

1 **Beyond Static Brain Atlases: AI-Powered Open Databasing**
2 **and Dynamic Mining of Brain-Wide Neuron Morphometry**

3 Shengdian Jiang^{1,+}, Lijun Wang^{1,+}, Zhixi Yun^{1,+}, Hanbo Chen^{2,+}, Jianhua Yao^{2,**}, Hanchuan
4 Peng^{1,3,*}

5
6 ¹ New Cornerstone Science Laboratory, Institute for Brain and Intelligence, Southeast
7 University, Nanjing, China.

8 ² Tencent AI, Shenzhen, China.

9 ³ Shanghai Academy of Natural Sciences, Shanghai, China.

10 ⁺ Equal contribution

11 ^{*} Correspondence: Hanchuan Peng <h@braintell.org>

12 ^{**} Co-Correspondence: Jianhua Yao <jianhuayao@tencent.com>

13

14 **Abstract**

15 We introduce NeuroXiv (neuroxiv.org), a large-scale, AI-powered database that provides
16 detailed 3D morphologies of individual neurons mapped to a standard brain atlas, designed to
17 support a wide array of dynamic, interactive neuroscience applications. NeuroXiv offers a
18 comprehensive collection of 175,149 atlas-oriented reconstructed morphologies of individual
19 neurons derived from more than 518 mouse brains, classified into 292 distinct types and
20 mapped into the Common Coordinate Framework Version 3 (CCFv3). Different from
21 conventional static brain atlases that are often limited to data-browsing, NeuroXiv allows
22 interactive analyses as well as uploading and databasing custom neuron morphologies, which
23 are mapped to the brain atlas for objective comparisons. Powered by a cutting-edge AI engine
24 (AIPOM), NeuroXiv enables dynamic, user-specific analysis and data mining. We specifically
25 developed a mixture-of-experts algorithm to harness the capabilities of multiple large language
26 models. We also developed a client program to achieve more than 10 times better performance
27 compared to a typical server-side setup. We demonstrate NeuroXiv's scalability, efficiency,
28 flexibility, openness, and robustness through various applications.

29

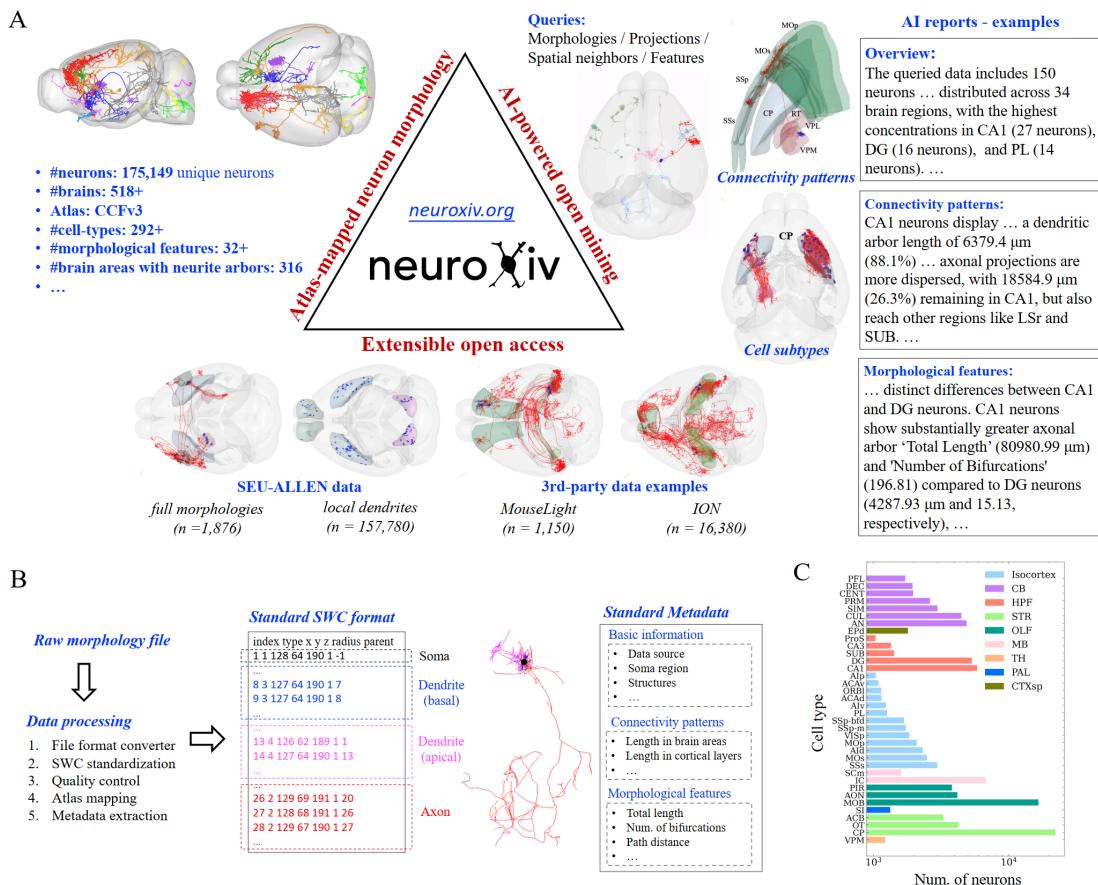
30 **Main**

31 Neuronal morphologies, characterized by their diverse branching patterns and anatomical
32 arborization, provide critical insights into cell types and brain functional networks (Zeng &
33 Sanes, 2017; Luo, 2021; Zeng, 2022). Recent advancements, including sparse labeling
34 techniques (Aransay et al., 2015; Karube et al., 2004; Rotolo et al., 2008), high-resolution brain
35 imaging (Economou et al., 2016; Gong et al., 2016; Zhong et al., 2021), terabyte-scale image
36 handling (Bria et al., 2016; Peng et al., 2017; Y. Wang et al., 2019), and neuron tracing methods
37 (Peng et al., 2010, 2011; Xiao & Peng, 2013; Feng et al., 2015; Jiang et al., 2022; Manubens-
38 Gil et al., 2023), have significantly enhanced our capability to digitize brain-wide neuron
39 morphologies. As a result, there has been a substantial increase in the volume of publicly
40 accessible neuron morphologies, which supports various quantitative analyses, including the
41 morphological characteristics of individual neurons (Peng et al., 2021; Winnubst et al., 2019),
42 dendritic microenvironments (Y. Liu et al., 2023), neuron typing (L. Liu et al., 2023; Xiong et
43 al., 2024), and the organizational principles of neuron projections (Gao et al., 2022, 2023; Jiao
44 et al., 2023; Qiu et al., 2024). However, a remarkable gap exists: how to harness these valuable
45 datasets from diverse sources for new knowledge discoveries while addressing dynamic needs
46 throughout the development process (Wilkinson et al., 2016; Martone, 2024).

47 Current data dissemination solutions for neuron morphology (Akram et al., 2018; Winnubst et
48 al., 2019; Kenney et al., 2024; Qiu et al., 2024) generally fall into two categories: browser-
49 based atlases, such as Neuron Browser (mouselight.janelia.org) and Digital Brain
50 (mouse.digital-brain.cn), and archiving platforms, including Brain Image Library
51 (brainimagelibrary.org) and NeuroMorpho.Org (neuromorpho.org). Archiving platforms
52 typically provide dataset downloads and offer a broader range of data from various sources,
53 while browser platforms furnish additional data exploration tools—such as visualization and
54 statistical analysis—but often limit access to data from their respective laboratories (**Extended**
55 **Data Table 1 and Supplementary Note 2**). Large-scale offline analyses, which require
56 aggregating datasets from multiple sources, introduce further neuroinformatics challenges.
57 These include dataset harmonization, alignment with common coordinate frameworks (CCFs),
58 extraction of key metadata such as morphological and anatomical features, and transforming
59 neuronal features and patterns into meaningful insights. Addressing these challenges demands
60 both domain-specific expertise and advanced coding skills (**Supplementary Note 1**).

61 We introduce the NeuroXiv platform (neuroxiv.org), currently hosted on Amazon Web
62 Services (AWS), designed to address challenges in databasing and mining brain-wide neuron
63 morphometry. Building upon the foundational work of the Allen Brain Atlas (Q. Wang et al.,
64 2020) and NeuroMorpho.Org, we have expanded efforts to establish a standardized atlas-

65 oriented database of neuron morphometry (**Fig. 1A and Methods**). To address challenges
 66 associated with large-scale analysis of neuron morphology, we have developed the AI-Powered
 67 Open Mining (AIPOM) engine. This engine offers functionalities such as searching and
 68 visualizing neuron morphometry, enabling analyses including data statistics, cell typing, and
 69 connectivity studies. Crucially, it incorporates advanced capabilities, such as generating AI-
 70 driven mining reports (**Fig. 1A and Fig. 2A**).



71

72 **Fig. 1 | Overview of NeuroXiv, an open, AI-assisted database for interactive brain-wide**
 73 **neuron analysis.** **A,** *NeuroXiv is founded on three pillars: atlas-mapped neuron morphology,*
 74 *AI-powered open mining, and extensible open access. AIPOM engine enables users to*
 75 *efficiently and flexibly retrieve neuronal data, explore neuron types and connectivity patterns,*
 76 *and offers an intelligent mining tool for generating comprehensive data mining reports.* **B,** *the*
 77 *dataset standardization process in NeuroXiv is performed server-side, which is crucial for*
 78 *ensuring data reusability and interoperability. This process involves formatting raw*
 79 *morphology files into the standard SWC format and storing them accordingly. Additionally, the*
 80 *data is mapped to the same atlas space, and rich metadata is extracted to enhance the dataset's*
 81 *utility.* **C,** *NeuroXiv has established the largest and most comprehensive dataset of neuron types,*

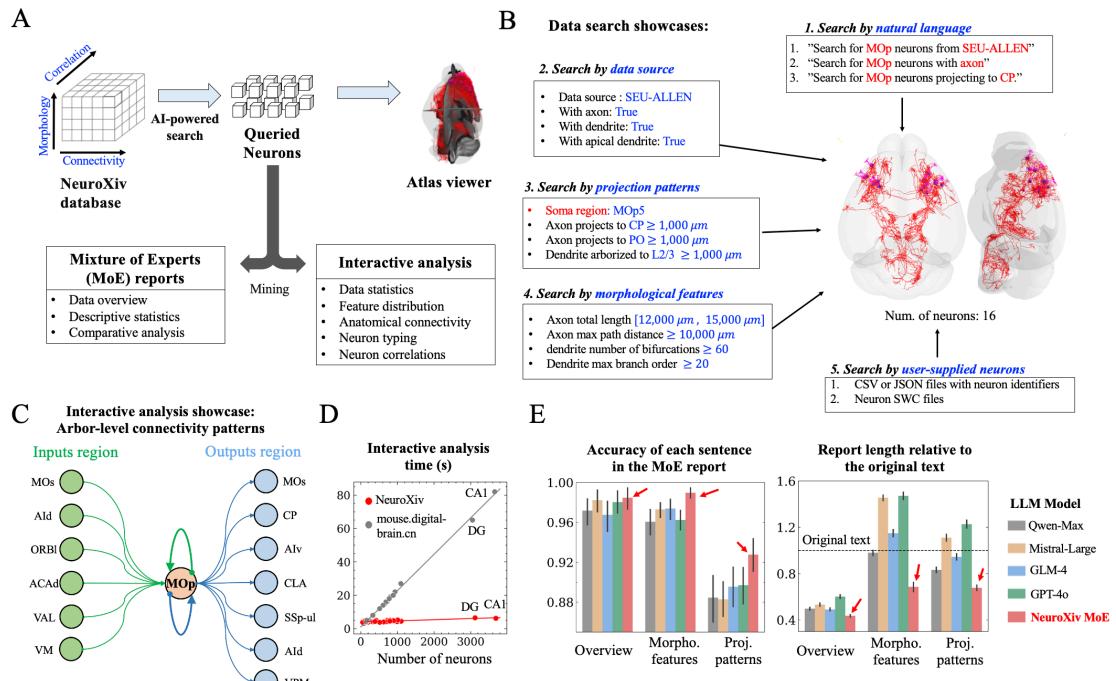
82 covering a wide range of brain regions including the TH, STR, Isocortex, HPF, and CB. A full
83 list of abbreviations for all brain structures in this study is provided in **Methods**.

84 We initiate by establishing a server-side data processing pipeline capable of continuously
85 aggregating publicly available datasets into our database (**Fig. 1B**). Interoperability is
86 maintained through standardization of neuron morphology in the widely used SWC format
87 (Mehta et al., 2023). Reusability is ensured by initially mapping neuron morphologies into the
88 Common Coordinate Framework Version 3 (CCFv3) (Q. Wang et al., 2020), a widely
89 recognized brain atlas for registering neuroanatomical data. We systematically document
90 comprehensive metadata, encompassing basic information, morphological features, and
91 anatomical arborization characteristics of neuron morphology (**Methods and Extended Data**
92 **Table 2**). All morphology data and their corresponding metadata are accessible and
93 downloadable through our web portal. Additionally, we have established an independent data
94 server (download.neuroxiv.org) on AWS to expedite the download of entire standardized
95 datasets.

