interrogate the cross talk  
131 between microbes and host cells. However, this task is associated with the major challenges of  
132 profiling the human gut microbiome including the vast inter-individual variability, differences  
133 in its composition along the intestinal tract as well as the dynamic changes of microbes over

134 time. As a result, integrating the gut microbiome in the HGCA is likely to be part of future  
135 efforts and will require close interaction with expert researches in this area. Another important  
136 aspect unique to the gastrointestinal (GI) tract is the exposure to and interaction with a wide  
137 range of food/nutritional components as well as antigens or potential toxins present in the daily  
138 diet. Robust evidence suggests that GI development, function, and predisposition to related  
139 disease are strongly correlated with dietary habits<sup>2,3</sup>. Therefore, documenting details on dietary  
140 habits will add critical information and further increase the value of molecular profiles  
141 generated. Importantly, epigenetic mechanisms can mediate the impact of exposure to  
142 environmental factors into stable cellular phenotypes. The possibility of simultaneously  
143 profiling multiple molecular layers at sc resolution including epigenetic, transcriptomic, and  
144 proteomic signatures will us to unravel interactions between intestinal host cells and our  
145 environment.

146

## 147 **Sampling strategies**

148 Acquiring human tissue in sufficient quantities and from a large number of donors is one of the  
149 major challenges associated with the generation of the HGCA. There are several considerations  
150 that apply to the GI tract. In general, there are three major strategies for sampling the alimentary  
151 tract including mucosal biopsies, live surgical resections, and resections from deceased donors.  
152 Each strategy and tissue source has unique advantages and disadvantages (Figure 2). Mucosal  
153 biopsies can be obtained during routine endoscopy, but access is restricted to the upper and  
154 lower GI tract, with jejunum and proximal ileum rarely sampled. Multiple biopsies can be taken  
155 from the same patient allowing comparisons of different anatomical regions. In contrast to  
156 mucosal biopsies, which can be obtained from healthy individuals, live resections are derived  
157 from patients with major GI pathology. Recent work identified transcriptional differences  
158 between healthy intestines and normal-appearing intestine proximal to chronic inflammation<sup>4</sup>.  
159 Hence, resections from patients with GI disease should be used with caution when mapping  
160 ‘healthy’ cells. Advanced endoscopic imaging technology such as confocal laser  
161 endomicroscopy can provide additional information and guide sampling strategies<sup>5</sup>. Tissue  
162 from deceased organ donors allows sampling of entire organs, providing the opportunity for  
163 widescale comparisons between regions. Another major advantage of gut resections is the  
164 provision of all tissue layers from luminal contents through to musculature and serosa. In  
165 contrast, forceps biopsies only capture the mucosa including the epithelium and immediate  
166 subepithelial cells. Size also varies drastically between biopsies and resections impacting on

167 the number and type of downstream analyses that can be performed. Finally, tissue quality may  
168 differ between strategies with the time from sampling to processing of tissue being of critical  
169 importance.

170 In summary, compared to other human organs, the GI tract is frequently being sampled during  
171 routine clinical procedures providing unique opportunities to obtain tissue from healthy  
172 individuals as well as many related gut diseases. However, for complete coverage of all  
173 anatomical regions and layers of the intestinal tube, the use of surgical resection materials as  
174 well as deceased donor tissue is essential.

175

## 176 **Documenting metadata**

177 The generation of a HGCA requires combining a vast number of different datasets generated  
178 by groups around the globe<sup>6,7</sup>. A fundamental aspect that determines the integration and  
179 comparability of datasets, is the documentation of detailed metadata<sup>8</sup>. The HCA provides  
180 extensive guidance on this topic<sup>9,10</sup>. In the following we briefly describe basic information  
181 required to allow future studies to be included in the HGCA and provide a template metadata  
182 table (Table 1). In line with recommendations provided by the HCA, required metadata can be  
183 divided into the following main aspects: a) study design, b) donor information, c) sample  
184 information, d) sample processing and e) data generation. In brief, a detailed description of  
185 study design including patient inclusion criteria as well as sampling strategy is of major  
186 importance. Donor information should include baseline demographic data, details on any  
187 known medical conditions, particularly those affecting the intestine, and medications can be of  
188 major benefit. Similarly, the impact of dietary habits on gut physiology is well established and  
189 information on dietary habits could be used to identify novel aspects of dietary factors involved  
190 in gut health and disease<sup>2,3</sup>. Baseline information on sample type include the method by which  
191 it was obtained, sample area covered and anatomic location. Depending on sample type, the  
192 exact anatomical location from which the sample was obtained can be difficult to determine.  
193 However, the level of detail provided has a major impact on the comparability of studies and  
194 interpretation of data generated. If available, representative haematoxylin & eosin (H&E)  
195 and/or other staining can be directly linked to samples processed for sc studies. Information on  
196 sampling procedures should include details on the time between sampling and processing, as  
197 well as storage duration. Furthermore, providing a detailed protocol on sample processing  
198 including tissue dissociation, cell viability, possible enrichment of individual cell types as well  
199 as equipment used is important. The development of standardised protocols is subject of

200 ongoing work within the Human Cell Atlas. Active engagement of the scientific community  
201 and use of available resources including recommended protocols will increase comparability  
202 of data generated by different groups.

203

## 204 **Protecting donor identity and engaging communities**

205 Documentation, storage and sharing of patient/donor information raises important ethical  
206 considerations<sup>11</sup>. Regulations vary considerably between different institutions and countries  
207 further complicating the sharing of tissues and data. In the European Union, regulations  
208 concerning the protection of personal data and the responsible sharing of such data have been  
209 formulated in the General Data Protection Regulation (GDPR) 2016/679<sup>12</sup>. In the United States  
210 there are several federal regulations that must be considered when using human tissue and data  
211 in the context of research including the Department of Health and Human Services (HHS) and  
212 the Food and Drug Administration (FDA)<sup>13</sup>. In addition to major differences in the rules and  
213 regulations between countries across the globe, exchanging tissue and data between countries  
214 poses further complexities and barriers. However, given the critical importance of metadata for  
215 the integration of different datasets, appropriate donor consent for sharing information, data  
216 and/or tissue is a prerequisite. In addition to explicit informed consent, pseudonymisation of  
217 data and donor privacy must be ensured. A critical aspect of achieving a broad utility for the  
218 HGCA is the inclusion and participation of ancestrally and geographically diverse populations.  
219 This requires close partnerships between researchers, funders, and potential patient/donor  
220 communities, particularly in the case of historically neglected or mistreated populations in  
221 research<sup>14</sup>. Ideally, studies providing data for the HGCA should include as many community  
222 stake-holders as feasible, with the long-term goal of having research questions formulated by  
223 researchers and participating populations working in partnership<sup>15</sup>.

224

## 225 **Available methodologies and their application**

226 A wide range of sc and single nuclei (sn) genomics tools have been developed and enabled the  
227 creation of high-resolution cellular atlases from human organs and tissues<sup>6 16 17</sup>. By and large,  
228 methodologies that will be used to generate the HGCA are the same as those used to profile  
229 other organs or systems (Figure 3). In addition to sc transcriptional profiling, key methods  
230 include spatial transcriptomics and sc profiling of other molecular layers such as the genotype  
231 (G&T-seq)<sup>18</sup> and epigenotype including chromatin states (sc Assay for Transposase-Accessible  
232 Chromatin sequencing (ATAC)-seq)<sup>19</sup> and DNA methylation (whole genome bisulfite



233 sequencing (scGBseq))<sup>20</sup>. Furthermore, several methods have been developed that allow  
234 simultaneous profiling of multiple molecular layers such as the 10x Genomics Multiome kit,  
235 which combines snATACseq and snRNAseq<sup>21</sup>. ScNOME-seq is capable of measuring DNA  
236 methylation and chromatin accessibility within single cells<sup>22</sup>, whilst scNMTseq captures  
237 transcript levels in addition to chromatin accessibility and DNA methylation<sup>23</sup>. Furthermore,  
238 the use of methods that enable enrichment for specific and/or rare cell types will be key to  
239 achieving complete coverage. For example, MIRACL-seq allows label free enrichment of rare  
240 cell types<sup>24</sup> whilst the Chromium Single Cell Immune Profiling assay provided by 10x  
241 Genomics allows detailed immune cell profiling including full-length, paired B-cell or T-cell  
242 receptor (BCR/TCR) sequences, surface protein expression, antigen specificity, and 5' gene  
243 expression. Another key challenge for the HGCA is the timely processing of fresh tissue  
244 samples. Isolation of tightly connected epithelial cells is associated with damage, highlighting  
245 the need for suitable dissociation protocols as well as computational tools to exclude dead or  
246 damaged cells. Furthermore, dissociation of cells at 37°C for prolonged time induces the  
247 expression of early-immediate genes, thereby disrupting nuanced cellular states. Recent  
248 advances have led to the development of snRNAseq protocols which can be applied to frozen  
249 tissue samples therefore easing the burden of sample processing. High quality nuclei can be  
250 readily isolated from frozen samples and subjected to snRNAseq<sup>25</sup>. Although fewer transcripts  
251 are recovered, these data capture most cell populations and vastly increase the availability of  
252 tissue sources. The decision on which tissue processing method is most appropriate will be  
253 largely determined by the individual setup including tissue availability. However, a detailed  
254 description on experimental procedures is of critical importance for the value of generated data  
255 for the HGCA.

256 As a result of methods requiring tissue dissociation prior to sc profiling, information on spatial  
257 arrangement of cells that are crucial for their function is lost entirely. Multiple spatial methods  
258 are used to capture cells in their anatomical context including Next Generation Sequencing  
259 (NGS)-based (e.g., 10x Genomics' Visium)<sup>26</sup> and imaging-based (e.g., MERFISH) assays<sup>27,28</sup>.  
260 Typically, frozen or Formalin-Fixed Paraffin-Embedded (FFPE) tissue sections are used as  
261 starting material and ideally are obtained from the same area from which tissue for sc assays  
262 was taken. Integration of spatial approaches with sc genomics provides both the cellular  
263 resolution as well as spatial organisation of cell combinations and states (functional tissue  
264 units) as an essential framework for a comprehensive atlas of the human gut. Several  
265 computational tools have been developed and will be discussed below.

