

# 1 Cell type ontologies of the Human Cell Atlas

2  
3 **Authors:** David Osumi-Sutherland<sup>1</sup>, Chuan Xu<sup>2</sup>, Maria Keays<sup>2</sup>, Adam P. Levine<sup>3</sup>, Peter V.  
4 Kharchenko<sup>4</sup>, Aviv Regev<sup>5</sup>, Ed Lein<sup>6</sup>, Sarah A. Teichmann<sup>2,7</sup>

## 5 6 **Affiliations:**

7 1: EMBL-European Bioinformatics Institute, Wellcome Genome Campus, Hinxton,  
8 Cambridge CB10 1SD, UK

9 2: Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA,  
10 UK

11 3: Research Department of Pathology, University College London, London WC1E 6DD, UK

12 4: Department of Biomedical Informatics, Harvard Medical School, Boston, Massachusetts  
13 02115, USA

14 5: Genentech, 1 DNA Way, South San Francisco, California 94080, USA

15 6: Allen Institute for Brain Science, Seattle, Washington 98109, USA

16 7: Cavendish Laboratory, University of Cambridge, JJ Thomson Ave, Cambridge CB3 0HE,  
17 UK

18 Correspondence should be addressed to S.A.T. (st9@sanger.ac.uk)

## 19 20 **Abstract**

21 Massive single-cell profiling efforts have accelerated our discovery of the cellular  
22 composition of the human body, while at the same time raising the need to formalise this  
23 new knowledge. Here, we discuss current efforts to harmonise and integrate different  
24 sources of annotations of cell types and states into a reference cell ontology. We illustrate  
25 with examples how a unified ontology can consolidate and advance our understanding of cell  
26 types across scientific communities and biological domains.

## 27 28 **Main**

29 With collaboration of over 2,000 scientists across more than 1,000 institutes from 76  
30 countries to date, the Human Cell Atlas (HCA) has generated comprehensive molecular  
31 profiles of tens of millions of single cells across 18 different organs and systems, which, in  
32 turn, are advancing our understanding of the definition of cell types and states<sup>1, 2</sup>.  
33 Technological advances in single-cell and spatial genomics are rapidly expanding the  
34 compendium of known cell types<sup>3</sup>, and accelerating discoveries of a large variety of novel  
35 cell populations.

36  
37 For instance, these efforts have been applied to system-level disciplines such as  
38 immunology and neuroscience, both of which require an understanding of vast networks of  
39 cells and tissues. In immunology, cell types have been historically recognised and well  
40 characterised. Yet, the number of discrete cell types and specific cell states identified from  
41 single-cell genomics has exceeded expectations, particularly with respect to the diversity of  
42 cell states derived from developmental dynamics<sup>4</sup>, tissue-resident phenotypes<sup>5</sup> and  
43 activation states<sup>6</sup>. For example, transcriptomic profiling identified three decidual natural killer  
44 cell populations at the maternal-fetal interface, which show varying levels of  
45 immunoregulatory properties and which modulate trophoblast invasion<sup>7</sup>. Transcriptomic and

46 genomic profiling has also captured an increasing variety of cell types and gene  
47 programmes in the central and peripheral nervous systems. Cell atlasing - i.e. the creation of  
48 a cell atlas - of mammalian brains has led to the discovery of previously uncharacterised cell  
49 types, including over a hundred cell types in one single region of the neocortex<sup>8</sup>, as well as  
50 of cellular diversity due to species-specific adaptations in the cortex<sup>8</sup>. A similar dramatic  
51 increase in diversity has been reported in the peripheral nervous system such as in the  
52 enteric nervous system<sup>9, 10</sup>.

53

54 This incredible progress takes us closer to answering a general question motivating stem  
55 and developmental cell biologists, as well as the HCA project: what is the complete cellular  
56 makeup of the human body? Annotating cells and gene programmes is crucial not only to  
57 address this question but also to fully exploit these data for biological discovery, including in  
58 pathological states. This can only be achieved by naming the entities we study in a  
59 consolidated way, such that findings can be related between studies and one study can build  
60 on findings from multiple previous ones as knowledge is accrued and expanded. However,  
61 most annotations of single-cell genomics datasets to date have used uncontrolled free text  
62 (i.e. arbitrary naming schemes) for cell type names, making cross-searching of annotations  
63 across separate datasets challenging and unreliable. In some cases, with a naming scheme  
64 absent, cells are described merely by a subset of their molecular characteristics and thus  
65 can be hard to match between studies.

66

67 To fully answer the question of what the cellular composition of the human body is, there is  
68 an urgent need to put new discoveries from the HCA into the context of classical cell biology  
69 and anatomy, as well as developmental biology, neurobiology, and pathology. Cell  
70 ontologies, a structured controlled vocabulary for cell types in animals, are a tremendously  
71 powerful way of formalising such knowledge, which in turn opens up opportunities for  
72 quantitative scientific interrogation of the HCA data in new and exciting ways.

73

74 In this Perspective, we discuss the utility and parts of cell ontologies, review the state of  
75 current cell ontologies, and conclude with ongoing efforts and how they can be applied for  
76 discovery over the coming years.

77

### 78 **Using cell ontology for knowledge integration and mining**

79 Biomedical ontologies originated in simple controlled vocabularies developed to supplement  
80 or replace the free text metadata in databases, clinical records and medical billing systems<sup>11</sup>.  
81 Standardising the text used to record, for example, diseases, gene functions, anatomical  
82 structures, and cell types within and between databases makes it possible to reliably search  
83 and group records referring to the same entities (diseases, cell types, etc.). However,  
84 controlled vocabularies are not sufficient for searching and grouping records with closely  
85 related contents. For example, a user searching a database for records relating to  
86 macrophages or liver sinusoid would not find records for Kupffer cells unless the data  
87 structures driving the search had some meaningful ways to relate the terms 'macrophage',  
88 'Kupffer cell' and 'liver sinusoid'. Cell ontologies provide mechanisms for this integration,  
89 allowing us to record a 'Kupffer cell' as a type of macrophage located in the liver sinusoid  
90 and then to enrich search results to take advantage of the classification and location  
91 relationships (Fig. 1).

