

cellular, anatomical, electrophysiological, behavioral, evolutionary, and computational. The experimental methods involve many spatial scales, from electron microscopy (EM) to whole-brain human neuroimaging, and time scales ranging from microseconds for ion channel gating to years for longitudinal studies of human development and aging. An increasing number of insights emerge from integration and synthesis across these spatial and temporal domains. However, such efforts face impediments related to the diversity of scientific subcultures and differing approaches to data acquisition, storage, description, and analysis, and even to the language in which they are described. It is often unclear how best to integrate the linear information of genetic sequences, the highly visual data of neuroanatomy, the time-dependent data of electrophysiology, and the more global level of analyzing behavior and clinical syndromes.

The great majority of neuroscientists carry out highly focused, hypothesis-driven research that can be powerfully framed in the context of known circuits and functions. Such efforts are complemented by a growing number of projects that provide large data sets aimed not at testing a specific hypothesis but instead enabling data-intensive discovery approaches by the community at large. Notable successes include gene expression atlases from the Allen Institute for Brain Sciences (4) and the Gene Expression Nervous System Atlas (GENSAT) Project (5), and disease-specific human neuroimaging repositories (6). However, the neuroscience community is not yet fully engaged in exploiting the rich array of data currently available, nor is it adequately poised to capitalize on the forthcoming data explosion.

Below we highlight several major endeavors that provide complementary perspectives on the challenges and opportunities in neuroscience data mining. One is a set of “connectome” projects that aim to comprehensively describe neural circuits at either the macroscopic or the microscopic level. Another, the Neuroscience Information Framework (NIF), encompasses all of neuroscience and provides access to existing knowledge and databases of many types. These and other efforts provide fresh approaches to the challenge of elucidating neural choreography.

Connectomes: Macroscopic and Microscopic

Brain anatomy provides a fundamental three-dimensional framework around which many types of neuroscience data can be organized and mined. Decades of effort have revealed immense amounts of information about local and long-distance connections in animal brains. A wide range of tools (such as immunohistochemistry and *in situ* hybridization) have characterized the biochemical nature of these circuits that are studied electrophysiologically, pharmacologically, and behaviorally (7). Several ongoing efforts aim to integrate anatomical information into searchable

resources that provide a backbone for understanding circuit biology and function (8–10). The challenge of integrating such data will dramatically increase with the advent of high-throughput anatomical methods, including those emerging from the nascent field of connectomics.

A connectome is a comprehensive description of neural connectivity for a specified brain region at a specified spatial scale (11, 12). Connectomics currently includes distinct subdomains for studying the macroconnectome (long-distance pathways linking patches of gray matter) and the microconnectome (complete connectivity within a single gray-matter patch).

The Human Connectome Project. Until recently, methods for charting neural circuits in the human brain were sorely lacking (13). This situation has changed dramatically with the advent

grants to two consortia (18). The consortium led by Washington University in St. Louis and the University of Minnesota (19) aims to characterize whole-brain circuitry and its variability across individuals in 1200 healthy adults (300 twin pairs and their nontwin siblings). Besides diffusion imaging and R-fMRI, task-based fMRI data will be acquired in all study participants, along with extensive behavioral testing; 100 participants will also be studied with magnetoencephalography (MEG) and electroencephalography (EEG). Acquired blood samples will enable genotyping or full-genome sequencing of all participants near the end of the 5-year project. Currently, data acquisition and analysis methods are being extensively refined using pilot data sets. Data acquisition from the main cohort will commence in mid-2012.