266 In summary, existing methodologies offer researchers a wide range of opportunities to address  
267 their research goals. Whilst integrating data derived from studies using different methodologies  
268 poses a challenge for the HGCA, the key factor determining the potential value is the quality  
269 of data generated combined with the documentation of detailed metadata. Additionally, we  
270 encourage researchers across the globe to engage with the HCA Gut Bionetwork prior to as  
271 well as during their sc studies in order to further maximise the use of generated data.

272

## 273 **Computational challenges and opportunities**

274 The HGCA aims to generate a resource for the scientific community that is reliable, easy to  
275 access and user friendly. Given the vast amount of diverse datasets generated by a wide range  
276 of research groups using different methodological approaches, the usage of existing as well as  
277 development of novel computational tools is of critical importance. Main challenges include  
278 adequate measures to identify studies with sufficient quality as well as combining and mapping  
279 of diverse datasets. Furthermore, the development of a user-friendly interface allowing the data  
280 to be explored by the scientific community is a key requirement. Whilst many of these tasks  
281 also apply to most if not all tissue mapping studies, there are several additional challenges and  
282 opportunities that are specific to the intestine. These include the integration of host cell  
283 molecular signatures with mucosa associated microbial profiles, mapping of datasets to specific  
284 anatomical locations along the cranio-caudal axis of the gut and alignment of healthy gut  
285 derived datasets with disease state and intestinal organoid models. With the rapid development  
286 of scRNAseq technology, numerous computational tools and pipelines have been developed to  
287 analyse sc data<sup>29-31</sup>. These include programmes/packages allowing pre-processing of data such  
288 as quality control (QC) checks, normalisation and batch correction, sequence alignment as well  
289 as detection and removal of cell doublets. A summary of relevant existing computational  
290 tools/packages is provided in **Table 2**. In the following, we briefly discuss the main  
291 computational challenges and opportunities related to the development of the HCGA.

292 ***Combining and integrating datasets across modalities:*** Combining and integrating studies  
293 performed by different groups using either the same or different methodological approaches  
294 poses a major challenge. An important first step is the application of stringent selection criteria  
295 and quality control measures to select high quality datasets. In addition to data quality, the  
296 provision of detailed metadata is a prerequisite for successful integration into the HGCA.  
297 Another major challenge is the removal of potentially confounding batch effects<sup>32</sup>. Available  
298 approaches including Seurat<sup>33</sup>, LIGER<sup>34</sup>, and Harmony<sup>35</sup>. The integration of different

299 molecular profiles forms another key computational challenge and analytical frameworks have  
300 been developed to integrate multiple data types in the same cells including GLUE<sup>36</sup>, MOFA+<sup>37</sup>  
301 and Cobolt<sup>38</sup>. Additionally, a statistical regression framework has been developed to integrate  
302 clinical metadata in sc profiling studies<sup>24</sup>.

303 **Cell type annotation:** Abundant and well characterised cell types can be reliably identified  
304 through unsupervised clustering algorithms followed by comparison of key marker gene  
305 expression profiles<sup>39</sup>. Generating a list of known marker genes based on existing literature  
306 forms a key aspect in the process of performing reliable cell annotation and requires the  
307 contribution of experts in the field. Indeed, combining expertise and large datasets provides a  
308 unique opportunity to develop an extensive cell marker gene list using both automated and  
309 supervised cell identification approaches (e.g. MACA<sup>40</sup>, singleR<sup>41</sup>, ScType<sup>42</sup>). For each  
310 confidently identified gut cell type, marker genes can then be inferred by distinguishing the  
311 known clusters (e.g. COMET<sup>43</sup>, COSG<sup>44</sup>) ultimately leading to a fully automated annotation  
312 procedure. Applying automated annotation approaches to larger datasets will also yield  
313 unknown/novel cell clusters, which must be subjected to additional validation studies and  
314 functional characterisations for example by using human gut organoid models as will be  
315 discussed below.

316 **Mapping cellular location:** Mapping individual cell types to their specific anatomical location  
317 within the human intestine is of critical importance. Combining sc sequencing technologies  
318 with spatial transcriptomics enables accurate profiling of cell type topography and several tools  
319 have been developed to address the computational aspects involved<sup>45</sup> including ASAP<sup>46</sup>,  
320 novoSpaRc<sup>47</sup>, and Cell2location<sup>48</sup>. Multi-omic spatial profiling efforts result in a gene, plus  
321 protein-by-cell matrix, annotated with spatial coordinates for the centroid position of each cell.  
322 This matrix can then be used as input to the recently developed open source for spatial analysis  
323 pipelines<sup>49</sup> allowing an unbiased identification of spatial patterns based on gene expression,  
324 distinct cellular neighborhoods and cell-to-cell interactions. Mapping the spatial network of  
325 cells in the human intestine will enable the identification of genes with coherent spatial  
326 expression patterns using methods such as BinSpect and SpatialDE<sup>50</sup>. Spatial domains with  
327 coherent gene expression patterns can be identified with a Hidden Markov random field  
328 models, that find spatial domains by comparing the gene expression patterns of each cell with  
329 its neighborhood to search for coherent patterns<sup>49</sup>. The prevalence of specific cell-to-cell  
330 interactions in the human intestine can be evaluated by the frequency that each pair of cell types  
331 is proximal to each other<sup>51</sup>. Finally, spatially interacting cell types can be analyzed to identify  
332 which known ligand-receptor pairs show significantly increased or decreased co-expression,

333 which could serve as a proxy for signaling activity, by creating a background distribution  
334 through spatially aware permutations<sup>49</sup>, which increases predictive power when compared to  
335 spatially unaware permutations. The inclusion of imaging data in the HGCA provides an  
336 additional approach towards mapping cellular location and will be discussed in more detail  
337 below.

338 ***Cellular and microbial cross talk:*** Investigating cellular cross-talk through the analyses of sc  
339 transcriptomic datasets is subject to active research in the field<sup>52</sup>. Amongst the most compelling  
340 existing approaches are CellPhoneDB, which enables interrogation of context specific  
341 crosstalk between different cell types based on an extensive database of known receptors and  
342 ligands<sup>53</sup>. Applying these algorithms to large and/or combined datasets is likely to yield novel  
343 insight into fundamental aspects of human intestinal physiology. Extending such analyses to  
344 include the crosstalk between human host cells and the gut microbiome represents another  
345 significant challenge. To this end, several computational strategies have been developed to  
346 identify cellular and microbial cross talk in the gut utilizing sc data from human gut and  
347 microbial samples<sup>54 55</sup>. Recent studies have provided evidence for the major value of  
348 combining in situ spatial-profiling technologies with single-cell sequencing to interrogate host-  
349 microbial interactions<sup>56</sup>.

350 ***Regulatory network inference:*** Successful integration of datasets and/or molecular profiles has  
351 the potential to unravel both known and novel cell type specific molecular networks. Several  
352 tools are available to enable the prediction of gene regulatory networks (GRNs) that control  
353 fundamental biological functions such as cellular differentiation and cell state transitions (e.g.  
354 SCENIC<sup>57</sup>, GRNBoost2<sup>58</sup>, PIDC<sup>59</sup>). Identifying the GRNs contributes towards our  
355 understanding of how coordinated expression of transcription factor networks drives the  
356 expression of their respective target genes and ultimately to shape and maintain gut cell  
357 identity<sup>60</sup>.

358 ***HGCA portal:*** A user-friendly and interactive web portal is essential for making the HGCA  
359 accessible to a broad research community<sup>61</sup>. There are a number of existing sc data portals  
360 providing a wide range of tools to interrogate datasets and a summary is provided below<sup>62-64</sup>.  
361 These portals provide an excellent starting point and we envision complete integration of the  
362 HGCA into the larger portals, which also enables gut specific tools to be developed.

363

364 **Development of a Common Coordinate Framework (CCF) for the HGCA**

365 A critical part of the HGCA is its ability to capture the location of cells in their anatomical and  
366 physiological context<sup>6</sup>. The concept of a Common Coordinate Framework (CCF) has been  
367 introduced to facilitate this capability. Rood et al. provide the following definition<sup>65</sup>: “*An*  
368 *underlying reference map of organs, tissues, or cells that allows new individual samples to be*  
369 *mapped to determine the relative location of structural regions between samples.*” Like any  
370 other computational model, a CCF provides an abstract representation of a real-world item.  
371 CCFs have been developed for several human organs including the lung<sup>66</sup>, the brain<sup>67</sup> and the  
372 use of vasculature as a CCF of the entire human body<sup>68</sup>. A first step towards the development  
373 of a HGCA CCF requires the generation of one-, two- and three-dimensional models of the  
374 intestinal tract. The one-dimensional conceptual model is based on the clinical view and  
375 biological organisation of the digestive tract in which distance from anatomical landmarks  
376 provides the critical location metric<sup>69</sup>. ‘Anatomograms’ are two-dimensional graphical models  
377 that include a simplified view of the gut for representational purposes thereby providing a basic  
378 framework for accurately capturing the anatomical location of the tissue component and cell  
379 types<sup>70</sup>. Three-dimensional models can be generated by integrating Computerized Tomography  
380 (CT) or Magnetic Resonance Imaging (MRI) data. Generating a framework also involves  
381 developing a system of associated coordinates that provide detailed specification of anatomical  
382 location as well as tissue and cell type. Hence, providing a detailed anatomical location of  
383 sampling sites is of critical importance. However, there are limitations to the accuracy of  
384 obtaining exact anatomical location during surgery or intestinal endoscopy. Stating  
385 proportional lengths and/or distances from reliable anatomical landmarks could improve  
386 accuracy and help to facilitate integration of generated datasets into the HGCA CCF. Another  
387 complementary strategy to localise cell types in 2D space is the combined use of scRNAseq  
388 with spatial transcriptomics as discussed above. This requires the development of  
389 computational tools and efforts to allow the integration of scRNAseq and spatial data into a  
390 CCF, and work in this area is ongoing.