96 We demonstrate the feasibility and scalability of the databasing method by consolidating data
97 from diverse sources, including our SEU-ALLEN datasets (Peng et al., 2021; Y. Liu et al., 2023)
98 and third-party examples (Winnubst et al., 2019; Gao et al., 2022, 2023; Qiu et al., 2024),
99 culminating in the largest database of brain-wide neuron morphologies to date (**Fig. 1A,**
100 **Methods and Extended Data Fig. 1**). The database features over 175,149 atlas-oriented
101 reconstructed morphologies of individual neurons derived from more than 518 mouse brains.
102 Each neuron reconstruction is characterized by its structural components (soma, axon, and
103 dendrite), alongside metadata documented using a common standardized description method
104 (**Supplementary Table 1**). The resource provides access to the most comprehensive atlas of
105 cell types based on soma anatomical locations, encompassing 12 major gray matter divisions
106 and including data from 292 out of 316 brain structures in the Allen Reference Atlas (ARA)
107 ontology (**Fig. 1C**).

108 Our database offers several advantages for neuron morphology research. By mapping neuron
109 morphologies from diverse sources onto a unified coordinate system, it enables rapid access to
110 neuronal data without the need to switch between different sources. This data aggregation
111 facilitates detailed, data-driven analyses of neuronal morphological characteristics and
112 enhances the study of brain connectivity at the single-cell level (L. Liu et al., 2023). Specifically,
113 we have identified a greater number of incoming neurons that extend their arbors into specific
114 brain regions (**Extended Data Fig. 3A**) and have uncovered additional projection combinations
115 of target regions formed by individual neurons (**Extended Data Fig. 3B**). Additionally, the
116 database's enhanced indexing system allows for improved retrieval of neurons based on spatial

117 proximity, morphological similarity, or shared arborization patterns (**Extended Data Fig. 2**
 118 and **Extended Data Fig. 3C-D**). This advanced indexing opens up new research opportunities,
 119 such as investigating whether spatially adjacent neurons consistently share similar
 120 morphological or projection characteristics (**Extended Data Fig. 2A-B**).



121

122 **Fig. 2 | AI-powered engine for open analysis and mining of neuron data.** **A**, the schematic
 123 diagram of data analysis and mining within NeuroXiv. The NeuroXiv database provides
 124 extensive morphology data, along with detailed morphological features, connectivity patterns,
 125 and inter-data correlations. The AI-powered search tool assists users in extracting relevant
 126 data of interest. Users can then interactively visualize the retrieved data on an atlas viewer and
 127 explore data features in depth. Additionally, NeuroXiv includes a Mixture of Experts (MoE)
 128 module that automatically generates reports to describe data characteristics and uncover
 129 patterns within the data. **B**, the showcases of data search in NeuroXiv demonstrate five common
 130 use cases. Users can search data using natural language queries, followed by searches based
 131 on the distribution of specific features provided by the user. Additionally, data can be retrieved
 132 using a user-supplied list of neurons, further enhancing the flexibility of data searches. **C**, the
 133 interactive analysis showcase in NeuroXiv presents an example where users can study arbor-
 134 level connectivity patterns. The input and output regions are determined by the spatial
 135 proximity of neurons within the database and the retrieved data, allowing for detailed
 136 exploration of neuronal connectivity. **D**, the comparison of interactive analysis time between
 137 the two platforms focuses on a shared cell type analysis scenario. The time measurement
 138 includes the entire process from data retrieval to the rendering of analysis charts. **E**, MoE

139 reports were evaluated against four LLMs (Qwen-Max, Mistral-Large, GLM-4, and GPT-4o)
140 for accuracy and text length across 100 random analysis cases. Our MoE showed higher
141 accuracy and conveyed the same information using shorter text.

142 The AIPOM engine streamlines the knowledge development workflow on the established
143 neuron morphometry database by integrating large language models (LLMs) (OpenAI et al.,
144 2024; Touvron et al., 2023; Yang et al., 2024), which have rapidly advanced in recent years
145 and demonstrated effectiveness in various domains (Bzdok et al., 2024; “Embedding AI in
146 Biology,” 2024) due to their robust text comprehension capabilities. With AIPOM, users can
147 flexibly define their data cohort of interest using natural language queries or a rule-based search
148 panel. The LLM-based mining tool automatically transforms queried neuron data into
149 comprehensive reports, including data overviews, descriptions of morphological features and
150 connectivity patterns, and comparative analyses among cell types. Simultaneously, NeuroXiv
151 provides an interactive analysis tool capable of generating a wide array of quantitative results
152 for queried neurons, including the distribution of data attributes, morphological characteristics,
153 and anatomical projection patterns through detailed visualizations (**Fig. 2A and**
154 **Supplementary Fig. 1**). Additionally, NeuroXiv integrates an enhanced visualization tool for
155 the interactive exploration of complex tree-like neuronal structures (**Methods, Extended Data**
156 **Fig. 8, and Supplementary Fig. 6-7**).

157 We demonstrate the board applicability and flexibility of data search tool through several
158 showcases of querying MOp neurons (**Fig. 2B**). We first illustrate that searches can be
159 conducted using an LLM-based method, enabling users to query for specific neuron types,
160 neurons with particular structures (such as axons or dendrites), or those exhibiting specific
161 projection patterns through natural language inputs (**Supplementary Fig. 2**). We then show
162 that precise searches can be performed by setting customized criteria based on neuron metadata
163 (**Extended Data Fig. 4A-C and Supplementary Fig. 3**). NeuroXiv further provides a database
164 interface enabling users to index specific neurons via an upload function, facilitating its use as
165 a downstream exploration tool following user-defined neuron classification (**Supplementary**
166 **Fig. 4**). Moreover, the search tool supports advanced capabilities such as similarity searches to
167 identify neurons with comparable morphological features and arborization patterns (**Extended**
168 **Data Fig. 2A-B and Supplementary Fig. 8**), as well as neighboring neuron queries to explore
169 arbor-level connectivity within the brain (**Extended Data Fig. 2C-D and Supplementary Fig.**
170 **9**).

171 We highlight the interactive analysis capabilities of AIPOM through two studies. In the first
172 study, users can identify which neuron types in the database provide input to a single neuron
173 and determine the brain regions that receive projections from that neuron (**Fig. 2C**). For the

174 second study, focusing on projection patterns, we use VPM neurons from our database as an
175 example (**Extended Data Fig. 9**). These VPM neurons, sourced from two datasets, exhibit
176 axonal arbors that extend across the CP into multiple cortical regions, including MOs, MOp,
177 SSp, and SSs, encompassing various projection subtypes (**Extended Data Fig. 9A**).
178 Additionally, our visualization tools allow users to observe the selectivity of different
179 projection subtypes across cortical regions and layers (**Extended Data Fig. 9B-C**), as well as
180 compare soma distribution and morphological features among these subtypes (**Extended Data**
181 **Fig. 9D-E**). These findings align with prior knowledge of VPM neuron projection patterns
182 (Peng et al., 2021; Y. Liu et al., 2023).

183 NeuroXiv also demonstrates high efficiency in online analyses (**Fig. 2D**). We perform
184 benchmark tests comparing NeuroXiv and the Digital Brain platform by analyzing the same
185 data categories and measuring the time to render results. Despite generating a broader range of
186 analyses than the Digital Brain platform, NeuroXiv completes most tasks within 4-5 seconds
187 and is largely unaffected by increases in data volume. In contrast, the Digital Brain platform
188 exhibits sensitivity to data size, with response times increasing linearly. For example,
189 generating results for the same number of CA1 neurons takes over 80 seconds on the Digital
190 Brain platform, nearly 20 times slower than NeuroXiv.

191 To improve the integration of LLMs into AIPOM, we implement two key optimizations. First,
192 to address the challenges of unpredictable outputs and occasional inaccuracies (Jin et al., 2024),
193 we developed an advanced Mixture of Experts (MoE) framework for more reliable mining
194 reports (**Fig. 2A and Extended Data Fig. 5-7**). This framework operates in three stages: first,
195 a program generates standardized reports that capture all relevant data details; second, multiple
196 LLM experts analyze and summarize these reports from a data scientist's perspective; and third,
197 a separate LLM reviews the outputs for accuracy and consistency, producing the final report.
198 This multi-expert approach allows MoE to deliver comprehensive data overviews while
199 effectively identifying morphological and projection differences. Our tests show that the MoE
200 framework yields higher accuracy and more concise reports compared to those generated by a
201 single LLM (**Fig. 2E and Methods**).

202 Second, to address the computational demands of server-side LLM deployment, we offer a
203 client-side solution using a natural language processing (NLP) model and a supervised decision
204 tree. This approach transforms natural language queries into actionable search operations within
205 2-3 seconds, achieving comparable accuracy with an 12.3-fold improvement in response time
206 compared to LLM-based server-side searches (**Methods and Extended Data Table 3**).

207 In summary, NeuroXiv offers neuroscientists worldwide access to the largest and most
208 comprehensive neuron morphometry resources. It aggregates publicly available neuron datasets
209 from diverse sources, standardizes them into the widely used SWC format, and maps them into
210 the CCFv3 atlas to enhance data reusability. Additionally, NeuroXiv integrates various tools,
211 including search, visualization, and analysis, to facilitate rapid knowledge development.
212 Leveraging advanced LLMs, the platform offers intuitive search functions and generates
213 mining reports, thereby streamlining the extraction of valuable insights from neuron
214 morphology data. To optimize LLM performance, AIPOM employs two approaches: the MoE
215 framework for enhanced report accuracy and a client-side deployment for faster query
216 responses. Together, the databasing method and AIPOM engine create an open, scalable,
217 efficient, and flexible platform for ongoing neuron data reuse in the neuroscience community.
218

219 **References**

- 220
- 221 AI, M. (n.d.). *Mistral AI | Frontier AI in your hands*. Retrieved August 16, 2024, from
222 <https://mistral.ai/>
- 223 Akram, M. A., Nanda, S., Maraver, P., Armañanzas, R., & Ascoli, G. A. (2018). An open
224 repository for single-cell reconstructions of the brain forest. *Scientific Data*, 5(1), Article
225 1. <https://doi.org/10.1038/sdata.2018.6>
- 226 Aransay, A., Rodríguez-López, C., García-Amado, M., Clascá, F., & Prensa, L. (2015). Long-
227 range projection neurons of the mouse ventral tegmental area: A single-cell axon tracing
228 analysis. *Frontiers in Neuroanatomy*, 9. <https://doi.org/10.3389/fnana.2015.00059>
- 229 Bria, A., Iannello, G., Onofri, L., & Peng, H. (2016). TeraFly: Real-time three-dimensional
230 visualization and annotation of terabytes of multidimensional volumetric images. *Nature
231 Methods*, 13(3), Article 3. <https://doi.org/10.1038/nmeth.3767>
- 232 Bzdok, D., Thieme, A., Levkovskyy, O., Wren, P., Ray, T., & Reddy, S. (2024). Data science
233 opportunities of large language models for neuroscience and biomedicine. *Neuron*,
234 112(5), 698–717. <https://doi.org/10.1016/j.neuron.2024.01.016>
- 235 Economo, M. N., Clack, N. G., Lavis, L. D., Gerfen, C. R., Svoboda, K., Myers, E. W., &
236 Chandrashekhar, J. (2016). A platform for brain-wide imaging and reconstruction of
237 individual neurons. *eLife*, 5, e10566. <https://doi.org/10.7554/eLife.10566>
- 238 Embedding AI in biology. (2024). *Nature Methods*, 21(8), 1365–1366.
239 <https://doi.org/10.1038/s41592-024-02391-7>
- 240 Feng, L., Zhao, T., & Kim, J. (2015). neuTube 1.0: A New Design for Efficient Neuron
241 Reconstruction Software Based on the SWC Format. *eNeuro*, 2(1).
242 <https://doi.org/10.1523/ENEURO.0049-14.2014>
- 243 Gao, L., Liu, S., Gou, L., Hu, Y., Liu, Y., Deng, L., Ma, D., Wang, H., Yang, Q., Chen, Z.,
244 Liu, D., Qiu, S., Wang, X., Wang, D., Wang, X., Ren, B., Liu, Q., Chen, T., Shi, X., ...
245 Yan, J. (2022). Single-neuron projectome of mouse prefrontal cortex. *Nature
246 Neuroscience*, 25(4), Article 4. <https://doi.org/10.1038/s41593-022-01041-5>
- 247 Gao, L., Liu, S., Wang, Y., Wu, Q., Gou, L., & Yan, J. (2023). Single-neuron analysis of
248 dendrites and axons reveals the network organization in mouse prefrontal cortex. *Nature
249 Neuroscience*, 26(6), 1111–1126. <https://doi.org/10.1038/s41593-023-01339-y>
- 250 GLM, T., Zeng, A., Xu, B., Wang, B., Zhang, C., Yin, D., Zhang, D., Rojas, D., Feng, G.,
251 Zhao, H., Lai, H., Yu, H., Wang, H., Sun, J., Zhang, J., Cheng, J., Gui, J., Tang, J.,
252 Zhang, J., ... Wang, Z. (2024). *ChatGLM: A Family of Large Language Models from
253 GLM-130B to GLM-4 All Tools* (arXiv:2406.12793). arXiv.
254 <https://doi.org/10.48550/arXiv.2406.12793>
- 255 Gong, H., Xu, D., Yuan, J., Li, X., Guo, C., Peng, J., Li, Y., Schwarz, L. A., Li, A., Hu, B.,
256 Xiong, B., Sun, Q., Zhang, Y., Liu, J., Zhong, Q., Xu, T., Zeng, S., & Luo, Q. (2016).
257 High-throughput dual-colour precision imaging for brain-wide connectome with
258 cytoarchitectonic landmarks at the cellular level. *Nature Communications*, 7(1), Article
259 1. <https://doi.org/10.1038/ncomms12142>