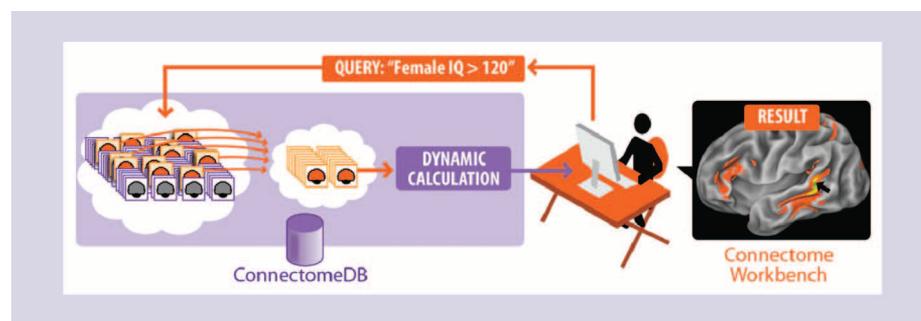


Fig. 1. Schematic illustration of online data mining capabilities envisioned for the HCP. Investigators will be able to pose a wide range of queries (such as the connectivity patterns of a particular brain region of interest averaged across a group of individuals, based on behavioral criteria) and view the search results interactively on three-dimensional brain models. Data sets of interest will be freely available for downloading and additional offline analysis.

of noninvasive neuroimaging methods. Two complementary modalities of magnetic resonance imaging (MRI) provide the most useful information about long-distance connections. One modality uses diffusion imaging to determine the orientation of axonal fiber bundles in white matter, based on the preferential diffusion of water molecules parallel to these fiber bundles. Tractography is an analysis strategy that uses this information to estimate long-distance pathways linking different gray-matter regions (14, 15). A second modality, resting-state functional MRI (R-fMRI), is based on slow fluctuations in the standard fMRI BOLD signal that occur even when people are at rest. The time courses of these fluctuations are correlated across gray-matter locations, and the spatial patterns of the resultant functional connectivity correlation maps are closely related but not identical to the known pattern of direct anatomical connectivity (16, 17). Diffusion imaging and R-fMRI each have important limitations, but together they offer powerful and complementary windows on human brain connectivity.

To address these opportunities, the National Institutes of Health (NIH) recently launched the Human Connectome Project (HCP) and awarded

Neuroimaging and behavioral data from the HCP will be made freely available to the neuroscience community via a database (20) and a platform for visualization and user-friendly data mining. This informatics effort involves major challenges owing to the large amounts of data (expected to be ~1 petabyte), the diversity of data types, and the many possible types of data mining. Some investigators will drill deeply by analyzing high-resolution connectivity maps between all gray-matter locations. Others will explore a more compact “parcellated connectome” among all identified cortical and subcortical parcels. Data mining options will reveal connectivity differences between subpopulations that are selected by behavioral phenotype (such as high versus low IQ) and various other characteristics (Fig. 1). The utility of HCP-generated data will be enhanced by close links to other resources containing complementary types of spatially organized data, such as the Allen Human Brain Atlas (21), which contains neural gene expression maps.

Microconnectomes. Recent advances in serial section EM, high-resolution optical imaging methods, and sophisticated image segmentation methods enable detailed reconstructions of the

microscopic connectome at the level of individual synapses, axons, dendrites, and glial processes (22–24). Current efforts focus on the reconstruction of local circuits, such as small patches of the cerebral cortex or retina, in laboratory animals. As such data sets begin to emerge, a fresh set of informatics challenges will arise in handling petabyte amounts of primary and analyzed data and in providing data mining platforms that enable neuroscientists to navigate complex local circuits and examine interesting statistical characteristics.

Micro- and macroconnectomes exemplify distinct data types within particular tiers of analysis that will eventually need to be linked. Effective interpretation of both macro- and microconnectomic approaches will require novel informatics and computational approaches that enable these two types of data to be analyzed in a common framework and infrastructure. Efforts such as the Blue Brain Project (25) represent an important initial thrust in this direction, but the endeavor will entail decades of effort and innovation.

Powerful and complementary approaches such as optogenetics operate at an intermediate (meso-connectome) spatial scale by directly perturbing neural circuits in vivo or in vitro with light-activated ion channels inserted into selected neuronal types (26). Other optical methods, such as calcium imaging with two-photon laser microscopy, enable analysis of the dynamics of ensembles of neurons in microcircuits (27, 28) and can lead to new conceptualizations of brain function (29). Such approaches provide an especially attractive window on neural choreography as they assess or perturb the temporal patterns of macro- or microcircuit activity.