92

93 Ontologies of cell types such as the Cell Ontology<sup>12</sup> and the Drosophila Anatomy Ontology<sup>13</sup>  
94 are increasingly used to annotate single-cell transcriptomic data. The use of ontology terms  
95 in dataset annotation relates annotated data back to hard-earned legacy knowledge,  
96 classical terminologies, and the accompanying understanding of cell types, anatomies, and  
97 development. Such annotation makes data cross-searchable, discoverable, integrable, and  
98 more accessible to general cell biologists. It facilitates cross-dataset analyses, allowing more  
99 quantitative analyses of similarities across thousands of individual cells, leading to more  
100 nuanced views of cell types, their classification, and their properties.

101

102 The Cell Ontology was first developed as a platform in 2004 to collect major cell types for  
103 humans and model organisms, and has been applied to various fields since then. For  
104 example, the Encyclopedia of DNA Elements (ENCODE) Consortium used the Cell Ontology  
105 to annotate its compendium of cell types, yielding a prioritised set of genetic and epigenetic  
106 elements<sup>14</sup>. Because the precise terms used for cell types, anatomical structures and  
107 diseases often vary greatly across sources, biomedical ontologies, including the cell  
108 ontology, typically use a bipartite system of universally resolvable IDs in the form of URLs for  
109 ontology terms, each linked to an official label. For example, the term with the primary label  
110 'Kupffer cell' in the Cell Ontology is identified by the persistent URL  
111 [http://purl.obolibrary.org/obo/CL\\_0000091](http://purl.obolibrary.org/obo/CL_0000091), which is further abbreviated to a compact form  
112 CL:0000091<sup>15</sup>. Critically, using resolvable IDs rather than labels to refer to cell types in  
113 database records allows associated metadata (labels, descriptions, and references) and  
114 their relationships (anatomy, development, functional and pathological relevance) to evolve  
115 over time with no cost for the databases and records that use IDs to refer to them (Fig. 1).

116

117 Ontologies can serve to link and integrate heterogeneous data types related to the same cell  
118 type across multiple modalities. For example, Virtual Fly Brain<sup>16, 17</sup> and the Fly Cell Atlas<sup>18</sup>  
119 use the same ontology terms to annotate 3D images of neurons (>70,000 images),  
120 connectomics data (>3.5 million pairwise connections), and single-cell transcriptomics data  
121 (~600,000 cells). Similarly, Cell Ontology terms, classifications and relationships are also  
122 increasingly used to define and classify terms in the Gene Ontology<sup>19</sup> (>750 terms) and in  
123 widely-used ontologies of phenotypes (730 terms in the Human Phenotype Ontology<sup>20</sup>) and  
124 diseases (>3,000 terms in the Mondo disease ontology<sup>21</sup>). These links make it possible to  
125 combine single-cell, phenotype, and disease data relating to the same cell types. With the  
126 advent of large-scale single-cell transcriptomic atlasing, community-driven nomenclature and  
127 ontology building projects have emerged and are coordinating with existing ontology building  
128 efforts (e.g. HCA Biological Networks<sup>2</sup>, HuBMAP<sup>22</sup>, BRAIN Initiative Cell Census Network  
129 (BICCN)<sup>23</sup> and Cell Annotation Platform (<http://celltype.info>)).

130

131 This is already impacting our ability to organise our knowledge of cell types for comparisons  
132 of datasets across individual laboratories, and notably, for effectively interpreting health and  
133 disease using the knowledge from both classical histopathology and single-cell genomics.  
134 For instance, ontological distinctions between fetal and mature cells in the kidney are  
135 mirrored by differences in their molecular signatures, which are critical to understanding the  
136 divergent origins of pediatric and adult kidney cancers, respectively<sup>24</sup>. Similarly, consistently  
137 annotated datasets allowed cross-tissue meta-analyses for COVID-19 that identified  
138 specialised nasal epithelial cells enriched for expression of SARS-CoV-2 entry factors<sup>25</sup>,  
139 identified covariates such as age, sex, and smoking status associated with the entry factor  
140 expression in lung and airway cells<sup>26</sup>, and compared cells in COVID-19 tissues from patient

141 autopsies to healthy and other disease conditions<sup>27</sup>, again highlighting the necessity and  
142 utility of establishing agreed-upon ontological classifications.

143

#### 144 **Considerations in the classification of human cell types**

145 Biologists have long recognised that the natural world lends itself to hierarchical systems of  
146 classification, which capture the underlying hierarchical processes driving biology, such as  
147 the phylogenetic classification of species by morphological and molecular observations.  
148 Similarly, cell types can be hierarchically classified and categorised in ever-increasing levels  
149 of resolution, from a general cell type like an endothelial cell, through more specialised types  
150 like a liver sinusoidal endothelial cell (LSEC), down to highly specialised types found in  
151 specific locations such as a periportal LSEC. As with a species' taxonomy, there are various  
152 kinds of observations informing the ultimate classification, and these different types of  
153 information are often used in concert to arrive at a particular cell type definition.

154

155 Take anatomical locations as an example: the Cell Ontology<sup>12</sup> imports information about  
156 anatomical structures and features from the Uber-anatomy Ontology (Uberon)<sup>28</sup> and relates  
157 them to the Cell Ontology terms using, for example, 'part of' to relate cell types to the tissues  
158 and organs, and 'located in' to relate cell types to cavities within structures. For example, the  
159 Cell Ontology definition of an LSEC includes a 'part of' relationship to 'hepatic sinusoid',  
160 which indicates that the liver sinusoidal endothelial cell forms part of the structure of the  
161 hepatic sinusoid as defined in Uberon, whereas the definition of Kupffer cells records that  
162 they are 'located in' (the lumen of) the hepatic sinusoid. In an anatomically higher hierarchy,  
163 the definition of hepatic sinusoid involves relations to the liver lobule and the liver overall,  
164 which is in turn defined by its structure, location and physiological role in the body. The  
165 LSEC is hence hierarchically defined relative to the whole organism down to its individual  
166 position in the specific tissue where it is found (Fig. 2a). Furthermore, since the Cell  
167 Ontology classifies cell types hierarchically from generic cell types down to more specialised  
168 types, an LSEC is also defined as a descendent of the general endothelial cell class in the  
169 Cell Ontology. The main LSEC class (officially 'endothelial cell of hepatic sinusoid') has its  
170 own descendent classes, representing further specialisations of LSECs: 'endothelial cell of  
171 periportal hepatic sinusoid' and 'endothelial cell of pericentral hepatic sinusoid'.

