

“The Brain is...”: A Survey of The Brain’s Many Definitions

Taylor Bolt¹, Lucina Q. Uddin^{1,2}

¹ Department of Psychiatry and Biobehavioral Sciences, University of California, Los Angeles, Los Angeles, CA, USA

² Department of Psychology, University of California, Los Angeles, Los Angeles, CA, USA

Abstract

A reader of the peer-reviewed neuroscience literature will often encounter expressions like the following: ‘the brain is a dynamic system’, ‘the brain is a complex network’, or ‘the brain is a highly metabolic organ’. These expressions attempt to define the essential functions and properties of the mammalian or human brain in a simple phrase or sentence, sometimes using metaphors. We sought to survey the most common phrases of the form ‘the brain is...’ in the biomedical literature to provide insights into current conceptualizations of the brain. Utilizing text analytic tools applied to a large sample of peer-reviewed full-text articles and abstracts, we extracted several thousand phrases of the form ‘the brain is...’ and identified over a dozen frequently appearing phrases. The most used phrases included metaphors (e.g., the brain as a ‘information processor’ or ‘prediction machine’) and descriptions of essential functions (e.g., ‘a central organ of stress adaptation’) or properties (e.g., ‘a highly vascularized organ’). The results of our analysis underscore the diversity of qualities and functions commonly attributed to the brain in the biomedical literature and suggest a range of conceptualizations that defy unification.

Introduction

The heart is a pump sending blood around the body, the primary organ of the circulatory system, and a muscular organ made up of cardiac muscle. While these definitions invariably fail to appreciate the complexity of the heart’s operations, they provide simple answers to a simple question: what sort of thing is a heart? They function as pedagogical or explanatory tools that convey the heart’s ‘essential’ functions and qualities. Some of these descriptions express analogies (e.g., ‘the heart is a pump’), others express part-whole relationships (e.g., the ‘primary organ of the circulatory system’), while others highlight the essential biochemical or cellular properties of the organ that enable its function (e.g., ‘made up of cardiac muscle’).

Not all organs of the human body can be so readily defined. Communicating the essential functions and qualities of the human or mammalian brain is a notoriously difficult task. For many scientists and philosophers, the brain’s essential qualities and functions are a matter of substantive scientific interest. At stake is not merely linguistic disagreements, but the guiding

interpretive framework for communicating neuroscientific results and potentially the favored methodologies by which we collect relevant data (Kelty-Stephen et al., 2022). For example, the computational or information-processing analogy has inspired a wealth of research findings and interpretive frameworks in cognitive, systems and computational neuroscience (Cobb, 2020). Favored definitions may be non-metaphorical as well, inspired by the brain's primary role in a physiological process (e.g., the brain's role in hormone regulation and stress).

To appreciate the diversity of definitions that scientists have ascribed to the brain in attempts to study it and determine whether unifying conceptualizations are evident, we conducted a survey of the biomedical literature (PubMed Central Open Access) using natural language processing (NLP) techniques. We searched for noun phrases following variations of the phrase 'The brain is a...'. Our survey targeted expressions where the subject ('the brain') is linked via a copular verb ('is') to a predicative expression consisting of a determiner ('the/a') followed by a noun phrase (e.g. 'complex computer'). An example expression that would fit this pattern is 'the brain is a complex computer'. NLP-based semantic embeddings and dimension-reduction techniques identified over a dozen commonly used copular expressions to describe the brain, many with quite different meanings. Some of these expressions were metaphors ('the brain is a computer'), others described an essential property ('the brain is a metabolically expensive organ'), while others described an essential function ('the brain is the key regulator of stress'). Our results underscore the diversity of attributes that scientists have ascribed to the brain in the biomedical literature.

Results

The biomedical corpus for analysis consisted of full text articles from the Pubmed Central Open Access Subset (N=4,993,411) and abstracts from leading neuroscience journals (N=253,022). Of the total corpus, 4,386 expressions from 895 peer-reviewed journals were found that matched expressions of the form 'the brain is ...'. Due to the disproportionate representation of large open-access journals in the corpus, journals such as *PLoS ONE* (N = 207), *International Journal of Molecular Sciences* (N = 187), *Frontiers in Neuroscience* (N = 172), *Psychology* (N = 129), and *Human Neuroscience* (N = 105), and *Scientific Reports* (N = 70) constituted the largest share of matched expressions. Articles from the current decade (≥ 2020) constituted the largest share of matched expressions (~56%). Of the 2,204 (51%) matched expressions identified from sections of text with 'standard' titles (e.g., titles containing 'Abstract', 'Introduction', 'Results', 'Discussion', 'Methods', etc.), 64% were found in introduction sections, 13% were found in discussions, and 10% were found in abstracts.