- 260 Guo, S., Zhao, X., Jiang, S., Ding, L., & Peng, H. (2022). Image enhancement to leverage the
261 3D morphological reconstruction of single-cell neurons. *Bioinformatics*, 38(2), 503–512.
262 <https://doi.org/10.1093/bioinformatics/btab638>
- 263 Jiang, S., Wang, Y., Liu, L., Ding, L., Ruan, Z., Dong, H.-W., Ascoli, G. A., Hawrylycz, M.,
264 Zeng, H., & Peng, H. (2022). Petabyte-Scale Multi-Morphometry of Single Neurons for
265 Whole Brains. *Neuroinformatics*, 20(2), 525–536. <https://doi.org/10.1007/s12021-022-09569-4>
- 266 Jiao, Z., Gao, T., Wang, X., Zhang, W., Biglari, N., Boxer, E. E., Steuernagel, L., Ding, X.,
267 Yu, Z., Li, M., Hao, M., Zhou, H., Cao, X., Li, S., Jiang, T., Qi, J., Jia, X., Feng, Z., Ren,
268 B., ... Xu, X. (2023). *Projectome-defined subtypes and modular intra-hypothalamic*
269 *subnetworks of peptidergic neurons* (p. 2023.05.25.542241). bioRxiv.
270 <https://doi.org/10.1101/2023.05.25.542241>
- 271 Jin, H., Zhang, Y., Meng, D., Wang, J., & Tan, J. (2024). *A Comprehensive Survey on*
272 *Process-Oriented Automatic Text Summarization with Exploration of LLM-Based*
273 *Methods* (arXiv:2403.02901). arXiv. <https://doi.org/10.48550/arXiv.2403.02901>
- 274 Karube, F., Kubota, Y., & Kawaguchi, Y. (2004). Axon branching and synaptic bouton
275 phenotypes in GABAergic nonpyramidal cell subtypes. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 24(12), 2853–2865.
276 <https://doi.org/10.1523/JNEUROSCI.4814-03.2004>
- 277 Kenney, M., Vasylyieva, I., Hood, G., Cao-Berg, I., Tuite, L., Laghaei, R., Smith, M. C.,
278 Watson, A. M., & Ropelewski, A. J. (2024). *The Brain Image Library: A Community-Contributed Microscopy Resource for Neuroscientists* (p. 2023.12.22.573024). bioRxiv.
279 <https://doi.org/10.1101/2023.12.22.573024>
- 280 Li, Y., Jiang, S., Ding, L., & Liu, L. (2023). NRRS: A re-tracing strategy to refine neuron
281 reconstruction. *Bioinformatics Advances*, 3(1), vbad054.
282 <https://doi.org/10.1093/bioadv/vbad054>
- 283 Li, Y., Wu, J., Lu, D., Xu, C., Zheng, Y., Peng, H., & Qu, L. (2022). mBrainAligner-Web: A
284 web server for cross-modal coherent registration of whole mouse brains. *Bioinformatics*,
285 38(19), 4654–4655. <https://doi.org/10.1093/bioinformatics/btac549>
- 286 Liu, L., Yun, Z., Manubens-Gil, L., Chen, H., Xiong, F., Dong, H.-W., Zeng, H., Hawrylycz,
287 M., Ascoli, G., & Peng, H. (2023). *Neuronal Connectivity as a Determinant of Cell*
288 *Types and Subtypes*. <https://doi.org/10.21203/rs.3.rs-2960606/v1>
- 289 Liu, Y., Jiang, S., Li, Y., Zhao, S., Yun, Z., Zhao, Z.-H., Zhang, L., Wang, G., Chen, X.,
290 Manubens-Gil, L., Hang, Y., Garcia-Forn, M., Wang, W., Rubeis, S. D., Wu, Z., Osten,
291 P., Gong, H., Hawrylycz, M., Mitra, P., ... Peng, H. (2023). Full-Spectrum Neuronal
292 Diversity and Stereotypy through Whole Brain Morphometry. *Research Square*, rs.3.rs-
293 3146034. <https://doi.org/10.21203/rs.3.rs-3146034/v1>
- 294 Luo, L. (2021). Architectures of neuronal circuits. *Science*, 373(6559), eabg7285.
295 <https://doi.org/10.1126/science.abg7285>
- 296 Manubens-Gil, L., Zhou, Z., Chen, H., Ramanathan, A., Liu, X., Liu, Y., Bria, A., Gillette, T.,
297 Ruan, Z., Yang, J., Radojević, M., Zhao, T., Cheng, L., Qu, L., Liu, S., Bouchard, K. E.,
298 Gu, L., Cai, W., Ji, S., ... Peng, H. (2023). BigNeuron: A resource to benchmark and

- 302 predict performance of algorithms for automated tracing of neurons in light microscopy
303 datasets. *Nature Methods*, 1–12. <https://doi.org/10.1038/s41592-023-01848-5>
- 304 Martone, M. E. (2024). The past, present and future of neuroscience data sharing: A
305 perspective on the state of practices and infrastructure for FAIR. *Frontiers in*
306 *Neuroinformatics*, 17. <https://doi.org/10.3389/fninf.2023.1276407>
- 307 Mehta, K., Ljungquist, B., Ogden, J., Nanda, S., Ascoli, R. G., Ng, L., & Ascoli, G. A.
308 (2023). Online conversion of reconstructed neural morphologies into standardized SWC
309 format. *Nature Communications*, 14(1), 7429. [https://doi.org/10.1038/s41467-023-42931-x](https://doi.org/10.1038/s41467-023-
310 42931-x)
- 311 OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida,
312 D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S.,
313 Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., ... Zoph, B. (2024). *GPT-4*
314 *Technical Report* (arXiv:2303.08774). arXiv. <https://doi.org/10.48550/arXiv.2303.08774>
- 315 Peng, H., Long, F., & Myers, G. (2011). Automatic 3D neuron tracing using all-path pruning.
316 *Bioinformatics*, 27(13), i239. <https://doi.org/10.1093/bioinformatics/btr237>
- 317 Peng, H., Ruan, Z., Atasoy, D., & Sternson, S. (2010). Automatic reconstruction of 3D
318 neuron structures using a graph-augmented deformable model. *Bioinformatics*, 26(12),
319 i38–i46. <https://doi.org/10.1093/bioinformatics/btq212>
- 320 Peng, H., Xie, P., Liu, L., Kuang, X., Wang, Y., Qu, L., Gong, H., Jiang, S., Li, A., Ruan, Z.,
321 Ding, L., Yao, Z., Chen, C., Chen, M., Daigle, T. L., Dalley, R., Ding, Z., Duan, Y.,
322 Feiner, A., ... Zeng, H. (2021). Morphological diversity of single neurons in molecularly
323 defined cell types. *Nature*, 598(7879), Article 7879. [https://doi.org/10.1038/s41586-021-03941-1](https://doi.org/10.1038/s41586-021-
324 03941-1)
- 325 Peng, H., Zhou, Z., Meijering, E., Zhao, T., Ascoli, G. A., & Hawrylycz, M. (2017).
326 Automatic tracing of ultra-volumes of neuronal images. *Nature Methods*, 14(4), Article
327 4. <https://doi.org/10.1038/nmeth.4233>
- 328 Qiu, S., Hu, Y., Huang, Y., Gao, T., Wang, X., Wang, D., Ren, B., Shi, X., Chen, Y., Wang,
329 X., Wang, D., Han, L., Liang, Y., Liu, D., Liu, Q., Deng, L., Chen, Z., Zhan, L., Chen,
330 T., ... Xu, C. (2024). Whole-brain spatial organization of hippocampal single-neuron
331 projectomes. *Science*, 383(6682), eadj9198. <https://doi.org/10.1126/science.adj9198>
- 332 Qu, L., Li, Y., Xie, P., Liu, L., Wang, Y., Wu, J., Liu, Y., Wang, T., Li, L., Guo, K., Wan,
333 W., Ouyang, L., Xiong, F., Kolstad, A. C., Wu, Z., Xu, F., Zheng, Y., Gong, H., Luo,
334 Q., ... Peng, H. (2022). Cross-modal coherent registration of whole mouse brains.
335 *Nature Methods*, 19(1), Article 1. <https://doi.org/10.1038/s41592-021-01334-w>
- 336 Rotolo, T., Smallwood, P. M., Williams, J., & Nathans, J. (2008). Genetically-Directed, Cell
337 Type-Specific Sparse Labeling for the Analysis of Neuronal Morphology. *PLOS ONE*,
338 3(12), e4099. <https://doi.org/10.1371/journal.pone.0004099>
- 339 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B.,
340 Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., & Lample, G.
341 (2023). *LLaMA: Open and Efficient Foundation Language Models* (arXiv:2302.13971).
342 arXiv. <https://doi.org/10.48550/arXiv.2302.13971>
- 343 *VTK - The Visualization Toolkit*. (n.d.). Retrieved August 16, 2024, from <https://vtk.org/>

- 344 Wang, Q., Ding, S.-L., Li, Y., Royall, J., Feng, D., Lesnar, P., Graddis, N., Naeemi, M.,
345 Facer, B., Ho, A., Dolbeare, T., Blanchard, B., Dee, N., Wakeman, W., Hirokawa, K. E.,
346 Szafer, A., Sunkin, S. M., Oh, S. W., Bernard, A., ... Ng, L. (2020). The Allen Mouse
347 Brain Common Coordinate Framework: A 3D Reference Atlas. *Cell*, 181(4), 936-
348 953.e20. <https://doi.org/10.1016/j.cell.2020.04.007>
- 349 Wang, Y., Li, Q., Liu, L., Zhou, Z., Ruan, Z., Kong, L., Li, Y., Wang, Y., Zhong, N., Chai,
350 R., Luo, X., Guo, Y., Hawrylycz, M., Luo, Q., Gu, Z., Xie, W., Zeng, H., & Peng, H.
351 (2019). TeraVR empowers precise reconstruction of complete 3-D neuronal morphology
352 in the whole brain. *Nature Communications*, 10(1), Article 1.
353 <https://doi.org/10.1038/s41467-019-11443-y>
- 354 Wilkinson, M. D., Dumontier, M., Aalbersberg, Ij. J., Appleton, G., Axton, M., Baak, A.,
355 Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., Bouwman, J.,
356 Brookes, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T.,
357 Finkers, R., ... Mons, B. (2016). The FAIR Guiding Principles for scientific data
358 management and stewardship. *Scientific Data*, 3(1), Article 1.
359 <https://doi.org/10.1038/sdata.2016.18>
- 360 Winnubst, J., Bas, E., Ferreira, T. A., Wu, Z., Economo, M. N., Edson, P., Arthur, B. J.,
361 Bruns, C., Rokicki, K., Schauder, D., Olbris, D. J., Murphy, S. D., Ackerman, D. G.,
362 Arshadi, C., Baldwin, P., Blake, R., Elsayed, A., Hasan, M., Ramirez, D., ...
363 Chandrashekhar, J. (2019). Reconstruction of 1,000 Projection Neurons Reveals New Cell
364 Types and Organization of Long-Range Connectivity in the Mouse Brain. *Cell*, 179(1),
365 268-281.e13. <https://doi.org/10.1016/j.cell.2019.07.042>
- 366 Xiao, H., & Peng, H. (2013). APP2: Automatic tracing of 3D neuron morphology based on
367 hierarchical pruning of a gray-weighted image distance-tree. *Bioinformatics (Oxford,*
368 *England)*, 29(11), 1448–1454. <https://doi.org/10.1093/bioinformatics/btt170>
- 369 Xiong, F., Xie, P., Zhao, Z., Li, Y., Zhao, S., Manubens-Gil, L., Liu, L., & Peng, H. (2024).
370 DSM: Deep sequential model for complete neuronal morphology representation and
371 feature extraction. *Patterns*, 5(1), 100896. <https://doi.org/10.1016/j.patter.2023.100896>
- 372 Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F.,
373 Dong, G., Wei, H., Lin, H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J., ... Fan,
374 Z. (2024). *Qwen2 Technical Report* (arXiv:2407.10671). arXiv.
375 <https://doi.org/10.48550/arXiv.2407.10671>
- 376 Zeng, H. (2022). What is a cell type and how to define it? *Cell*, 185(15), 2739–2755.
377 <https://doi.org/10.1016/j.cell.2022.06.031>
- 378 Zeng, H., & Sanes, J. R. (2017). Neuronal cell-type classification: Challenges, opportunities
379 and the path forward. *Nature Reviews Neuroscience*, 18(9), Article 9.
380 <https://doi.org/10.1038/nrn.2017.85>
- 381 Zhao, Z.-H., Liu, L., & Liu, Y. (2024). NIEND: Neuronal image enhancement through noise
382 disentanglement. *Bioinformatics*, 40(4), btae158.
383 <https://doi.org/10.1093/bioinformatics/btae158>
- 384 Zhong, Q., Li, A., Jin, R., Zhang, D., Li, X., Jia, X., Ding, Z., Luo, P., Zhou, C., Jiang, C.,
385 Feng, Z., Zhang, Z., Gong, H., Yuan, J., & Luo, Q. (2021). High-definition imaging

386 using line-illumination modulation microscopy. *Nature Methods*, 18(3), 309–315.
387 <https://doi.org/10.1038/s41592-021-01074-x>
388
389

390 **Methods**

391 **Nomenclature and abbreviations of brain regions**

392 The 12 “major divisions” of gray matter in the Allen Reference Atlas (ARA) ontology:
393 Isocortex, Olfactory areas (OLF), Hippocampal formation (HPF), Cortical subplate (CTXsp),
394 Striatum (STR), Pallidum (PAL), Thalamus (TH), Hypothalamus (HY), Midbrain (MB), Pons
395 (P), Medulla (MY), and Cerebellum (CB).