172

173 Sources of information contributing to a cell type categorisation include morphological  
174 features, developmental origins, and functional profiles. Ontologies attempt to capture all  
175 terms that are used by different scientific communities to refer to the same cell type, as well  
176 as alternative names that may not be commonly used. Historically, different fields in biology  
177 have focused on different aspects of cells to drive their naming. For example, many immune  
178 cells have been classified according to which cell surface protein(s) they express<sup>29-36</sup>,  
179 whereas cells of the nervous system have been named according to a combination of  
180 features including morphologies, physiologies, connectivities and the roles they play in the  
181 neuronal circuitry<sup>37</sup>. In some systems, such as the retina<sup>38</sup>, there is strong evidence that cell  
182 types can be classified consistently regardless of the features used to classify them. In these  
183 cases, classically defined cell types typically align well with those identified by analysis of  
184 single-cell transcriptomic data, making cell annotation straightforward. In other cases,  
185 different features could in principle lead to different cell type classifications, making  
186 consistent annotation more challenging. Formal ontologies are able to support multiple  
187 overlapping classification schemes, and thus can potentially help reconcile different  
188 classification schemes, at least at the level of more generally grouped classes.

189

190 Cell ontologies also represent developmental lineages and, to a more limited extent, cell  
191 states such as activation, cycling, morphological changes and stresses (Fig. 2b) - either  
192 directly or through extensions of existing annotations. Cell-cycle states, for example, can be  
193 represented in the annotation system by combining a Cell Ontology term with a term from  
194 the Gene Ontology Cell Cycle Phase terms. Developmental or actively regenerating tissues  
195 present particular challenges to cell ontology development, as a plethora of intermediate  
196 states and continuous branching lineages can be partitioned. In such a setting, cell  
197 annotation needs to emphasise the relative ordering of states, or their positions on a  
198 continuous differentiation path. There are also striking examples of developmental  
199 convergence (developmental homoplasy). Somatosensory neurons, for example, can be of  
200 mixed origin, from the neural crest or sensory placodes<sup>39</sup>. Similarly, dermal fibroblasts in  
201 different parts of the trunk or face are derived from distinct embryonic lineages, despite  
202 molecular and phenotypic likeness<sup>40</sup>. Nevertheless, cell ontologies record gross lineage  
203 relationships, with limited temporal resolution between developing/progenitor and mature cell  
204 types using specific relations where these relationships are stereotyped and consistent. To  
205 date, the Cell Ontology records lineage and differentiation relationships for more than 1,900  
206 cell types, connecting developing cell types to developing tissues and stages via links to  
207 Uberon.

208

209 Many processes driving cell diversifications, including ontogeny (cell differentiation),  
210 morphogenesis (often driven by continuous gradients), and the dual impact of a cell's  
211 differentiation history and tissue context, are imprinted in a cell's molecular properties and  
212 can be captured by hierarchical representations. Therefore, molecular features can serve as  
213 the basis for robust cell type classification, reflecting these underlying processes (even when  
214 the process is not explicitly known). Currently, cell types and states can be elucidated from  
215 single-cell transcriptomic, epigenomic and proteomic expression profiles, using different  
216 software such as SCCAF<sup>41</sup>. Further complemented by morphological, physiological,  
217 developmental, and functional properties, this data-driven framework makes cell annotations  
218 comparable across independent ontology efforts and the inferred cell types understandable  
219 across different communities. Of note, while these inferences are unbiased, it is important to  
220 reconcile them with conventional biological and clinical understanding and terminologies.

221

## 222 **Current state of ontologies**

223 First developed as the platforms to integrate cross-species ontology information, the Cell  
224 Ontology and Uberon are now species-neutral ontologies with a strong focus on mammalian  
225 cell types and anatomies with standard mechanisms for recording the species applicability of  
226 terms. To date, the Cell Ontology has 2,401 terms covering all major cell types. The  
227 granularity of this coverage is variable, with the greatest coverage currently for the immune  
228 system (>500 cell types). Uberon defines over 14,000 types of anatomical structures and  
229 records many types of relationships between them. Practically, the Cell Ontology and  
230 Uberon are tightly integrated with each other. Almost 2,000 cell types in the Cell Ontology  
231 are linked by 'part of' relationships to the anatomical structures defined in Uberon. Further  
232 combining the Cell Ontology with newly discovered cell populations from HCA data, we are  
233 beginning to extensively cover major organs and cell types in the human body (Table 1).

234

235 The human-applicable components of the Cell Ontology and Uberon are under active  
236 development as part of multiple collaborative efforts. For human data, terms are being added

237 in a coordinated fashion to both ontology platforms in response to the requests of individual  
238 labs, as well as to the annotation needs of atlasing projects including HCA's Data  
239 Coordination Platform<sup>2</sup> (<https://data.humancellatlas.org>), and the Cambridge Cell Atlas portal  
240 ([www.cambridgecellatlas.org](http://www.cambridgecellatlas.org)). Editing of the Cell Ontology and Uberon is coordinated by a  
241 team of researchers drawn from a growing number of collaborating projects including the  
242 Human Cell Atlas (Chan Zuckerberg Initiative), HuBMAP (NIH), the Monarch Initiative (NIH)  
243 and the Cell Annotation Platform (a collaborative effort funded by Schmidt Futures). This  
244 team of editing researchers runs regular open training sessions, and anyone trained to edit  
245 the ontology can join the editing team. Edits are coordinated and reviewed on GitHub  
246 (<https://github.com/obophenotype/cell-ontology>), with all changes and releases subject to  
247 automated quality-control tests prior to approval. Issues not resolved after discussion on  
248 open tickets are coordinated via monthly editor video conferences, which also coordinate the  
249 general focus of Cell Ontology and Uberon efforts. These calls frequently feature guest  
250 speakers with a particular interest in extending the Cell Ontology or Uberon in specific areas.  
251 Cell Ontology and Uberon are both members of the Open Biological and Biomedical  
252 Ontology (OBO) Foundry group of ontologies<sup>15</sup>, a loose alliance of ontologies committed to  
253 adopting common standards and aligning semantics and ontology infrastructure. All these  
254 endow the Cell Ontology and Uberon with the ability to continuously evolve with inputs from  
255 various projects and perspectives and to supply formalised ontology information back to the  
256 projects (Table 2). Examples of the co-evolution of the Cell Ontology and human cell  
257 ontology-building efforts are listed below.