391 Several CCFs have been developed and are accessible through existing data portals. For  
392 example, as part of the Human BioMolecular Atlas Program (HuBMAP), 3D models of both  
393 the human large and small intestine have been developed. The output is a series of surface  
394 models of the gut representing the outline within a 3D context allowing placement of tissue or  
395 cell types at the macro-level. This data can be accessed online through a CCF Registration User  
396 Interface (CCF-RUI, [Figure 4A](#)) and mapped data can be explored using the Exploration User  
397 Interface (CCF-EUI). Furthermore, as part of a collaborative project funded by the Leona M.  
398 and Harry B. Helmsley Charitable Trust (Gut Cell Atlas)<sup>71</sup>, researchers from The European

399 Bioinformatics Institute (EBI Cambridge) and Edinburgh have developed a 1D linear  
400 conceptual CCF model for the human large and small intestine that is linked to the 2D  
401 anatomogram<sup>69,70</sup>. These mappings coupled with the inverse transform from the 3D and 2D  
402 spaces back to the 1D linear model allow spatial interoperability between all representations  
403 and therefore the capability to compare and query data registered to any CCF using a web-  
404 based visualisation platforms (Figure 4B, C and D)<sup>69</sup>.

405 In addition to providing a framework for accurately mapping single datasets, CCFs are also  
406 capable of integrating other related datasets and can be managed across geographically  
407 dispersed data repositories. Importantly, the value of CCFs is directly related to the level of  
408 standardisation achieved as it determines the degree of potential interoperability across the  
409 datasets. Amongst the ways to increase standardisation are the development of consensus  
410 guidelines on how anatomical location can be most accurately documented for the purpose of  
411 integrating sc studies into the HGCA CCF. A starting point can be the use of a standard  
412 metadata template (Table 1) and the adherence to the minimum information standard for the  
413 description of cell/tissue sources in the gut. Ultimately, successful development of a CCF will  
414 greatly enhance the value and broad application of the HGCA.

415

## 416 **Summary of relevant existing datasets, studies and portals**

417 The number of published studies reporting on sc profiling of the human gut has increased  
418 substantially since the founding of the HCA in 2016 (Figure 5 and Supplementary Table 1).  
419 Successful integration of these studies provides a strong foundation for the development of the  
420 HGCA. Indeed, several publicly accessible portals have already been established that contain  
421 intestinal sc datasets and enable access in a user-friendly way. Furthermore, a substantial  
422 number of studies have profiled intestinal tissues obtained from patients with gut diseases  
423 and/or human intestinal organoids (Figure 5). In this section, we provide a summary of the  
424 main existing datasets and portals, highlight examples of novel findings derived from sc  
425 research, and provide recommendations for future work.

426

### 427 **In utero development and healthy gut datasets**

428 The high throughput scalability intrinsic to most sc technologies has provided unprecedented  
429 advances to the field. One area that benefited from this aspect is GI organogenesis, since its  
430 progress has been hindered by the scarcity of samples. Indeed, datasets derived from scRNAseq  
431 spanning from 6 to 25 post-conception weeks have yielded several novel insights. For example,

432 during early fetal development, the crypt-villus axis begins to emerge in the small intestine  
433 concomitant to the appearance of *FOXL1*<sup>+</sup> mesenchyme cells co-expressing *PDGFRA* and  
434 *F3*<sup>72</sup>. In other sc studies a population of mesenchymal cells displaying a similar transcriptional  
435 signature were found to co-express *NRG1*, which has been demonstrated to support  
436 differentiation of *LGR5*<sup>+</sup> stem cells into mature intestinal epithelial cells<sup>73,74</sup>. Additionally,  
437 distinct clusters of sc transcriptomes identified during early stages of human intestinal  
438 development provided novel insight into the processes of regionalization during early intestinal  
439 development<sup>75,74</sup>. There are numerous examples of novel findings based on scRNAseq datasets  
440 derived from healthy adult gut samples. An assessment of epithelial cells from ileum, colon,  
441 and rectum revealed a high degree of functional diversity between the small and large intestine,  
442 as reflected by different nutrient absorption preferences<sup>76</sup>. The proposed existence of Paneth-  
443 like cells in the large intestine based on a cluster of colorectal epithelial cells co-expressing  
444 *LYZ*, *CA4*, *CA7*, and *SPIB* was later attributed to a new absorptive epithelial cell type  
445 expressing *BEST4*<sup>+</sup><sup>16,77,78</sup>. Recently, *BEST4*<sup>+</sup> epithelial cells have also been reported to vary  
446 in abundance and transcriptional signature across different regions of the gut<sup>79</sup>. Together, these  
447 findings illustrate the major benefit of combining and comparing datasets from different studies  
448 ultimately reaching reliable insight into healthy gut physiology and cellular function. Besides  
449 overcoming the challenges inherent to scarcity of material, datasets generated by scRNAseq  
450 have aided the characterization of the cellular diversity and transcriptional signatures of rare  
451 cell types. For instance, the characteristics of the human enteric nervous system has remained  
452 elusive until recently. By employing MIRACL-seq, a novel method designed to enrich samples  
453 for rare cell types, 1,445 enteric neurons were recovered from the human colon and found to  
454 cluster into 14 subsets based on their transcriptional signatures<sup>24</sup>. Interstitial Cells of Cajal  
455 (ICC) are rare entities critical for GI peristalsis through both the generation of slow wave  
456 pacemaker activity to smooth muscle cells, and mediation of neurotransmission from enteric  
457 neurons. By applying the same strategy, transcriptional signatures of 1,103 ICCs from the  
458 human colon have been generated<sup>24</sup>. Another sc database, generated by immunophenotyping  
459 and fluorescence-activated cell sorting, enriched ICCs from gastric resections and provided a  
460 comprehensive characterization of pathways and channels participating in their pacemaker  
461 activity<sup>80</sup>. Chemosensory cells such as tuft cells (TF) and enteroendocrine cells (EEC) are rare  
462 intestinal epithelial cells operating as an interface for signal transduction between the intestinal  
463 lumen and the body, relaying diet- and microbiota-derived signals through the release of  
464 numerous peptide hormones, neurotransmitters, and cytokines. In addition, TF cells are  
465 essential for mounting Th2 immune responses against parasites. Based on sc studies, TF cells

466 were demonstrated to interact with the innate and adaptive immune systems through previously  
467 unreported receptors<sup>79</sup>, including IgG receptors<sup>16</sup>. EECs sense intestinal content and release  
468 hormones to regulate gastrointestinal activity, systemic metabolism, and food intake. By using  
469 an organoid-based platform wherein EEC differentiation was induced by transient expression  
470 of *NEUROG3*, and hormones were tagged with gene reporters, the authors generated a  
471 comprehensive dataset of EEC subtypes derived from the small intestine and colon<sup>81</sup>.

## 472 **Intestinal diseases and gut organoids**

473 Generating a complete HGCA in health provides unique opportunities to study pathogenesis  
474 of related diseases. Furthermore, sc transcriptional profiles of primary tissue samples can be  
475 utilised as a valuable reference map allowing validation of existing and future organoid models.  
476 It is therefore of critical importance to take steps towards ensuring that datasets generated from  
477 related sample and patient cohorts can be integrated into the HGCA (Figure 5). Amongst the  
478 related gut diseases that have been investigated using sc profiling approaches are colorectal  
479 cancer (CRC)<sup>24,76,82-92</sup>, the inflammatory bowel diseases (IBD)<sup>93</sup>, Crohn's disease (CD)<sup>72,94-97</sup>  
480 and ulcerative colitis (UC)<sup>98-102</sup> as well as celiac disease<sup>103</sup>.

481 Examples of major findings in IBD include the identification of distinct immune cell signatures  
482 in both UC and CD<sup>104</sup>, a pathogenic cellular module associated with resistance to anti-TNF  
483 therapy<sup>95</sup>, inference of genetic risk genes to sc function<sup>98</sup> as well as the reactivation of fetal  
484 intestinal epithelial transcriptional profiles in childhood onset CD<sup>72</sup>. Similarly, application of  
485 sc molecular profiling methods to colonic tissue obtained from CRC patients has led to major  
486 advances in our understanding of disease pathogenesis. Specifically, recent sc insights into the  
487 stem and metaplastic origins of human pre-cancers<sup>84</sup> have led to reclassification of the  
488 consensus molecular subtypes of CRC by their intrinsic features<sup>105</sup>. Transition of benign lesions  
489 into malignancy is accompanied by tumor cell acquisition of stem characteristics<sup>84,92</sup>, and  
490 reorganization of the microenvironment into suppressive immune-stromal hubs that can  
491 potentially be therapeutically targeted<sup>83,106</sup>. Although a comprehensive summary of all relevant  
492 available sc studies in IBD and CRC is beyond the scope of this manuscript, it is important to  
493 highlight that the ability to integrate and compare studies performed on disease tissues at  
494 different stages of progression is of major benefit. Hence, ensuring compatibility with the  
495 HGCA remains a key priority for every study. Although integration of common gut related  
496 conditions for which extensive datasets are already available will be prioritised in the first  
497 phase of data integration, it is important to emphasise the major value of investigating and  
498 ultimately integrating rarer conditions. Indeed, combining datasets from less commonly  
499 profiled conditions represents another unique opportunity of the HGCA by increasing



500 computational power and allowing validation of key findings. Furthermore, applying a variety  
501 of methodologies (e.g. single multi-ome profiling) to the same condition will further improve  
502 the value of the HGCA and the chances of gaining novel insight in diseases pathogenesis.