Following extraction, the text in each matched expression was converted to a vector space via a pretrained embedding model (Deka et al., 2022). To identify commonly used expressions, we dimension-reduced the embedding space to two dimensions using UMAP (McInnes et al., 2020) and clustered the phrases using a hierarchical density-based clustering algorithm (HDBSCAN; (Campello et al., 2013) (**Figure 1**). Of the 29 clusters extracted by the HDBSCAN algorithm (see *Methods and Materials*), 24 were found to constitute semantically coherent groups of expressions – i.e., groups of expressions that express similar meaning. Two

pairs of clusters (9 & 10, 22 & 25) were found to contain similar meanings and were merged, leaving a total of 22 semantically coherent clusters. Labels for each cluster were generated from manual inspection. In terms of overall organization of the semantic embedding space, the phrases were organized along a single dominant dimension, such that cellular/biochemical phrases were concentrated in the top-left, to more abstract/metaphorical phrases in the bottom-right.

Commonly used noun phrases following the phrase ‘the brain is...’ span multiple levels of organization, from the biochemical (e.g., ‘a lipid rich organ’) to the structural (e.g., ‘a highly vascularized organ’); as well as different types of expressions, including metaphors (e.g., ‘a prediction machine’ and ‘a control system’). Of the 24 semantically coherent clusters, the top clusters were (in descending order), the brain is ‘an energy demanding organ’ (N = 426), ‘a complex network’ (N = 412), ‘a heterogenous organ (N = 248), ‘a control system’ (N = 160), and ‘an immune-privileged organ’ (N = 157).

The occurrence of each expression was unequally distributed across journals (**Supplementary Table 1**), with expressions clustering into journals with distinct disciplinary concentrations. Some of these were predictable – e.g., biochemical expressions, such as ‘insulin-sensitive’ and ‘cholesterol-’ and ‘lipid-rich’ tended to appear more frequently in journals with a focus in biochemistry and molecular biology (e.g. *Oxidative Medicine and Cellular Longevity, Molecules*), and mentions of the brain as a ‘common site of metastasis’ appear more frequently in oncology journals (e.g. *Cancers, Frontiers in Oncology*).

Expressions of metaphors tended to appear more frequently in psychology and human neuroscience journals. For example, the majority of the expressions found in the journal *Frontiers in Psychology* were the brain as a ‘prediction machine’ or ‘computer/information processor’. The expressions of the brain as a ‘non-linear’ or ‘complex, dynamic’ system and ‘complex network’ tended to occur in human neuroscience/neuroimaging journals, including *Human Brain Mapping, Frontiers in Human Neuroscience, Neuroimage, Brain and Behavior*, and *Network Neuroscience*.

