

# Towards a pragmatic definition of cell type

This manuscript ([permalink](#)) was automatically generated from [lubianat/technotype@136df16](#) on January 22, 2021.

## Authors

---

- **Tiago Lubiana**

 [0000-0003-2473-2313](#) ·  [lubianat](#)

Department of Clinical and Toxicological Analyses, School of Pharmaceutical Sciences, University of São Paulo, São Paulo, Brazil

- **Helder I Nakaya**

 [0000-0001-5297-9108](#)

Department of Clinical and Toxicological Analyses, School of Pharmaceutical Sciences, University of São Paulo, São Paulo, Brazil

## Abstract

---

The concept of cell type is key for modeling biology. Recent technological advances are prompting us to rethink what we understand by cell type and how we classify them. There is currently no consensus for a definition of cell type, which makes it hard to integrate knowledge across life sciences. We propose here that a cell type should represent any class of cell that (1) is explicitly defined; (2) is identifiable within a taxon; and (3) is theoretically useful. We also present four classes of cell types: *sensu stricto* cell types, archetypes, infratypes, and technotypes. They respectively specify cell type concepts applied to a single species, multiple species, populations below the species level, and particular experiments. The flexible and rigorous framework we propose can base annotation of single-cell omics datasets, and reconcile knowledge about cells across all different domains of science.

# Introduction

One of the basic subjects in any undergraduate major in life sciences is histology. The students are required to identify cell types across various tissues and look for color and shape patterns in hematoxylin-eosin stains. Textbooks, like Junqueira's Basic Histology [1] work as manuals that perpetuate the paradigms (in the Kuhnian sense) [2] of what we know about a few hundred cell types.

Our concept of "cell type" is still based on centuries-old histochemical techniques, such as the Golgi-stains of neurons immortalized by Ramon y Cajal [3]. The histological influence is noticeable even in the names given to cell types, such as "erythrocytes", "eosinophils", "basophils", and "oxyphilic cells of the thyroid". The concepts we use are drawn from studies of microanatomy. This connection with anatomy leads us to think about cell types as anatomical entities as if they are dissectible and fixed in an organism. The limits of resolution perpetuated by the histological-anatomical view may be why attempts to quantify cell types use the scale of "hundreds" of human cell types [4,5].

New techniques have challenged this anatomy based conceptualization. From flow cytometry to patch clamping, to single-cell RNA-seq, we saw a burst of new categories, and novel cell "subtypes" and "families" popped up in the literature. The bursting intensified in the past few years, with the rise of projects to characterize *all* human cell types, like the Human Cell Atlas and HUBMAP [4,5].

The advances in biology require us to find better answers for how to define a cell type. Such a concept might not even have a "true" meaning, in a philosophical-realistic sense. Nevertheless, we can strive to find nominal, pragmatic definitions for the real challenges of large-scale biology. Otherwise, how can we precisely label single-cell data? How can we formalize the discovery of new cell types? How can we integrate the knowledge from millions of published scientific articles?

The need for a conceptual advance is being perceived by the community, and new perspectives are rising [6,7,8,9,10,11,12,13,14,15,16]. In an opinion article published in Cell Systems in 2017, a series of researchers presented their views on the conceptual definition of 'cell type' in the context of a mature organism [6]. Many of the scientists believed that cell functions have a core role in defining cell types, which is a slippery road, as the very meaning of "function" in biology is elusive [17]. The opinions were varied, and no consensus was achieved.

One core line of thought is based on the cell type as an evolutionary unit defined by a Core Regulatory Complex (CoRC) of transcription factors. That definition enables the drawing of parallels, from the evolution of other biological entities (such as genes, proteins, and species) to the evolution of cell types. Models of how multicellular life works greatly benefit from concepts such as "sister types" (cell types that diverged from a single ancestor), "cell type homology" (cell types in different species that share a common evolutionary origin), and "cell type convergence" (cell types that execute similar functions but which are not directly evolutionarily related) [18,19]

However, as much as different concepts of species coexist [20], our quest to define cell types may take various forms. The challenge of representing cell types in the context of evolution is conceptually different from the challenge of representing cell types in biomedical experimentation. In that second direction, the groundwork of the Cell Ontology [21,22,23] and the contributions of the International Workshop on Cells in Experimental Life Sciences series [24,25] are notable. Their contributions base much of the views here and will be discussed in detail throughout the article.

We chose to use the term "cell type" to emphasize the focus on types as classes (or "kinds") in contrast to real-world objects. The similar term "cell state" is used both to describe classes (e.g. activated T-cell) and real-world observations (e.g. the current state of a particular cell). Other similar notions, such as a

“cell set”, “cell population” and “cell cluster” can also reminisce of a specific, countable group of cells, frequently from the same experiment.

The term “cell class” is also used in the literature, and is a suitable synonym for our notion of cell type, as the main goal here is to refine the human-based theoretical classes. Classes that we can instantiate, i.e. assign to an observation of any real cell, in the same way we assign the class *Homo sapiens* to each and every human. The term “cell identity” has also been suggested for avoiding the cell type/cell state dilemma [26], but the notion of identity is slightly different from the idea of class. We opted to frame our work around the term “cell type” due to its historical usage and familiarity for the life sciences community.

The conceptual quest addressed by this work is one of research synthesis and is summarized in the following question: Which cell type definition can be crafted for rigorously describing biomedical experiments?

Towards that goal, the body of the article is divided into 4 parts. In Part 1, we propose a set of rules that are necessary and sufficient for defining cell types. Part 2 offers a small set of names for differentiating the main classes of cell types. In Part 3, we address the logical consequences of the proposed definitions, while Part 4 is a discussion of the pragmatic challenges envisaged in employing such definitions.

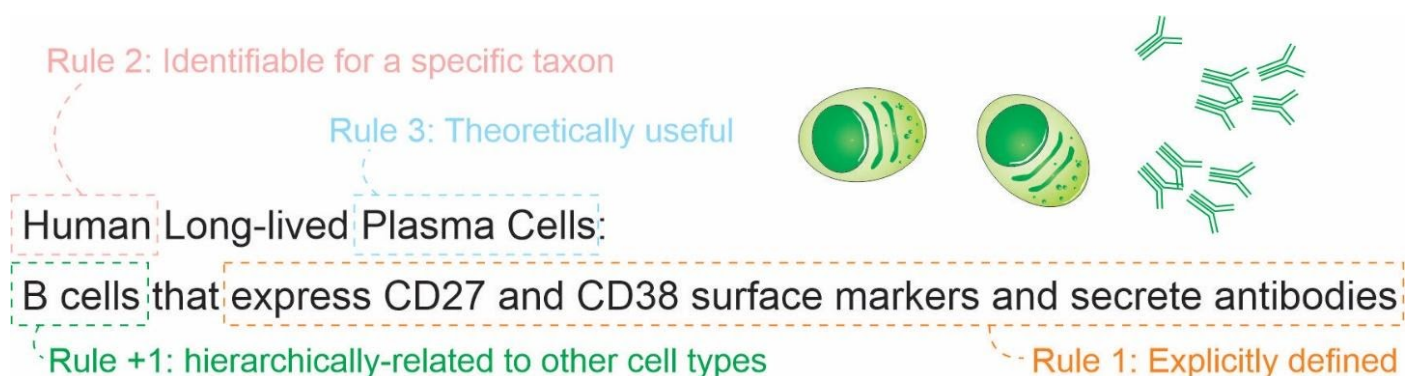
# 1. A set of 3 + 1 rules for defining a cell type

Our pragmatic definition of cell type (for eukaryotic, multicellular organisms) consists of 3 + 1 simple rules (Figure 1). A cell type is a class of cells that must be:

1. Explicitly defined
2. Identifiable for a defined taxon
3. Theoretically useful

And that should be:

4. Hierarchically related to other cell types



**Figure 1:** The set of 3 + 1 rules for defining a cell type.

Here, “must” represents an absolute requirement, whereas “should” suggests that “there may exist valid reasons in particular circumstances to ignore a particular item” (as per RFC 2119 [27]).

For rule 1, we mean that the cell type needs to be followed by a clear definition that would allow rational judgments of whether a singular cell belongs to the type or not. Such definitions should be complete, [28] and provide necessary and sufficient criteria for classification. An example is a cell type

defined by “expression of the proteins CD3 and CD4, but lacking CD8.” Even though there is still some ambiguity (see [29,30] for longer discussions), it already states clear and reasonable criteria. The degree of rigorousness cannot be decided a priori, as we still do not have a rigorous framework for representing biological knowledge, but we should strive to make definitions as rigorous as possible. Other examples of what could be explicit definitions are as follows:

- “Big cell” is a class of cells that have a length of more than 50 micrometers on any axis.
- “Human cortical neuron” is a class of cells in human cortex that are capable of producing an action potential.
- “Leukocyte” is a class of cells found in animal blood which are achromatic cells.