503 Development of human intestinal organoid culture models has transformed many  
504 aspects of gut-related research providing researchers with unprecedented opportunities to study  
505 fundamental aspects of intestinal biology<sup>107</sup>. Performing transcriptional profiling of patient  
506 derived intestinal organoids on a sc level is of great value as reflected in the increasing numbers  
507 of studies reporting novel findings. Main benefits include the ability to validate cellular  
508 composition of organoids to further improve culture conditions as well as evaluating to what  
509 extent disease associated cellular alterations are retained in vitro. A recent example is a study  
510 published by He and colleagues, who applied scRNAseq to human small intestinal  
511 organoids<sup>108</sup>. Exposure of organoids to IL-22 resulted in increased expression of antimicrobial  
512 peptides suggesting the presence of Paneth cells<sup>108</sup>. Furthermore, sc profiling of primary  
513 human fetal gut and organoids derived from the same tissue sample revealed in vitro maturation  
514 of fetal organoids highlighting their value to investigate early stages of human intestinal  
515 epithelial cell development<sup>72</sup>. Work by Ishikawa and colleagues combined scRNAseq of  
516 human colonic epithelium cells with genetically modified human intestinal organoids leading  
517 to the identification of quiescent LGR5+ stem cells in the human colon<sup>109</sup>. Similarly,  
518 combining cancer organoid assays and scRNAseq of biopsies from patients with CRC  
519 identified a novel role of a  $\beta$ -hydroxybutyrate (BHB)-triggered pathway in regulating intestinal  
520 tumorigenesis<sup>110</sup>. As for studies based on profiling primary tissue, the provision of metadata  
521 including clinical donor details, sampling sites, biobanking as well as experimental procedures  
522 are equally important to organoid related studies as they will determine their future  
523 compatibility and ensure successful integration into the HGCA.

### 524 525 526 **Existing Data portals**

527 In an effort to consolidate the growing number of sc datasets that have been generated in the  
528 recent past, several web-based data repositories have emerged to provide access in a user-  
529 friendly format. While most portals share common features of storing data, curating datasets  
530 by different levels of metadata, and offering a variety of analysis and visualization tools,  
531 distinct aspects render each portal valuable. Here we summarize key features of the main  
532 existing data portals relevant to the HGCA.

533 **Single Cell Expression Atlas**<sup>70,111</sup>: The Single Cell Expression Atlas (SCEA) is the sc  
534 component of EMBL-EBI's Expression Atlas. This is an added-value resource that enables  
535 simple gene and meta-data queries, allowing users to answer questions such as “where is my  
536 gene of interest expressed”, “how does its expression level change in a disease”. SCEA collects  
537 expression data from all species, annotates their metadata with appropriate ontology terms to  
538 allow nested across studies and crucially, re-analyses all datasets using standardised analysis  
539 pipelines to enable comparison across studies. The SCEA works in close collaboration with the  
540 HGCA, releasing the datasets as they become publicly available and generating the first full  
541 gut 2D anatomogram, that will enable easy visual exploration of gene and marker gene  
542 expression across the cell types in the different anatomical sub-structures of the gut.

543 **HubMAP**<sup>112</sup>: The Human BioMolecular Atlas Program (HuBMAP) is a NIH Common Fund  
544 effort to integrate and map diverse biological data across the healthy human body. Three main  
545 features enable: 1) analysis of scRNAseq experiments with *Azimuth*, a web application that  
546 uses reference datasets to automate annotation and interpretation of data; 2) spatial sc data  
547 visualization with *Vitessce*; and 3) navigation of the healthy human cells with the CCF to  
548 interact with the virtual human body to focus on anatomical structures, cell types, and  
549 biomarkers.

550 **Tabula Sapiens**<sup>113</sup>: Like HuBMAP, the Tabula Sapiens focuses on healthy human subjects,  
551 and has created a first draft HCA of 24 organs from 15 different human subjects. Tabula  
552 Sapiens is funded by the Chan Zuckerberg Initiative (CZI) and is unique in that the sc data sets  
553 derived for each organ are from the same human subject, controlling for inter-individual  
554 factors. This also allows comparison of cell types that are shared between different organs. The  
555 web portal offers the ability to peruse the all the sc data combined (currently at 500,000 cells)  
556 or curated cell subsets (*i.e. endothelial, epithelial, immune, and stromal*) in an easy-to-navigate  
557 graphical user interface.

558 **UCSC Cell Browser**<sup>114</sup>: The University of California at Santa Cruz (UCSC) *Cell Browser* is a  
559 sub-project under supervision of the UCSC Genome Browser project that was developed and  
560 is maintained by a cross-departmental team in the UCSC Genomics Institute. The Cell Browser  
561 is an interactive viewer where users can interrogate sc data sets from a wide variety of species,  
562 organs, and tissues from a menu list. A unique feature is the broad scientific focus of the  
563 datasets generated from human, mouse, fly, and sponges, and curated as conventional Atlases  
564 or analysed in the context of development and evolution. The datasets are converted for  
565 compatibility by the UCSC Cell Browser Group and placed in an open-source portal. The  
566 analysis and viewer package can also be downloaded and installed locally.

567 **Broad Institute Single Cell Portal**<sup>115</sup>: The Single Cell Portal is hosted by the Broad Institute  
568 and was created with the idea to make sc data easy to share and access. The portal is a deposit  
569 site where investigators can create their own collection as a body of work or contribute to the  
570 growing list of curated datasets. The site is easy to use and has succinct overviews describing  
571 study design and experimental conditions. Like other portals, the Single Cell Portal has  
572 conventional visualization tools that are interactive in a simple graphical user interface that  
573 allows the user to filter several parameters in the data sets and sub-sampling data.

574 **Gut Cell Atlas, Wellcome Sanger Institute**<sup>116</sup>: This resource provides a detailed gut cell  
575 survey by combining single cell data generated from a range of human gut tissues. Datasets  
576 include studies that profiled various gut segments obtained from fetal, pediatric and adult  
577 donors as well as patients diagnosed with CD. In addition to direct access to raw datasets, the  
578 site provides an interactive viewer that allows basic analyses and interrogation of datasets such  
579 as exploring single cell expression profiles of individual cell lineages, gut regions, or  
580 comparisons between age groups.

581

582

### 583 **Future opportunities including incorporating datasets**

584 In many aspects, the current collection of high quality published sc studies achieve good  
585 coverage of the human intestinal tract and therefore provide a solid foundation for the  
586 generation of the HGCA. Areas that have been explored using scRNAseq technology include  
587 intestinal development<sup>72,75,117</sup>, profiling of cell types, states and tissue composition<sup>81</sup>, mapping  
588 regional differences across intestinal tissues and along the intestinal tube<sup>81,118</sup>, as well as  
589 comparisons between healthy individuals and IBD<sup>79,119</sup>, CRC<sup>89,90</sup>, or other related  
590 diseases<sup>92,120</sup>. However, there are several less well explored areas that require careful  
591 consideration in future studies. For example, to achieve a comprehensive and truly global  
592 HGCA, greater effort must be made to include samples obtained from underrepresented and  
593 ethnically diverse communities to reflect potential cellular differences across human  
594 populations. Currently only a few studies have sampled multiple intestinal regions within the  
595 same individual. However, such cross-tissue studies are of major value as they enable the  
596 identification of common/distinct cell types across intestinal regions as well as studying  
597 migratory patterns of specific immune cell populations<sup>16</sup>. Furthermore, greater emphasis  
598 should be placed on profiling rare cell types particularly in less frequently sampled gut regions.  
599 This can be achieved by applying methods designed to enrich for such cell types prior to  
600 scRNAseq and/or the application of spatial transcriptomics. Additionally, analysing single

601 nuclei rather than single cells is a useful strategy to profile cells that cannot be readily recovered  
602 using standard dissociation protocols<sup>24</sup>. However, advantages and disadvantages of these  
603 methods must be carefully considered in the context of each particular research study<sup>121,122</sup>.  
604 Finally, there has been a major bias with regards to anatomical sampling sites with a large  
605 proportion focussing on profiling the colon and distal small bowel. This is likely to be partly  
606 caused by ease of access during endoscopic procedures as well as the major relevance to related  
607 GI diseases such as IBD and CRC. Future studies should also include less frequently sampled  
608 gut segments, such as the mid small intestine, oesophagus, stomach as well as other organs  
609 involved in the digestive process including the liver and pancreas. Last but not least, profiling  
610 the intestine at different developmental stages by obtaining tissue from donors of all age groups  
611 will provide further insight in the pathophysiology of related diseases including those occurring  
612 in specific age groups.

613

## 614 **Conclusions**

615 Generating a complete map of the human intestine at the sc level will improve our  
616 understanding of gut health and disease. The inherent complexities and scale of this task  
617 requires a coordinated effort led by experts in this rapidly evolving field. This manuscript forms  
618 a key part of our strategy by providing a detailed roadmap to the scientific community. Broad  
619 distribution and constructive discussion of our proposal is essential to achieving our goal the  
620 generation of a HGCA.

621

622

## 623 **References:**

624

- 625 1 Rooks, M. G. & Garrett, W. S. Gut microbiota, metabolites and host immunity.  
626 *Nature Reviews Immunology* **16**, 341-352, doi:10.1038/nri.2016.42 (2016).
- 627 2 Makki, K., Deehan, E. C., Walter, J. & Backhed, F. The Impact of Dietary Fiber on Gut  
628 Microbiota in Host Health and Disease. *Cell Host Microbe* **23**, 705-715,  
629 doi:10.1016/j.chom.2018.05.012 (2018).
- 630 3 Khalili, H. *et al.* The role of diet in the aetiopathogenesis of inflammatory bowel  
631 disease. *Nat Rev Gastroenterol Hepatol* **15**, 525-535, doi:10.1038/s41575-018-0022-  
632 9 (2018).
- 633 4 Smillie, C. S. *et al.* Intra- and Inter-cellular Rewiring of the Human Colon during  
634 Ulcerative Colitis. *Cell* **178**, 714-730.e722, doi:10.1016/j.cell.2019.06.029 (2019).
- 635 5 Pilonis, N. D., Januszewicz, W. & di Pietro, M. Confocal laser endomicroscopy in  
636 gastro-intestinal endoscopy: technical aspects and clinical applications. *Transl*  
637 *Gastroenterol Hepatol* **7**, 7, doi:10.21037/tgh.2020.04.02 (2022).