The NIF

Connectome-related projects illustrate ways in which neuroscience as a field is evolving at the level of neural circuitry. Other discovery efforts include genome-wide gene expression profiling [for example, (30)] or epigenetic analyses across multiple brain regions in normal and diseased brains. This wide range of efforts results in a sharp increase in the amount and diversity of data being generated, making it unlikely that neuroscience will be adequately served by only a handful of centralized databases, as is largely the case for the genomics and proteomics community (31). How, then, can we access and explore these resources more effectively to support the data-intensive discovery envisioned in *The Fourth Paradigm* (32)?

Tackling this question was a prime motivation behind the NIF (33). The NIF was launched in 2005 to survey the current ecosystem of neuroscience resources (databases, tools, and materials) and to establish a resource description framework and search strategy for locating, accessing, and using digital neuroscience-related resources (34).

The NIF catalog, a human-curated registry of known resources, currently includes more than 3500 such resources, and new ones are added

daily. Over 2000 of these resources are databases that range in size from hundreds to millions of records. Many were created at considerable effort and expense, yet most of them remain underused by the research community.

Clearly, it is inefficient for individual researchers to sequentially visit and explore thousands of databases, and conventional online search engines are inadequate, insofar as they do not effectively index or search database content. To promote the discovery and use of online databases, the NIF created a portal through which users can search not only the NIF registry but the content of multiple databases simultaneously. The current NIF federation includes more than 65 databases accessing ~30 million records (35) in major domains of relevance to neuroscience (Fig. 2). Besides very large genomic collections, there are nearly 1 million antibody records, 23,000 brain connectivity records, and >50,000 brain activation coordinates. Many of these areas are covered by multiple databases, which the NIF knits together into a coherent view. Although impressive, this represents only the tip of the iceberg. Most individual databases are underpopulated because of insufficient community contributions. Entire domains of neuroscience (such as electrophysiology and behavior) are underrepresented as compared to genomics and neuroanatomy.

Ideally, NIF users should be able not only to locate answers that are known but to mine available data in ways that spur new hypotheses regarding what is not known. Perhaps the single biggest roadblock to this higher-order data mining is the lack of standardized frameworks for organizing neuroscience data. Individual investigators often use terminology or spatial coordinate systems customized for their own particular analysis approaches. This customization is a substantial barrier to data in-

tegration, requiring considerable human effort to access each resource, understand the context and content of the data, and determine the conditions under which they can be compared to other data sets of interest.

To address the terminology problem, the NIF has assembled an expansive lexicon and ontology covering the broad domains of neuroscience by synthesizing open-access community ontologies (36). The Neurolex and accompanying NIFSTD (NIF-standardized) ontologies provide definitions of over 50,000 concepts, using formal languages to represent brain regions, cells, subcellular structures, molecules, diseases, and functions, and the relations among them. When users search for a concept through the NIF, it automatically expands the query to include all synonymous or closely related terms. For example, a query for “striatum” will include “neostriatum, dorsal striatum, caudoputamen, caudate putamen” and other variants.

Neurolex terms are accessible through a wiki (37) that allows users to view, augment, and modify these concepts. The goal is to provide clear definitions of each concept that can be used not only by humans but by automated agents, such as the NIF, to navigate the complexities of human neuroscience knowledge. A key feature is the assignment of a unique resource identifier to make it easier for search algorithms to distinguish among concepts that share the same label.

For example, nucleus (part of cell) and nucleus (part of brain) are distinguished by unique IDs. Using these identifiers in addition to natural language to reference concepts in databases and publications, although conceptually simple, is an especially powerful means for making data maximally discoverable and useful.