258  
259 The Brain Data Standards Initiative, part of the NIH BRAIN Initiative Cell Census Network, is  
260 extending the Cell Ontology with terms for cortical cell types defined by single-cell  
261 transcriptomics, with a current focus on the primary motor cortex of human, marmoset, and  
262 mouse<sup>42</sup>. This work leverages existing efforts on nomenclature standards<sup>43</sup>, but importantly  
263 aims to use the quantitative hierarchical cell type classification from single-cell genomics as  
264 a data-driven foundation for ontological definitions. Different data types about these cell  
265 types are integrated at different levels of the hierarchy, including their spatial tissue  
266 distributions, morphological and physiological properties, and axonal projection targets.  
267 Ultimately such a data-driven approach may be used across the entire human body,  
268 providing a common metric in gene usage to measure similarities and potential common  
269 developmental origins across organs.

270  
271 The ASCT+B effort<sup>44</sup> presented as an accompanying Perspective in this issue is a  
272 HuBMAP/HTAN/HCA community-wide project to build tables representing the human  
273 anatomy and cell type terminology needed for annotating scRNA-seq datasets, and to record  
274 expert-approved lists of markers for cell types. Entries in these tables are mapped to existing  
275 Cell Ontology or Uberon terms where possible or turned into term requests for these data  
276 resources, when new terms are needed. The relationships between cell types and  
277 anatomical structures encoded in these tables are validated against the Cell Ontology and  
278 Uberon. The results of this validation are relayed to improve the tables, Uberon, and Cell  
279 Ontology based on discussions and agreement with experts. For example, the ASCT+B  
280 project is building an expert-validated ontological model of the human vasculature that is  
281 feeding hundreds of new terms and relationships back into Uberon. One important outcome  
282 of this work will be a curated subset of Cell Ontology and Uberon terms for reliably  
283 annotating human scRNA-seq data, both for the healthy HCA data as well as disease  
284 samples.

285

286 As part of the human cell-focused Sanger-EBI (European Bioinformatics Institute)  
287 Cambridge Cell Atlas portal (<https://www.cambridgecellatlas.org>), an effort to make results  
288 from human single-cell gene expression experiments easily accessible to a broad  
289 community of users including clinicians, the Cell Ontology is being enriched and extended  
290 based on contributions from pathologists and clinicians. This will introduce human cell types  
291 annotated with details of specific immunohistochemical markers that are in routine clinical  
292 use in diagnostic pathology. This ontology can then be integrated into the search  
293 functionality of the Cambridge Cell Atlas platform to enable searching based on a specific  
294 immunohistochemical marker or panel of markers, allowing for the identification of the  
295 normal cell type(s) (and potentially pathogenic cell types as well) that express the marker(s).  
296 This functionality could be useful to pathologists in interpreting and contextualising the range  
297 of cell types stained by different immunohistochemical markers on histological sections,  
298 cytological preparations or by flow cytometry, and in understanding perturbations in staining  
299 patterns in pathological states.

300

### 301 **Applications of a cell ontology**

302 Cell ontologies provide a single place to look up cell types for the community. Through this,  
303 knowledge can be aggregated and standardised in an encyclopaedic sense. First, cross-  
304 modal data integration can reinforce or refine the identity of a cell type. For example, the  
305 survey on the mammalian neocortex revealed the correspondence of various cellular  
306 properties when overlapping imaging, electrophysiology and connectivity with transcriptomic  
307 profiles<sup>37</sup>. Second, mining of an ontological classification system can reveal major trends  
308 with respect to shared cell types across organ-specific atlases (e.g. immune, stromal and  
309 endothelial cells) versus specialised types (e.g. goblet cell in the gut and lung), emphasising  
310 the concept of a tissue being the collective of its cells operating in concert in a specific 3D  
311 organisation.

312

313 Importantly, with more single-cell resources employing the cell and anatomy ontologies,  
314 including but not limited to the Fly Cell Atlas, EBI's Single Cell Expression Atlas and Sanger-  
315 EBI Cambridge Cell Atlas, cell ontologies can link scientific and medical communities  
316 through common nomenclatures and markers for human cell biology, pathology and disease.  
317 This link, in a broader sense, represents cross-community research where a common cell  
318 type reference can be referred. For example, a well-defined cell type classification of human  
319 head and neck tumors, which covered major immune and non-immune cell populations, was  
320 utilised as the reference to interrogate the cellular signals contributing to bulk samples of  
321 head and neck squamous cell carcinoma from The Cancer Genome Atlas (TCGA), revealing  
322 the association of tumor-infiltrating regulatory T cells with improved survival in head and  
323 neck cancer<sup>45</sup>.

324

325 At the same time, immunohistochemical markers in routine clinical use (such as those listed  
326 by Pathology Outlines, <https://www.pathologyoutlines.com/stains.html>), which are linked to  
327 the non-pathological cell types by the Cambridge Cell Atlas project, could also be curated  
328 and further linked to pathological tissues and cell states that express them. This would  
329 provide hundreds of antibodies to link cell types and anatomical structures with the Cell  
330 Ontology and Uberon, albeit with a focus on pathological states (of course the Cell Ontology  
331 and Uberon currently focus on healthy homeostatic states).

332

333 The application of cell ontologies will be most pertinent in the context of interactive and  
334 automated systems for the interpretation and annotation of single-cell genomic datasets. A  
335 number of efforts to design such systems are under way, including automated cell  
336 annotation projection pipelines<sup>46-52</sup>. For example, as part of the HCA initiative, the Cell  
337 Annotation Platform (CAP) aims to provide a general repository for cell annotations of  
338 different datasets, in combination with interactive tools for annotating new datasets. For a  
339 cell of interest, CAP user interfaces will suggest the appropriate ontology terms based on  
340 text search, learned synonyms, and eventually molecular signatures themselves. Where no  
341 appropriate term is available from the Cell Ontology, free text annotation will be used as the  
342 basis for new term addition to the Cell Ontology. Similarly, the HuBMAP data portal assigns  
343 cell annotations to scRNA-seq datasets with an Azimuth-based label transfer procedure<sup>49</sup>  
344 based on a vocabulary of cell types from the Cell Ontology, aiming at assessing cellular  
345 diversities at different levels of resolution. With an initial focus on immune cells, CellTypist  
346 uses an expandable cross-tissue cell reference before predicting cell identities with a logistic  
347 regression-based label transfer pipeline, with all derived cell types directly interpretable by  
348 the Cell Ontology<sup>48</sup>. Conversely, the resulting knowledgebase of commonly used annotation  
349 terms and associated molecular signatures will provide a useful resource to extend  
350 ontologies as well as to train and optimise machine learning models that automate the  
351 annotation task. In parallel to these efforts, data-driven ontology development is advancing  
352 community engagement in specific research domains such as NeMO Analytics for the brain,  
353 <https://nemoanalytics.org>, and gEAR for the ear<sup>53</sup>.