- 638 6 Regev, A. *et al.* The Human Cell Atlas. *Elife* **6**, doi:ARTN e27041  
639 10.7554/eLife.27041 (2017).
- 640 7 Rozenblatt-Rosen, O., Stubbington, M. J. T., Regev, A. & Teichmann, S. A. The Human  
641 Cell Atlas: from vision to reality. *Nature* **550**, 451-453, doi:DOI 10.1038/550451a  
642 (2017).
- 643 8 Skinnider, M. A., Squair, J. W. & Courtine, G. Enabling reproducible re-analysis of  
644 single-cell data. *Genome Biol* **22**, doi:ARTN 215  
645 10.1186/s13059-021-02422-y (2021).
- 646 9 The Human Cell Atlas - Metadata.
- 647 10 Füllgrabe, A. *et al.* Guidelines for reporting single-cell RNA-seq experiments. *Nat*  
648 *Biotechnol* **38**, 1384-1386, doi:10.1038/s41587-020-00744-z (2020).
- 649 11 Lee, S. S. J. The Ethics of Consent in a Shifting Genomic Ecosystem. *Annual Review of*  
650 *Biomedical Data Science, Vol 4* **4**, 145-164, doi:10.1146/annurev-biodatasci-030221-  
651 125715 (2021).
- 652 12 EUR-Lex, T. E. P. a. t. C. o. t. E. U. *General Data Protection Regulation (EU) 2016/679*  
653 *(GDPR)*, <<https://eur-lex.europa.eu/eli/reg/2016/679/oj>> (2016).
- 654 13 Bledsoe, M. J. & Grizzle, W. E. Use of human specimens in research: the evolving  
655 United States regulatory, policy, and scientific landscape. *Diagn Histopathol (Oxf)* **19**,  
656 322-330, doi:10.1016/j.mpdhp.2013.06.015 (2013).
- 657 14 Shore, N. *et al.* Understanding community-based processes for research ethics  
658 review: a national study. *Am J Public Health* **101 Suppl 1**, S359-364,  
659 doi:10.2105/ajph.2010.194340 (2011).
- 660 15 Mikesell, L., Bromley, E. & Khodyakov, D. Ethical community-engaged research: a  
661 literature review. *Am J Public Health* **103**, e7-e14, doi:10.2105/ajph.2013.301605  
662 (2013).
- 663 16 Elmentaite, R. *et al.* Cells of the human intestinal tract mapped across space and  
664 time. *Nature* **597**, 250-+, doi:10.1038/s41586-021-03852-1 (2021).
- 665 17 Camp, J. G., Platt, R. & Treutlein, B. Mapping human cell phenotypes to genotypes  
666 with single-cell genomics. *Science* **365**, 1401-+, doi:10.1126/science.aax6648 (2019).
- 667 18 Macaulay, I. C. *et al.* G&T-seq: parallel sequencing of single-cell genomes and  
668 transcriptomes. *Nat Methods* **12**, 519-522, doi:10.1038/nmeth.3370 (2015).
- 669 19 Lareau, C. A. *et al.* Droplet-based combinatorial indexing for massive-scale single-cell  
670 chromatin accessibility. *Nat Biotechnol* **37**, 916-924, doi:10.1038/s41587-019-0147-6  
671 (2019).
- 672 20 Smallwood, S. A. *et al.* Single-cell genome-wide bisulfite sequencing for assessing  
673 epigenetic heterogeneity. *Nat Methods* **11**, 817-820, doi:10.1038/nmeth.3035  
674 (2014).
- 675 21 Hao, Y. H. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-+,  
676 doi:10.1016/j.cell.2021.04.048 (2021).
- 677 22 Pott, S. Simultaneous measurement of chromatin accessibility, DNA methylation, and  
678 nucleosome phasing in single cells. *Elife* **6**, doi:10.7554/eLife.23203 (2017).
- 679 23 Clark, S. J. *et al.* scNMT-seq enables joint profiling of chromatin accessibility DNA  
680 methylation and transcription in single cells. *Nat Commun* **9**, doi:ARTN 781  
681 10.1038/s41467-018-03149-4 (2018).
- 682 24 Drokhyansky, E. *et al.* The Human and Mouse Enteric Nervous System at Single-Cell  
683 Resolution. *Cell* **182**, 1606-+, doi:10.1016/j.cell.2020.08.003 (2020).

- 684 25 Slyper, M. *et al.* A single-cell and single-nucleus RNA-Seq toolbox for fresh and frozen  
685 human tumors. *Nat Med* **26**, 792-802, doi:10.1038/s41591-020-0844-1 (2020).
- 686 26 Wu, S. Z. *et al.* A single-cell and spatially resolved atlas of human breast cancers. *Nat*  
687 *Genet* **53**, 1334-1347, doi:10.1038/s41588-021-00911-1 (2021).
- 688 27 Xia, C., Babcock, H. P., Moffitt, J. R. & Zhuang, X. Multiplexed detection of RNA using  
689 MERFISH and branched DNA amplification. *Scientific Reports* **9**, 7721,  
690 doi:10.1038/s41598-019-43943-8 (2019).
- 691 28 Black, S. *et al.* CODEX multiplexed tissue imaging with DNA-conjugated antibodies.  
692 *Nat Protoc* **16**, 3802-3835, doi:10.1038/s41596-021-00556-8 (2021).
- 693 29 Hwang, B., Lee, J. H. & Bang, D. Single-cell RNA sequencing technologies and  
694 bioinformaticspipelines. *Experimental & Molecular Medicine* **50**, 1-14,  
695 doi:10.1038/s12276-018-0071-8 (2018).
- 696 30 Zappia, L., Phipson, B. & Oshlack, A. Exploring the single-cell RNA-seq analysis  
697 landscape with the scRNA-tools database. *PLOS Computational Biology* **14**,  
698 e1006245, doi:10.1371/journal.pcbi.1006245 (2018).
- 699 31 Moreno, P. *et al.* User-friendly, scalable tools and workflows for single-cell RNA-seq  
700 analysis. *Nat. Methods* **18**, 327-328, doi:10.1038/s41592-021-01102-w (2021).
- 701 32 Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell  
702 RNA sequencing data. *Genome Biology* **21**, 12, doi:10.1186/s13059-019-1850-9  
703 (2020).
- 704 33 Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-  
705 3587.e3529, doi:10.1016/j.cell.2021.04.048 (2021).
- 706 34 Liu, J. *et al.* Jointly defining cell types from multiple single-cell datasets using LIGER.  
707 *Nature protocols* **15**, 3632-3662, doi:10.1038/s41596-020-0391-8 (2020).
- 708 35 Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with  
709 Harmony. *Nat. Methods* **16**, 1289-1296, doi:10.1038/s41592-019-0619-0 (2019).
- 710 36 Cao, Z. J. & Gao, G. Multi-omics single-cell data integration and regulatory inference  
711 with graph-linked embedding. *Nat Biotechnol*, doi:10.1038/s41587-022-01284-4  
712 (2022).
- 713 37 Argelaguet, R. *et al.* MOFA+: a statistical framework for comprehensive integration  
714 of multi-modal single-cell data. *Genome Biology* **21**, 111, doi:10.1186/s13059-020-  
715 02015-1 (2020).
- 716 38 Gong, B., Zhou, Y. & Purdom, E. Cobolt: integrative analysis of multimodal single-cell  
717 sequencing data. *Genome Biology* **22**, 351, doi:10.1186/s13059-021-02556-z (2021).
- 718 39 Jones, R. C. *et al.* The Tabula Sapiens: A multiple-organ, single-cell transcriptomic  
719 atlas of humans. *Science* **376**, 711-+, doi:10.1126/science.abl4896 (2022).
- 720 40 Xu, Y., Baumgart, S. J., Stegmann, C. M. & Hayat, S. MACA: Marker-based automatic  
721 cell-type annotation for single cell expression data. *Bioinformatics*,  
722 doi:10.1093/bioinformatics/btab840 (2021).
- 723 41 Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a  
724 transitional profibrotic macrophage. *Nature Immunology* **20**, 163-172,  
725 doi:10.1038/s41590-018-0276-y (2019).
- 726 42 Ianevski, A., Giri, A. K. & Aittokallio, T. Fully-automated and ultra-fast cell-type  
727 identification using specific marker combinations from single-cell transcriptomic  
728 data. *Nature Communications* **13**, 1246, doi:10.1038/s41467-022-28803-w (2022).
- 729 43 Delaney, C. *et al.* Combinatorial prediction of marker panels from single-cell  
730 transcriptomic data. *Mol Syst Biol* **15**, e9005, doi:10.15252/msb.20199005 (2019).