The recognition of multiple valid characteristics to define types is not new. The first Cell Ontology article, in 2005, explicitly acknowledged criteria based on function, histology, lineage, and ploidy. [21] These features were combined in the definitions of “species-neutral” cell types, arguably useful for integrating databases or for teaching biology. [23] Gradually, we are acknowledging that we might need more specific classes to characterize experimental biology, leading to the definition of species-specific types defined by granular characteristics. [22,31]

Rule 2 is an explicit criterion that must be followed while discussing cell types scientifically; we need to define the taxa for which a given cell type is expected to manifest. The scope is not only a taxonomic constraint (in the sense used in the Gene Ontology [32]); it states that the cell type needs to be discoverable in any individual of the taxon (or taxa) of interest, given the appropriate conditions (e.g., stage of life and biological sex). The set of taxa covered by a cell type is called here a taxonomic scope (or just scope) of the cell type. Note that, as cell types can be defined by function and functions can converge, the taxonomic scope is not restricted to monophyletic taxa (clades). The definition of taxon used here is liberal and applies to any class of organisms that any researcher identifies explicitly as a unit.

Knowing the scope is important to avoid the pitfalls of extrapolation. A recurrent theme is that theories corroborated by mouse experiments are valid for human cell types. Such extrapolation is an instance of the classic problem of induction, which is discussed thoroughly in “The Logic of Scientific Discovery”. [33] The taxonomic scope allows us, researchers, to be clear regarding our claims, and better discern what we claim to be true for a strain, a species or any other class of organisms.

Rule 3, regarding usefulness, deals with a practical concern. Rigorously, there is an infinite number of explicit definitions that any scientist might come up with. One simple proof of this infinitude is that size-based cell definitions (as for “big cell” above) may alone consider any of the infinite real numbers. Thus, a cell type “bigger than 7.835 micrometers” might fit the first two rules, but will likely fail rule 3. If we, as a research community, want to characterize *all* human cell types, we may need to have a finite number of cell types. Rule 3 could be paraphrased as: a valid cell type is a class of cells that any researcher rationally finds useful for a theoretical perspective of reality. For example, a recent study used single-cell RNA-seq experiments to assign 275,000 *Drosophila* cells into 200 cell types. [34] Since these 200 cell types were useful for Özel and colleagues when describing the world, they automatically satisfied rule 3.

Rule 4 is a practical extension of the usefulness rule: a cell type has to be hierarchically-related to other cell types for increased usefulness. This means that a definition of a cell class is (for research synthesis concerns) less useful if it cannot be considered a “subclass” of another cell type. For practical concerns, all imaginable mammalian cell types are subclasses of a “eukaryotic cell” (defined as any cell of an eukaryotic organism) and likely can be subclasses of more specific cell types. The rule 4 is presented as a recommendation instead of a requirement as, in practice, it might be an overhead and not strictly necessary for tasks like claiming the discovery of a new cell type.

Ontological organization is important for integrating knowledge across studies. A cell type that is based on its transcriptome is not the same as one based on its electrophysiology. They can, nevertheless, be connected by a superclass that matches either one or the other criterion. For example, the green-OFF bipolar cells of the retina and the Syt2<sup>-</sup>/NK3R<sup>+</sup> cells of the retina are considered to be the same cell type. [35]

However, as these features are often measured separately, we have, in fact, two individual classes for which knowledge is produced. These classes, then, can be combined in the superclass “(green-OFF) OR (Syt2<sup>-</sup>/NK3R<sup>+</sup>) cells” for the integration of claims across domains. Practically, when describing a cell type, one should make an effort to insert it into the universe of interrelated cell types, even if that implies creating new superclasses.

The consequences of this set of criteria will be discussed further in the sections 3 and 4.

## 2. Naming classes of cell types

To facilitate communication among life scientists, we propose a set of naming conventions for different classes of cell types. Much of the literature mixes cell types in one species (e.g., when dealing with a cell type as an evolutionary unit) or in multiple species (e.g., in the Cell Ontology). It is useful to distill these different concepts into names. Given the importance of the concept of species in biological classification [36], we derive a species-centric view on the naming of classes of cell types. The four classes (Figure 2) we propose are as follows:

- archetypes, for when the taxonomic scope of the type is beyond the level of species; for example, “mammal neutrophils.”
- *sensu stricto* cell types, for when the taxonomic scope of the type corresponds to a single species; for example, *Mus musculus* neutrophils.”
- infratypes, for when the taxonomic scope is below the level of species; for example, considering the mouse strain “C57BL/6J”, “neutrophils from C57BL/6J mice”.
- technotypes, for specific, experimentally defined cell types that harbor in their definition the precise conditions of the cells sampled; “2-month-old male C57BL/6J, Ly-6G<sup>+</sup> CD11b<sup>+</sup> M-CSF R<sup>-</sup> CD244<sup>-</sup> neutrophils”.

archetypes



mammal  
neutrophils

*sensu stricto* cell types



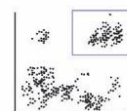
mouse  
neutrophils

infratypes



neutrophils from  
C57BL/6J mice

technotypes



2-month-old  
male C57BL/6J,  
Ly-6G<sup>+</sup> CD11b<sup>+</sup>  
M-CSF R<sup>-</sup> CD244<sup>-</sup>  
neutrophils

**Figure 2:** Names for classes of cell types.

By adopting a precise vocabulary, we can avoid misunderstandings and communicate more clearly. At the level of individual scientific experiments, scientists rarely reach the *sensu stricto* cell type level; the samples come only from a subpopulation of the species of interest and cannot be assumed to be randomly sampled from all individuals of the species. This has important practical considerations to, once again, avoid failing implicitly at the problem of induction.

Besides, in individual experiments, we work with cells of very specific classes. They are not only infratypes but very specific infratypes defined by non-random research setups and pragmatic choices. For example, we might call “CD4 T cells” what are CD3<sup>+</sup>, CD4<sup>+</sup>, CD8<sup>-</sup> cells from the axillary lymph node of 2-month-old chow-fed female C57BL6/J mice from the mouse-house of the Institute of Biochemistry of the University of São Paulo collected on several mornings around 10 pm. Although quite specific, all the mentioned facets (markers, anatomical location, age, diet, biological sex, strain, housing conditions and circadian clock) are known to alter what we know about cell types (see Tabansky et al [13] for a longer discussion ). Thus, we benefit from using a name for the experimentally-constrained cell classes: technotypes.

Even if it is specific, a technotype is still a class. Unless a study used only one single-cell, it likely contained some sampling method. Samples are from a specific population for which hypotheses are tested. This is the most granular cell type, in our considered view, for research synthesis. This is the type that can be strictly annotated in single-cell RNA-seq datasets, for example.

Single claims are made and tested for technotypes, and the claims can be logically combined in “upper” ontological levels for reaching a higher degree of universality. The propagation of knowledge to upper levels cannot be implicit (see Yarkoni 2020 for an analogous problem in the psychological sciences [37]). As Popper defends, knowledge should travel “quasi-inductionally” by fostering hypotheses with higher degrees of generality, which can then be tested for the more universal class. [33]

### 3. Logical consequences of the definition of a cell type

One notable logical consequence of the proposed set of criteria is that the definition of a cell state is left as a subclass for cell type. For the pragmatic purpose adopted here, we avoid the dissection of the differences between persistent classes of cells (often called “cell types”) or the transient, fugacious classes of cells (often called “cell states”) (see “Definition of cell identity” section in [38] for an example). We also consider only the cell as it was observed in an experiment, not necessarily the future conditions of any cell (i.e. the “cell fate”). [39]

Even though such a distinction is an important topic for theoretical research, it is not a requirement for representing biomedical experiments.

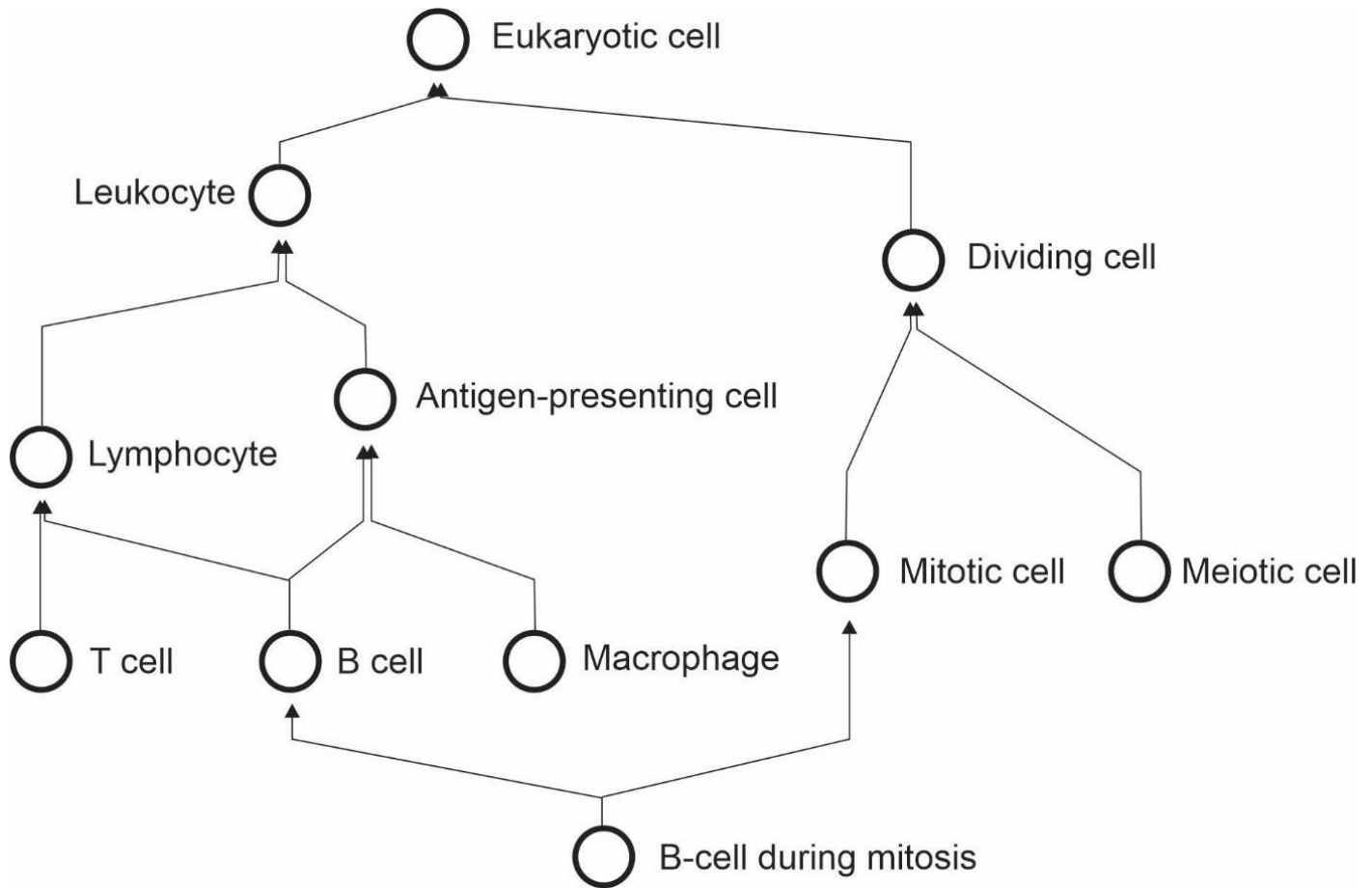
One example of this entailment is that the class “human cells in metaphase of mitosis” can be considered a cell type, as it can be explicitly defined and restricted to a taxon. Even though “metaphase” itself is a biological process, we can describe all cells executing this process as a singular cell class.