These efforts to develop and deploy a semantic framework for neuroscience, spearheaded by

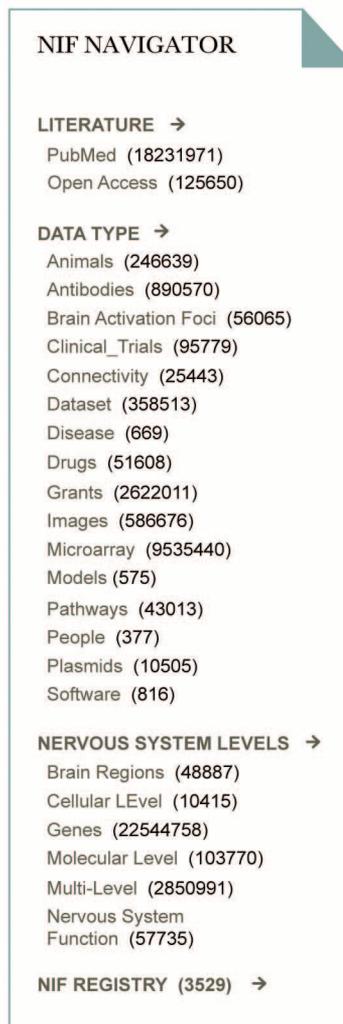


Fig. 2. Current contents of the NIF. The NIF navigator displays the current contents of the NIF data federation, organized by data type and level of the nervous system. The number of records in each category is displayed in parentheses.

the NIF and the International Neuroinformatics Coordinating Facility (38), are complemented by projects related to brain atlases and spatial frameworks (39–41), providing tools for referencing data to a standard coordinate system based on the brain anatomy of a given organism.

Neuroinformatics as a Prelude to New Discoveries

How might improved access to multiple tiers of neurobiological data help us understand the brain? Imagine that we are investigating the neurobiology of bipolar disorder, an illness in which moods are normal for long periods of time, yet are labile and sometimes switch to mania or depression without an obvious external trigger. Although highly heritable, this disease appears to be genetically very complex and possibly quite heterogeneous (42). We may discover numerous genes that impart vulnerability to the illness. Some may be ion channels, others synaptic proteins or transcription factors. How will we uncover how disparate genetic causes lead to a similar clinical phenotype? Are they all affecting the morphology of certain cells; the dynamics of specific microcircuits, for example, within the amygdala; the orchestration of information across regions, for example, between the amygdala and the prefrontal cortex? Can we create genetic mouse models of the various mutated genes and show a convergence at any of these levels? Can we capture the critical changes in neuronal and/or glial function (at any of the levels) and find ways to prevent the illness? Discovering the common thread for such a disease will surely benefit from tools that facilitate navigation across the multiple tiers of data—genetics, gene expression/epigenetics, changes in neuronal activity, and differences in dynamics at the micro and macro levels, depending on the mood state. No single focused level of analysis will suffice to achieve a satisfactory understanding of the disease. In neural choreography terms, we need to identify the dancers, define the nature of the dance, and uncover how the disease disrupts it.

Recommendations

Need for a cultural shift. To meet the grand challenge of elucidating neural choreography, we need increasingly powerful scientific tools to study brain activity in space and in time, to extract the key features associated with particular events, and to do so on a scale that reveals commonalities and differences between individual brains. This requires an informatics infrastructure that has built-in flexibility to incorporate new types of data and navigate across tiers and domains of knowledge.

The NIF currently provides a platform for integrating and systematizing existing neuroscience knowledge and has been working to define

best practices for those producing new neuroscience data. Good planning and future investment are needed to broaden and harden the overall framework for housing, analyzing, and integrating future neuroscience knowledge. The International Neuroinformatics Coordinating Facility (INCF) plays an important role in coordinating and promoting this framework at a global level.

But can neuroscience evolve so that neuroinformatics becomes integral to how we study the brain? This would entail a cultural shift in the field regarding the importance of data sharing and mining. It would also require recognition that neuroscientists produce data not just for consumption by readers of the conventional literature, but for automated agents that can find, relate, and begin to interpret data from databases as well as the literature. Search technologies are advancing rapidly, but the complexity of scientific data continues to challenge. To make neuroscience data maximally interoperable within a global neuroscience information framework, we encourage the neuroscience community and the associated funding agencies to consider the following set of general and specific suggestions:

1) Neuroscientists should, as much as is feasible, share their data in a form that is machine-accessible, such as through a Web-based database or some other structured form that benefits from increasingly powerful search tools.