354

### 355 **Summary and outlook**

356 Resolving the cellular makeup of the human body warrants the categorisation of cells in a  
357 standardised framework. The Cell Ontology offers one such avenue to consolidating this  
358 knowledge in an encyclopaedic manner, with applications from cell and tissue biology all the  
359 way to the clinic. Despite potential cell classification ambiguities and transient cellular states,  
360 each facet of a cell ranging from morphological to molecular features can be taken into  
361 account, until a defining status is reached and recognised by the community.

362

363 Many HCA-related resources, such as cellxgene<sup>54</sup>, have been using the Cell Ontology for *de*  
364 *novo* cell annotation. Cell Ontologies serve other sources of data by retrieving or delivering  
365 ontology-level information. We anticipate the synergy between the HCA project and the Cell  
366 Ontology will continue to grow over the coming years and beyond the completion of HCA,  
367 with dimensions of human genetic variation, ageing and disease on the horizon. HCA single-  
368 cell omics data provide a foundation for the development of cell ontologies, which are  
369 powerful resources to define cell types that are universal across the entire body or specific to  
370 subsets of tissues and which will facilitate future research. This will become more pressing  
371 and clearer as the number of HCA studies of individual tissues and organs increases. The  
372 HCA Biological Networks will provide nucleation points for expert community efforts to  
373 achieve gold standard, consensus cell annotations with cell ontology terms. With such a  
374 quantitative approach, common phenotypes and developmental origins of cell types will  
375 become understandable through shared gene usage, and functional similarities will be  
376 revealed in gene patterns. Whole-body consequences of disease will be understandable  
377 through differential gene usage in differently located cells. This will thus create opportunities  
378 for a new and different kind of quantitative data-driven framework extending and potentially  
379 transforming existing ontology efforts.

380



381 **Acknowledgements**

382 We are grateful to Jana Eliasova (scientific illustrator) for support with the figure, to Roser  
383 Vento-Tormo for comments on the figure and texts, and to the following clinicians and  
384 researchers for information on standard pathology markers for tissues and cells: Lia  
385 Campos, Andrew Dean, Luiza Moore, Neil Sebire, Teresa Brevini, Muzlifah Haniffa, Jing  
386 Eugene Kwa, James McCaffrey, and Alexandra Kreins. We also thank all members of the  
387 Cell Ontology and Uberon editorial teams including Chris Mungall, Nicolas Matentzoglou,  
388 Alexander Diehl, Nicole Washington, Shawn Tan, Paola Roncaglia, Tiago Lubiana and  
389 Damien Goutte-Gattat. Research reported in this publication was supported by the Wellcome  
390 Trust Grant 108413/A/15/D, the Office Of The Director, National Institutes Of Health of the  
391 National Institutes of Health under Award Number OT2OD026682', grants from the CZI  
392 (Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community Foundation),  
393 and Schmidt Futures (Grant 74). This publication is part of the Human Cell Atlas -  
394 [www.humancellatlas.org/publications/](http://www.humancellatlas.org/publications/).

395

396 **Competing interests**

397 Since January 2019, S.A.T. has been remunerated for consulting and SAB membership by  
398 Foresite Labs, GlaxoSmithKline, Biogen, Roche and Genentech, and is a founder and equity  
399 holder of Transition Bio. A.R. is a founder and equity holder in Celsius Therapeutics, an  
400 equity holder in Immunitas Therapeutics, and was a scientific advisory board member for  
401 ThermoFisher Scientific, Syros Pharmaceuticals and Neogene Therapeutics until August 1,  
402 2020. From August 1, 2020, A.R. is an employee of Genentech. A.R. is a named inventor on  
403 several patents and patent applications filed by the Broad Institute in the area of single cell  
404 and spatial genomics. Other authors declare no competing interests.