- 731 44 Dai, M., Pei, X. & Wang, X.-J. Accurate and fast cell marker gene identification with  
732 COSG. *Briefings in Bioinformatics* **23**, doi:10.1093/bib/bbab579 (2022).
- 733 45 Andersson, A. *et al.* Single-cell and spatial transcriptomics enables probabilistic  
734 inference of cell type topography. *Communications Biology* **3**, 565,  
735 doi:10.1038/s42003-020-01247-y (2020).
- 736 46 David, F. P. A., Litovchenko, M., Deplancke, B. & Gardeux, V. ASAP 2020 update: an  
737 open, scalable and interactive web-based portal for (single-cell) omics analyses.  
738 *Nucleic Acids Research* **48**, W403-W414, doi:10.1093/nar/gkaa412 (2020).
- 739 47 Moriel, N. *et al.* NovoSpaRc: flexible spatial reconstruction of single-cell gene  
740 expression with optimal transport. *Nature protocols* **16**, 4177-4200,  
741 doi:10.1038/s41596-021-00573-7 (2021).
- 742 48 Kleshchevnikov, V. *et al.* Cell2location maps fine-grained cell types in spatial  
743 transcriptomics. *Nat. Biotechnol.*, doi:10.1038/s41587-021-01139-4 (2022).
- 744 49 Dries, R. *et al.* Giotto: a toolbox for integrative analysis and visualization of spatial  
745 expression data. *Genome Biol* **22**, 78, doi:10.1186/s13059-021-02286-2 (2021).
- 746 50 Svensson, V., Teichmann, S. A. & Stegle, O. SpatialDE: identification of spatially  
747 variable genes. *Nat Methods* **15**, 343-346, doi:10.1038/nmeth.4636 (2018).
- 748 51 Goltsev, Y. *et al.* Deep Profiling of Mouse Splenic Architecture with CODEX  
749 Multiplexed Imaging. *Cell* **174**, 968-981 e915, doi:10.1016/j.cell.2018.07.010 (2018).
- 750 52 Armingol, E., Officer, A., Harismendy, O. & Lewis, N. E. Deciphering cell-cell  
751 interactions and communication from gene expression. *Nature Reviews Genetics* **22**,  
752 71-88, doi:10.1038/s41576-020-00292-x (2021).
- 753 53 Efremova, M., Vento-Tormo, M., Teichmann, S. A. & Vento-Tormo, R. CellPhoneDB:  
754 inferring cell-cell communication from combined expression of multi-subunit ligand-  
755 receptor complexes. *Nat Protoc* **15**, 1484-1506, doi:10.1038/s41596-020-0292-x  
756 (2020).
- 757 54 Kang, R., Park, B., Eady, M., Ouyang, Q. & Chen, K. J. Single-cell classification of  
758 foodborne pathogens using hyperspectral microscope imaging coupled with deep  
759 learning frameworks. *Sensors and Actuators B-Chemical* **309**, doi:ARTN 127789  
760 10.1016/j.snb.2020.127789 (2020).
- 761 55 Chattopadhyay, P. K., Roederer, M. & Bolton, D. L. A deadly dance: the choreography  
762 of host-pathogen interactions, as revealed by single-cell technologies. *Nat Commun*  
763 **9**, doi:ARTN 4638  
764 10.1038/s41467-018-06214-0 (2018).
- 765 56 Galeano Niño, J. L. *et al.* Effect of the intratumoral microbiota on spatial and cellular  
766 heterogeneity in cancer. *Nature* **611**, 810-817, doi:10.1038/s41586-022-05435-0  
767 (2022).
- 768 57 Aibar, S. *et al.* SCENIC: single-cell regulatory network inference and clustering. *Nat.*  
769 *Methods* **14**, 1083-1086, doi:10.1038/nmeth.4463 (2017).
- 770 58 Moerman, T. *et al.* GRNBoost2 and Arboreto: efficient and scalable inference of gene  
771 regulatory networks. *Bioinformatics* **35**, 2159-2161,  
772 doi:10.1093/bioinformatics/bty916 (2019).
- 773 59 Chan, T. E., Stumpf, M. P. H. & Babbie, A. C. Gene Regulatory Network Inference from  
774 Single-Cell Data Using Multivariate Information Measures. *Cell Syst* **5**, 251-267.e253,  
775 doi:10.1016/j.cels.2017.08.014 (2017).

- 776 60 Pratapa, A., Jalihal, A. P., Law, J. N., Bharadwaj, A. & Murali, T. M. Benchmarking  
777 algorithms for gene regulatory network inference from single-cell transcriptomic  
778 data. *Nat. Methods* **17**, 147-154, doi:10.1038/s41592-019-0690-6 (2020).
- 779 61 Cakir, B. *et al.* Comparison of visualization tools for single-cell RNAseq data. *NAR*  
780 *Genomics and Bioinformatics* **2**, doi:10.1093/nargab/lqaa052 (2020).
- 781 62 Megill, C. *et al.* chanzuckerberg/cellxgene: Release 0.15.0.  
782 doi:10.5281/ZENODO.3710410 (2020).
- 783 63 Elmentaite, R. *et al.* Single-Cell Sequencing of Developing Human Gut Reveals  
784 Transcriptional Links to Childhood Crohn's Disease. *Developmental Cell* **55**, 771-  
785 783.e775, doi:<https://doi.org/10.1016/j.devcel.2020.11.010> (2020).
- 786 64 Moreno, P. *et al.* Expression Atlas update: gene and protein expression in multiple  
787 species. *Nucleic Acids Research* **50**, D129-D140, doi:10.1093/nar/gkab1030 (2021).
- 788 65 Rood, J. E. *et al.* Toward a Common Coordinate Framework for the Human Body. *Cell*  
789 **179**, 1455-1467, doi:10.1016/j.cell.2019.11.019 (2019).
- 790 66 Schiller, H. B. *et al.* The Human Lung Cell Atlas: A High-Resolution Reference Map of  
791 the Human Lung in Health and Disease. *American Journal of Respiratory Cell and*  
792 *Molecular Biology* **61**, 31-41, doi:10.1165/rcmb.2018-0416TR (2019).
- 793 67 Eze, U. C., Bhaduri, A., Haeussler, M., Nowakowski, T. J. & Kriegstein, A. R. Single-cell  
794 atlas of early human brain development highlights heterogeneity of human  
795 neuroepithelial cells and early radial glia. *Nature Neuroscience* **24**, 584-594,  
796 doi:10.1038/s41593-020-00794-1 (2021).
- 797 68 Weber, G. M., Ju, Y. N. & Borner, K. Considerations for Using the Vasculature as a  
798 Coordinate System to Map All the Cells in the Human Body. *Frontiers in*  
799 *Cardiovascular Medicine* **7**, doi:ARTN 29  
800 10.3389/fcvm.2020.00029 (2020).
- 801 69 Burger, A. *et al.* Towards a Clinically-based Common Coordinate Framework for the  
802 Human Gut Cell Atlas - The Gut Models. *bioRxiv*, 2022.2012.2008.519665,  
803 doi:10.1101/2022.12.08.519665 (2022).
- 804 70 Moreno, P. *et al.* Expression Atlas update: gene and protein expression in multiple  
805 species. *Nucleic Acids Res* **50**, D129-D140, doi:10.1093/nar/gkab1030 (2022).
- 806 71 *Gut Cell Atlas, an Initiative Supported by the Helmsley Charitable Trust*  
807 [<https://helmsleytrust.org/our-focus-areas/crohns-disease/crohns-disease-  
808 therapeutics/gut-cell-atlas/>](https://helmsleytrust.org/our-focus-areas/crohns-disease/crohns-disease-therapeutics/gut-cell-atlas/) (
- 809 72 Elmentaite, R. *et al.* Single-Cell Sequencing of Developing Human Gut Reveals  
810 Transcriptional Links to Childhood Crohn's Disease. *Developmental Cell* **55**, 771-+,  
811 doi:10.1016/j.devcel.2020.11.010 (2020).
- 812 73 Holloway, E. M. *et al.* Mapping Development of the Human Intestinal Niche at Single-  
813 Cell Resolution. *Cell Stem Cell* **28**, 568-+, doi:10.1016/j.stem.2020.11.008 (2021).
- 814 74 Yu, Q. H. *et al.* Charting human development using a multi-endodermal organ atlas  
815 and organoid models. *Cell* **184**, 3281-+, doi:10.1016/j.cell.2021.04.028 (2021).
- 816 75 Gao, S. *et al.* Tracing the temporal-spatial transcriptome landscapes of the human  
817 fetal digestive tract using single-cell RNA-sequencing (vol 20, pg 721, 2018). *Nature*  
818 *Cell Biology* **20**, 1227-1227, doi:10.1038/s41556-018-0165-5 (2018).
- 819 76 Wang, Y. L. *et al.* Single-cell transcriptome analysis reveals differential nutrient  
820 absorption functions in human intestine. *Journal of Experimental Medicine* **217**,  
821 doi:10.1084/jem.20191130 (2020).



- 822 77 Parikh, K. *et al.* Colonic epithelial cell diversity in health and inflammatory bowel  
823 disease. *Nature* **567**, 49-+, doi:10.1038/s41586-019-0992-y (2019).
- 824 78 Busslinger, G. A. *et al.* Human gastrointestinal epithelia of the esophagus, stomach,  
825 and duodenum resolved at single-cell resolution. *Cell Reports* **34**, doi:ARTN 108819  
826 10.1016/j.celrep.2021.108819 (2021).
- 827 79 Burclaff, J. *et al.* A Proximal-to-Distal Survey of Healthy Adult Human small Intestine  
828 and Colon Epithelium by Single-Cell Transcriptomics. *Cellular and Molecular*  
829 *Gastroenterology and Hepatology* **13**, 1554-1589, doi:10.1016/j.jcmgh.2022.02.007  
830 (2022).
- 831 80 Foong, D. *et al.* Single-cell RNA sequencing predicts motility networks in purified  
832 human gastric interstitial cells of Cajal. *Neurogastroenterology and Motility* **34**,  
833 doi:ARTN e14303  
834 10.1111/nmo.14303 (2022).
- 835 81 Beumer, J. *et al.* High-Resolution mRNA and Secretome Atlas of Human  
836 Enteroendocrine Cells. *Cell* **181**, 1291-+, doi:10.1016/j.cell.2020.04.036 (2020).
- 837 82 Lee, H. O. *et al.* Lineage-dependent gene expression programs influence the immune  
838 landscape of colorectal cancer. *Nat Genet* **52**, 594-603, doi:10.1038/s41588-020-  
839 0636-z (2020).
- 840 83 Pelka, K. *et al.* Spatially organized multicellular immune hubs in human colorectal  
841 cancer. *Cell* **184**, 4734-4752.e4720, doi:10.1016/j.cell.2021.08.003 (2021).
- 842 84 Chen, B. *et al.* Differential pre-malignant programs and microenvironment chart  
843 distinct paths to malignancy in human colorectal polyps. *Cell* **184**, 6262-6280.e6226,  
844 doi:10.1016/j.cell.2021.11.031 (2021).
- 845 85 Li, H. *et al.* Reference component analysis of single-cell transcriptomes elucidates  
846 cellular heterogeneity in human colorectal tumors. *Nat Genet* **49**, 708-718,  
847 doi:10.1038/ng.3818 (2017).
- 848 86 Qian, J. *et al.* A pan-cancer blueprint of the heterogeneous tumor microenvironment  
849 revealed by single-cell profiling. *Cell Res* **30**, 745-762, doi:10.1038/s41422-020-0355-  
850 0 (2020).
- 851 87 Zhang, L. *et al.* Single-Cell Analyses Inform Mechanisms of Myeloid-Targeted  
852 Therapies in Colon Cancer. *Cell* **181**, 442-459.e429, doi:10.1016/j.cell.2020.03.048  
853 (2020).
- 854 88 Domanska, D. *et al.* Single-cell transcriptomic analysis of human colonic  
855 macrophages reveals niche-specific subsets. *J Exp Med* **219**,  
856 doi:10.1084/jem.20211846 (2022).
- 857 89 Uhlitz, F. *et al.* Mitogen-activated protein kinase activity drives cell trajectories in  
858 colorectal cancer. *Embo Molecular Medicine* **13**, doi:ARTN e14123  
859 10.15252/emmm.202114123 (2021).
- 860 90 Zhang, L. *et al.* Lineage tracking reveals dynamic relationships of T cells in colorectal  
861 cancer. *Nature* **564**, 268-+, doi:10.1038/s41586-018-0694-x (2018).
- 862 91 Che, L. H. *et al.* A single-cell atlas of liver metastases of colorectal cancer reveals  
863 reprogramming of the tumor microenvironment in response to preoperative  
864 chemotherapy. *Cell Discov* **7**, 80, doi:10.1038/s41421-021-00312-y (2021).
- 865 92 Becker, W. R. *et al.* Single-cell analyses define a continuum of cell state and  
866 composition changes in the malignant transformation of polyps to colorectal cancer.  
867 *Nature Genetics* **54**, 985-+, doi:10.1038/s41588-022-01088-x (2022).