2) Databases spanning a growing portion of the neuroscience realm need to be created, populated, and sustained. This effort needs adequate support from federal and other funding mechanisms.

3) Because databases become more useful as they are more densely populated (43), adding to existing databases may be preferable to creating customized new ones. The NIF, INCF, and other resources provide valuable tools for finding existing databases.

4) Data consumption will increasingly involve machines first and humans second. Whether creating database content or publishing journal articles, neuroscientists should annotate content using community ontologies and identifiers. Coordinates, atlas, and registration method should be specified when referencing spatial locations.

5) Some types of published data (such as brain coordinates in neuroimaging studies) should be reported in standardized table formats that facilitate data mining.

6) Investment needs to occur in interdisciplinary research to develop computational, machine-learning, and visualization methods for synthesizing across spatial and temporal information tiers.

7) Educational strategies from undergraduate through postdoctoral levels are needed to ensure that neuroscientists of the next generation are proficient in data mining and using the data-sharing tools of the future.

8) Cultural changes are needed to promote widespread participation in this endeavor. These ideas are not just a way to be responsible and collaborative; they may serve a vital role in attaining a deeper understanding of brain function and dysfunction.

With such efforts, and some luck, the machinery that we have created, including powerful computers and associated tools, may provide us with the means to comprehend this “most unaccountable of machinery” (44), our own brain.

References and Notes

- World Health Organization, *Mental Health and Development: Targeting People With Mental Health Conditions As a Vulnerable Group* (World Health Organization, Geneva, 2010).
- F. A. Azevedo *et al.*, *J. Comp. Neurol.* **513**, 532 (2009).
- B. Pakkenberg *et al.*, *Exp. Gerontol.* **38**, 95 (2003).
- www.alleninstitute.org/
- www.gensat.org/
- http://adni.loni.ucla.edu
- A. Bjorklund, T. Hokfelt, Eds., *Handbook of Chemical Neuroanatomy Book Series*, vols. 1 to 21 (Elsevier, Amsterdam, 1983–2005).
- http://cocomac.org/home.asp
- http://brancusi.usc.edu/bkms/
- http://brainmaps.org/
- http://en.wikipedia.org/wiki/Connectome
- O. Sporns, G. Tononi, R. Kotter, *PLoS Comput. Biol.* **1**, e42 (2005).
- F. Crick, E. Jones, *Nature* **361**, 109 (1993).
- H. Johansen-Berg, T. E. J. Behrens, *Diffusion MRI: From Quantitative Measurement to in-vivo Neuroanatomy* (Academic Press, London, ed. 1, 2009).
- H. Johansen-Berg, M. F. Rushworth, *Annu. Rev. Neurosci.* **32**, 75 (2009).
- J. L. Vincent *et al.*, *Nature* **447**, 83 (2007).
- D. Zhang, A. Z. Snyder, J. S. Shimony, M. D. Fox, M. E. Raichle, *Cereb. Cortex* **20**, 1187 (2010).
- http://humanconnectome.org/consortia
- http://humanconnectome.org
- D. S. Marcus, T. R. Olsen, M. Ramaratnam, R. L. Buckner, *Neuroinformatics* **5**, 11 (2007).
- http://human.brain-map.org/
- K. L. Briggman, W. Denk, *Curr. Opin. Neurobiol.* **16**, 562 (2006).
- J. W. Lichtman, J. Livet, J. R. Sanes, *Nat. Rev. Neurosci.* **9**, 417 (2008).
- S. J. Smith, *Curr. Opin. Neurobiol.* **17**, 601 (2007).
- H. Markram, *Nat. Rev. Neurosci.* **7**, 153 (2006).
- K. Deisseroth, *Nat. Methods* **8**, 26 (2011).
- B. F. Grewe, F. Helmchen, *Curr. Opin. Neurobiol.* **19**, 520 (2009).
- B. O. Watson *et al.*, *Front. Neurosci.* **4**, 29 (2010).
- G. Buzsaki, *Neuron* **68**, 362 (2010).
- R. Bernard *et al.*, *Mol. Psychiatry*, 13 April 2010 (e-pub ahead of print).
- M. E. Martone, A. Gupta, M. H. Ellisman, *Nat. Neurosci.* **7**, 467 (2004).
- The Fourth Paradigm: Data Intensive Scientific Discovery*, T. Hey, S. Tansler, K. Tolle, Eds. (Microsoft Research Publishing, Redmond, WA, 2009).
- http://neuiinfo.org
- D. Gardner *et al.*, *Neuroinformatics* **6**, 149 (2008).
- A. Gupta *et al.*, *Neuroinformatics* **6**, 205 (2008).
- W. J. Bug *et al.*, *Neuroinformatics* **6**, 175 (2008).
- http://neurolex.org
- http://infc.org
- http://www.brain-map.org
- http://wholebraincatalog.org
- http://infc.org/core/programs/atlasimg
- H. Akil *et al.*, *Science* **327**, 1580 (2010).