However, does a dividing fibroblast stop being a fibroblast, even if temporarily? Again, we do not aim to answer this in a philosophical-ontological sense. Pragmatically, if the explicit definition used for fibroblast (e.g., expression of a marker) still holds during duplication, this cell can be assigned to two classes that are not hierarchically related: “fibroblasts” and “doubling cells”. If cells can be assigned to multiple classes that are not hierarchically related, it is not possible to annotate cell types with a single identifier using a taxonomic tree, in which each concept is represented by a single node with one (and



only one) direct parent node. This is in conflict with attempts to classify cell-types using single hierarchies in the form of a tree [40] [41] [42].

Cell types need to be represented ontologically with multiple inheritance, which can be thought of as multiple, intertwining trees that take into account different ways of classifying cells (Figure 3).



**Figure 3:** The cell type hierarchy is not a tree - it requires multiple inheritance for completeness.

Another logical consequence of the definition is that the concept of subtype becomes redundant with the concept of cell type. The notion of subtype, then, only makes sense when discussing classes with different degrees of universality. Thus, claims to discovery of new cell “subtypes” or “types” differ only stylistically and can be considered indistinguishable from the perspective of research synthesis.

## 4. Practical consequences of the definition of a cell type

In the previous section, we discussed the logical entailments of accepting the proposed rules as valid. Here, we extend the pragmatic considerations on using such a system for real-world applications. In a recent attempt to define cell types for single-cell RNA-Seq, Aevermann et al came up with a set of needs: “The minimum set of necessary and sufficient marker genes selectively expressed by the cell type”, “A parent cell class in the CL (Cell Ontology)”, and “A specimen source description (anatomic structure þ species).” [43] Their approach has great merit in defining clear guidelines for marking a cell type. The requirement of markers is reasonable for the field of single-cell RNA-seq, where marker information is abundant. The Cell Ontology has used markers for defining cell types, an approach employed in particular for immune cells [22,29,30].

The use of markers, however, leaves us with a conceptual problem – definitions of cell type used by electrophysiologists, or even in the manuals of histology courses, are not based on markers.

Rigorously adopted, this requirement would leave aside an entire segment of what we consider biomedical knowledge. Moreover, gene markers are not defined for cell types that span multiple species, a problem already discussed in the Cell Ontology report of 2011 [22]. Thus, our set of rules was crafted to accommodate the different ways that people classify cells.

In fact, with so many different takes on the field, vast amounts of data, and loose definitions of cell type, it becomes uncannily easy to claim a new cell type. Our set of rules may contribute to formalizing cell type discovery.

If one explicitly claims to have discovered a new *sensu stricto* cell type, one should provide enough evidence that cells from this class are identifiable across all individuals of a species (given the constraints as age and biological sex). A claim of an archetype would require evidence of existence in more than one species. Consequently, experiments that only use a specific strain of mice have a more robust claim if the expectation is limited to the infratype.

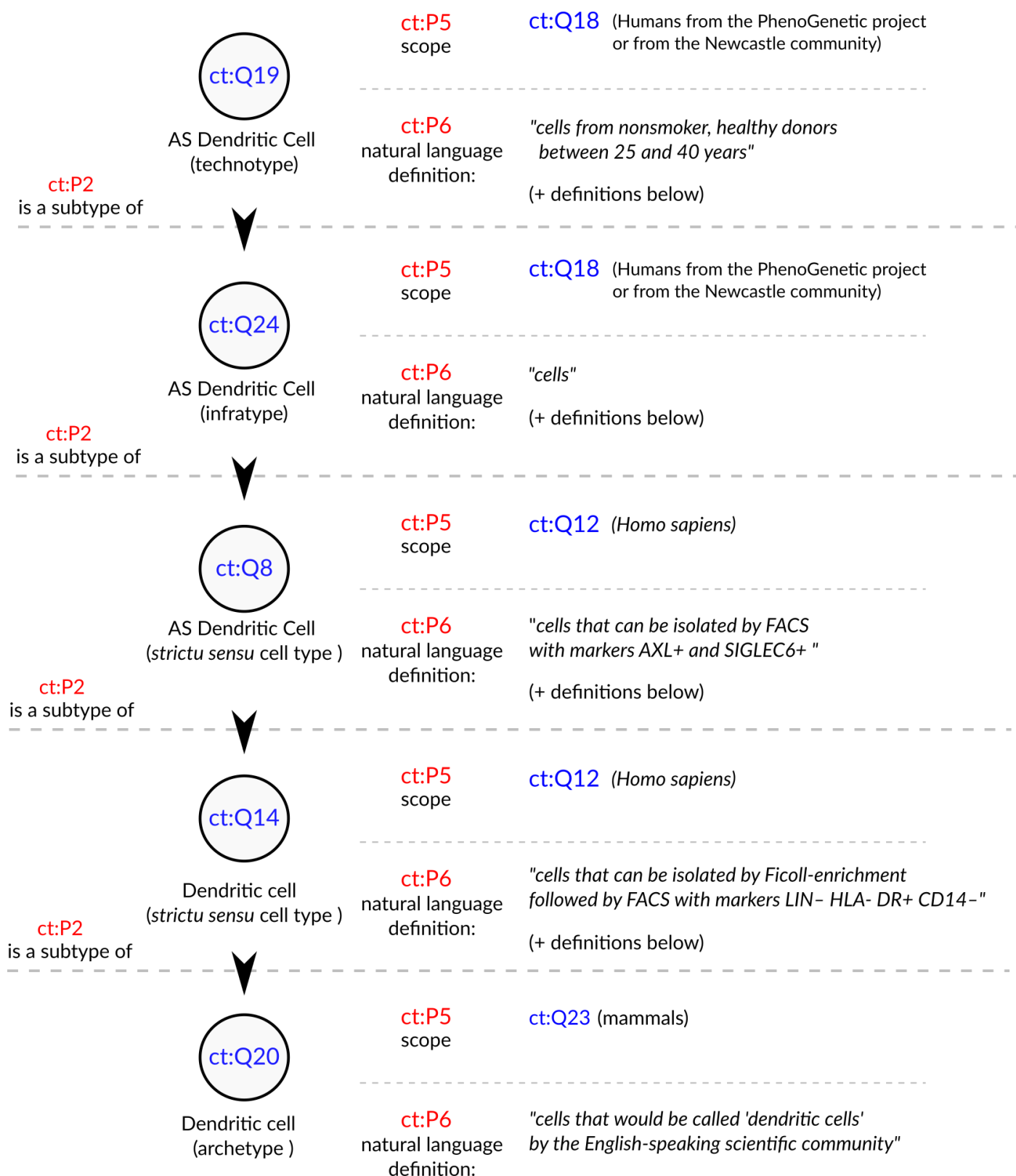
An example of the discovery of a new archetype is the pair of articles published in Nature in 2018 [44,45] about the newly found “ionocyte”, a class of cells in the trachea enriched for the expression of genes homologous to the human *CFTR* gene. Both studies displayed evidence for such a class in both mouse and human samples, corroborating the existence of an archetype. This discovery of an archetype has been denominated by both articles as a discovery of a new cell type.

Another example of cell type discovery is found in a pioneering article by Villani et al [46]. The authors describe subclasses of monocytes and dendritic cells in humans and pragmatically use markers for their definition. The patients were recruited from “the Boston-based PhenoGenetic project (...) and the Newcastle community.” Arguably, they did not have a random sample of humanity, and the observed results might not hold for different populations. This discovery of infratypes has also been described as the discovery of a new cell type.