- 868 93 Bolton, C. *et al.* An Integrated Taxonomy for Monogenic Inflammatory Bowel  
869 Disease. *Gastroenterology* **162**, 859-876, doi:10.1053/j.gastro.2021.11.014 (2022).
- 870 94 Kanke, M. *et al.* Single-Cell Analysis Reveals Unexpected Cellular Changes and  
871 Transposon Expression Signatures in the Colonic Epithelium of Treatment-Naive  
872 Adult Crohn's Disease Patients. *Cell Mol Gastroenterol Hepatol* **13**, 1717-1740,  
873 doi:10.1016/j.jcmgh.2022.02.005 (2022).
- 874 95 Martin, J. C. *et al.* Single-Cell Analysis of Crohn's Disease Lesions Identifies a  
875 Pathogenic Cellular Module Associated with Resistance to Anti-TNF Therapy. *Cell*  
876 **178**, 1493-1508 e1420, doi:10.1016/j.cell.2019.08.008 (2019).
- 877 96 Jaeger, N. *et al.* Single-cell analyses of Crohn's disease tissues reveal intestinal  
878 intraepithelial T cells heterogeneity and altered subset distributions. *Nat Commun*  
879 **12**, 1921, doi:10.1038/s41467-021-22164-6 (2021).
- 880 97 Uniken Venema, W. T. *et al.* Single-Cell RNA Sequencing of Blood and Ileal T Cells  
881 From Patients With Crohn's Disease Reveals Tissue-Specific Characteristics and Drug  
882 Targets. *Gastroenterology* **156**, 812-815 e822, doi:10.1053/j.gastro.2018.10.046  
883 (2019).
- 884 98 Smillie, C. S. *et al.* Intra- and Inter-cellular Rewiring of the Human Colon during  
885 Ulcerative Colitis. *Cell* **178**, 714-+, doi:10.1016/j.cell.2019.06.029 (2019).
- 886 99 Uzzan, M. *et al.* Ulcerative colitis is characterized by a plasmablast-skewed humoral  
887 response associated with disease activity. *Nat Med* **28**, 766-779,  
888 doi:10.1038/s41591-022-01680-y (2022).
- 889 100 Chen, E. *et al.* Inflamed Ulcerative Colitis Regions Associated With MRGPRX2-  
890 Mediated Mast Cell Degranulation and Cell Activation Modules, Defining a New  
891 Therapeutic Target. *Gastroenterology* **160**, 1709-1724,  
892 doi:10.1053/j.gastro.2020.12.076 (2021).
- 893 101 Devlin, J. C. *et al.* Single-Cell Transcriptional Survey of Ileal-Anal Pouch Immune Cells  
894 From Ulcerative Colitis Patients. *Gastroenterology* **160**, 1679-1693,  
895 doi:10.1053/j.gastro.2020.12.030 (2021).
- 896 102 Corridoni, D. *et al.* Single-cell atlas of colonic CD8(+) T cells in ulcerative colitis. *Nat*  
897 *Med* **26**, 1480-1490, doi:10.1038/s41591-020-1003-4 (2020).
- 898 103 Atlasy, N. *et al.* Single cell transcriptomic analysis of the immune cell compartment in  
899 the human small intestine and in Celiac disease. *Nat Commun* **13**, 4920,  
900 doi:10.1038/s41467-022-32691-5 (2022).
- 901 104 Mitsialis, V. *et al.* Single-Cell Analyses of Colon and Blood Reveal Distinct Immune  
902 Cell Signatures of Ulcerative Colitis and Crohn's Disease. *Gastroenterology* **159**, 591-  
903 608 e510, doi:10.1053/j.gastro.2020.04.074 (2020).
- 904 105 Joanito, I. *et al.* Single-cell and bulk transcriptome sequencing identifies two  
905 epithelial tumor cell states and refines the consensus molecular classification of  
906 colorectal cancer. *Nat Genet* **54**, 963-975, doi:10.1038/s41588-022-01100-4 (2022).
- 907 106 Qi, J. J. *et al.* Single-cell and spatial analysis reveal interaction of FAP(+) fibroblasts  
908 and SPP1(+) macrophages in colorectal cancer. *Nat Commun* **13**, doi:ARTN 1742  
909 10.1038/s41467-022-29366-6 (2022).
- 910 107 Gunther, C., Winner, B., Neurath, M. F. & Stappenbeck, T. S. Organoids in  
911 gastrointestinal diseases: from experimental models to clinical translation. *Gut* **71**,  
912 1892-1908, doi:10.1136/gutjnl-2021-326560 (2022).

- 913 108 He, G. W. *et al.* Optimized human intestinal organoid model reveals interleukin-22-  
914 dependency of paneth cell formation. *Cell Stem Cell* **29**, 1333-1345 e1336,  
915 doi:10.1016/j.stem.2022.08.002 (2022).
- 916 109 Ishikawa, K. *et al.* Identification of Quiescent LGR5(+) Stem Cells in the Human Colon.  
917 *Gastroenterology*, doi:10.1053/j.gastro.2022.07.081 (2022).
- 918 110 Dmitrieva-Posocco, O. *et al.* beta-Hydroxybutyrate suppresses colorectal cancer.  
919 *Nature* **605**, 160-165, doi:10.1038/s41586-022-04649-6 (2022).
- 920 111 *Single Cell Expression Atlas*, <<https://www.ebi.ac.uk/gxa/sc/home>> (  
921 112 *Human BioMolecular Atlas Program (HuBMAP) Data Portal*,  
922 <<https://portal.hubmapconsortium.org>> (  
923 113 *Tabula Sapiens Data Portal*, <<https://tabula-sapiens-portal.ds.czbiohub.org>> (  
924 114 *UCSC Cell Browser*, <<https://cells.ucsc.edu>> (  
925 115 *Broad Institute Single Cell Portal*, <[https://singlecell.broadinstitute.org/single\\_cell](https://singlecell.broadinstitute.org/single_cell)> (  
926 116 *Gut Cell Survey*, <<https://www.gutcellatlas.org>> (  
927 117 Fawcner-Corbett, D. *et al.* Spatiotemporal analysis of human intestinal development  
928 at single-cell resolution. *Cell* **184**, 810-+, doi:10.1016/j.cell.2020.12.016 (2021).
- 929 118 James, K. R. *et al.* Distinct microbial and immune niches of the human colon. *Nat*  
930 *Immunol* **21**, 343-+, doi:10.1038/s41590-020-0602-z (2020).
- 931 119 Kinchen, J. *et al.* Structural Remodeling of the Human Colonic Mesenchyme in  
932 Inflammatory Bowel Disease. *Cell* **175**, 372-+, doi:10.1016/j.cell.2018.08.067 (2018).
- 933 120 Nowicki-Osuch, K. *et al.* Molecular phenotyping reveals the identity of Barrett's  
934 esophagus and its malignant transition. *Science* **373**, 760-+,  
935 doi:10.1126/science.abd1449 (2021).
- 936 121 Nguyen, Q. H., Pervolarakis, N., Nee, K. & Kessenbrock, K. Experimental  
937 Considerations for Single-Cell RNA Sequencing Approaches. *Frontiers in Cell and*  
938 *Developmental Biology* **6**, doi:ARTN 108  
939 10.3389/fcell.2018.00108 (2018).
- 940 122 Haque, A., Engel, J., Teichmann, S. A. & Lonnberg, T. A practical guide to single-cell  
941 RNA-sequencing for biomedical research and clinical applications. *Genome Medicine*  
942 **9**, doi:ARTN 75  
943 10.1186/s13073-017-0467-4 (2017).

944  
945  
946

## 947 **Acknowledgements**

948 This publication is part of the Human Cell Atlas  
949 (HCA)- [www.humancellatlas.org/publications](http://www.humancellatlas.org/publications). The HCA initiative receives funding from The  
950 Wellcome Trust, the UK Research and Innovation Medical Research Council, EU Horizon  
951 2020, INSERM (HuDeCA), and the Knut and Alice Wallenberg and Erling-Persson  
952 foundations. We thank the HCA Executive Office for their support. The Gut Cell Atlas is  
953 organised by The Leona M. and Harry B. Helmsley Charitable Trust and provides funding for  
954 members in form of project grants. MZ was supported by an MRC New Investigator research  
955 grant (MR/T001917/1) and a project grant from the Great Ormond Street Hospital Children's

956 Charity, Sparks (V4519); KSL was supported by NIDDK R01DK103831, and The Helmsley  
957 Trust - G-1903-03793. STM received funding from (STM) National Institutes of Health  
958 USA R01DK115806 & P30DK034987. TS was supported by the Japanese Science and  
959 Technology (JST) forest, and the Japanese Society for the Promotion of Science (JSPS)  
960 (21K18272). LAC and KTW. were supported by The Helmsley Charitable Trust - G-1903-  
961 03793. KTW was also supported by NIDDK R01DK128200. LAC was supported by a  
962 Veterans Affairs Merit Award 1I01BX004366. MK was supported by the National Research  
963 Foundation, South Africa grant no: 129356.