396 Isocortex: primary motor area (MOp), secondary motor area (MOs), primary somatosensory
397 area (SSp), supplemental somatosensory area (SSs), gustatory area (GU), visceral area (VISC),
398 dorsal auditory area (AUDd), primary auditory area (AUDp), posterior auditory area (AUDpo),
399 ventral auditory area (AUDv), primary visual area (VISp), anterior cingulate area, dorsal part
400 (ACAd), anterior cingulate area, ventral part (ACAv), prelimbic area (PL), infralimbic area
401 (ILA), orbital area, lateral part (ORBl), orbital area, medial part (ORBm), orbital area,
402 ventrolateral part (ORBvl), agranular insular area, dorsal part (AId), agranular insular area,
403 posterior part (Alp), agranular insular area, ventral part (AIv), retrosplenial area, ventral part
404 (RSPv), temporal association area (TEa).

405 Olfactory areas (OLF): piriform area (PIR).

406 Hippocampal formation (HPF): hippocampal region (HIP), fields CA1, CA2, CA3, dentate
407 gyrus (DG), entorhinal area, lateral part (ENTl), entorhinal area, medial part (ENTm),
408 parasubiculum (PAR), postsubiculum (POST), presubiculum (PRE), subiculum (SUB),
409 prosubiculum (ProS).

410 Cortical subplate (CTXsp): claustrum (CLA).

411 Cerebral nuclei (CNU): striatum (STR), caudoputamen (CP), nucleus accumbens (ACB),
412 globus pallidus, external segment (GPe), globus pallidus, internal segment (GPi).

413 Thalamus (TH): ventral anterior-lateral complex (VAL), ventral medial nucleus (VM), ventral
414 posterolateral nucleus (VPL), ventral posterolateral nucleus, parvicellular part (VPLpc), ventral
415 posteromedial nucleus (VPM), ventral posteromedial nucleus, parvicellular part (VPMpc),
416 medial geniculate complex, dorsal part (MGd), lateral geniculate complex, dorsal part (LGd),
417 lateral posterior nucleus (LP), posterior complex (PO), anteromedial nucleus (AM),
418 mediodorsal nucleus (MD), submedial nucleus (SMT), paraventricular nucleus (PVT), reticular
419 nucleus (RT).

420 Hypothalamus (HY): subthalamic nucleus (STN), zona incerta (ZI).

421 Midbrain (MB): substantia nigra, reticular part (SNr), midbrain reticular nucleus (MRN).

422 **NeuroXiv platform**

423 The architecture of NeuroXiv (neuroxiv.org) is designed to support large-scale analysis of
424 brain-wide neuron morphometry, utilizing a cohesive and highly integrated technology stack.

425 *Frontend*

426 The frontend of NeuroXiv is developed using *Vue.js* (*v2.6.12*), a progressive JavaScript
427 framework known for its efficiency in building dynamic and responsive single-page
428 applications. To create a user-friendly and visually engaging interface, *Element Plus* (*v2.7.3*),
429 a Vue 3-based component library, is employed. Additionally, *Three.js* (*v0.134.0*) is integrated
430 into the frontend to handle the rendering of complex 3D visualizations, including brain regions
431 and neuron reconstructions. This powerful WebGL-based library allows for detailed and
432 interactive 3D models, providing users with an immersive experience in exploring
433 neuroanatomical data.

434 *Backend*

435 The backend is constructed with *Python* (*v3.9.12*) and *Flask* (*v3.0.0*), a lightweight WSGI web
436 application framework. Flask serves as the backbone of the server-side architecture, enabling
437 seamless communication between the frontend and the database. *SQLite* (*v3.38.2*) is used as the
438 database engine, offering a self-contained, serverless solution for efficient data storage and
439 retrieval. This setup ensures that the platform remains agile and capable of handling the
440 substantial datasets inherent to neuroinformatics research.

441 To manage web traffic and optimize performance, *Nginx* (*v1.24.0*) is deployed as a reverse
442 proxy server. Nginx efficiently distributes incoming requests across backend processes,
443 enhancing the platform's ability to support a high volume of concurrent users while maintaining
444 fast response times and secure connections.

445 NeuroXiv is hosted on Amazon Web Services (AWS) with 4 CPUs, 32 GB of memory, and
446 12.5 Gbps network bandwidth, leveraging AWS's scalable and resilient cloud infrastructure to
447 provide reliable access for users worldwide. This deployment strategy ensures that the platform
448 remains highly available and capable of scaling in response to increasing user demand, thereby
449 offering a stable and responsive environment for researchers. And we also noticed that the
450 current performance of NeuroXiv on AWS is somewhat compromised and is expected to
451 perform better on a server with the addition of more CPUs and acceleration through GPU
452 devices.

453 **Datasets**

454 Currently, the NeuroXiv database reports the integration of several brain-wide neuron
455 morphology datasets shared by the community. Each dataset will be described in detail in the
456 following sections. In the future, we plan to continuously add new mouse brain datasets and
457 encourage users to contribute their own datasets to the NeuroXiv platform. Additionally, in
458 upcoming updates, we plan to incorporate the BigNeuron Project (Manubens-Gil et al., 2023)—
459 a community-contributed resource for benchmarking neuron morphology auto-tracing
460 algorithms—into NeuroXiv, providing ongoing support for users worldwide.

461 *SEU-ALLEN full dataset*

462 This dataset (Peng et al., 2021; Y. Liu et al., 2023) was initially generated using a semi-
463 automatic annotation pipeline with 1,741 neurons and has been expanded to 1,876 neurons with
464 improved quality (Li et al., 2023). Each neuron includes fully traced axonal and dendritic arbors,
465 with 512 apical dendrites additionally annotated. The data mainly covers neurons in the VPM
466 (389 neurons), CP (324 neurons), and lots of cortical regions.

467 *SEU-ALLEN local dataset*

468 This dataset was produced using an automatic tracing method described in our previous work
469 (Y. Liu et al., 2023). Initially, image volumes centered on the cell body (soma) were extracted
470 from whole-brain image data, with a size greater than 200 μm in each dimension, sufficient to
471 capture most of the neuron's dendritic arbor. These image volumes were then processed using
472 image enhancement algorithms (Guo et al., 2022) to improve image quality. Automatic
473 reconstructions were generated and cross-validated using the APP2 (Peng et al., 2011; Xiao &
474 Peng, 2013) and NeuTube (Feng et al., 2015) algorithms, followed by neurite fiber pruning to
475 remove extraneous signals (Zhao et al., 2024). In NeuroXiv, we retained only the neurite
476 segments within 100 μm of the soma to ensure consistency. This dataset contains 155,743
477 neurons, which are extensively distributed across various brain regions such as CP, MOB, OT,
478 AON, and PIR.

479 *ION datasets*

480 Currently, NeuroXiv integrates two datasets (Gao et al., 2022, 2023; Qiu et al., 2024) from ION:
481 a prefrontal cortex dataset comprising 6,357 neurons (Xiaofei Wang. (2023). Single-neuron
482 projectome of mouse prefrontal cortex (with dendrite). Brain Science Data Center, Chinese
483 Academy of Sciences. <https://cstr.cn/33145.11.BSDC.1689837400.1681922768243666945>
484 and <https://doi.org/10.12412/BSDC.1690164952.20001>), and a hippocampus dataset

485 consisting of 10,100 neurons (Xiaofei Wang. (2024). Single-neuron datasets for mouse
486 hippocampus. Brain Science Data Center, Chinese Academy of Sciences.
487 <https://cstr.cn/33145.11.BSDC.1667284058.1585980235450376194> and
488 <https://doi.org/10.12412/BSDC.1667278800.20001>). During data integration, 77 neurons with
489 indeterminate soma locations were excluded, resulting in a final dataset comprising 16,380 fully
490 reconstructed axons and 6,106 fully reconstructed dendrites. The neurons are primarily
491 distributed across brain regions such as CA1 (3,657 neurons), DG-sg (2,618 neurons), SUB
492 (934 neurons), and CA3 (887 neurons).

493 *MouseLight dataset*

494 The MouseLight project (Winnubst et al., 2019) currently publishes data on 1,200 neurons
495 available at MouseLight NeuronBrowser (<http://ml-neuronbrowser.janelia.org>). During data
496 integration, 50 neurons with somas located in fiber tracts were excluded, resulting in 1,150
497 neurons from MouseLight being included in NeuroXiv. This dataset contains 1,150 fully
498 reconstructed axons and 1,138 fully reconstructed dendrites. The neurons are distributed across
499 various brain regions, such as MOs, SUB, PRE, VAL, DG-mo, and VPM, with some overlap
500 with data from SEU-ALLEN and ION.

501 **Data aggregation**

502 Data aggregation in NeuroXiv involves collecting datasets from our own datasets (SEU-
503 ALLEN) and third-party sources like ION and MouseLight, and processing the data to convert
504 it into a consolidated format.

505 *Data Format Conversion*

506 SEU-ALLEN datasets have already been processed into the standardized SWC format (Mehta
507 et al., 2023) and registered to the CCFv3 atlas. Therefore, our focus here is on processing the
508 datasets from ION and MouseLight. The ION and MouseLight datasets had different format
509 issues. We standardized the ION datasets into the SWC format, aligning the structure domain
510 types, for example, soma (type label = 1), axon (type label = 2), basal dendrite (type label = 3),
511 and apical dendrite (type label = 4). We also converted the neuron reconstruction data from the
512 MouseLight dataset from JSON files into SWC files.

513 *Quality Control*

514 We first performed a quality screening process to ensure data usability, filtering out non-
515 compliant data. The specific steps included:

- 516 1. *Single Connected Tree*: Ensuring that all nodes have only one parent node, tracing back to
517 a single root node (soma).
- 518 2. *Root Node (Soma) Labeling*: Verifying that there is exactly one node labeled as type=1
519 with parent=-1.
- 520 3. *Structure Domain Type Correctness*: Confirming that type attributes 1-4 are valid and that
521 the type attribute remains consistent when tracing from terminal nodes back to the root.
- 522 4. *SWC Tree Structure Integrity*: Checking for the presence of loops and trifurcations in the
523 SWC tree structure.

524 *Atlas Mapping*

525 Using mBrainAligner (Li et al., 2022; Qu et al., 2022), we mapped all data points to the CCFv3
526 atlas. We then resampled the atlas-oriented reconstruction data, ensuring that the distance
527 between parent and child nodes was set to 1 μm .

528 *Data Curation*

529 All data points were renamed to follow a standardized format:
530 “<resource_name>_<full/local>_<brainid>_<neuronid>_..._<atlas>”. For example: "SEU-
531 ALLEN_full_17302_00001_CCFv3".

532 *Metadata Extraction*

533 We first extracted basic information such as the soma region for each neuron in the dataset.
534 Then, using the atlas annotation template, we calculated the arborization strength for axons and
535 dendrites across different brain regions based on neurite length. Finally, we extracted
536 morphological features for axons and dendrites (**Extended Data Table 2**).

537 We further generated a list of morphology similar neurons for each neuron based on the
538 distances between their morphological features. Additionally, we created a list of projection
539 similar neurons by calculating the distance between point clouds formed by key axonal nodes
540 (soma, bifurcation nodes, and terminal nodes). We also defined two types of neighboring
541 neurons based on the overlap between neuron arbors: axon neighboring neurons and dendrite
542 neighboring neurons. All metadata and proximity or similarity tables between data points have
543 been imported into an SQL database for easy user access.