43. G. A. Ascoli, *Nat. Rev. Neurosci.* 7, 318 (2006).
44. N. Nicolson, J. Trautmann, Eds., *The Letters of Virginia Woolf, Volume V: 1932–1935* (Harcourt, Brace, New York, 1982).
45. Funded in part by grants from NIH (5P01-DA021633-02), the Office of Naval Research (ONR-N00014-02-1-0879),

and the Pritzker Neuropsychiatric Research Consortium (to H.A.); and the HCP (1U54MH091657-01) from the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research (to D.C.V.E.). The NIF is supported by NIH Neuroscience Blueprint contract HHSN271200800035C via the National

Institute on Drug Abuse to M.E.M. We thank D. Marcus, R. Poldrack, S. Curtiss, A. Bandrowski, and S. J. Watson for discussions and comments on the manuscript.

10.1126/science.1199305

PERSPECTIVE

The Disappearing Third Dimension

Timothy Rowe¹ and Lawrence R. Frank²

Three-dimensional computing is driving what many would call a revolution in scientific visualization. However, its power and advancement are held back by the absence of sustainable archives for raw data and derivative visualizations. Funding agencies, professional societies, and publishers each have unfulfilled roles in archive design and data management policy.

Three-dimensional (3D) image acquisition, analysis, and visualization are increasingly important in nearly every field of science, medicine, engineering, and even the fine arts. This reflects rapid growth of 3D scanning instrumentation, visualization and analysis algorithms, graphic displays, and graduate training. These capabilities are recognized as critical for future advances in science, and the U.S. National Science Foundation (NSF) is one of many funding agencies increasing support for 3D imaging and computing. For example, one new initiative aims at digitizing biological research collections, and NSF’s Earth Sciences program will soon announce a new target in cyberinfrastructure, with a spotlight on 3D imaging of natural materials.

Many consider the advent of 3D imaging a “scientific revolution” (1), but in many ways the revolution is still nascent and unfulfilled. Given the increasing ease of producing these data and rapidly increased funding for imaging in non-medical applications, a major unmet challenge to ensure maximal advancement is archiving and managing the science that will be, and even has already been, produced. To illustrate the problems, we focus here on one domain, volume elements or “voxels,” the 3D equivalents of pixels, but the argument applies more broadly across 3D computing. Useful, searchable archiving requires infrastructure and policy to enable disclosure by data producers and to guarantee quality to data consumers (2). Voxel data have generally lacked a coherent archiving and dissemination policy, and thus the raw data behind thousands of published reports are not released or available for validation and reuse. A solution requires new

infrastructure and new policy to manage ownership and release (3).