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

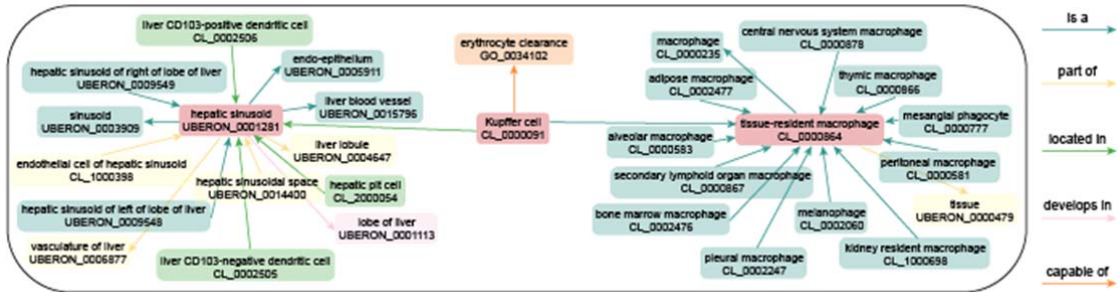
421

422

423

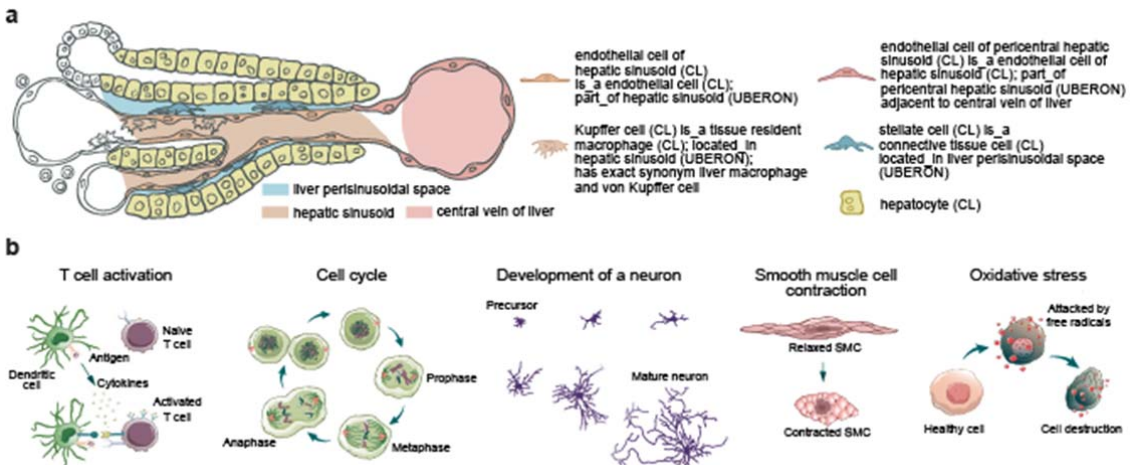
424

425



426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437

**Fig. 1: A graph representation of a portion of the Cell Ontology centred around the term Kupffer cell.** Graph showing the relationships between terms for anatomical structures (e.g. hepatic sinusoid), cell types (e.g. macrophage), and functional roles (e.g. erythrocyte clearance). Relationships shown include 'is a' which records the classification, 'part of' which relates cells to their tissues and organs, 'located in' which relates cells to spaces such as the hepatic sinusoid, 'develops in' which records the developmental origin, and 'capable of' which records the function.



438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452

**Fig. 2: A cell ontology links human cell types with anatomy and cell state transition.** **a**, The Cell Ontology (CL) has terms for a variety of cell types associated with the hepatic sinusoid (UBERON:0001281). The classification of these cell types allows them to be grouped with other cells from the same location (e.g. Kupffer cells (CL:0000091) can be grouped with other tissue-resident macrophages or with cells of the hepatic sinusoid). **b**, Ontologies can be used to encode transitions through diverse cell states. Examples include T cell activation following antigen recognition, cell cycling, neuron development and maturation, smooth muscle cell contraction and relaxation, and cell destruction after oxidative stress.

453

454

**Table 1 Current status of cell type enumerations in the Cell Ontology and HCA data.**

455

Summary of cell type numbers in the Cell Ontology and HCA data.

<b>Tissue</b>	<b>No. cell types (Cell Ontology version:2021-04- 22)</b>	<b>No. cell types as per HCA Ref</b>	<b>HCA Ref</b>
Kidney	127	33 (mature)/44 (fetal)	Stewart et al., 2019 <sup>55</sup>
Lymph node	12	19	James et al., 2020 <sup>56</sup>
Small and large intestine	125	132	Elmentaite et al., 2021 <sup>10</sup>
Lung	27	21; 58	Vieira Braga et al., 2019 <sup>57</sup> ; Travaglini et al., 2020 <sup>58</sup>
Liver	19	21; 39	Ramachandran et al., 2019 <sup>59</sup> ; Aizarani et al., 2019 <sup>60</sup>
Muscle	31	19	Litviňuková et al., 2020 <sup>61</sup>
Esophagus	11	18	Madisson et al., 2019 <sup>5</sup>
Heart	54	67	Litviňuková et al., 2020 <sup>61</sup>
Thymus	55	44	Park et al., 2020 <sup>62</sup>
Brain (primary motor cortex)	133	127	Bakken et al., 2020 <sup>42</sup>
Bone marrow and blood	515	48	HCA Data Portal
Skin	71	34	Reynolds et al., 2021 <sup>63</sup>
Endometrium and decidua	5	14; 11	Garcia-Alonso et al., 2021 <sup>64</sup> ; Vento-Tormo et al., 2018 <sup>7</sup>
Placenta	10	5	Vento-Tormo et al., 2018 <sup>7</sup>

456

457

**Table 2 Projects using and contributing to the Cell Ontology (CL).**

<b>Project</b>	<b>Description</b>	<b>CL Use</b>	<b>URL</b>
Cell Annotation Platform	An open annotation platform for scRNA-seq data	Uses CL and free text for cell type annotation	<a href="http://celltype.info">http://celltype.info</a>
EBI Single Cell Expression Atlas & Cambridge Cell Atlas	Open public repository for exploration of single cell gene expression data	Uses CL to annotate samples and cell types in tertiary analysis	<a href="https://www.ebi.ac.uk/gxa/sc">https://www.ebi.ac.uk/gxa/sc</a> and <a href="https://www.cambridgecellatlas.org">https://www.cambridgecellatlas.org</a>
HCA/DCP	Community generated, multi-omic, open data processed by standardized pipelines	Uses CL to annotate samples and cell types in tertiary analysis	<a href="https://data.humancellatlas.org">https://data.humancellatlas.org</a>
HuBMAP/CCF ASCT+B tables	Expert curated tables of human cell types, their markers and anatomical context	Maps all cell types to CL	<a href="https://hubmapconsortium.github.io/ccf-asct-reporter">https://hubmapconsortium.github.io/ccf-asct-reporter</a>
cellxgene	An open annotation platform requiring annotation with ontology terms	Uses CL to annotate samples and cell types in tertiary analysis	<a href="https://chanzuckerberg.github.io/cellxgene">https://chanzuckerberg.github.io/cellxgene</a>
Tabula Muris	Curated whole mouse scRNA-seq atlas	Uses CL to annotate gross cell types, extending definitions with free text and markers	<a href="https://tabula-muris.ds.czbiohub.org">https://tabula-muris.ds.czbiohub.org</a>
Monarch Initiative	A resource building ontologies of phenotypes and disease and using these to build an integrated collection of phenotype/disease to gene/variant associations	Defines cellular phenotypes and diseases	<a href="https://monarchinitiative.org">https://monarchinitiative.org</a>
Gene Ontology	The world's largest source of information on the function and location of gene products	Defines cell type-specific organelles and biological processes	<a href="http://geneontology.org">http://geneontology.org</a>