544

545 **Visualization**

546 *Brain Atlas Visualization*

547 We use the Visualization Toolkit (VTK) to render brain atlases and neuron morphologies on
548 the website. For brain atlas visualization, we start with annotation template files that record
549 spatial coordinates for different brain regions, based on the ARA brain region table and brain
550 templates. We apply the marching cubes algorithm to obtain the mesh contours for each brain
551 region and generate corresponding mesh models. These models are then smoothed using the
552 Laplacian smoothing algorithm. To optimize performance, we reduce the number of triangles
553 in the mesh using the progressive mesh decimation algorithm while preserving the geometric
554 information and other attributes. As a result, VTK files for 838 brain regions were generated
555 for the CCFv3 atlas.

556 *Neuron Morphology Visualization*

557 This includes two components: a thumbnail view of each neuron morphology for the Neuron
558 Browser and an interactive visualization in 3D atlas viewer for detailed neuron morphology.
559 The thumbnails are designed to provide a quick overview of neuron morphology and assist
560 users in locating neurons. To achieve this, we resample neuron morphologies with a step size
561 of 100 μm . For generating 2D projection thumbnails, we first use Principal Component
562 Analysis (PCA) to transform the neuron morphology coordinates, ensuring that the projection
563 retains as much structural information as possible. For 3D visualization of neuron morphology,
564 including axons, dendrites and arbors, we convert the neuron morphological structures into
565 renderable line objects (OBJ files) for VTK. Additionally, soma visualization is implemented
566 using Three.js, rendering a sphere with a radius of 50 μm .

567 **Mixture of Experts (MoE)**

568 We have developed a Mixture of Experts (MoE) framework that leverages four large language
569 models (LLMs), each containing trillions of parameters. This system is specifically designed to
570 collaboratively mitigate errors and hallucinations that are commonly associated with LLM-
571 generated content, thereby producing reliable, accurate, and coherent data analysis reports. The
572 MoE framework operates in three distinct stages:

- 573 1. Descriptive reports generation: Initially, data retrieved from the database is
574 programmatically organized into a standardized data description format. This ensures
575 consistency and facilitates accurate analysis by the models.
- 576 2. LLM export reports: The organized data is then independently analyzed by three models—
577 Qwen-Max-0428 (Yang et al., 2024), Mistral-Large-2407 (AI, n.d.), and GLM-4-0520
578 (GLM et al., 2024). Each model is tasked with generating an analysis report based on the
579 following prompts:

580 *2.1 Prompt for Overview*

581 Objective: To provide a concise summary that enhances readability and clarity, with a
582 focus on accurately representing significant numerical values.
583 Methodology: The model is instructed to prioritize larger statistics while summarizing
584 key findings in a coherent paragraph without bullet points.
585 Data Input: "Original statistical data: {data}"
586
587 *2.2 Prompt for Morphological Features Mining*
588 Objective: To analyze neuronal morphology data, particularly focusing on critical features
589 such as 'Total Length' and 'Number of Bifurcations.'
590 Methodology: The model generates a comparative summary that emphasizes the
591 importance of these features, ensuring numerical accuracy throughout.
592 Data Input: "Original neuronal morphology data: {data}"
593
594 *2.3 Prompt for Projection Pattern Mining*
595 Objective: To analyze neuronal projection data, with a specific focus on axon and dendrite
596 projections, and their implications for neuronal connectivity.
597 Methodology: The model produces a summary that highlights the key points related to
598 projection length and strength of connectivity, maintaining numerical precision and
599 coherence.
600 Data Input: "Original neuronal projection data: {data}"
601 3. Report confirmation: The GPT-4o-2024-05-13 model (OpenAI et al., 2024) serves as the
602 final synthesis expert. This model evaluates the analysis reports generated by the three
603 previous models against the original data and synthesizes them into a comprehensive,
604 refined analysis report. The process follows a structured evaluation as outlined below:
605 *3.1 Prompt for Overview*
606 Objective: To assess the precision of three summaries relative to the original statistical
607 dataset insights.
608 Methodology: The model ensures that numerical data in the summaries aligns with the
609 source material. The most accurate summary is then refined into a new summary that
610 enhances readability, brevity, and consistency.
611 Data Input: "Original text: {origin_input}"
612 *3.2 Prompt for Morphological Features Mining*
613 Objective: To meticulously assess the accuracy of three summaries in relation to an
614 original text detailing neuronal morphology data.
615 Methodology: The model compares numerical values, particularly those related to 'Total
616 Length', 'Number of Bifurcations', 'Max Path Distance', and 'Center Shift', ensuring
617 accuracy and consistency in the summaries.

618 Data Input: "Original text: {origin_input}"
619 *3.3 Prompt for Projection Pattern Mining*
620 Objective: To evaluate the precision of summaries concerning neuronal projection
621 characteristics, particularly focusing on axon and dendrite projections as indicators of
622 connectivity strength.
623 Methodology: The model confirms numerical congruity and validates the logical
624 consistency of comparisons in the summaries, generating a final, coherent summary.
625 Data Input: "Original text: {origin_input}"

626 MoE evaluation

627 The evaluation methodology is centered on assessing the accuracy and logical consistency of
628 text summaries by comparing them against a source text. This process is implemented through
629 a custom Python script that systematically evaluates key aspects of the summaries, particularly
630 focusing on numerical data accuracy and logical consistency.

631 *Data Accuracy Evaluation*

632 The evaluation begins by extracting numerical data from both the source text and the generated
633 summaries. A custom function utilizes natural language processing (NLP) tools, such as *spaCy*,
634 to identify numbers within their contextual surroundings. These extracted numbers are then
635 compared between the source text and the summaries to determine how accurately the
636 numerical data has been represented.

637 A data accuracy score is calculated by examining the occurrence and contextual integrity of
638 each number in the summaries relative to the source text. This score reflects the proportion of
639 correctly matched numerical values, providing a quantitative measure of how faithfully the
640 summaries represent the original data.

641 *Logic Consistency Verification*

642 Beyond numerical accuracy, the script also evaluates the logical consistency of the summaries.
643 This involves verifying whether the statements in the summaries logically follow from the
644 information provided in the source text.

645 The script employs a large language model (LLM) to perform this verification. It generates a
646 prompt that includes both the source text and the summary statement in question, asking the
647 model to determine whether the summary statement can be logically and numerically inferred
648 from the source. The output from the LLM is then parsed to decide whether the summary is

649 logically consistent. The logic accuracy score is derived by calculating the percentage of
650 summary sentences that were deemed logically valid.

651 *Comprehensive Evaluation*

652 The script integrates the results from both the data accuracy and logic consistency assessments
653 to provide a comprehensive evaluation of the summaries. By quantifying the alignment of
654 numerical data and logical coherence, the evaluation method offers a robust approach to
655 determining the quality and reliability of text summaries in capturing the essence of the source
656 material.

657 *Documentation and Reporting*

658 The results of the evaluation process, including both data accuracy and logic consistency scores,
659 are meticulously recorded. This documentation includes relevant metadata, such as the models
660 used and the specific instances evaluated, ensuring that the evaluation process is both
661 transparent and reproducible for further analysis and refinement.

662 **AI-powered natural language search**

663 Our framework integrates multiple components to achieve accurate and context-aware natural
664 language understanding and data retrieval.

- 665 1. Entity Recognition and Intent Classification: The core of our Natural Language Processing
666 (NLP) framework is built on a combination of machine learning models and rule-based
667 systems. A supervised decision tree classifier, trained on a specialized dataset of
668 neuroscience-related queries, is used to recognize key entities such as neuron types, brain
669 regions, and projection relationships. The classifier works alongside rule-based
670 components that handle domain-specific terminology variations, ensuring a robust
671 response to user queries.
- 672 2. Semantic Parsing and Contextual Understanding: The framework employs semantic
673 parsing techniques to accurately extract and interpret user intent from natural language
674 input. It detects complex phrases related to neuroscience, such as neuron classifications
675 and brain region relationships. Using contextual analysis, the system discerns detailed
676 query intents (e.g., "projection from region X to region Y"), allowing precise and relevant
677 data to be retrieved.
- 678 3. Dynamic Mapping and Knowledge Integration: The framework integrates domain-specific
679 structured schemas to map both full terminologies and their corresponding abbreviations
680 into a standardized format compatible with database queries. This dynamic mapping

681 ensures consistency and accuracy by aligning user input with the system's structured
682 knowledge base. This capability enhances the system's flexibility and robustness in
683 providing comprehensive and relevant responses.

684 4. Multi-Stage Query Processing Pipeline: The NLP module operates through a multi-stage
685 query processing pipeline, encompassing tokenization, entity extraction, context
686 recognition, and result formulation. Each stage is designed to maximize the understanding
687 of user input and generate accurate database queries, providing users with precise and
688 comprehensive results.

689 *Front-End Deployment and Benefits*

690 The NLP framework is deployed on the front end using Vue.js, which brings two significant
691 advantages:

- 692 1. Protecting User Privacy: By processing queries directly on the client side, the framework
693 ensures that user inputs remain private and are not exposed to external servers. This
694 approach is particularly beneficial in sensitive research settings where data privacy is
695 paramount.
- 696 2. Improved Query Speed and Responsiveness: Client-side processing significantly reduces
697 latency by eliminating unnecessary server round trips. This results in faster response times
698 and a more interactive user experience, enabling researchers to explore neuroscience data
699 efficiently.

700 *Implementation and Model Training*

701 The implementation leverages JavaScript-based libraries combined with tailored AI algorithms
702 optimized for the neuroscience domain. The decision tree model is trained on a diverse set of
703 domain-specific queries to ensure robust performance and generalization.

704 **Data availability**

705 Atlas-mapped neuronal morphology data and discovery results—including metadata, figures,
706 and mining text—are available for direct download via the web portal. Additionally, we have
707 set up an AWS server (<https://download.neuroxiv.org>) to facilitate easy access to standardized
708 datasets. VTK files for various brain regions in CCFv3 are also accessible through our GitHub
709 repository (<https://github.com/SEU-ALLEN-codebase/NeuroXiv/VTK>). All the materials
710 available from NeuroXiv should only be used exclusively for academic purposes and must
711 adhere to the CC-BY NC license (<https://creativecommons.org/licenses/by-nc/4.0/legalcode>).

712 **Code availability**

713 The source code of the NeuroXiv project can be obtained from our GitHub repository
714 (<https://github.com/SEU-ALLEN-codebase/NeuroXiv>), where detailed deployment
715 instructions are provided. User manuals and video-tutorials are available at our website
716 (<https://neuroxiv.org>). Data processing utilizes the plugin system provided by the Vaa3D
717 platform (version 4.001, available at <https://github.com/Vaa3D>). The source code for quality
718 control procedures can be found at
719 https://github.com/Vaa3D/vaa3d_tools/tree/master/hackathon/shengdian/NeuroMorphoLib.
720 The source code for mBrainAlinger is accessible at
721 https://github.com/Vaa3D/vaa3d_tools/tree/master/hackathon/mBrainAligner. Atlas mapping
722 can also be conducted via the mBrainAlinger web portal (<http://mbrainaligner.ahu.edu.cn>).

723 **Supplementary Data**

724 **The Supplementary Notes** (*NeuroXiv_supplementary_notes_and_figures.pdf*),
725 **Supplementary Tables** (*Supplementary_Table1.csv*) and **Supplementary Figures**
726 (*NeuroXiv_supplementary_notes_and_figures.pdf*) can be found along with the submission
727 files of this manuscript.

728 **Acknowledgements**

729 This work was mainly supported by several initiatives of neuroscience and a New Cornerstone
730 grant awarded to H.P.. The Southeast University team was also supported by a STI2030-Major
731 Projects Grant No. 2022ZD0205200/2022ZD0205204 awarded to Lijuan Liu. We thank Brain
732 Science Data Center, Chinese Academy of Sciences (<https://braindatacenter.cn/>) for the open
733 sharing of ION datasets. We thank MouseLight project (<http://ml-neuronbrowser.janelia.org>)
734 for the open sharing of its dataset. We extend our gratitude to Lijuan Liu for her contributions
735 to the project design and discussions surrounding the first-generation platform, and to Xuan
736 Zhao for early involvement in the project's development. We also thank Sujun Zhao for
737 assisting with the use of the auto-arbor tool to generate arbor data for all axon-containing
738 neurons, and Xiaoxuan Jiang for testing NeuroXiv's functionalities and drafting the NeuroXiv
739 user manual.