Technological advancements in both medical and industrial scanners have increased the resolution of 3D volumetric scanners to the point that they can now volumetrically digitize structures from the size of a cell to a blue whale, with exquisite sensitivity toward scientific targets ranging from tissue properties to the composition of meteorites. Voxel data sets are generated by rapidly diversifying instruments that include x-ray computed tomographic (CT) scanners (Fig. 1), magnetic resonance imaging (MRI) scanners (Fig. 2), confocal microscopes, synchrotron light sources, electron-positron scanners, and other innovative tools to digitize entire object volumes.

CT and MRI have evolved across the widest range of applications. CT is sensitive to density. Its greatest strength is imaging dense materials like rocks, fossils, the bones in living organisms and, to a lesser extent, soft tissue. The first clinical CT scan, made in 1971, used an 80 by 80 matrix of 3 mm by 3 mm by 13 mm voxels, each slice measuring 6.4 Kb and taking up to 20 min to acquire. Each slice image took 7 min to reconstruct on a mainframe computer (4). In 1984, the first fossil, a 14-cm-long skull of the extinct mammal *Stenoposchoerus*, was scanned in its entirety (5), signaling CT’s impact beyond the clinical setting and in digitizing entire object volumes. The complete data set measured 3.28 Mb. With a mainframe computer, rock matrix was removed to visualize underlying bone, and surface models were reconstructed in computations taking all night. By 1992, industrial adaptations brought CT an order-of-magnitude higher resolution (high-resolution x-ray computed tomography, HRXCT) to inspect much smaller, denser objects. A fossil skull of the stem-mammal *Thrinaxodon* (68 mm long) was scanned in 0.2-mm-thick slices measuring 119 Kb each. Scanning the entire volume took 6 hours, and the complete raw data set occupied 18.4 Mb (6). The scans revealed all internal details reported earlier, from destructive mechanical serial sectioning of a different specimen, and pushed

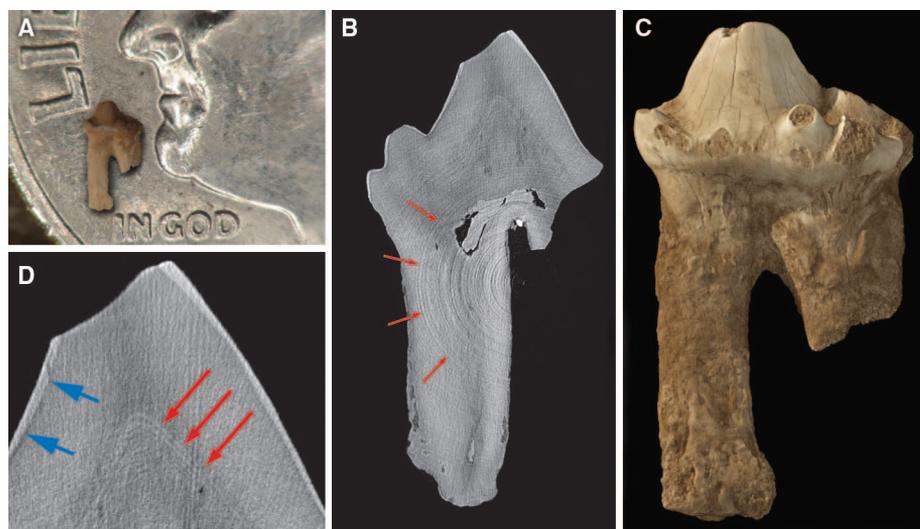


Fig. 1. (A) Photomicrograph of fossil tooth (*Morganucodon* sp.). (B) MicroXCT slice, 3.2- μ m voxel size; arrows show ring artifact, a correctable problem with raw data but an interpretive challenging with compressed data. (C) Digital 3D reconstruction. (D) Slice through main cusp; red arrows show growth bands in dentine, and blue arrows mark the enamel-dentine boundary, both observed in *Morganucodon* for the first time in these scans.

¹Jackson School of Geosciences, The University of Texas, Austin, TX 78712, USA. ²Department of Radiology, University of California, San Diego, CA 92093, USA.

*To whom correspondence should be addressed. E-mail: rowe@mail.utexas.edu (T.R.); lfrank@ucsd.edu (L.R.F.)