964

965 **Author contributions** MZ, KTW, AS, MPS and RJX are coordinators of the HCA  
966 Gut Biological Network. MZ conceived the idea, co-ordinated the writing process, wrote parts  
967 of the paper and edited all sections. SP was section lead for methodologies. ZL was section  
968 lead for computational tools and challenges. MK was section lead of the integration of HGCA  
969 with gut diseases and organoids. KJ was section lead for the summary of existing datasets and  
970 portals. AGB was section lead for the Common Coordinate Framework (CCF). AH helped to  
971 design and create the figures. All other authors wrote parts of the paper, contributed to critical  
972 discussions and provided feedback on the entire manuscript.

973

974 **Competing interests:** SAT: In the past three years, SAT. has consulted or been a member of  
975 scientific advisory boards at Roche, Genentech, Biogen, GlaxoSmithKline, Qiagen and  
976 ForeSite Labs and is an equity holder of Transition Bio. GM has received grant funding from  
977 Boehringer Ingelheim. A.R. is a co-founder and equity holder of Celsius Therapeutics, an  
978 equity holder in Immunitas, and was an SAB member of ThermoFisher Scientific, Syros  
979 Pharmaceuticals, Neogene Therapeutics and Asimov until 31 July 2020. Since 1 August 2020,  
980 A.R. has been an employee of Genentech and has equity in Roche. A.R. is an inventor on  
981 patents and patent applications filed at the Broad related to single cell genomics.

982

983

984

985

986

987

988 **Key points:**

- 989
- The number of studies applying single cell sequencing methods to human intestinal tissue has been rapidly increasing providing a unique opportunity to generate a complete map of the human intestine.
  - The generation of a Human Gut Cell Atlas (HGCA) requires the coordinated efforts of groups across the globe and the integration of various datasets followed by their computational analyses.
  - This article provides a roadmap for the generation of the HGCA based on the expertise and recommendations of the Human Cell Atlas (HCA) Gut Biological Network.
  - The Human Gut Cell Atlas (HGCA) will provide a unique and highly valuable reference map enhancing research in intestinal health and disease.
- 990
- 991
- 992
- 993
- 994
- 995
- 996
- 997
- 998
- 999

1000

1001

## 1002 **Figure legends**

1003

1004 **Figure 1: Profiling the human digestive tract:** A complete map of the human intestinal tract requires the inclusion of the entire intestinal tube and associated organs (left panel). Profiling must also consider the impact of developmental stage, the gut microbiome, and the potential impact of environmental factors such as diet, toxins and medication. In addition to transcriptional profiling, capturing the underlying genome sequence and epigenetic programme will yield critical information (right panel).

1005

1006

1007

1008

1009

1010

1011 **Figure 2: Summary of main tissue types and sampling strategies available.** Main advantages and limitations are illustrated.

1012

1013

1014 **Figure 3: Data integration, processing and analysing strategies:** 1) Generation of the HGCA requires successful integration of various datasets, datatypes, and associated patient/donor metadata. 2 and 3) Successful integration will allow a range of queries to be performed and outputs generated.

1015

1016

1017

1018

1019 **Figure 4: Current applications for the Human Gut Cell Atlas Common Coordinate Framework: A:** HuBMAP Registration User Interface (CCF-RUI) showing the gut visibility controls on the LH panel, interactive block registration interface in the centre panel and the

1020

1021

1022 block specification controls in the RH panel. **B:** Gut Atlas CCF browser interface showing the  
 1023 1D conceptual model for the full large and small intestines with zoom panel and ontology  
 1024 listings at the top with 2D and 3D interactive displays for all of the available models. The  
 1025 position of the current ROI is displayed in each CCF view and is fully interoperable in the  
 1026 sense that position selection in any view will be updated immediately in all the other views. **C:**  
 1027 EBI SCEA anatomogram for the large and small intestines allowing selection of any structure  
 1028 identified in the ASCT tables. **D:** Expanded view of the anal canal to show relevant cell-types  
 1029 and tissue organisation of that region.

1030

1031 **Figure 5: HGCA in health and disease: A)** Integration of datasets generated from intestinal  
 1032 tissues obtained from healthy individuals and patients with related diseases including IBD and  
 1033 colorectal cancer will enhance the future value of the HGCA. Single cell profiling of intestinal  
 1034 organoids followed by their integration in the HGCA will provide unique opportunities for  
 1035 translational research including regenerative medicine, drug testing and development. **B)**  
 1036 Summary of existing studies that have used sc profiling methods to human tissue samples  
 1037 obtained from healthy gut, colorectal cancer, Inflammatory Bowel Diseases and patient derived  
 1038 intestinal organoids.

1039

1040

1041 **Table 1: Template metadata for gut-related single-cell studies**

|            |  |
|------------|--|
| Sample ID  |  |
| Patient ID |  |

1042

| Study design               |                              |
|----------------------------|------------------------------|
| Summary                    | <i>(brief study summary)</i> |
| Patient numbers            |                              |
| Sample numbers             |                              |
| Patient inclusion criteria |                              |
| Sampling strategy          |                              |

| <b>Donor information</b>                             |                                |
|--|--------------------------------|
| Age  |                                |
| Gender   |                                |
| Ethnicity  |                                |
| Medical condition                                    | Gastro-intestinal (GI) related |
|  | Non-GI related                 |
| Medication   | Gastro-intestinal (GI) related |
|  | Non-GI related                 |
| Dietary habits (if applicable)                       |                                |
| <b>Sample Information</b>                            |                                |
| <i>Sample type</i>                                   |                                |
| Mucosal biopsy                                       |                                |
| Surgical resection                                   |                                |
| Deceased donor                                       |                                |
| <i>Macroscopical appearance of sample</i>            |                                |
| Normal   |                                |
| Pathological<br>(Provide details)                    |                                |
| <i>Microscopic appearance of sample</i>              |                                |
| Normal   |                                |
| Pathological<br>(Provide details)                    |                                |
| Representative<br>H&E stain available                | Yes/No                         |
| <i>(Provide details of H&amp;E and other stains)</i> |                                |
| <b>Anatomical Sample location</b>                    |                                |
| Oral cavity  |                                |

|                              |            |                  |  |
|------------------------------|------------|------------------|--|
| Oesophagus                   | Proximal   |                  |  |
|                              | Mid        |                  |  |
|                              | Distal     |                  |  |
| Stomach                      | Antrum     |                  |  |
|                              | Body       |                  |  |
|                              | Fundus     |                  |  |
|                              | Pylorus    |                  |  |
| Small bowel*                 | Duodenum   | D1               |  |
|                              |            | D2               |  |
|                              |            | D3               |  |
|                              |            | D4               |  |
|                              | Jejunum    | Proximal         |  |
|                              |            | Mid              |  |
|                              |            | Distal           |  |
|                              | Ileum      | Proximal         |  |
|                              |            | Mid              |  |
| Terminal Ileum               |            |                  |  |
| Large Bowel*                 | Appendix   |                  |  |
|                              | Caecum     |                  |  |
|                              | Ascending  |                  |  |
|                              | Transverse |                  |  |
|                              | Descending |                  |  |
|                              | Sigmoid    |                  |  |
|                              | Rectum     |                  |  |
| <b>CCF Sample Location *</b> |            |                  |  |
| Proximal Landmark            |            | Distance (mm/cm) |  |
| Distal Landmark              |            | Distance (mm/cm) |  |
| <b>Sample processing</b>     |            |                  |  |



|                                       |  |
|---------------------------------------|--|
| <b><i>Sample storage</i></b>          |  |
| 4°C                                   |  |
| -80°C                                 |  |
| Liquid nitrogen                       |  |
| <b><i>Storage solution</i></b>        |  |
| <b><i>Starting material</i></b>       |  |
| Fresh tissue                          |  |
| Frozen tissue                         |  |
| <b><i>Tissue dissociation</i></b>     |  |
| Whole tissue                          |  |
| Cell type purification                |  |
| Cell viability (%)                    |  |
| <b><i>Downstream applications</i></b> |  |
| Cell culture                          |  |
| Organoids generation                  |  |
| -omics processing                     |  |
| <b><i>Data generation</i></b>         |  |
| <b><i>omics approach used</i></b>     |  |
| RNA-seq                               |  |
| scRNA-seq                             |  |
| ATAC-seq                              |  |
| Proteomics                            |  |
| Whole genome sequencing               |  |
| <b><i>Library preparation</i></b>     |  |
| <b><i>Sequencing platform</i></b>     |  |

1043 \*If available provide distance (in mm/cm) from nearest anatomical landmark: gastroduodenal  
 1044 junction, ileocaecal valve, hepatic flexure, splenic flexure, or anus.

1045

1046 **Table 2 Summary of currently available computational packages.**

1047

1048

| Main categories                                      | Task of analysis                 | Available tools/algorithms   |
|--|----------------------------------|--|
| Combining and integrating datasets across modalities | General workflow                 | Seurat, Scanpy   |
|  | Trajectory inference             | Monocle, PAGA, Slingshot, Velocity, scVelo                                 |
|  | Imputation                       | MAGIC, scVI, SAVERX  |
|  | Batch effect correction          | LIGER, Harmony, BBKNN, scVI  |
|  | Cell type matching and searching | FR-Match, Cell BLAST   |
|  | Multi-modal integration          | GLUE, MOFA+, Cobolt, MultiMAP, ArchR, MultiVI                              |
|  | Metadata integration             | (Drokhlyansky et al., 2020)  |
| Cell type annotation                                 | Automatic cell type annotation   | MACA, singleR, ScType, Celltypist  |
|  | Novel markers                    | COMET, COSG  |
| Mapping spatial location                             | Spatial location inference       | ASAP, novoSpaRc, Cell2location, BinSpect, SpatialDE                        |
| Cellular and microbial cross talk                    | Cellular cross talk              | CellChat, scConnect, CellComm, CellPhoneDB                                 |
|  | Microbial cross talk             | (Kang et al., 2020)<br>(Chattopadhyay et al., 2018)<br>(Liao et al., 2020) |
| Regulatory network inference                         | TF regulatory network            | SCENIC, GRNBoost2, PIDC  |

1049

1050

1051

1052