740 **Author contributions**

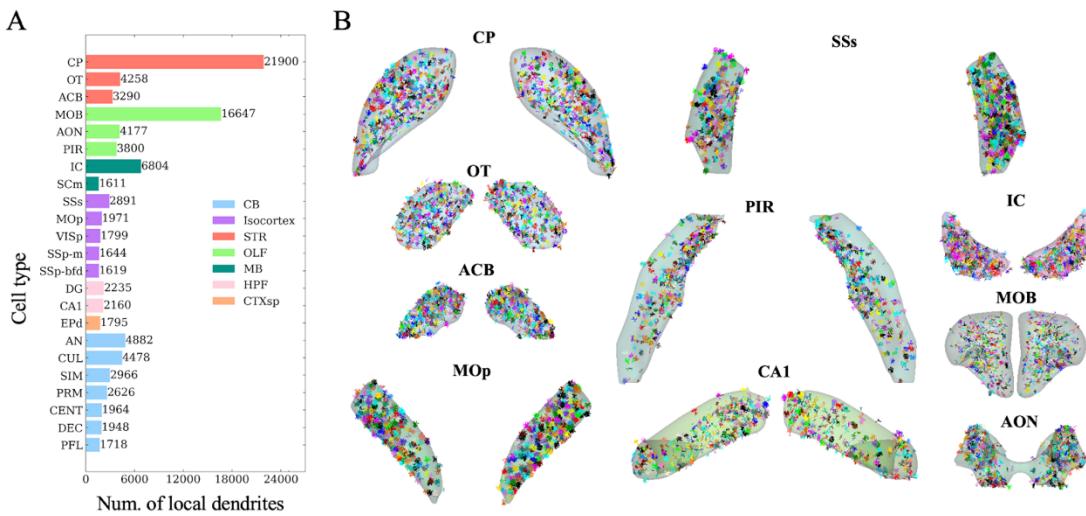
741 H.P. conceptualized and managed this study, and invented AIPOM. Z.Y., H.C., and J.Y.
742 designed the initial version of the NeuroXiv (neuroxiv.net) platform, hosted originally on

743 Tencent Cloud server. Z.Y. was responsible for backend development of the first-generation
744 platform, while H.C. handled frontend coding. S.J. and L.W. co-developed the new version of
745 the NeuroXiv (neuroxiv.org) platform and migrated the servers to the AWS cloud platform.
746 L.W. undertook the majority of website development tasks, including project deployment,
747 backend, and frontend development. S.J. managed the data aggregation tasks, including data
748 standardization processing, metadata generation, and the production of required atlases and
749 reconstruction files (obj files) for the website. H.P. and S.J. wrote the manuscript with
750 assistance of all authors, who reviewed and revised the manuscript.

751 **Competing interests**

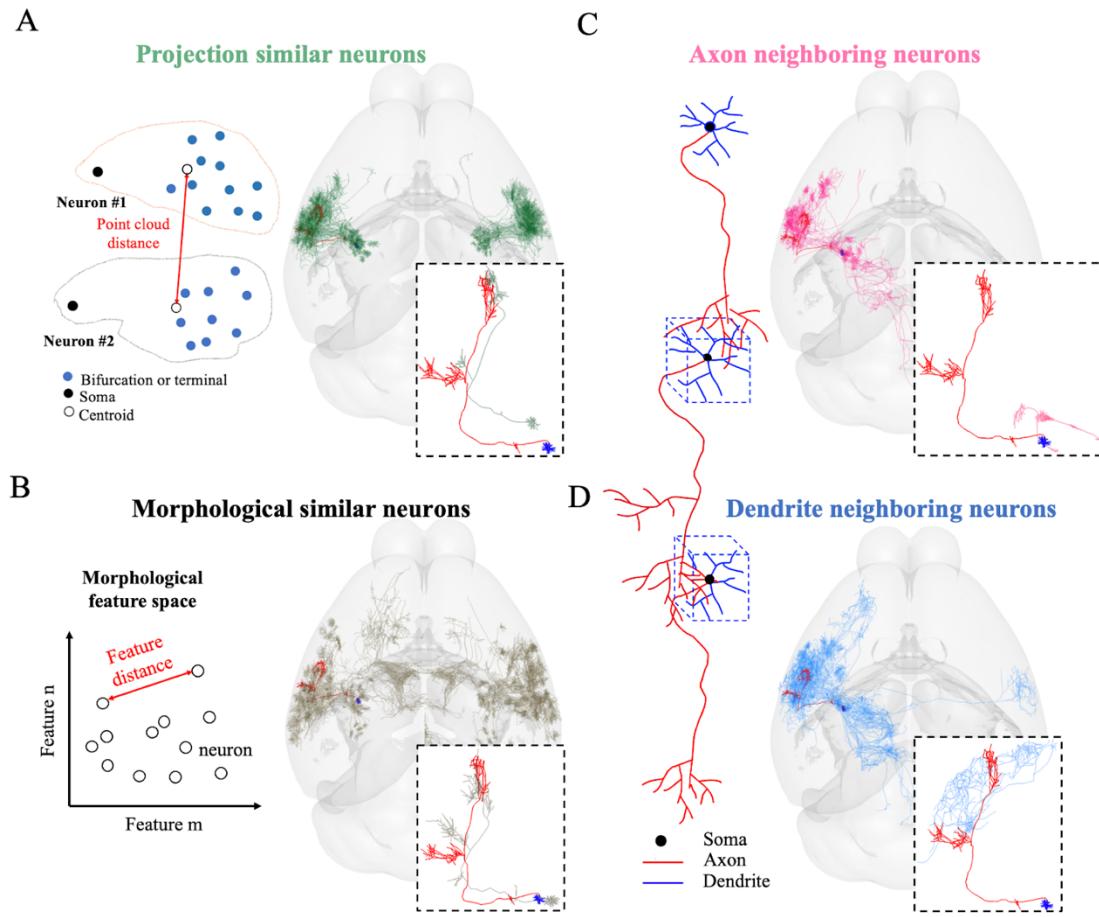
752 The authors declare no competing interests.
753

754 **Extended Data Figures**



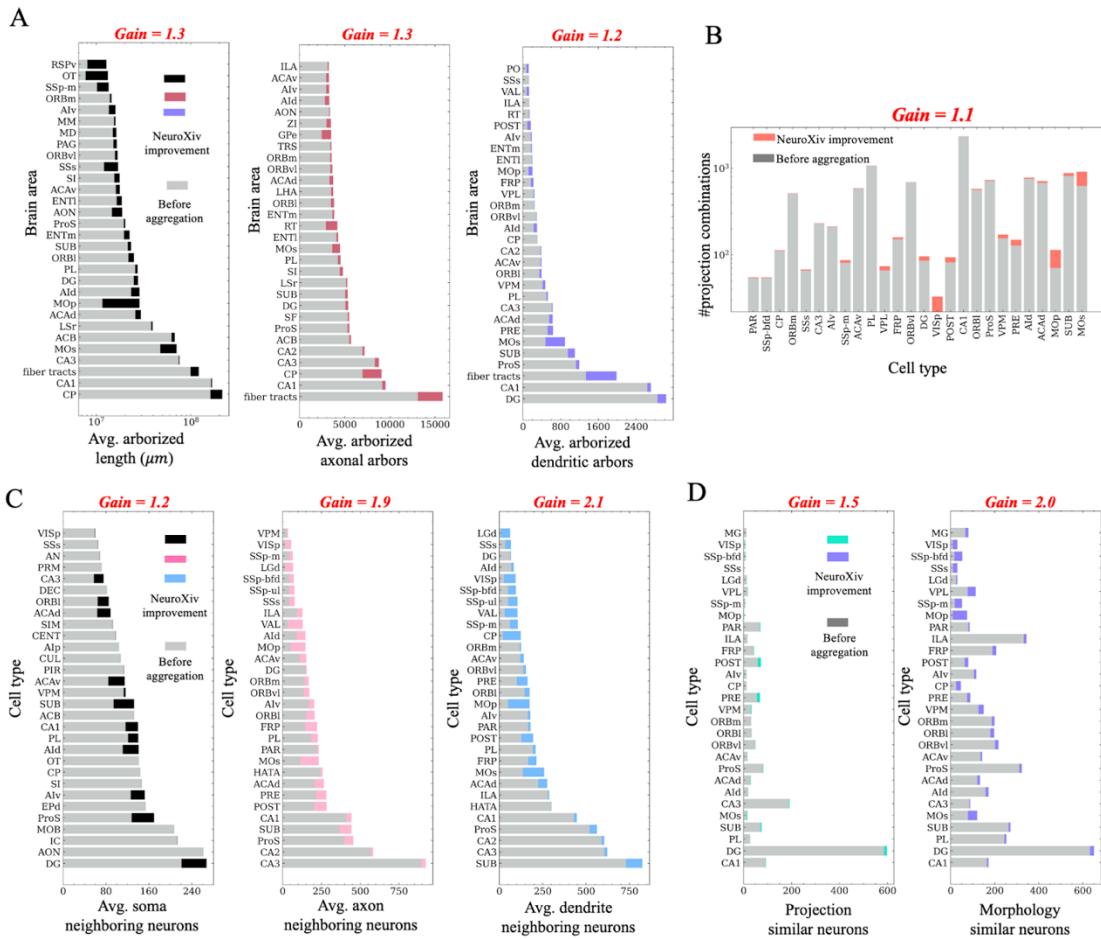
755

756 **Extended Data Fig. 1 | Local dendrites in the NeuroXiv database.** *A*, the statistics of the
757 number of local dendrites in different brain regions within the CCFv3 atlas. *B*, Visualization
758 of local dendritic data across various brain regions. Local dendrites are rendered in distinct
759 colors to enhance differentiation.



760

761 **Extended Data Fig. 2 | Illustrative diagrams depicting the definition of correlated neuron**
 762 **data on the NeuroXiv platform. A, projection similarity neurons are defined by the distance**
 763 **measured between key axonal nodes of the neurons, including the soma, bifurcation points, and**
 764 **terminal points. B, morphological similarity neurons are defined by calculating the distance**
 765 **between neuron pairs in morphological feature space. In the database, we rank similarity based**
 766 **on the distances, with closer distances indicating greater similarity. C and D, Axon neighboring**
 767 **neurons are those where the axonal arbor of neurons in the database spatially overlaps with**
 768 **the dendritic arbor of the subject neuron. Conversely, dendrite neighboring neurons are those**
 769 **where the dendritic arbor of neurons in the database spatially overlaps with the axonal arbor**
 770 **of the subject neuron. In the database, we rank neighboring neurons based on the length of the**
 771 **overlapping arbor regions, with greater lengths indicating higher proximity.**