An example from the Villani et al article is the discovery of the “AXL+ SIGLEC6+ AS Dendritic cell”. This and other cell types are presented in the article as part of a “Human dendritic cell atlas”, generalizing the theory for the whole of humanity. The jump from technotype (which takes into consideration also descriptors like “healthy” and “age between 25 and 40 years”) to infratype (“all humans in this population scope”) to cell type *sensu stricto* (all humans) is depicted in Figure 4 and exemplifies the logical flow in our proposed framework.

Of note, “dendritic cells” are one of the cell types most thoroughly modeled by the Cell Ontology. [47,48] The current definition of the dendritic cell (CL 0000451) is coupled to the definition of leukocyte (CL 0000738), which defines it as “An achromatic cell of the myeloid or lymphoid lineages capable of ameboid movement”. This definition is not reconcilable with the “dendritic cells” studied by Villani et al. We have no way of knowing if the cells in their work are achromatic or capable of ameboid movement. That might sound pedantic and might, unfortunately, be so, but the logical requirements of computational systems lead to both biocurators and computers being seen as pedantic. This high level of precision is necessary to accurately depict not only the complexities of cell types but also of research settings.





**Figure 4:** Conceptualization of a set of the cell types in Villani et al, 2017 [46]. The depicted cell types were manually curated from the article, where they are either implicitly or explicitly mentioned. Identifiers for cell types are written in pseudocode based on the Turtle serialization for RDF (<https://www.w3.org/TR/turtle/>) and represent valid URIs (described in the database [https://celltypes.wiki.opencura.com/wiki/Main\\_Page](https://celltypes.wiki.opencura.com/wiki/Main_Page)). URI: Universal Resource Identifier; RDF: Resource Description Framework; ct: <http://celltypes.wiki.opencura.com/entity/> .

Even if we are not able to represent all the aspects that go into a cell type definition using ontologies, we can use an explicit “natural language definition” property to define cell types. As David Osumi-Sutherland puts in his 2017 article about cell type classification: there is a “*mismatch between quantified logic, which records assertions about all members of a class, and the messy, noisy reality of biology and the data we collect about it.*” [31]. Luckily, we do not need to have all the biology formalized before we deal with cell types. Taking the example in Figure 4, all cell types treated as

“dendritic cells” in the literature are valid subclasses of the dendritic cell archetype ([ct:Q20](#)). To reach a middleground between natural languages and complete axiomatization, we can use less imprecise Controlled Natural Languages [[49](#)] for better clarity. Such system might still lack the power for full computational reasoning, but it could already provide a coherent scaffold for representing experimental data (e.g., from single-cell transcriptomics) and allow logically robust data integration. Such a restricted system is more feasible than an “ontology of everything” and would suit the creation of ontologies in accordance to the principle of minimal ontological commitment [[28](#)].

The commitment to logical coherence will require us to deal with many more types than we are used to. Given the variety of species on Earth, the complexity of multicellular life, and the diversity of research settings, a count of cell types may far exceed the mark of one million. Sabina Leonelli stated that the challenges thrown up by big data in biology require the advancement of our philosophical theories [[50](#)]. We agree and argue that the converse is also true: to advance the theoretical foundations of modern biology, we need to harness the power of computational tools. Computational ontologies provide a solution for dealing with complex concepts. Classes in ontologies can have alpha-numeric identifiers. We can, thus, assign each technotype a Unique Resource Identifier, a URI, similar to the Cell Ontology (CL)[[21,22,23](#)] or in the knowledge graph of Wikidata [[51](#)]. The power of using knowledge graphs for integrating knowledge about cell types is gaining momentum [[11](#)], and they rely heavily on the precise usage of unique identifiers.

The classification of cells into types and the naming of cell types are parallel tasks. While there has been progress on rules for naming cell types (particularly in neuroscience [[52,53,54](#)]), nomenclature is outside the scope of this article. Using identifiers/URIs without semantic sense already suffices for our purposes. Semantically void identifiers also help us to steer away from the Aristotelian essentialist view upon cell types, as discussed by Rowe and Stone in 1977 [[55](#)]. Identifiers can have labels that can be freely changed, while keeping a persistent URI. Our effort to refine the logical aspects of cell-type definitions can be combined with any commonly agreed naming/labeling system.

The URIs at the level of technotype allow precise labeling of cell types in real-world experiments. The technotype annotation empowers researchers to craft their cell type of interest, and connect this cell type to a common network of knowledge. Several single-cell transcriptomics tools try to assign labels to cells. While some approaches avoid ontologies [[40,56](#)], others utilize the Cell Ontology [[57,58,59](#)] or MeSH IDs [[58,60](#)] to identify the most likely cell type label for each cell or cell cluster. Different studies, however, almost always study different technotypes. Thus, the task of finding the exact type of cells in a given experiment, algorithms could try to find where the new technotypes should be inserted in an ontological network. For example, instead of claiming that cells from study A and study B are myeloid dendritic cells, we can claim that both cell types belong to the myeloid dendritic cell branch. By embracing these real differences between studies (and cells), the precise metadata of the study will enable a precise statement of the cell type. This will ultimately allow the coherent reuse of publicly available data.

This flexible, yet rigorous, framework for defining cell types can help us to deal with the challenges of varying resolution levels of interest and the scaling large datasets. [[doi:10.1186/s13059-020-1926-6](#)] The need of an identification routine for cell-type taxonomies is acknowledge for more than 45 years [[61](#)], and still is a core challenge of human cell-type atlases [[doi:10.1186/s13059-020-1926-6](#)]. The quest for data-driven cell classification is at least as old [[61](#)]. The framework here proposed provides ideas to improve our solutions to both tasks.

## Final remarks

In this article, we have proposed a set of three rules (explicit definition, taxon scope restriction, and theoretical usefulness) and one recommendation (hierarchical linking) to be followed when defining

cell types. We have also proposed four types of naming to clarify discussions on the topic: archetypes (a class with a scope above species level), *sensu stricto* cell types (a class with scope equal to one species), infratypes (a class with scope below the species level) and technotypes (the exact cell type defined for an experimental setup). The concept of the “technotype” can be harnessed as the unit for classifying cells, in a manner analogous to how the “species” is the conventional unit for classifying organisms into higher-order taxa. We have dissected some logical entailments of such definition, which admittedly might conflict with current views on defining cell types. We do not aim to solve such conflicts or negate the other perspectives but only to propose a unique way of organizing our knowledge on cell types. This article clarifies some of the meanings and provides directions for the future development of the theoretical basis of a cell type definition. The discussion on cell types’ definition is still in its infancy, and we need human power to tackle these huge theoretical challenges. Biologists, philosophers, and computer scientists ought to distill the details of defining cell types, powering the Human Cell Atlas, and the life sciences research enterprise of this century.

## Acknowledgments

We would like to express our gratitude and acknowledge the researches that dedicated time specifically to help us to discuss and refine the basis for developing concepts here presented. Nominally, we thank Kleber Neves, Gabriel Lovate, Cesar Prada, Diógenes Saulo Lima, Lucas Cardozo, Juliane Fernandes, Pedro Medeiros, Érika Molina, Antonio Pedro Vieira, João Vitor Cavalcante, Maria Fernanda Forni, Diorge Souza, Jean Bezerra, Gabriel Sato, Roberta Andrejew and Dimitrius Pramio. Part of this work was supported by grant [#2019/26284-1, São Paulo Research Foundation \(FAPESP\)](#).

## Additional thoughts

A set of additional thoughts on the pertaining matters, but which do not fit the article (neither as part of the body or the supplements).

This part was not and will not be submitted with this article in its final form or to preprint servers.

## What to do when two researchers disagree on a definition?

---

Cell type names are notoriously ambiguous and one definition might collide with an other, specially regarding the natural language name used to described. There are many different, equally valid definitions of a “dendritic cell.” we do not aim to solve this problem from a societal standpoint. However, from a computational-ontology standpoint, there is one simple solution: split the concept.

This approach is similar to King Solomon’s solution in a famous bible story, called the [Judgement of Solomon](#). In a dispute between two women that claimed to be the mothers of a child, the solution of the king was simple: split the baby. However, babies are notoriously indivisible, and the true mother did not really like the idea.

It may be that some scientists are attached to their names, as mothers are to their babies. However, unlike babies, namings can be divided. Each of the scientists gets to name their specific conceptualization however they choose. Many names might “collide” in that way, and that is okay. Under the hood, however, the names refer to different identifiers. Computationally there would be no ambiguity. Then, it is just a matter of the researcher to respect the choice of their peers of calling something by the *wrong* name, as long as the identifier is correct.

Splitting concepts upon conflicts in the end is more the multiplication of bread and fish in the [Feeding the multitude](#) episode, and everyone gets to eat.

But ontologies are different from ordinary babies and magical fish. The splitting of concepts would not only create new concepts, but leave a trace. They would be immediate subclasses of their conjunction. An equally valid superclass that can be defined by “a cell containing characteristics of any of their subclasses”.

In a parallel with text-book mitosis, the concept gets divided in two new, equally real concepts. And as we can trace cells in an animal to a single zygote, we can keep track of concepts while they keep dividing, whenever a new conflict pops up.

## The big assumption of continuity in time

---

One assumption that underlies the validity of the models proposed here is that taxons preserve their characteristics throughout time.