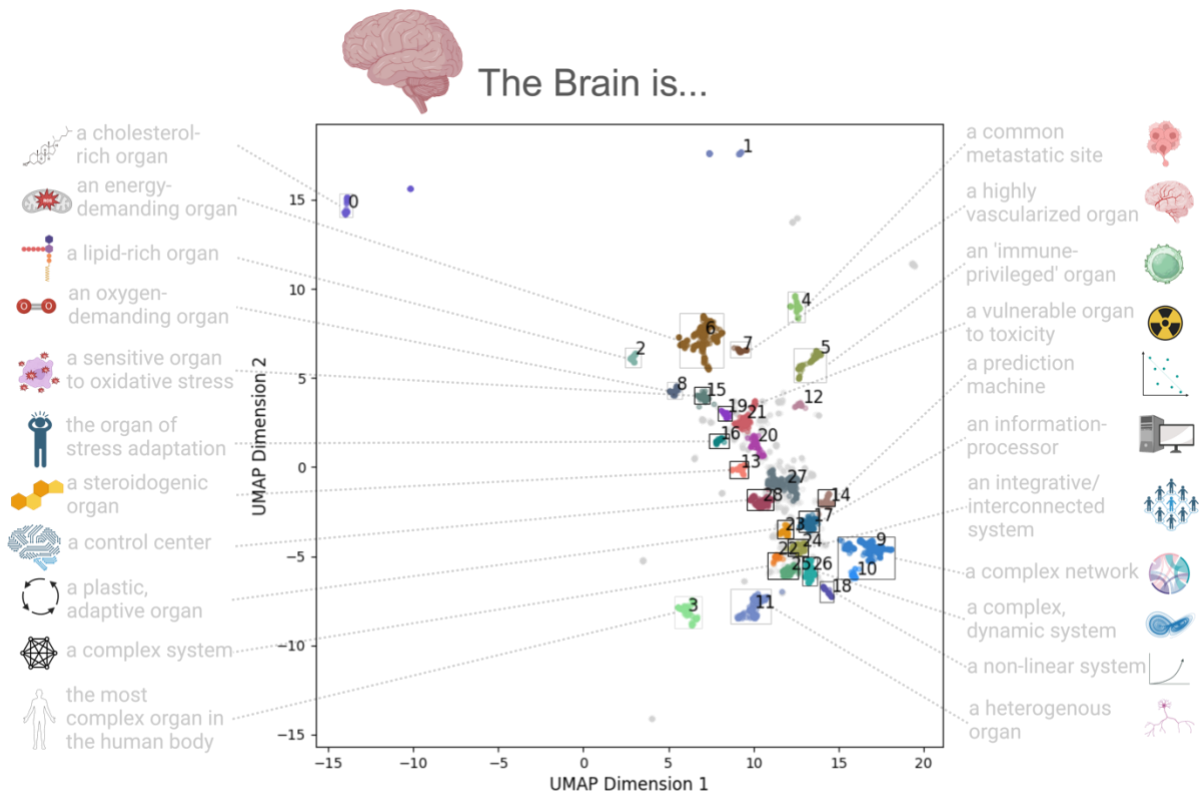


Figure 1. **'The brain is...'** Expressions in a Two-Dimensional Embedding Space. Expressions matching the form 'The brain is...' embedded into a two-dimensional space via a dimension-reduction (UMAP) applied to their semantic embeddings. The distance between points in this space reflect the semantic similarity between the expressions – i.e., expressions (points) in this space that are closer together reflect similar meanings. Expressions are color-coded according to their cluster assignment from the HDBSCAN clustering algorithm. Each semantically coherent cluster (N=24) is labeled by a manual interpretation of the expressions in the cluster (note that two pairs of clusters were merged together, forming 12 clusters).

Discussion

What is the brain? The results of our analysis underscore the diversity of qualities and functions commonly attributed to the human or mammalian brain in the biomedical literature. For scientists in some disciplines, the brain is understood from an abstract level, in the form of analogy – a computer/information processor, a prediction (Bayesian) machine, or a dynamic system. For others, the brain is a site of remarkable functional capacities, including its role in stress adaptation, its status as a relatively 'immune-privileged' organ, and its high metabolic activity. For others still the brain is defined in terms its biochemical properties as a lipid- and cholesterol-rich organ or its vulnerability to oxidative stress.

The appearance of this diversity of expression in the biomedical literature raises the question of implications for scientific communication. In many cases, it's clear that these varied

expressions do serve a purpose: expressions such as the ‘the brain is a highly vascularized organ’ or ‘the brain is an oxygen-demanding organ’ communicates an attribute of the brain that enables its unique functions or demarcates it from other organs in the body. Others attempt to communicate the ‘essential’ function or property of the brain that makes it what it is – e.g., ‘the brain is an information processor/computer’ or ‘the brain is a dynamic, non-linear system’. In many cases, these analogical expressions function to justify the author’s research approach or perspective, sometimes in contrast to a rival approach or perspective.

Importantly, these expressions often serve to introduce a research or review article, as evidenced by the fact they often appear in the Introduction section of articles. They are not propositions to be assessed by empirical observations or experiments reported in the research article. Rather, they serve to introduce more concrete and testable hypotheses, detailed frameworks and theories, and/or methodological approaches.

We agree with and reiterate a common expression in the biomedical literature: the brain is the most complex organ of the human body. This complexity is not only found in its biological structure, but also in the predicates that scientists use to describe it. The fact that a range of conceptualizations exist that defy unification has implications for science communication, neuroscience education, and society more broadly.