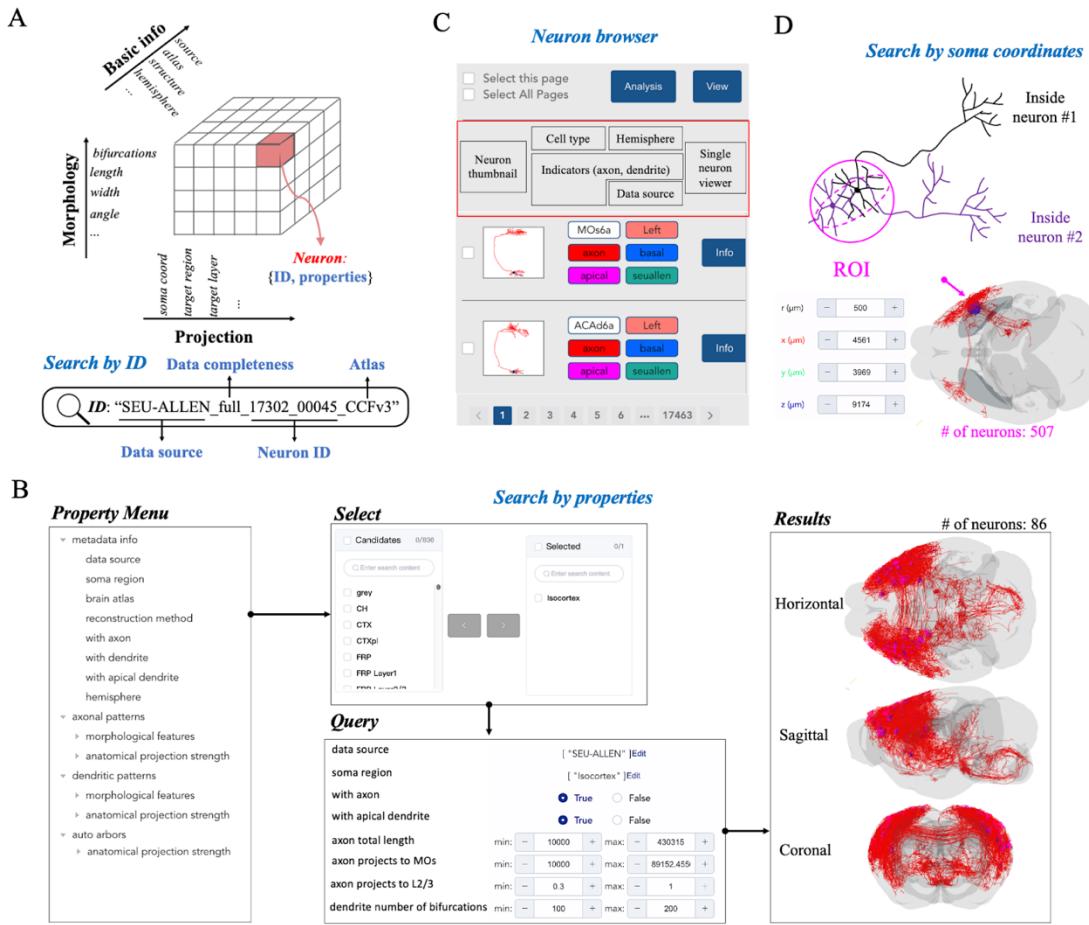


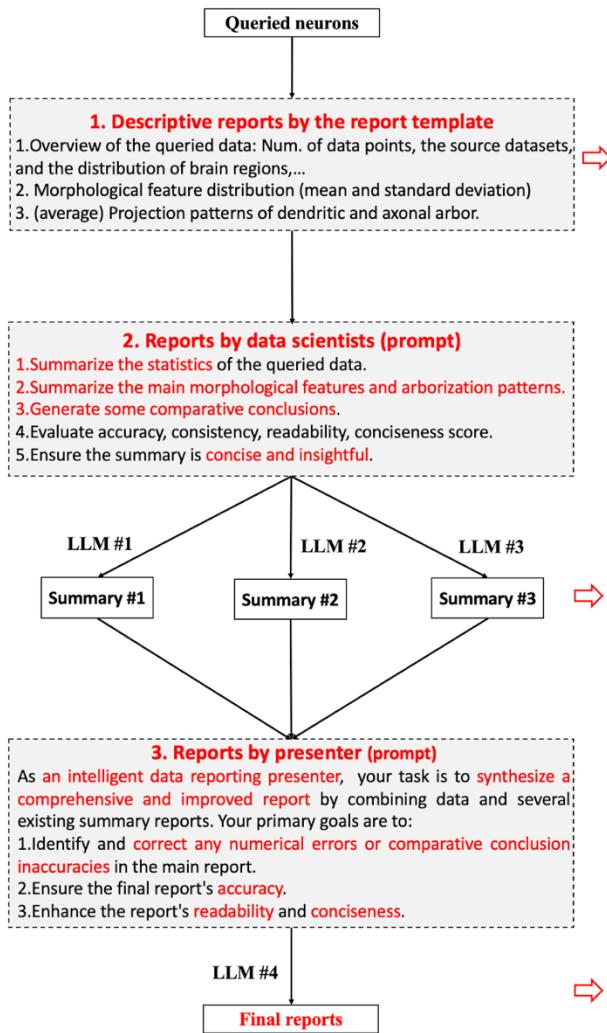
772

773 **Extended Data Fig. 3 | Gains in improvements resulting from data aggregation.** *A*, we
 774 aggregated three data sources to formulate the NeuroXiv database. Compared to the
 775 information obtainable from a single data source, this approach allows us to capture more
 776 neurite length across various regions. The gray bar represents the maximum neurite length that
 777 can be obtained from one data source, while the other colored bars represent the improvements
 778 gained through aggregation. *B*, the NeuroXiv database offers a greater number of projection
 779 combinations, where neurons extend into brain regions with lengths exceeding 1000 μm . *C* and
 780 *D*, we can identify more neighboring neurons, including soma neighboring, axon neighboring,
 781 and dendrite neighboring neurons. Additionally, we can find more similar neurons, including
 782 those that are projection similar and morphology similar. Definitions of neighboring and
 783 similar neurons can be found in **Extended Data Fig. 2**.

784

785 **Extended Data Fig. 4 | An illustrative diagram of data retrieval and filtering on the**
 786 **NeuroXiv platform.** **A**, each neuron reconstruction in the NeuroXiv database is assigned a
 787 *a unique ID and includes three types of metadata: basic information, morphological features,*
 788 *and projection characteristics.* **B**, users can customize their search strategies using either the
 789 *ID or the metadata.* **C**, to facilitate efficient data filtering, we have implemented a neuron
 790 *browser on the web portal. This tool displays the morphology of each neuron (neuron*
 791 *thumbnail) and key information, and includes an entry point for navigating to detailed data*
 792 *pages.* **D**, users can define regions of interest (ROI) and then retrieve neurons with soma
 793 *located within these ROI (Supplementary Fig. 5).*





The queried data comprises 150 neurons extracted from 3 datasets: ION (132 neurons), SEU-ALLEN (12 neurons) and MouseLight (6 neurons). This selection encompasses neuron structures, including axons (150), basal dendrites (67), apical dendrites (1), and local dendrites (0). The queried data locates in left hemisphere (49) and right hemisphere (101). The queried data is distributed across 26 brain regions, detailed as follows: CA1 (34 neurons), DG (29 neurons), ACAv (11 neurons), Specifically, there are 46 neurons in cortical layers, including L5 (30 neurons), L2/3 (10 neurons), and L6a (6 neurons).

The dataset consists of 150 neurons obtained from three sources: ION with 132 neurons, SEU-ALLEN with 12, and MouseLight with 6. These neurons encompass various structures, primarily 150 axons, alongside 67 basal dendrites, a single apical dendrite, and no local dendrites. The data is divided between the left hemisphere (49 neurons) and the right hemisphere (101 neurons). Spanning 26 brain regions, the most densely represented areas are CA1 with 34 neurons, followed by DG with 29, and ACAv with 11. Notably, 46 neurons belong to cortical layers, predominantly L5 with 30 neurons, L2/3 with 10, and L6a with 6. Other regions host fewer neurons, ranging from 8 in CA3 to single representations in areas like CLA, VPL, CP, and more. Overall, the dataset offers a diverse neurological scope for analysis.

The dataset comprises 150 neurons from three sources: ION (132 neurons), SEU-ALLEN (12 neurons), and MouseLight (6 neurons). These neurons include 150 axons, 67 basal dendrites, 1 apical dendrite, and no local dendrites. They are mainly situated in the right hemisphere (101 neurons) versus the left hemisphere (49 neurons). Spanning 26 brain regions, notable concentrations are found in CA1 (34 neurons), DG (29 neurons), and ACAv (11 neurons). Additionally, 46 neurons are in cortical layers, mostly in L5 (30 neurons), L2/3 (10 neurons), and L6a (6 neurons). The dataset offers a diverse range for analysis.

794

795 Extended Data Fig. 5 | A schematic illustration of the mixture of experts (MoE) system.

796 *The MoE-based report generation process can be divided into three distinct stages. 1) Descriptive reports generation: A program generates reports in a fixed format, capturing all relevant details of the retrieved data. 2) LLM Expert reports: Multiple LLM Experts analyze and summarize the descriptive reports from a data scientist's perspective. Although three experts are shown, the process can involve one or more. 3) Report confirmation: A different LLM Expert evaluates the previous reports for accuracy, readability, and coherence, and refines the final report accordingly. An actual case is shown on the far right, with red arrows pointing to the reports generated at each stage.*

A Overview of the queried data (CA neurons with reconstructed axon)

The dataset encompasses 4632 neurons sourced from three datasets: ION (4590 neurons), MouseLight (29 neurons), and SEU-ALLEN (13 neurons). These neurons include 4632 axons, 1987 basal dendrites, and 3 apical dendrites, with no local dendrites recorded. There are 1023 neurons in the left hemisphere and 3609 in the right hemisphere. The neurons are distributed across three brain regions: CA1 (3691 neurons), CA3 (895 neurons), and CA2 (46 neurons). All the neurons are located in cortical layers, with the specific layer distribution not specified.

B The distribution of morphological features

The analysis of neuronal morphology data from ION, MouseLight, and SEU-ALLEN sources reveals significant variability in the axonal and dendritic features of CA1 and CA3 neurons.

CA3 neurons consistently exhibit greater Total Length (around 159,000 to 189,000 µm) and Number of Bifurcations (310 to 341) compared to CA1 neurons (approximately 52,000 to 74,000 µm in length and 171 to 215 bifurcations), indicating more extensive and complex connectivity in CA3 neurons. Furthermore, CA3 neurons generally show higher Max Path Distance, suggesting broader reach. The measures for Center Shift are relatively comparable across neuron types and sources, indicating similar spatial distribution balance.

Dendritic arbor data mirrors these patterns, with CA3 neurons typically having higher Total Length and Number of Bifurcations than CA1 neurons. This highlights the functional diversity tied to the morphological complexity of these neurons. However, substantial variations exist among the different data sources, emphasizing the need to consider data provenance when interpreting neuronal morphology.

In summary, the observed differences in Total Length, Number of Bifurcations, Max Path Distance, and Center Shift underscore the morphological intricacy of CA1 and CA3 neurons, with CA3 neurons displaying notably more extensive and complex arborization patterns.

C Projection / Arborized patterns

The analysis of neuronal projection data from CA1 and CA3 neurons across sources (ION, MouseLight, and SEU-ALLEN) reveals complex and diverse connectivity patterns.

Dendritic Arbor Data:

- CA1 neurons: Displayed high proportions of their dendritic length within the CA1 region across all sources (72.1% to 87.7%), indicating strong local signal reception.
- CA3 neurons: Similarly exhibited significant dendritic lengths within CA3 (78.9% to 82.0%) in the ION and MouseLight datasets. However, the SEU-ALLEN data showed more distribution with 41.6% in CA3 and additional projections to the DG (39.9%).

Axonal Arbor Data:

- CA1 neurons: Showed extensive axonal projections both locally within CA1 and to other regions like SUB and ProS. The ION source highlighted 28.1% of axonal length within CA1 itself, while MouseLight showed dominant projections to SUB (14.0%) and ProS (11.9%).
- CA3 neurons: Exhibited major axonal projections to CA1 and retained significant intra-regional connections within CA3. The ION dataset noted 40.4% of CA3 axons targeting CA1, whereas SEU-ALLEN data indicated a dominant 71.2% of projections remaining within CA3.

In conclusion, CA1 and CA3 neurons exhibit robust intraregional connectivity and significant strategic long-range projections, primarily between CA3 and CA1. These intricate patterns underscore the brain's complex neural communication and the importance of comprehensive comparative data to fully understand neuronal connectivity.

804

805 **Extended Data Fig. 6 | MoE report of the queried CA neurons.** *The report is structured into*
806 *three sections: A) an overview of the queried data, B) the distribution of morphological features*
807 *within the queried data, and C) the projection patterns observed in the queried data. Key*
808 *statistical points and comparative analyses are highlighted in red throughout the report.*

A Overview of the queried data (MOs neurons with reconstructed axon)

The dataset comprises 991 neurons from three datasets: ION (661 neurons), MouseLight (278 neurons), and SEU-ALLEN (52 neurons). These neurons include 991 axons, 594 basal dendrites, 43 apical dendrites, and no local dendrites. The data is distributed between the left hemisphere (319 neurons) and the right hemisphere (672 neurons). The neurons are located in five brain regions: MOs5 (483 neurons), MOs2/3 (307 neurons), MOs6 (149 neurons), MOs1 (46 neurons), and MOs (6 neurons). Within cortical layers, the distribution is L5 (483 neurons), L2/3 (307 neurons), L6a (149 neurons), L1 (46 neurons), and an unspecified layer with 6 neurons.

B The distribution of morphological features

The neuronal morphology dataset from ION, MouseLight, and SEU-ALLEN repositories reveals significant variations in key features of MOs5 and MOs2/3 neurons, shedding light on their structural complexity.

For MOs5 neurons:

- ION source: Total Length: 166229.11 μm , Number of Bifurcations: 328.43
- MouseLight: Total Length: 129225.81 μm , Number of Bifurcations: 236.61
- SEU-ALLEN: Total Length: 111741.34 μm , Number of Bifurcations: 304.21

For MOs2/3 neurons:

- ION: Total Length: 116075.31 μm , Number of Bifurcations: 194.33
- MouseLight: Total Length: 90393.69 μm , Number of Bifurcations: 166.37
- SEU-ALLEN: Total Length: 103496.7 μm , Number of Bifurcations: 253.23

Max Path Distance remains consistent, averaging around 9,000-11,000 μm for axons, and about 600 μm for dendrites. Center Shift indicates spatial distribution balance, ranging from 1,751.33 to 2,827.01 for axons, and much lower for dendrites. Overall, MOs5 neurons exhibit higher values and variability in these key features compared to MOs2/3 neurons.

C Projection / Arborized patterns

The analysis of neuronal projection data from different sources (ION, MouseLight, SEU-ALLEN) reveals consistent and significant patterns in dendritic and axonal arborization lengths for MOs5 and MOs2/3 neurons across various brain regions, indicating their connectivity strength and functional roles.