CellTypist	An open source tool for automated cell type annotations as well as a work group in charge of curating models and ontologies	Maps all cell types to CL	<a href="https://www.celltypist.org">https://www.celltypist.org</a>
Human Immunology Project Consortium (HIPC)	A comprehensive, centralised research resource with the goal of facilitating a comprehensive understanding of the human immune system and its regulation	Works with CL to improve the representation of human immune cell types for use in data annotation	<a href="https://www.immune-profiling.org/hipc">https://www.immune-profiling.org/hipc</a>

460

461 **References**

462

- 463 1. Regev, A. *et al.* The Human Cell Atlas. *Elife* **6** (2017).
- 464 2. Rozenblatt-Rosen, O., Stubbington, M.J.T., Regev, A. & Teichmann, S.A. The  
465 Human Cell Atlas: from vision to reality. *Nature* **550**, 451-453 (2017).
- 466 3. Aldridge, S. & Teichmann, S.A. Single cell transcriptomics comes of age. *Nat*  
467 *Commun* **11**, 4307 (2020).
- 468 4. Popescu, D.M. *et al.* Decoding human fetal liver haematopoiesis. *Nature* **574**, 365-  
469 371 (2019).
- 470 5. Madisson, E. *et al.* scRNA-seq assessment of the human lung, spleen, and  
471 esophagus tissue stability after cold preservation. *Genome Biol* **21**, 1 (2019).
- 472 6. Hagai, T. *et al.* Gene expression variability across cells and species shapes innate  
473 immunity. *Nature* **563**, 197-202 (2018).
- 474 7. Vento-Tormo, R. *et al.* Single-cell reconstruction of the early maternal-fetal interface  
475 in humans. *Nature* **563**, 347-353 (2018).
- 476 8. Hodge, R.D. *et al.* Conserved cell types with divergent features in human versus  
477 mouse cortex. *Nature* **573**, 61-68 (2019).
- 478 9. Drokhyansky, E. *et al.* The Human and Mouse Enteric Nervous System at Single-  
479 Cell Resolution. *Cell* **182**, 1606-1622 e1623 (2020).
- 480 10. Elmentaite, R. *et al.* Cells of the human intestinal tract mapped across space and  
481 time. *Nature* **597**, 250-255 (2021).
- 482 11. Bodenreider, O. & Stevens, R. Bio-ontologies: current trends and future directions.  
483 *Brief Bioinform* **7**, 256-274 (2006).
- 484 12. Diehl, A.D. *et al.* The Cell Ontology 2016: enhanced content, modularization, and  
485 ontology interoperability. *J Biomed Semantics* **7**, 44 (2016).
- 486 13. Costa, M., Reeve, S., Grumblin, G. & Osumi-Sutherland, D. The Drosophila  
487 anatomy ontology. *J Biomed Semantics* **4**, 32 (2013).
- 488 14. Consortium, E.P. An integrated encyclopedia of DNA elements in the human  
489 genome. *Nature* **489**, 57-74 (2012).
- 490 15. Smith, B. *et al.* The OBO Foundry: coordinated evolution of ontologies to support  
491 biomedical data integration. *Nat Biotechnol* **25**, 1251-1255 (2007).
- 492 16. Milyaev, N. *et al.* The Virtual Fly Brain browser and query interface. *Bioinformatics*  
493 **28**, 411-415 (2012).
- 494 17. Osumi-Sutherland, D., Costa, M., Court, R. & O'Kane, C.J. Virtual Fly Brain - Using  
495 OWL to support the mapping and genetic dissection of the Drosophila brain. *CEUR*  
496 *Workshop Proc* **1265**, 85-96 (2014).
- 497 18. Li, H. *et al.* Fly Cell Atlas: a single-cell transcriptomic atlas of the adult fruit fly.  
498 *bioRxiv*, 2021.2007.2004.451050 (2021).
- 499 19. Gene Ontology, C. The Gene Ontology resource: enriching a GOLD mine. *Nucleic*  
500 *Acids Res* **49**, D325-D334 (2021).
- 501 20. Mungall, C.J. *et al.* The Monarch Initiative: an integrative data and analytic platform  
502 connecting phenotypes to genotypes across species. *Nucleic Acids Res* **45**, D712-  
503 D722 (2017).
- 504 21. Jacqz, E., Branch, R.A., Heidemann, H. & Aujard, Y. [Prevention of nephrotoxicity of  
505 amphotericin B during the treatment of deep candidiasis]. *Ann Biol Clin (Paris)* **45**,  
506 689-693 (1987).
- 507 22. Hu, B.C. The human body at cellular resolution: the NIH Human Biomolecular Atlas  
508 Program. *Nature* **574**, 187-192 (2019).
- 509 23. Ecker, J.R. *et al.* The BRAIN Initiative Cell Census Consortium: Lessons Learned  
510 toward Generating a Comprehensive Brain Cell Atlas. *Neuron* **96**, 542-557 (2017).
- 511 24. Young, M.D. *et al.* Single-cell transcriptomes from human kidneys reveal the cellular  
512 identity of renal tumors. *Science* **361**, 594-599 (2018).
- 513 25. Sungnak, W. *et al.* SARS-CoV-2 entry factors are highly expressed in nasal epithelial  
514 cells together with innate immune genes. *Nat Med* **26**, 681-687 (2020).