| | adaptive plastic | cholesterol rich | complex dynamic | complex network | complex system | control system | energy demanding | heterogenous organ | high vascularization | immune privileged | information processor | integrated system | lipid rich | metastasis | most complex organ | non linear system | oxidative stress | oxygen demanding | prediction machine | steroidogenic | stress adaptation | toxicity |
|--|------------------|------------------|-----------------|-----------------|----------------|----------------|------------------|--------------------|----------------------|-------------------|-----------------------|-------------------|------------|------------|--------------------|-------------------|------------------|------------------|--------------------|---------------|-------------------|----------|
| Aging (Albany NY) | 0 | 2 | 0 | 0 | 1 | 0 | 6 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| Antioxidants | 0 | 2 | 0 | 0 | 0 | 1 | 11 | 1 | 0 | 1 | 0 | 0 | 2 | 0 | 1 | 0 | 7 | 6 | 0 | 0 | 0 | 2 |
| BioMed Research International | 0 | 0 | 1 | 1 | 1 | 1 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 1 | 1 | 1 | 0 | 0 | 1 | 2 |
| Biology | 0 | 0 | 1 | 1 | 4 | 2 | 1 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| Biomolecules | 0 | 1 | 0 | 0 | 0 | 0 | 6 | 2 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 1 |
| Brain Sciences | 4 | 0 | 0 | 5 | 8 | 3 | 0 | 1 | 0 | 0 | 0 | 3 | 0 | 1 | 3 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| Brain and Behavior | 1 | 0 | 0 | 12 | 1 | 1 | 0 | 2 | 1 | 0 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Cancers | 0 | 0 | 0 | 2 | 0 | 0 | 2 | 1 | 1 | 5 | 0 | 0 | 0 | 9 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Cells | 0 | 0 | 0 | 2 | 0 | 1 | 6 | 1 | 1 | 4 | 0 | 1 | 1 | 1 | 3 | 0 | 2 | 0 | 0 | 2 | 1 | 2 |
| Current Neuropharmacology | 0 | 1 | 0 | 0 | 1 | 1 | 6 | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Current opinion in neurology | 1 | 0 | 0 | 1 | 2 | 0 | 2 | 0 | 0 | 4 | 2 | 1 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 0 |
| Entropy | 0 | 0 | 1 | 1 | 6 | 2 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 3 | 0 | 0 | 1 | 0 | 0 | 0 |
| Frontiers in Aging Neuroscience | 0 | 3 | 2 | 17 | 2 | 3 | 17 | 1 | 1 | 0 | 0 | 0 | 2 | 1 | 2 | 0 | 3 | 3 | 0 | 3 | 0 | 0 |
| Frontiers in Cell and Developmental Biology | 0 | 1 | 0 | 0 | 0 | 2 | 3 | 3 | 0 | 1 | 1 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Frontiers in Cellular Neuroscience | 1 | 2 | 0 | 2 | 0 | 2 | 4 | 7 | 0 | 4 | 0 | 0 | 1 | 0 | 4 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| Frontiers in Computational Neuroscience | 2 | 0 | 3 | 11 | 5 | 2 | 2 | 0 | 0 | 1 | 3 | 1 | 0 | 0 | 1 | 2 | 0 | 0 | 2 | 0 | 0 | 0 |
| Frontiers in Endocrinology | 0 | 1 | 1 | 1 | 1 | 0 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 1 | 0 |
| Frontiers in Human Neuroscience | 4 | 1 | 15 | 22 | 6 | 1 | 1 | 0 | 0 | 0 | 5 | 1 | 0 | 0 | 0 | 4 | 0 | 0 | 9 | 1 | 0 | 0 |
| Frontiers in Immunology | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 3 | 0 | 12 | 0 | 0 | 0 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 0 |
| Frontiers in Molecular Neuroscience | 0 | 3 | 0 | 2 | 0 | 0 | 3 | 4 | 1 | 2 | 0 | 0 | 0 | 1 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| Frontiers in Neurology | 0 | 4 | 6 | 7 | 3 | 1 | 11 | 3 | 2 | 4 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 2 | 0 | 0 | 0 |
| Frontiers in Neuroscience | 5 | 0 | 9 | 40 | 9 | 0 | 10 | 4 | 1 | 2 | 0 | 0 | 4 | 0 | 1 | 5 | 3 | 3 | 3 | 2 | 1 | 1 |
| Frontiers in Oncology | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 2 | 0 | 0 | 3 | 9 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Frontiers in Pharmacology | 0 | 0 | 1 | 0 | 0 | 1 | 5 | 3 | 1 | 0 | 0 | 1 | 2 | 0 | 1 | 0 | 1 | 3 | 1 | 0 | 0 | 1 |
| Frontiers in Physiology | 0 | 0 | 1 | 5 | 1 | 4 | 6 | 1 | 4 | 1 | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 1 | 0 | 0 | 1 | 2 |
| Frontiers in Psychiatry | 1 | 0 | 5 | 11 | 1 | 2 | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 0 | 2 | 1 | 2 | 2 | 4 | 0 | 3 | 0 |
| Frontiers in Psychology | 3 | 0 | 2 | 6 | 0 | 2 | 0 | 1 | 0 | 0 | 21 | 5 | 0 | 0 | 4 | 1 | 0 | 0 | 10 | 0 | 1 | 0 |
| Frontiers in Systems Neuroscience | 1 | 0 | 2 | 9 | 3 | 2 | 1 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 3 | 0 | 0 | 0 |
| Human Brain Mapping | 1 | 0 | 0 | 16 | 4 | 1 | 0 | 1 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| I.J. of Environmental Research and Public Health | 0 | 2 | 0 | 1 | 1 | 0 | 2 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 2 | 3 |
| International Journal of Molecular Sciences | 1 | 11 | 0 | 1 | 4 | 5 | 10 | 0 | 2 | 11 | 0 | 1 | 9 | 2 | 4 | 0 | 11 | 5 | 0 | 4 | 2 | 1 |
| Metabolites | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Molecules | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 5 | 2 | 0 | 0 | 1 | 1 |
| Nature Communications | 0 | 0 | 1 | 1 | 0 | 0 | 4 | 4 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Network Neuroscience | 2 | 0 | 2 | 13 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Neural Plasticity | 0 | 2 | 4 | 5 | 3 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| NeuroImage | 1 | 1 | 6 | 10 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| NeuroImage : Clinical | 2 | 0 | 3 | 9 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Nutrients | 0 | 0 | 1 | 0 | 0 | 2 | 14 | 3 | 1 | 1 | 0 | 0 | 5 | 0 | 3 | 0 | 4 | 0 | 0 | 0 | 2 | 1 |
| Oxidative Medicine and Cellular Longevity | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 2 | 1 | 0 | 0 | 0 | 4 | 0 | 0 | 0 | 7 | 4 | 0 | 1 | 3 | 1 |
| PLoS Computational Biology | 0 | 0 | 1 | 7 | 2 | 0 | 0 | 3 | 0 | 0 | 4 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 0 | 0 |
| PLoS ONE | 4 | 1 | 5 | 30 | 9 | 11 | 17 | 11 | 0 | 0 | 0 | 3 | 0 | 2 | 7 | 3 | 0 | 1 | 1 | 3 | 1 | 3 |
| Scientific Reports | 1 | 0 | 2 | 9 | 3 | 3 | 6 | 6 | 0 | 3 | 2 | 1 | 2 | 0 | 1 | 0 | 0 | 1 | 2 | 3 | 0 | 1 |
| Sensors (Basel, Switzerland) | 0 | 2 | 2 | 8 | 1 | 7 | 1 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 2 | 3 | 0 | 0 | 0 | 0 | 1 | 0 |
| eCAM | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 2 | 0 | 0 |
| eLife | 0 | 1 | 1 | 3 | 1 | 1 | 3 | 4 | 0 | 0 | 0 | 1 | 0 | 0 | 3 | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