In Popper’s Logic of Scientific Research, he states that he has a metaphysical faith on the continuity of laws of nature through time.

We have no way of testing this metaphysical faith, and it is absolutely necessary for the scientific endeavour as we understand.

While in physics this assumption seems to be reasonable, evolution makes biology quite more complicated. Statements that we have about the human species, for example, might be valid today, but were not valid 2000 years ago, and vice-versa.

When we talk about sub-species taxons, which might be a local population of a town, for example, this unit is not immutable. The population of Newcastle, as per the example, might change in time, with immigration and emigration, mutation, natural selection, neutral evolution and the many forms of modifications of a gene pool.

It is, thus, and heuristic, to call “population of Newcastle” a class. We could specify a period in time for which we expect the information to be valid. For example, we may say we are sampling from “the Newcastle population in the years 2019-2020.” This would be a valid statement, but it would not be falsifiable, as by 1st of January 2021, no independent tests of the theory can be done.

It is technically possible to have a technotype so precise as to have a scope with a time constraint. In fact, that might be the right way of representing information, if we want to compare experiments done in evolving populations.

While evolutionary definitions take this dimension into account, they are fit to theoretical research, but still lack the rigour for explaining real world experiments.

All research that uses human samples are subject to strong influence of time.

Thus, the explicit assumption here is that taxons are consistent in time. And, for what we know, it is blatantly false.

This is a great flaw of the model, and, maybe, of a great part of biomedical research. The logical consequences are so dire, that it merits a separate, dedicated work.

# Clusters are not cells

---

In the era of large-scale omics, we are starting to see declarations of cell types that are not based on pre-selected criteria, but derived from unsupervised clustering followed by labelling.

This is a powerful exploratory approach, which, as mentioned in the main text, has led to discoveries of ionocytes and new classes of dendritic cells, for example.

Many algorithms and “expert-based” annotation protocols focus on labeling *clusters* instead of labeling *cells*.

Cells in a cluster are arbitrarily similar (as determined by the clustering algorithm) and so they will, by definition, differ from other cells in the sample.

For single-cell RNA-seq, one usually checks which genes are differentially expressed when comparing the cells in a cluster with cells in other clusters. These genes are called “markers” and used for labeling a cell cluster.

What does it mean to label a cell cluster, though? Does it mean that *all* cells there conform to the cell type? Does it mean that *most* cells there conform to the cell type? Does it mean that cells from other clusters in the dataset *definitely do not* conform to the cell type?

So far, we haven’t seen a clear, explicit, coherent definition for a cluster label. Not even once.

Marker-based definitions are assumed for the group as a whole, but in current pipelines, nothing blocks one cell in a cluster to lack the expression of a “name-giver” marker.

The classification scheme proposed here works to classify *cells*, but is not sufficient for labeling unsupervisedly-defined *cell clusters*.

What is possible, though, is to use clustering for data exploration. From then on, strict patterns can be decided (ex: a cell that expresses A and B, but not C) and then apply this pattern to the whole dataset. For clusters with consistent markers, this approach should be roughly equivalent to the previously described.

Using such “regular expressions” might lead to a cell being assigned to multiple clusters. Even if we assume that the sample is free of doublets, that cannot be a problem. Cells may have multiple functions. As argued in the main text, each cell can be labeled by multiple standards.

We may avoid multiple labeling if we really need in practice, though, and make preferential claims (if a cell matches definitions X and Y, it is assigned only to X, for example).

By having explicit “regular expression” patterns for cell definitions in single-cell datasets, the “cell-type assessment” problem becomes trivial: a cell in a new dataset is of the *exact* same type if (and only if) it matches the *exact* definition.

When that is not the case, current algorithms for reconciling single cell datasets can still be successfully employed. But instead of propagating a label, it would propagate a parent class, looking for cells of a similar, sister class.

## What this work is not

---

This is not an attempt to substitute the Cell Ontology (CL) or contradict it in any way. CL is an amazing resource, built by a community of wonderful researchers. Its relation with CL is coexistential, and topics discussed here might be or might not be of interest of CL, and that is OK.

This is not an attempt to create an ontology itself, or a system that allows reasoning. It is a set of suggestions that can be taken into consideration for building a coherent ontology. The [Cell Type Wikibase](#) is an experimental ontology, and far from ready for professional use.

This is not an attempt to have a one-size-fits-all definition of cell type. It is built as a theoretical solution for one cell-type related task. Similar to species definitions, we need an ecosystem of cell-type and cell state definitions that better suit different areas.

This is not an attempt to claim anything about the “true” nature of cell types, in the biological sense. It is a proposal of practical guidelines to represent research data.

This is not an attempt to solve *all* problems for cell type data annotation. It is the introduction of alternatives that need to be further developed and discussed.

## Why are we doing this?

---

A series of examples of why a pragmatic definition is useful ### Immune Epithope DB

The Immune Epitope Database and Analysis Resource (IEDB) announced it in 2006 [doi:10.1371/journal.pcbi.0020125]:

“the goal of the IEDB is to present as much information as possible without subjective interpretation, we can never presume any information, but rather we must try to capture the data exactly as presented in the reference, while maintaining the conclusions of the reference in a uniform manner. For example, if all experiments are performed with a whole cell population, but the authors attribute the response to a particular cell type without any evidence, we must capture the effector cells as the entire population.”

That is the kind of challenge that the “technotype” solves in theory, as it gets the objective population sampled in any article.

## Nanopublications

Nanopublications [62] are semantically-rich publications of single assertions. If we want to have computable assertions, it is a good idea to have pragmatic definitions.

This is the case for having better representations of primary datasets representation of primary data and datasets (as explained in [63] and exemplified in [64]).

Some argue (eloquently) for immutable, citable datasets, which could be shared via trustworthy URIs, for example, which contain hash values in the URI that prevent silent changes [65]. It is of our interest that these immutable, core datasets also use very solid concepts when representing information. Ideally, not only using unique identifiers [66] but pointing to concepts backed up by a solid theory. As Thomas Gruber argues, “an ontology is only a specification, and the utility of an ontology ultimately depends on the utility of the theory it represents.” [28]

## Phenetic Species Concept

---



The Phenetic Species Concept states that species are the smallest groups that are consistently and persistently distinct and distinguishable by ordinary means. Maybe that is crosslinkable.

The link between the phenetic species concept and the classification of cell types has been noted before. [67] The ideas of explicit definitions and hierarchical relations resound with the goals of this article.

In this work[67], they mention 5 facets of phenetic systematics: - Group (that is, 'type') membership should be based on multiple criteria rather than on a single so-called 'essential' feature that the investigator favours. - The criteria for group membership should be rule-based, explicit and quantitative. - Groupings should be hierarchical rather than flat to acknowledge the validity of both coarse and fine divisions. - Groupings generated by this approach should be viewed as hypotheses to be tested rather than inflexible rules. - Classification should focus on discontinuities between groups and ignore parameters that vary continuously.

## Calculation of theoretical maxima and minima of human cell types

---

<https://twitter.com/lubianat/status/1295923945770823682>

## References

---

### 1. Junqueiras Basic Histology: Text And Atlas

Luiz Carlos Uchôa Junqueira, Anthony L. Mescher

McGraw-Hill Education (2018-01-01) <https://www.wikidata.org/wiki/Q102104590>

ISBN: [9781260026184](https://www.wikidata.org/wiki/Q102104590)

### 2. The Structure of Scientific Revolutions

Thomas Kuhn

University of Chicago Press (1962-01-01) <https://www.wikidata.org/wiki/Q951060>

### 3. Histologie du système nerveux de l'homme & des vertébrés.

Santiago Ramón y Cajal

(1909-01-01) <https://www.wikidata.org/wiki/Q51488921>

### 4. The Human Cell Atlas

Aviv Regev, Sarah A Teichmann, Eric S Lander, Ido Amit, Christophe Benoist, Ewan Birney, Bernd Bodenmiller, Peter Campbell, Piero Carninci, Menna Clatworthy, ... Human Cell Atlas Meeting Participants

eLife (2017-12-05) <https://doi.org/gcnzcv>

DOI: [10.7554/elife.27041](https://doi.org/10.7554/elife.27041) · PMID: [29206104](https://pubmed.ncbi.nlm.nih.gov/29206104/) · PMCID: [PMC5762154](https://pubmed.ncbi.nlm.nih.gov/PMC5762154/)

### 5. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program

HuBMAP Consortium

Nature (2019-10-09) <https://doi.org/gf92ht>

DOI: [10.1038/s41586-019-1629-x](https://doi.org/10.1038/s41586-019-1629-x) · PMID: [31597973](https://pubmed.ncbi.nlm.nih.gov/31597973/) · PMCID: [PMC6800388](https://pubmed.ncbi.nlm.nih.gov/PMC6800388/)

### 6. What Is Your Conceptual Definition of "Cell Type" in the Context of a Mature Organism?

Paul Blainey, Hans Clevers, Cole Trapnell, Ed Lein, Emma Lundberg, Alfonso Martinez Arias, Joshua R. Sanes, Jay Shendure, James Eberwine, Junhyong Kim, ... Mathias Uhlén

*Cell systems* (2017-03-01) <https://www.wikidata.org/wiki/Q87649649>  
DOI: [10.1016/j.cels.2017.03.006](https://doi.org/10.1016/j.cels.2017.03.006)