- 515 26. Muus, C. *et al.* Single-cell meta-analysis of SARS-CoV-2 entry genes across tissues  
516 and demographics. *Nat Med* **27**, 546-559 (2021).
- 517 27. Delorey, T.M. *et al.* COVID-19 tissue atlases reveal SARS-CoV-2 pathology and  
518 cellular targets. *Nature* (2021).
- 519 28. Mungall, C.J., Torniai, C., Gkoutos, G.V., Lewis, S.E. & Haendel, M.A. Uberon, an  
520 integrative multi-species anatomy ontology. *Genome Biol* **13**, R5 (2012).
- 521 29. Bernard, A., Institut national de la santé et de la recherche médicale (France), World  
522 Health Organization. & IUIS. *Leucocyte typing : human leucocyte differentiation*  
523 *antigens detected by monoclonal antibodies : specification, classification,*  
524 *nomenclature = Typage leucocytaire : antigènes de différenciation leucocytaire*  
525 *humains révélés par les anticorps monoclonaux.* (Springer-Verlag, Berlin ; New York;  
526 1984).
- 527 30. Reinherz, E.L. *Leukocyte typing II.* (Springer-Verlag, New York; 1986).
- 528 31. McMichael, A.J. *Leucocyte typing III : white cell differentiation antigens.* (Oxford  
529 University Press, Oxford ; New York; 1987).
- 530 32. Knapp, W. *Leucocyte typing IV : white cell differentiation antigens.* (Oxford University  
531 Press, Oxford ; New York; 1989).
- 532 33. Schlossman, S.F. *Leucocyte typing V : white cell differentiation antigens :  
533 proceedings of the fifth international workshop and conference held in Boston, USA,*  
534 *3-7 November, 1993.* (Oxford University Press, Oxford ; New York; 1995).
- 535 34. Kishimoto, T. *Leucocyte typing VI : white cell differentiation antigens : proceedings of  
536 the sixth international workshop and conference held in Kobe, Japan, 10-14  
537 November 1996.* (Garland Pub., New York; 1998).
- 538 35. Mason, D. *Leucocyte typing VII : white cell differentiation antigens : proceedings of  
539 the Seventh International Workshop and Conference held in Harrogate, United  
540 Kingdom.* (Oxford University Press, Oxford; 2002).
- 541 36. Zola, H. *Leukocyte and stromal cell molecules : the CD markers.* (Wiley-Liss,  
542 Hoboken, N.J.; 2007).
- 543 37. Yuste, R. *et al.* A community-based transcriptomics classification and nomenclature  
544 of neocortical cell types. *Nat Neurosci* **23**, 1456-1468 (2020).
- 545 38. Shekhar, K. *et al.* Comprehensive Classification of Retinal Bipolar Neurons by  
546 Single-Cell Transcriptomics. *Cell* **166**, 1308-1323 e1330 (2016).
- 547 39. Vermeiren, S., Bellefroid, E.J. & Desiderio, S. Vertebrate Sensory Ganglia: Common  
548 and Divergent Features of the Transcriptional Programs Generating Their Functional  
549 Specialization. *Front Cell Dev Biol* **8**, 587699 (2020).
- 550 40. Driskell, R.R. & Watt, F.M. Understanding fibroblast heterogeneity in the skin. *Trends  
551 Cell Biol* **25**, 92-99 (2015).
- 552 41. Miao, Z. *et al.* Putative cell type discovery from single-cell gene expression data. *Nat  
553 Methods* **17**, 621-628 (2020).
- 554 42. Bakken, T.E. *et al.* Evolution of cellular diversity in primary motor cortex of human,  
555 marmoset monkey, and mouse. *bioRxiv*, 2020.2003.2031.016972 (2020).
- 556 43. Miller, J.A. *et al.* Common cell type nomenclature for the mammalian brain. *Elife* **9**  
557 (2020).
- 558 44. Börner, K. *et al.* Anatomical Structures, Cell Types, and Biomarkers Tables Plus 3D  
559 Reference Organs in Support of a Human Reference Atlas. *bioRxiv*,  
560 2021.2005.2031.446440 (2021).
- 561 45. Qi, Z. *et al.* Single-Cell Deconvolution of Head and Neck Squamous Cell Carcinoma.  
562 *Cancers (Basel)* **13** (2021).
- 563 46. Kiselev, V.Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell RNA-seq data  
564 across data sets. *Nat Methods* **15**, 359-362 (2018).
- 565 47. Kimmel, J.C. & Kelley, D.R. Semi-supervised adversarial neural networks for single-  
566 cell classification. *Genome Res* (2021).
- 567 48. Domínguez, C.C. *et al.* Cross-tissue immune cell analysis reveals tissue-specific  
568 adaptations and clonal architecture across the human body. *bioRxiv*,  
569 2021.2004.2028.441762 (2021).

- 570 49. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587  
571 e3529 (2021).
- 572 50. Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a  
573 transitional profibrotic macrophage. *Nat Immunol* **20**, 163-172 (2019).
- 574 51. Bernstein, M.N., Ma, Z., Gleicher, M. & Dewey, C.N. CellO: comprehensive and  
575 hierarchical cell type classification of human cells with the Cell Ontology. *iScience* **24**,  
576 101913 (2021).
- 577 52. Hou, R., Denisenko, E. & Forrest, A.R.R. scMatch: a single-cell gene expression  
578 profile annotation tool using reference datasets. *Bioinformatics* **35**, 4688-4695  
579 (2019).
- 580 53. Orvis, J. *et al.* gEAR: Gene Expression Analysis Resource portal for community-  
581 driven, multi-omic data exploration. *Nat Methods* **18**, 843-844 (2021).
- 582 54. Megill, C. *et al.* cellxgene: a performant, scalable exploration platform for high  
583 dimensional sparse matrices. *bioRxiv*, 2021.2004.2005.438318 (2021).
- 584 55. Stewart, B.J. *et al.* Spatiotemporal immune zonation of the human kidney. *Science*  
585 **365**, 1461-1466 (2019).
- 586 56. James, K.R. *et al.* Distinct microbial and immune niches of the human colon. *Nat*  
587 *Immunol* **21**, 343-353 (2020).
- 588 57. Vieira Braga, F.A. *et al.* A cellular census of human lungs identifies novel cell states  
589 in health and in asthma. *Nat Med* **25**, 1153-1163 (2019).
- 590 58. Travaglini, K.J. *et al.* A molecular cell atlas of the human lung from single-cell RNA  
591 sequencing. *Nature* **587**, 619-625 (2020).
- 592 59. Ramachandran, P. *et al.* Resolving the fibrotic niche of human liver cirrhosis at  
593 single-cell level. *Nature* **575**, 512-518 (2019).
- 594 60. Aizarani, N. *et al.* A human liver cell atlas reveals heterogeneity and epithelial  
595 progenitors. *Nature* **572**, 199-204 (2019).
- 596 61. Litvinukova, M. *et al.* Cells of the adult human heart. *Nature* **588**, 466-472 (2020).
- 597 62. Park, J.E. *et al.* A cell atlas of human thymic development defines T cell repertoire  
598 formation. *Science* **367** (2020).
- 599 63. Reynolds, G. *et al.* Developmental cell programs are co-opted in inflammatory skin  
600 disease. *Science* **371** (2021).
- 601 64. Garcia-Alonso, L. *et al.* Mapping the temporal and spatial dynamics of the human  
602 endometrium *in vivo* and *in vitro*. *bioRxiv*,  
603 2021.2001.2002.425073 (2021).
- 604