Supplementary Table 1. **Counts of Expressions from Each Cluster by Journal Title.** A cross-tabulation table containing the counts of expressions in each cluster (column) by journal title (row). Only journals containing ≥ 20 expressions were included in the table. Clusters are labeled with the same labels as appears in **Figure 1** (abbreviated for space). Cells of the table are color-coded according to their row-wise normalized (min-max normalization) counts, with lighter colors corresponding to more frequent counts relative to the total counts in each journal (row).

| Journal | Abstract Count |
|---|----------------|
| Journal of Neuroscience | 37967 |
| Journal of Cognitive Neuroscience | 4227 |
| Neuroscience Letters | 35108 |
| Neuropsychologia | 8996 |
| Neuroimage | 22060 |
| Nature Reviews Neuroscience | 1319 |
| Journal of Neurophysiology | 19269 |
| Trends in Cognitive Science | 2316 |
| Neuroscience and Biobehavioral Reviews | 5492 |
| Human Brain Mapping | 5911 |
| Psychophysiology | 4178 |
| Cortex | 4875 |
| Experimental Brain Research | 14804 |
| Nature Reviews Neurology | 888 |
| European Journal of Neuroscience | 13157 |
| Neuron | 11945 |
| The Neuroscientist | 1083 |
| Annual Review of Neuroscience | 696 |
| Molecular Neurobiology | 5933 |
| Brain | 8470 |
| Social Cognitive and Affective Neuroscience | 1980 |
| Nature Neuroscience | 5054 |
| Journal of Neurology | 9330 |
| Behavioral and Brain Sciences | 3465 |
| Progress in Neurobiology | 1921 |
| Biological Psychiatry | 9314 |
| Trends in Neuroscience | 3003 |
| Current Opinion in Neurology | 2644 |
| Cerebral Cortex | 7617 |

Supplementary Table 2. **Neuroscience Journal Abstracts Included in Corpus.** The count of abstracts of neuroscience journals included in corpus for analysis.