For MOs5 Neurons (dendritic arbor):

The MOs region consistently shows the highest dendritic arborization across sources: - ION: 8161.6 μm (45.8%) - MouseLight: 7717.8 μm (48.5%) - SEU-ALLEN: 8233.7 μm (42.4%)

For MOs2/3 Neurons (dendritic arbor):

The MOs region also dominates in dendritic arborization for MOs2/3 neurons: - ION: 6415.0 μm (47.7%) - MouseLight: 6716.2 μm (49.3%) - SEU-ALLEN: 7789.0 μm (44.6%)

For MOs5 Neurons (axonal arbor):

- ION Source: - CP: 47505.8 μm (15.8%) - MOs: 31232.2 μm (12.5%)
- MouseLight Source: - MOs: 27780.4 μm (16.0%) - CP: 37482.4 μm (15.0%)
- SEU-ALLEN Source: - MOs: 28349.8 μm (20.0%) - CP: 9900.6 μm (6.2%)

For MOs2/3 Neurons (axonal arbor):

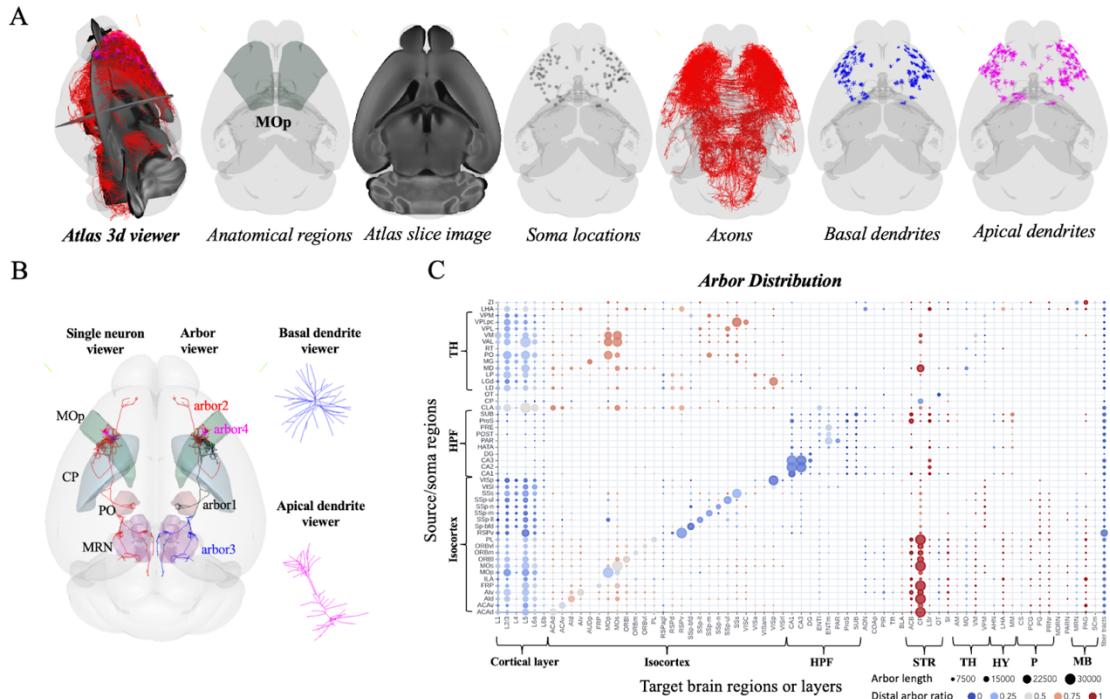
- ION Source: - CP: 38746.1 μm (20.0%) - MOs: 29712.3 μm (18.8%)
- MouseLight Source: - MOs: 29181.1 μm (21.9%) - CP: 24252.6 μm (15.0%)
- SEU-ALLEN Source: - MOs: 45005.7 μm (26.8%) - MOp: 10683.7 μm (6.6%)

Conclusion

The MOs region is a central hub for both dendritic and axonal arborization in MOs5 and MOs2/3 neurons, highlighting its crucial role in signal reception and transmission. Axonal projections of these neurons show a considerable distribution across the CP and MOp regions, emphasizing their involvement in motor function and coordination. Notable variations between sources underscore the importance of comprehensive analysis for understanding connectivity patterns in neural networks.

809

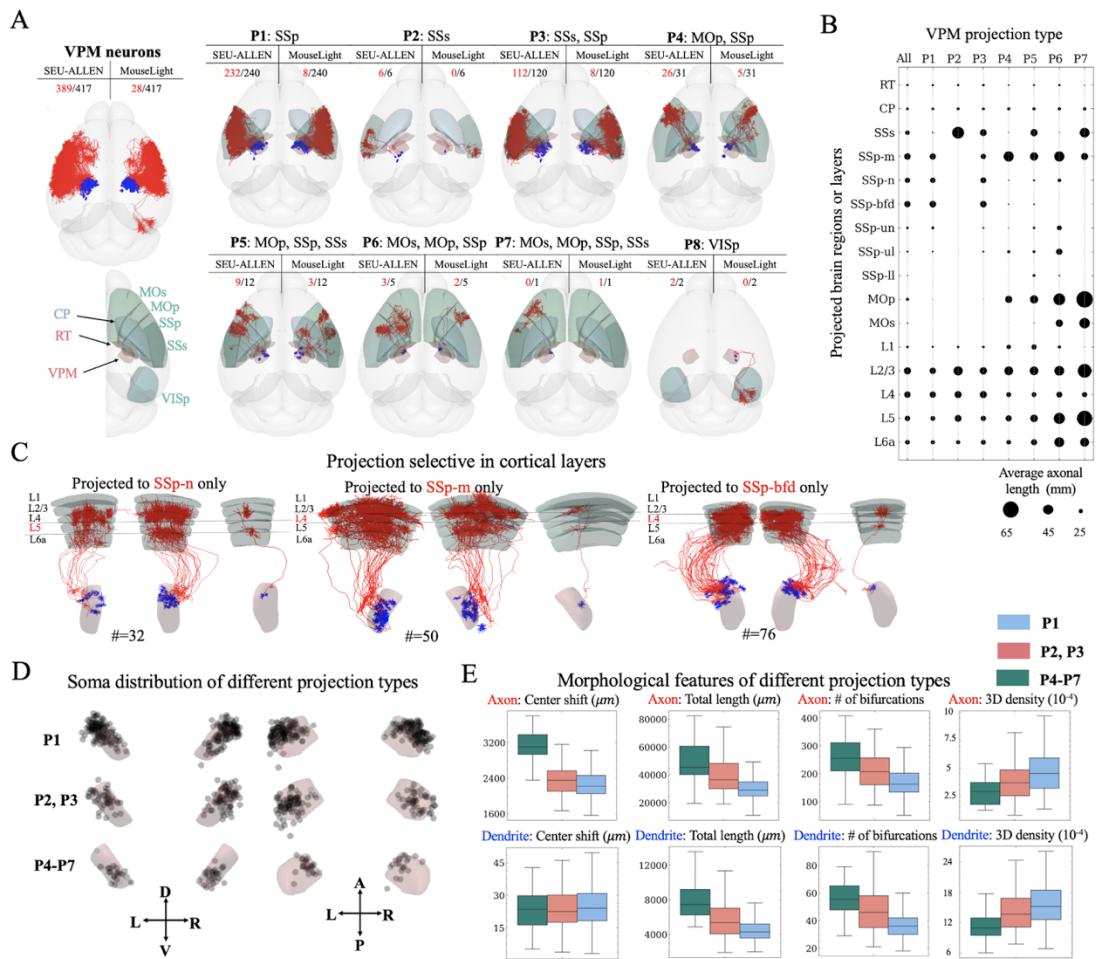
810 **Extended Data Fig. 7 | MoE report of the queried MOs neurons.** The report is structured
811 into three sections: A) an overview of the queried data, B) the distribution of morphological
812 features within the queried data, and C) the projection patterns observed in the queried data.
813 Key statistical points and comparative analyses are highlighted in red throughout the report.



814

815 **Extended Data Fig. 8 | Interactive visualization diagram of NeuroXiv.** *A*, we have
 816 implemented a 3D viewer on the web platform, allowing users to visualize the atlas and neurons
 817 or neuronal structures of interest. *B*, visualization of individual neuron data, featuring an
 818 embedded arbor viewer with zoom-in views of basal and apical dendrites. *C*, arbor distribution
 819 of major cell types across various brain regions. We also visualize arbor distribution across
 820 different cortical layers.

821
822
823
824
825
826
827
828
829
830
831
832
833
834
835



836 **Extended Data Tables**

837 **Extended Data Table 1 | Comparison of alternative neuron morphology web platforms**
 838 (*Extended_Data_Table1.pdf*).

	Name	NeuroXiv	MouseLight neuron browser	Digital Brain (ION)	NeuroMorpho
	Website	neuroxiv.org	mouselight.janelia.org	mouse.digitalbrain.cn	neuromorpho.org
Database	number of neurons (mouse)	177,186	1,227	16,788	137,457 (unknown completeness)
	number of cell types (anatomy)	294	126	25	72
	structural domains	soma, axon, (basal and apical) dendrite, arbor	soma, axon, dendrite	soma, axon, dendrite	soma, axon, (basal and apical) dendrite
	number of data sources	3	1	1	13
Data query	data format and standard	standard SWC format	JSON file containing metadata and neuron structures	SWC format file with non-standard and inconsistent structural domain identifiers	standard SWC format (no quality control)
	reference atlases	CCFv3	CCFv2.5, CCFv3	CCFv3	Not exist one for all neurons
Data visualization	neuron browser	Yes (neuron thumbnail and metadata)	Yes (id based)	Yes (id based)	Yes (id based)
	by soma region	Yes	Yes	Yes	Yes
	by spatial coordinates	Yes	Yes	Yes	
	by projected regions	Yes	Yes	Yes	
	by morphological features	Yes	Yes	Yes	
	by natural language	Yes			Yes
Data analysis	by user-supplied neuron id (s)	Yes			
	interactive visualization	Yes	Yes	Yes	2D snapshot only
	statistics	Yes		Yes	Yes
Data mining	projection patterns	Yes			
	morphological feature distribution	Yes			measures of one neuron
	AI report	Yes			
Open access	morphology similar neurons	Yes			
	projection similar neurons	Yes			
	spatial neighboring neurons	Yes			
	no authentication	Yes	Yes		Yes
	(meta) data	Yes	Yes		Yes
	queried (meta) data	Yes	Yes		Yes
	discovery data	Yes	Yes (limited number : 20)	morphology data only Yes (limited number : 500) only metadata of part of figures	
	processing tools / pipelines	Yes			Yes
	processing required before data reuse			file format conversion and standardization, quality control, atlas mapping, metadata extraction	

839

840

841

842 **Extended Data Table 2 | The metadata of neurons and their descriptions in the NeuroXiv**
843 **database (*Extended_Data_Table2.pdf*).**

844 This table can be found along with the submission files of this manuscript.

845

846 **Extended Data Table 3 | Comparison of the performance of two types of natural language**
 847 **query methods in AIPOM** (*Extended_Data_Table3.pdf*). To achieve this, we defined three
 848 common retrieval scenarios: querying neuron types, querying neurons with particular structures,
 849 and querying neurons with specific projection patterns, each comprising 10 test cases. To
 850 ensure fair testing, both methods used the same computer configuration, eliminating disparities
 851 due to varying computational power. The server-side method involved setting up a local
 852 NeuroXiv server on the test computer, while the client-side method accessed the NeuroXiv
 853 server hosted on AWS directly. As a result, compared to the server-side approach, the client-
 854 side method demonstrated an average response time improvement of 12.3.

	query items	server-side method (ms)	client-side method (ms)	Gain
query neuron types	search CA1 neurons	16,599.4	1,687.2	9.8
	search DG neurons	14,286.5	1,569.9	9.1
	search PL neurons	19,531.8	1,975.4	9.9
	search MOs neurons	12,846.5	1,048.0	12.3
	search CA3 neurons	21,297.8	1,482.3	14.4
	search Ald neurons	24,380.3	3,622.0	6.7
	search SUB neurons	18,534.9	3,340.7	5.5
	search ACAd neurons	19,674.8	1,163.7	16.9
	search VPM neurons	13,016.1	1,435.2	9.1
	search CP neurons	12,664.5	2,642.9	4.8
query neurons with particular structure	search CA1 neurons with axon	24,787.2	1,630.7	15.2
	search DG neurons with axon	22,434.6	2,105.8	10.7
	search PL neurons with axon	24,884.9	1,138.6	21.9
	search MOs neurons with axon	27,150.4	1,178.4	23.0
	search CA3 neurons with axon	21,155.0	1,277.4	16.6
	search Ald neurons with axon	18,351.7	2,532.8	7.2
	search SUB neurons with axon	25,905.8	2,792.3	9.3
	search ACAd neurons with axon	17,380.2	3,053.6	5.7
	search VPM neurons with axon	37,543.3	2,326.7	16.1
	search CP neurons with axon	12,683.2	2,251.7	5.6
query neurons with specific projection patterns	search CA1 neurons projecting to ACB	15,865.6	2,568.2	6.2
	search DG neurons projecting to CA3	33,939.7	4,821.4	7.0
	search PL neurons projecting to CP	31,327.7	2,796.8	11.2
	search MOs neurons projecting to SSs	59,691.9	2,008.7	29.7
	search CA3 neurons projecting to LSr	35,424.7	2,596.1	13.6
	search Ald neurons projecting to MOs	31,140.7	2,357.7	13.2
	search SUB neurons projecting to MM	24,500.8	2,155.4	11.4
	search ACAd neurons projecting to CP	49,261.0	2,418.1	20.4
	search VPM neurons projecting to MOp	26,355.9	2,323.7	11.3
	search CP neurons projecting to SNr	31,316.8	2,203.5	14.2
Average		24,797.8	2,216.8	12.3

855