## 7. A periodic table of cell types

Bo Xia, Itai Yanai

*Development* (2019-06-15) <https://doi.org/ggctwf>

DOI: [10.1242/dev.169854](https://doi.org/10.1242/dev.169854) · PMID: [31249003](https://pubmed.ncbi.nlm.nih.gov/31249003/) · PMCID: [PMC6602355](https://pubmed.ncbi.nlm.nih.gov/PMC6602355/)

## 8. Exciting times to study the identity and evolution of cell types

Maria Sachkova, Pawel Burkhardt

*Development* (2019-09-19) <https://doi.org/ghdb9v>

DOI: [10.1242/dev.178996](https://doi.org/10.1242/dev.178996) · PMID: [31537583](https://pubmed.ncbi.nlm.nih.gov/31537583/)

## 9. The Human Cell Atlas: from vision to reality.

Orit Rozenblatt-Rosen, Michael J. T. Stubbington, Aviv Regev, Sarah Teichmann

*Nature* (2017-10-01) <https://www.wikidata.org/wiki/Q47565008>

DOI: [10.1038/550451a](https://doi.org/10.1038/550451a)

## 10. Human Cell Atlas and cell-type authentication for regenerative medicine

Yulia Panina, Peter Karagiannis, Andreas Kurtz, Glyn N. Stacey, Wataru Fujibuchi

*Experimental and Molecular Medicine* (2020-09-15) <https://www.wikidata.org/wiki/Q99418657>

DOI: [10.1038/s12276-020-0421-1](https://doi.org/10.1038/s12276-020-0421-1)

## 11. A community-based transcriptomics classification and nomenclature of neocortical cell types

Rafael Yuste, Michael J. Hawrylycz, Nadia Aalling, Argel Aguilar-Valles, Detlev Arendt, Rubén

Armañanzas, Giorgio A. Ascoli, Concha Bielza, Vahid Bokharaie, Tobias B. Bergmann, ... Ed S. Lein

*Nature Neuroscience* (2020-08-24) <https://www.wikidata.org/wiki/Q98665291>

DOI: [10.1038/s41593-020-0685-8](https://doi.org/10.1038/s41593-020-0685-8)

## 12. The evolving concept of cell identity in the single cell era

Samantha A. Morris

*Development* (2019-06-27) <https://www.wikidata.org/wiki/Q93086971>

DOI: [10.1242/dev.169748](https://doi.org/10.1242/dev.169748)

## 13. Implications of Epigenetic Variability within a Cell Population for “Cell Type” Classification

Inna Tabansky, Joel Stern, Donald W. Pfaff

*Frontiers in Behavioral Neuroscience* (2015-12-16) <https://www.wikidata.org/wiki/Q26770736>

DOI: [10.3389/fnbeh.2015.00342](https://doi.org/10.3389/fnbeh.2015.00342)

## 14. Geometry of the Gene Expression Space of Individual Cells

Yael Korem, Pablo Szekely, Yuval Hart, Hila Sheftel, Jean Hausser, Avi Mayo, Michael E. Rothenberg, Tomer Kalisky, Uri Alon

*PLOS Computational Biology* (2015-07-10) <https://www.wikidata.org/wiki/Q35688096>

DOI: [10.1371/journal.pcbi.1004224](https://doi.org/10.1371/journal.pcbi.1004224)

## 15. Evolution of Cellular Differentiation: From Hypotheses to Models

Pedro Márquez-Zacarías, Rozenn M. Pineau, Marcella Gomez, Alan Veliz-Cuba, David Murrugarra, William C. Ratcliff, Karl J. Niklas

*Trends in Ecology & Evolution* (2020-08-20) <https://www.wikidata.org/wiki/Q98633613>

DOI: [10.1016/j.tree.2020.07.013](https://doi.org/10.1016/j.tree.2020.07.013)

**16. An era of single-cell genomics consortia**

Yoshinari Ando, Andrew Tae-Jun Kwon, Jay W. Shin

*Experimental and Molecular Medicine* (2020-09-15) <https://www.wikidata.org/wiki/Q99418649>

DOI: [10.1038/s12276-020-0409-x](https://doi.org/10.1038/s12276-020-0409-x)

**17. The meanings of “function” in biology and the problematic case of de novo gene emergence**

Diane Marie Keeling, Patricia Garza, Charisse Michelle Nartey, Anne-Ruxandra Carvunis

*eLife* (2019-11-01) <https://doi.org/ggjnmv>

DOI: [10.7554/elife.47014](https://doi.org/10.7554/elife.47014) · PMID: [31674305](https://pubmed.ncbi.nlm.nih.gov/31674305/) · PMCID: [PMC6824840](https://pubmed.ncbi.nlm.nih.gov/PMC6824840/)

**18. The evolution of cell types in animals: emerging principles from molecular studies.**

Detlev Arendt

*Nature reviews. Genetics* (2008-11) <https://www.ncbi.nlm.nih.gov/pubmed/18927580>

DOI: [10.1038/nrg2416](https://doi.org/10.1038/nrg2416) · PMID: [18927580](https://pubmed.ncbi.nlm.nih.gov/18927580/)

**19. The origin and evolution of cell types**

Detlev Arendt, Jacob M. Musser, Clare V. H. Baker, Aviv Bergman, Connie Cepko, Douglas H. Erwin, Mihaela Pavlicev, Gerhard Schlosser, Stefanie Widder, Manfred D. Laubichler, Günter P. Wagner

*Nature Reviews Genetics* (2016-11-07) <https://doi.org/f9b62x>

DOI: [10.1038/nrg.2016.127](https://doi.org/10.1038/nrg.2016.127) · PMID: [27818507](https://pubmed.ncbi.nlm.nih.gov/27818507/)

**20. Species Concepts and Species Delimitation**

Kevin De Queiroz

*Systematic Biology* (2007-12) <https://doi.org/c34kzf>

DOI: [10.1080/10635150701701083](https://doi.org/10.1080/10635150701701083) · PMID: [18027281](https://pubmed.ncbi.nlm.nih.gov/18027281/)

**21. An ontology for cell types**

Jonathan Bard, Seung Y. Rhee, Michael Ashburner

*Genome Biology* (2005-01-01) <https://www.wikidata.org/wiki/Q21184168>

DOI: [10.1186/gb-2005-6-2-r21](https://doi.org/10.1186/gb-2005-6-2-r21)

**22. Logical Development of the Cell Ontology**

Terrence F Meehan, Anna Maria Masci, Amina Abdulla, Lindsay G Cowell, Judith A Blake, Christopher J Mungall, Alexander D Diehl

*BMC Bioinformatics* (2011-01-05) <https://doi.org/c7kw6x>

DOI: [10.1186/1471-2105-12-6](https://doi.org/10.1186/1471-2105-12-6) · PMID: [21208450](https://pubmed.ncbi.nlm.nih.gov/21208450/) · PMCID: [PMC3024222](https://pubmed.ncbi.nlm.nih.gov/PMC3024222/)

**23. The Cell Ontology 2016: enhanced content, modularization, and ontology interoperability**

Alexander D. Diehl, Terrence F. Meehan, Yvonne M. Bradford, Matthew H. Brush, Wasila M. Dahdul, David S. Dougall, Yongqun He, David Osumi-Sutherland, Alan Ruttenberg, Sirarat Sarntivijai, ... Christopher J. Mungall

*Journal of Biomedical Semantics* (2016-07-04) <https://doi.org/gg99b9>

DOI: [10.1186/s13326-016-0088-7](https://doi.org/10.1186/s13326-016-0088-7) · PMID: [27377652](https://pubmed.ncbi.nlm.nih.gov/27377652/) · PMCID: [PMC4932724](https://pubmed.ncbi.nlm.nih.gov/PMC4932724/)

**24. Cells in experimental life sciences - challenges and solution to the rapid evolution of knowledge**

Sirarat Sarntivijai, Alexander D. Diehl, Yongqun He

*BMC Bioinformatics* (2017-12-21) <https://doi.org/gg99b7>

DOI: [10.1186/s12859-017-1976-2](https://doi.org/10.1186/s12859-017-1976-2) · PMID: [29322916](https://pubmed.ncbi.nlm.nih.gov/29322916/) · PMCID: [PMC5763506](https://pubmed.ncbi.nlm.nih.gov/PMC5763506/)

**25. Cells in Experimental Life Sciences (CELLS-2018): capturing the knowledge of normal and diseased cells with ontologies**

Sirarat Sarntivijai, Yongqun He, Alexander D. Diehl

*BMC Bioinformatics* (2019-04-25) <https://doi.org/gg99b8>  
DOI: [10.1186/s12859-019-2721-9](https://doi.org/10.1186/s12859-019-2721-9) · PMID: [31272374](https://pubmed.ncbi.nlm.nih.gov/31272374/) · PMCID: [PMC6509796](https://pubmed.ncbi.nlm.nih.gov/PMC6509796/)

**26. Current best practices in single-cell RNA-seq analysis: a tutorial**

Malte D. Luecken, Fabian J. Theis

*Molecular Systems Biology* (2019-06-19) <https://www.wikidata.org/wiki/Q64974172>

DOI: [10.15252/msb.20188746](https://doi.org/10.15252/msb.20188746)

**27. Key words for use in RFCs to Indicate Requirement Levels**

Scott Bradner

(1997-01-01) <https://www.wikidata.org/wiki/Q104060055>

**28. Toward principles for the design of ontologies used for knowledge sharing?**

Thomas R. Gruber

*International Journal of Human-Computer Studies* (1995-11-01)