Methods and Materials

Dataset

The corpus used for analysis consisted of all full-text articles from the PubMed Central (PMC) Open Access subset (<https://www.ncbi.nlm.nih.gov/pmc/tools/openftlist/>) and the PMC Author Manuscript Dataset as of July 23, 2023 (N=4,993,411). The PMC Open Access Subset provides access to full texts from open access peer-reviewed journals. The PMC Author Manuscript Dataset provides access to full texts of manuscripts made available in PMC by authors in compliance with the NIH Public Access Policy. Both sources form part of PMC's Open Access Collection (<https://www.ncbi.nlm.nih.gov/pmc/tools/textmining/>). Bulk downloads of the full Open Access Collection articles were conducted using the PMC FTP service. To supplement our corpus with scientific text outside the PMC Open Access Collection, we pulled abstracts from the PubMed database from select neuroscience journals ($N_{\text{journals}} = 29$, $N_{\text{articles}} = 253,022$; **Supplementary Table 2**) as of August 26, 2023 via the EFetch utility (Sayers, 2022). Overall, a total of approximately 5 million articles were downloaded and screened for our analysis. All code for preprocessing and analysis are provided at https://github.com/tsb46/the_brain_is.

Preprocessing of PMC Open Access Collection

Full-text articles from the PMC Open Access Collection were downloaded in XML format. To parse the article XML files into structured data for analysis we used the Pubmed Parser package in Python (Achakulvisut et al., 2020). As an initial filtering step to reduce the number of articles in the next stage of preprocessing, the text of each section in the XML article files were searched for mentions of the entity 'brain' (case insensitive and lemmatized) via the entity extraction pipeline in ScispaCy (Neumann et al., 2019; *en_core_sci_sm* model; v0.5.2). Note, other related entity terms were searched for in this filtering process but were not included in the manuscript. This initial filtering step reduced the size of the PMC Open Access Collection to 561,836 articles.

Detection of Copular Expressions of the Form 'The Brain is...'

To find sentences containing phrases/expressions of the form 'the brain is...' in our corpus, we used the *PhraseMatcher* utility in spaCy package (CITE). The phrase matcher detects sequence of tokens in text with prespecified linguistic properties, such as string matching, part-of-speech tags, dependency tree labels. The phrase we sought to detect in the corpus has a specific linguistic structure: the token 'brain' as the nominal subject of the clause (case insensitive match) preceded by a determiner (the/a) and followed by the copular verb ('is'). In addition, this phrase (e.g. 'The brain is...') is followed by another determiner (the/a) that precedes the noun phrase. Hypothetical phrases that match this pattern would be 'the brain is a

complex system' or 'the brain is the most metabolically active organ in the body'. To capture adjectives that modify or further describe the noun 'brain', such as the 'human brain' or 'mammalian brain', we also allowed for matches with an adjective token preceding the token 'brain'. Tokenization, dependency parsing and part-of-speech tagging were performed with the *en_core_sci_lg* model in ScispaCy (v0.5.2; Neumann et al., 2019). The *PhraseMatcher* pipeline applied to the filtered PMC Open Access Collection articles (N = 561,836) and neuroscience journal abstracts (N=253,022) identified 4,836 sentences containing a sequence of tokens that matched the token sequence template described above.

Information Extraction of Matched Expressions

The *PhraseMatcher* pipeline detected the noun phrase and copular verb ('is') that forms the subject of the sentence, the phrase 'The brain is...'. To extract the entire expression, including the full noun phrase(s) that follow the phrase 'The brain is...', further text parsing was performed. For accurate quantification of semantic similarity, it was important to extract the relevant tokens of the noun phrases following the phrase 'The brain is...', or irrelevant tokens/words would contribute the semantic similarity estimates. One possibility is extract the entire sentence containing the matched expressions, but this was found to include too much irrelevant information. For example, extracting the entire sentence: 'The brain is a complex system, and this is the inspiration for our analysis in the manuscript' would include the irrelevant second clause ('and this is the inspiration...'). We developed a custom rule-based algorithm based on traversal of the dependency tree extracted from the sentence by the ScispaCy model (<https://spacy.io/api/dependencyparser>). Briefly, the algorithm consisted of the following primary steps: 1) start with the immediate rightward children tokens of the copular verb 'is', 2) loop through the children tokens and detect any instance of a nominal modifier, conjunction, or relative clause modifier (in that order) syntactic relationship between the head and the child token. 3) If one token with any of these syntactic relationships is found, extract the immediate rightward children of that token, and repeat the process. Minor modifications of this algorithm were included based on repeated trial-and-error experiments on the matched expressions. The full list of the 4,836 expressions extracted from this algorithm are available at (URL).

Phrase Embedding and Clustering

The primary goal of our analysis was to estimate the common types of 'The brain is...' expressions in the biomedical literature. To identify different types of expressions, we grouped the 4,836 expressions into groups with identical or very similar semantic meaning using a clustering approach. First, the expressions (extracted from the information extraction pipeline described above) were transformed into a semantic embedding space ($N_{\text{dimensions}} = 768$) via a pretrained neural network model, previously trained on semantic similarity tasks for scientific text (Deka et al., 2022; <https://huggingface.co/pritamdeka/S-Scibert-snli-multinli-stsb>). The expressions in the embedding vector space were then fed to the uniform manifold approximation algorithm (UMAP; McInnes et al., 2020) for dimension reduction to two dimensions. Two central parameters – the number of nearest neighbors in initial graph

construction and the minimum distance points are allowed to be apart - control the coordinates and global structure of the semantic expressions in the dimension-reduced space. Based on visualization of all expressions in the two-dimensional space, the following parameters were found to produce Euclidean distances between expressions that most closely followed the semantic similarity/dissimilarity in the expressions ($N_{Neighbors} = 15$; $min_{dist} = 0$).

The two-dimensional coordinates of each expression in the dimension-reduced space were then fed to a hierarchical density-based clustering algorithm (HDBSCAN; McInnes et al., 2017) to isolate groups of expressions with very similar or identical semantic meanings. The number of clusters extracted from the HDBSCAN algorithm is controlled indirectly by setting the minimum cluster size. Manual inspection of cluster solutions across varying values of this parameter revealed that a minimum cluster size of 50 ($N_{clusters} = 29$) maximized semantic interpretability.

References

- Achakulvisut, T., Acuna, D. E., & Kording, K. (2020). Pubmed Parser: A Python Parser for PubMed Open-Access XML Subset and MEDLINE XML Dataset XML Dataset. *Journal of Open Source Software*, 5(46), 1979. <https://doi.org/10.21105/joss.01979>
- Campello, R. J. G. B., Moulavi, D., & Sander, J. (2013). Density-Based Clustering Based on Hierarchical Density Estimates. In J. Pei, V. S. Tseng, L. Cao, H. Motoda, & G. Xu (Eds.), *Advances in Knowledge Discovery and Data Mining* (pp. 160–172). Springer. https://doi.org/10.1007/978-3-642-37456-2_14
- Cobb, M. (2020). *The Idea of the Brain: The Past and Future of Neuroscience* (Illustrated edition). Basic Books.
- Deka, P., Jurek-Loughrey, A., & Deepak. (2022). Unsupervised Keyword Combination Query Generation from Online Health Related Content for Evidence-Based Fact Checking. *The 23rd International Conference on Information Integration and Web Intelligence*, 267–277. <https://doi.org/10.1145/3487664.3487701>

Kelty-Stephen, D. G., Cisek, P. E., De Bari, B., Dixon, J., Favela, L. H., Hasselman, F., Keijzer, F.,

Raja, V., Wagman, J. B., Thomas, B. J., & Mangalam, M. (2022). *In search for an alternative to the computer metaphor of the mind and brain* (arXiv:2206.04603). arXiv.

<https://doi.org/10.48550/arXiv.2206.04603>

McInnes, L., Healy, J., & Astels, S. (2017). hdbscan: Hierarchical density based clustering. *Journal of Open Source Software*, 2(11), 205. <https://doi.org/10.21105/joss.00205>

McInnes, L., Healy, J., & Melville, J. (2020). *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction* (arXiv:1802.03426). arXiv.