<https://www.wikidata.org/wiki/Q47303277>

DOI: [10.1006/ijhc.1995.1081](https://doi.org/10.1006/ijhc.1995.1081)

**29. Reporting and connecting cell type names and gating definitions through ontologies**

James A. Overton, Randi Vita, Patrick Dunn, Julie G. Burel, Syed Ahmad Chan Bukhari, Kei-Hoi Cheung, Steven H. Kleinstein, Alexander D. Diehl, Bjoern Peters

*BMC Bioinformatics* (2019-04-25) <https://doi.org/ghbk9r>

DOI: [10.1186/s12859-019-2725-5](https://doi.org/10.1186/s12859-019-2725-5) · PMID: [31272390](https://pubmed.ncbi.nlm.nih.gov/31272390/) · PMCID: [PMC6509839](https://pubmed.ncbi.nlm.nih.gov/PMC6509839/)

**30. flowCL: ontology-based cell population labelling in flow cytometry**

Mélanie Courtot, Justin Meskas, Alexander D. Diehl, Radina Droumeva, Raphael Gottardo, Adrin Jalali, Mohammad Jafar Taghiyar, Holden T. Maecker, J. Philip McCoy, Alan Ruttenberg, ... Ryan R. Brinkman

*Bioinformatics* (2015-04-15) <https://doi.org/f7cc46>

DOI: [10.1093/bioinformatics/btu807](https://doi.org/10.1093/bioinformatics/btu807) · PMID: [25481008](https://pubmed.ncbi.nlm.nih.gov/25481008/) · PMCID: [PMC4393520](https://pubmed.ncbi.nlm.nih.gov/PMC4393520/)

**31. Cell ontology in an age of data-driven cell classification**

David Osumi-Sutherland

*BMC Bioinformatics* (2017-12-21) <https://doi.org/ghcbdk>

DOI: [10.1186/s12859-017-1980-6](https://doi.org/10.1186/s12859-017-1980-6) · PMID: [29322914](https://pubmed.ncbi.nlm.nih.gov/29322914/) · PMCID: [PMC5763290](https://pubmed.ncbi.nlm.nih.gov/PMC5763290/)

**32. Formalization of taxon-based constraints to detect inconsistencies in annotation and ontology development.**

Jennifer I. Deegan née Clark, Emily C. Dimmer, Chris Mungall

*BMC Bioinformatics* (2010-10-25) <https://www.wikidata.org/wiki/Q33727235>

DOI: [10.1186/1471-2105-11-530](https://doi.org/10.1186/1471-2105-11-530)

**33. Logik der Forschung**

Karl Popper

(1934-01-01) <https://www.wikidata.org/wiki/Q1868040>

ISBN: [9783161484100](https://www.wikidata.org/wiki/Q1868040)

**34. Neuronal diversity and convergence in a visual system developmental atlas**

Mehmet Neset Özel, Félix Simon, Shadi Jafari, Isabel Holguera, Yen-Chung Chen, Najate Benhra, Rana Naja El-Danaf, Katarina Kapuralin, Jennifer Amy Malin, Nikolaos Konstantinides, Claude Desplan

*Nature* (2020-11-04) <https://www.wikidata.org/wiki/Q101226729>

DOI: [10.1038/s41586-020-2879-3](https://doi.org/10.1038/s41586-020-2879-3)

35. **The neuronal organization of the retina.**  
Richard H. Masland  
*Neuron* (2012-10-17) <https://www.wikidata.org/wiki/Q34307217>  
DOI: [10.1016/j.neuron.2012.10.002](https://doi.org/10.1016/j.neuron.2012.10.002)
36. **International Code of Phylogenetic Nomenclature (PhyloCode)**  
Philip D. Cantino, Kevin de Queiroz  
(2020) <http://phylonames.org/code/>
37. **The generalizability crisis**  
Tal Yarkoni  
*PsyArXiv* (2019-11-22) <https://psyarxiv.com/jqw35>
38. **The Human Cell Atlas: Technical approaches and challenges.**  
Chung Chau Hon, Jay W. Shin, Piero Carninci, Michael J. T. Stubbington  
*Briefings in functional genomics* (2017-10-28) <https://www.wikidata.org/wiki/Q48563763>  
DOI: [10.1093/bfpg/elx029](https://doi.org/10.1093/bfpg/elx029)
39. **Theory of cell fate**  
Michael J. Casey, Patrick S. Stumpf, Ben D. MacArthur  
*Wiley interdisciplinary reviews. Systems biology and medicine* (2019-12-12)  
<https://www.wikidata.org/wiki/Q91908361>  
DOI: [10.1002/wsbm.1471](https://doi.org/10.1002/wsbm.1471)
40. **CHETAH: a selective, hierarchical cell type identification method for single-cell RNA sequencing**  
Jurrian K de Kanter, Philip Lijnzaad, Tito Candelli, Thanasis Margaritis, Frank CP Holstege  
*Nucleic Acids Research* (2019-09-19) <https://doi.org/gg99dp>  
DOI: [10.1093/nar/gkz543](https://doi.org/10.1093/nar/gkz543) · PMID: [31226206](https://pubmed.ncbi.nlm.nih.gov/31226206/) · PMCID: [PMC6895264](https://pubmed.ncbi.nlm.nih.gov/PMC6895264/)
41. **Mapping the transcriptional diversity of genetically and anatomically defined cell populations in the mouse brain**  
Ken Sugino, Erin Clark, Anton Schulmann, Yasuyuki Shima, Lihua Wang, David L Hunt, Bryan M Hooks, Dimitri Tränkner, Jayaram Chandrashekar, Serge Picard, ... Sacha B Nelson  
*eLife* (2019-04-12) <https://doi.org/ghbc3p>  
DOI: [10.7554/elife.38619](https://doi.org/10.7554/elife.38619) · PMID: [30977723](https://pubmed.ncbi.nlm.nih.gov/30977723/) · PMCID: [PMC6499542](https://pubmed.ncbi.nlm.nih.gov/PMC6499542/)
42. **How Single-Cell Genomics Is Changing Evolutionary and Developmental Biology**  
John C. Marioni, Detlev Arendt  
*Annual Review of Cell and Developmental Biology* (2017-10-06) <https://doi.org/ggb632>  
DOI: [10.1146/annurev-cellbio-100616-060818](https://doi.org/10.1146/annurev-cellbio-100616-060818) · PMID: [28813177](https://pubmed.ncbi.nlm.nih.gov/28813177/)
43. **Cell type discovery using single-cell transcriptomics: implications for ontological representation.**  
Brian D Aeversmann, Mark Novotny, Trygve Bakken, Jeremy A Miller, Alexander D Diehl, David Osumi-Sutherland, Roger S Lasken, Ed S Lein, Richard H Scheuermann  
*Human molecular genetics* (2018-05-01) <https://www.ncbi.nlm.nih.gov/pubmed/29590361>  
DOI: [10.1093/hmg/ddy100](https://doi.org/10.1093/hmg/ddy100) · PMID: [29590361](https://pubmed.ncbi.nlm.nih.gov/29590361/) · PMCID: [PMC5946857](https://pubmed.ncbi.nlm.nih.gov/PMC5946857/)
44. **A revised airway epithelial hierarchy includes CFTR-expressing ionocytes**  
Daniel T. Montoro, Adam L. Haber, Moshe Biton, Vladimir Vinarsky, Brian Lin, Susan E. Birket, Feng Yuan, Sijia Chen, Hui Min Leung, Jorge Villoria, ... Jayaraj Rajagopal

*Nature* (2018-08-01) <https://doi.org/gdwskh>  
DOI: [10.1038/s41586-018-0393-7](https://doi.org/10.1038/s41586-018-0393-7) · PMID: [30069044](https://pubmed.ncbi.nlm.nih.gov/30069044/) · PMCID: [PMC6295155](https://pubmed.ncbi.nlm.nih.gov/PMC6295155/)

**45. A single-cell atlas of the airway epithelium reveals the CFTR-rich pulmonary ionocyte**

Lindsey W. Plasschaert, Rapolas Žilionis, Rayman Choo-Wing, Virginia Savova, Judith Knehr, Guglielmo Roma, Allon M. Klein, Aron B. Jaffe

*Nature* (2018-08-01) <https://doi.org/gdwsjZ>

DOI: [10.1038/s41586-018-0394-6](https://doi.org/10.1038/s41586-018-0394-6) · PMID: [30069046](https://pubmed.ncbi.nlm.nih.gov/30069046/) · PMCID: [PMC6108322](https://pubmed.ncbi.nlm.nih.gov/PMC6108322/)

**46. Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors**

Alexandra-Chloé Villani, Rahul Satija, Gary Reynolds, Siranush Sarkizova, Karthik Shekhar, James Fletcher, Morgane Griesbeck, Andrew Butler, Shiwei Zheng, Suzan Lazo, ... Nir Hacohen

*Science* (2017-04-20) <https://doi.org/f94x5t>

DOI: [10.1126/science.aah4573](https://doi.org/10.1126/science.aah4573) · PMID: [28428369](https://pubmed.ncbi.nlm.nih.gov/28428369/) · PMCID: [PMC5775029](https://pubmed.ncbi.nlm.nih.gov/PMC5775029/)