<https://doi.org/10.48550/arXiv.1802.03426>

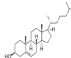


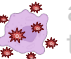




Neumann, M., King, D., Beltagy, I., & Ammar, W. (2019). ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. *Proceedings of the 18th BioNLP Workshop and Shared Task*, 319–327. <https://doi.org/10.18653/v1/W19-5034>

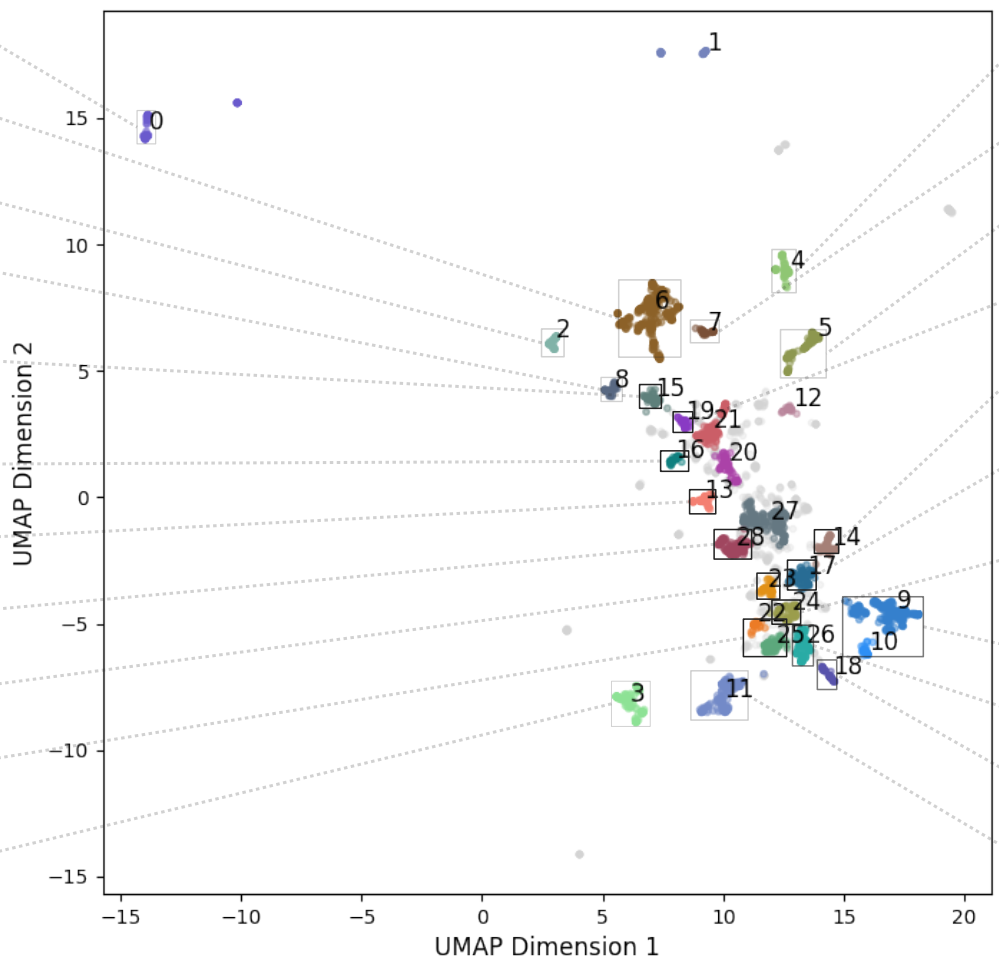
Sayers, E. (2022). The E-utilities In-Depth: Parameters, Syntax and More. In *Entrez Programming Utilities Help [Internet]*. National Center for Biotechnology Information (US).





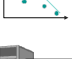





<https://www.ncbi.nlm.nih.gov/books/NBK25499/>



The Brain is...

-  a cholesterol-rich organ
-  an energy-demanding organ
-  a lipid-rich organ
-  an oxygen-demanding organ
-  a sensitive organ to oxidative stress
-  the organ of stress adaptation
-  a steroidogenic organ
-  a control center
-  a plastic, adaptive organ
-  a complex system
-  the most complex organ in the human body



-  a common metastatic site
-  a highly vascularized organ
-  an 'immune-privileged' organ
-  a vulnerable organ to toxicity
-  a prediction machine
-  an information-processor
-  an integrative/interconnected system
-  a complex network
-  a complex, dynamic system
-  a non-linear system
-  a heterogenous organ