**47. An improved ontological representation of dendritic cells as a paradigm for all cell types**

Anna Masci, Cecilia N Arighi, Alexander D Diehl, Anne E Lieberman, Chris Mungall, Richard H Scheuermann, Barry Smith, Lindsay G Cowell

*BMC Bioinformatics* (2009) <https://doi.org/cpxdhs>

DOI: [10.1186/1471-2105-10-70](https://doi.org/10.1186/1471-2105-10-70) · PMID: [19243617](https://pubmed.ncbi.nlm.nih.gov/19243617/) · PMCID: [PMC2662812](https://pubmed.ncbi.nlm.nih.gov/PMC2662812/)

**48. Hematopoietic cell types: Prototype for a revised cell ontology**

Alexander D. Diehl, Alison Deckhut Augustine, Judith A. Blake, Lindsay G. Cowell, Elizabeth S. Gold, Timothy A. Gondré-Lewis, Anna Maria Masci, Terrence F. Meehan, Penelope A. Morel, Anastasia Nijnik, ... Christopher J. Mungall

*Journal of Biomedical Informatics* (2011-02) <https://doi.org/c6dmmh>

DOI: [10.1016/j.jbi.2010.01.006](https://doi.org/10.1016/j.jbi.2010.01.006) · PMID: [20123131](https://pubmed.ncbi.nlm.nih.gov/20123131/) · PMCID: [PMC2892030](https://pubmed.ncbi.nlm.nih.gov/PMC2892030/)

**49. A Survey and Classification of Controlled Natural Languages**

Tobias Kuhn

*Computational Linguistics* (2014-03-01) <https://www.wikidata.org/wiki/Q57402167>

DOI: [10.1162/coli\\_a\\_00168](https://doi.org/10.1162/coli_a_00168)

**50. The challenges of big data biology**

Sabina Leonelli

*eLife* (2019-04-05) <https://doi.org/gfzw8q>

DOI: [10.7554/elife.47381](https://doi.org/10.7554/elife.47381) · PMID: [30950793](https://pubmed.ncbi.nlm.nih.gov/30950793/) · PMCID: [PMC6450665](https://pubmed.ncbi.nlm.nih.gov/PMC6450665/)

**51. Wikidata as a knowledge graph for the life sciences**

Andra Waagmeester, Gregory Stupp, Sebastian Burgstaller-Muehlbacher, Benjamin M Good, Malachi Griffith, Obi L Griffith, Kristina Hanspers, Henning Hermjakob, Toby S Hudson, Kevin Hybiske, ... Andrew I Su

*eLife* (2020-03-17) <https://doi.org/gggqc6>

DOI: [10.7554/elife.52614](https://doi.org/10.7554/elife.52614) · PMID: [32180547](https://pubmed.ncbi.nlm.nih.gov/32180547/) · PMCID: [PMC7077981](https://pubmed.ncbi.nlm.nih.gov/PMC7077981/)

**52. Common Cell type Nomenclature for the mammalian brain: A systematic, extensible convention**

Jeremy A. Miller, Nathan W. Gouwens, Bosiljka Tasic, Forrest Collman, Cindy T. J. van Velthoven, Trygve E. Bakken, Michael J. Hawrylycz, Hongkui Zeng, Ed S. Lein, Amy Bernard

*arXiv* (2020-11-13) <https://www.wikidata.org/wiki/Q104247451>



53. **Neuron Names: A Gene- and Property-Based Name Format, With Special Reference to Cortical Neurons**  
Gordon M. Shepherd, Luis Marenco, Michael L. Hines, Michele Migliore, Robert A. McDougal, Nicholas T. Carnevale, Adam J. H. Newton, Monique Surles-Zeigler, Giorgio A. Ascoli  
*Frontiers in Neuroanatomy* (2019-01-01) <https://www.wikidata.org/wiki/Q64065346>  
DOI: [10.3389/fnana.2019.00025](https://doi.org/10.3389/fnana.2019.00025)
54. **Name-calling in the hippocampus (and beyond): coming to terms with neuron types and properties.**  
D. J. Hamilton, D. W. Wheeler, C. M. White, C. L. Rees, A. O. Komendantov, M. Bergamino, Giorgio A. Ascoli  
*Brain informatics* (2016-06-09) <https://www.wikidata.org/wiki/Q37656899>  
DOI: [10.1007/s40708-016-0053-3](https://doi.org/10.1007/s40708-016-0053-3)
55. **Naming of neurones. Classification and naming of cat retinal ganglion cells.**  
Rowe MH, Stone J  
*Brain, Behavior and Evolution* (1977-01-01) <https://www.wikidata.org/wiki/Q41052480>  
DOI: [10.1159/000125660](https://doi.org/10.1159/000125660)
56. **Probabilistic gene expression signatures identify cell-types from single cell RNA-seq data**  
Isabella N. Grabski, Rafael A. Irizarry  
*bioRxiv* (2020-01-23) <https://www.wikidata.org/wiki/Q104371272>  
DOI: [10.1101/2020.01.05.895441](https://doi.org/10.1101/2020.01.05.895441)
57. **ontoProc: processing of ontologies of anatomy, cell lines, and so on**  
Vince Carey  
*Bioconductor* (2020) <https://www.bioconductor.org/packages/release/bioc/html/ontoProc.html>
58. **Searching large-scale scRNA-seq databases via unbiased cell embedding with Cell BLAST**  
Zhi-Jie Cao, Lin Wei, Shen Lu, De-Chang Yang, Ge Gao  
*Nature Communications* (2020-07-10) <https://doi.org/gg4mm3>  
DOI: [10.1038/s41467-020-17281-7](https://doi.org/10.1038/s41467-020-17281-7) · PMID: [32651388](https://pubmed.ncbi.nlm.nih.gov/32651388/) · PMCID: [PMC7351785](https://pubmed.ncbi.nlm.nih.gov/PMC7351785/)
59. **Unifying single-cell annotations based on the Cell Ontology**  
Sheng Wang, Angela Oliveira Pisco, Aaron McGeever, Maria Brbic, Marinka Zitnik, Spyros Darmanis, Jure Leskovec, Jim Karkanias, Russ B. Altman  
(2019-10-20) <https://www.wikidata.org/wiki/Q104057222>
60. **CellMeSH: Probabilistic Cell-Type Identification Using Indexed Literature**  
Shunfu Mao, Yue Zhang, Georg Seelig, Sreeram Kannan  
(2020-05-31) <https://www.wikidata.org/wiki/Q104371393>
61. **The naming of neurons: applications of taxonomic theory to the study of cellular populations.**  
Tyner CF  
*Brain, Behavior and Evolution* (1975-01-01) <https://www.wikidata.org/wiki/Q34065481>  
DOI: [10.1159/000124141](https://doi.org/10.1159/000124141)
62. **The anatomy of a nanopublication**  
Paul Groth, Andrew Gibson, Jan Velterop  
*Information Services & Use* (2010-09-21) <https://www.wikidata.org/wiki/Q57011346>  
DOI: [10.3233/isu-2010-0613](https://doi.org/10.3233/isu-2010-0613)

**63. The value of data**

Barend Mons, Herman van Haagen, Christine Chichester, Peter-Bram 't Hoen, Johan T. den Dunnen, Gertjan van Ommen, Erik M. van Mulligen, Bharat Singh, Rob Hooft, Marco Roos, ... Erik Schultes

*Nature Genetics* (2011-03-29) <https://www.wikidata.org/wiki/Q22676713>

DOI: [10.1038/ng0411-281](https://doi.org/10.1038/ng0411-281)

**64. Publishing DisGeNET as nanopublications**

Núria Queralt Rosinach, Tobias Kuhn, Christine Chichester, Michel Dumontier, Ferran Sanz, Laura I. Furlong

*Semantic Web: interoperability, usability, applicability* (2016-06-23)

<https://www.wikidata.org/wiki/Q31194033>

DOI: [10.3233/sw-150189](https://doi.org/10.3233/sw-150189)

**65. Trusty URIs: Verifiable, Immutable, and Permanent Digital Artifacts for Linked Data**

Tobias Kuhn, Michel Dumontier

*Lecture Notes in Computer Science* (2014-01-01) <https://www.wikidata.org/wiki/Q56915510>

DOI: [10.1007/978-3-319-07443-6\\_27](https://doi.org/10.1007/978-3-319-07443-6_27)

**66. Identifiers for the 21st century: How to design, provision, and reuse persistent identifiers to maximize utility and impact of life science data**

Julie A. McMurry, Julie A. McMurry, Nick Juty, Niklas Blomberg, Tony Burdett, Tom Conlin, Nathalie Conte, Melanie Courtot, John Deck, Michel Dumontier, ... Helen Parkinson

*PLOS Biology* (2017-06-29) <https://www.wikidata.org/wiki/Q33037209>

DOI: [10.1371/journal.pbio.2001414](https://doi.org/10.1371/journal.pbio.2001414)

**67. Neuronal cell-type classification: challenges, opportunities and the path forward**

Hongkui Zeng, Joshua R. Sanes

*Nature Reviews Neuroscience* (2017-08-03) <https://doi.org/10.1038/nrn.2017.85>

DOI: [10.1038/nrn.2017.85](https://doi.org/10.1038/nrn.2017.85) · PMID: [28775344](https://pubmed.ncbi.nlm.nih.gov/28775344/)