

Matplotlib - Foundation of Visualization in Python

July 26, 2019

1 Goals

Matplotlib (<http://matplotlib.org>) is the fundamental data visualization library for the scientific Python stack of tools, providing a visual interface used by over a million¹ scientists in conjunction with other foundational tools like Numpy and SciPy (<https://www.scipy.org/about.html>). It is used across a wide spectrum of fields, including bio-medical imaging, microscopy, and genomics [CJL⁺06, WAT18, vdWSNI⁺14] [SIW⁺11, GHWB09, HWSY12, KR12, CRDZ⁺14, LRA⁺14, JSC⁺15, APE⁺14], and we expect the user base to continue to grow as Python is adopted by more scientists in the life sciences. There are many downstream packages that build on Matplotlib to implement domain-specific plotting tools.

Matplotlib has been actively developed and maintained by a vibrant, primarily volunteer, community over the last 16 years. However, given the scale, scope, and importance of the project, we are at the limit of what can be developed and maintained with a primarily volunteer effort.

This proposal supports 2.4 person years of developer support for Matplotlib, support that is crucial to set the groundwork and shepherd the project through the next 16 years:

- a) Maintenance of the library, including curating new and existing Issues and Pull Requests.
- b) Developing a comprehensive plan to evolve the core architecture of Matplotlib.
- c) Developing the tools, documentation, and community to foster a rich eco-system of domain-specific plotting tools built on Matplotlib.

1.1 Revitalizing the developer community

Matplotlib is a community driven project, however we are large enough that we need supported developers who have the time organize, plan and make timely decisions. Currently, New Issues and Pull Requests are being submitted faster than our volunteers can review them. Over the past few years we have merged PRs and closed Issues at about [X per month], but about [Y] new issues and PRs are opened monthly; currently we have about 1200 open issues and 300 open PRs. Among the latter are useful contributions and bug fixes that, possibly with additional attention and polish, could improve Matplotlib for direct users and downstream packages. The backlog is discouraging for new or occasional contributors, and distracting for core developers.

To maintain Matplotlib's health we need to do the following:

- Curate the current backlog of Issues and Pull Requests in terms of topic, difficulty, and urgency.
- Label and review newly opened Issues and Pull Requests promptly.

¹Estimated from `pypi` download numbers, `conda` download numbers, and the number of unique monthly visitors to the documentation website

- Fix critical bugs and regressions promptly.
- On-board new contributors to the project team; critical to sustaining and diversifying our developer community.
- Maintain backward compatibility. If we do break API, ensure it is intentional and well documented.
- Manage discussions about proposed enhancements, features, and breaking API changes.

None of this is to demote the importance of the volunteer contributors, but instead to better co-ordinate and nurture their efforts, with the goal of growing and sustaining a diverse community of expert contributors, both volunteer and paid.

1.2 Road-map and Architecture

Matplotlib needs sustained attention to map out the library’s architecture for the next decade. The current architecture² of Matplotlib was developed 15 years ago [Hun07]. That it is still in use is a testament to its initial design; but that design does not reflect recent developments in data structures, software design, and visualization. Matplotlib does not natively know how to exploit structured (e.g. `pandas` or `xarray`) or streaming data. Some of the design choices about how to store user data internally and where the main plotting name space lives are starting to be constraints as Matplotlib scales to more domains. While there are many downstream domain-specific libraries that are built on top of Matplotlib, interoperability between them is problematic.

1.2.1 Homogenization of the Application Programming Interface (API)

The library has grown organically over time through the contributions of many people (approx. 900 individuals) and the code has accumulated many small inconsistencies in the API. Similar methods have different argument order, e.g., `ax.text(x, y, s)` vs `ax.annotation(s, (x, y))`, and some keyword arguments can be singular or plural, e.g., `color` vs `colors`. These subtle issues add friction for users, but are hard to fix without breaking existing code somewhere in our large user base. Our goal is to minimize breakage, but still unify the API. Taking into account **all** of the APIs, we will carefully consider which to leave as they are, which to deprecate, and which to replace.

1.2.2 API generalization

Currently Matplotlib has two main user-facing APIs: the `pyplot` API and the `Object Oriented (OO)` API. The `pyplot` API closely follows MATLAB and is built around the concept of a global “current Axes”. While convenient for quick interactive usage, it frequently produces surprising and undesired coupling between in the code when used in libraries. On the other hand the OO interface is more explicit and flexible but more verbose. However, the main name space for plotting methods is on the `Axes` object which leads to three issues with the API. First, third party domain-specific packages can never feel “First Class” as they will not be implemented as `Axes` methods like the “built in” plots. Second, some visualizations require putting `Artists` on multiple `Axes` which can not be naturally expressed as an `Axes` method. Lastly, there are over 250 methods on the `Axes`, all of the plotting methods and some additional `Axes` specific methods, which makes it extremely hard to discover if the method you need exists via tab-completion.

To address these issues we will move to top-level functions that take in data, style, and `Axes` as the primary API. During this refactoring we will use consistent naming and call conventions,

²<https://www.aosabook.org/en/matplotlib.html>

as discussed in 1.2.1, that can be used by downstream libraries. These functions will return the rich composites discussed in section 1.2.3 and consume the data objects discussed in section 1.2.4. From the large pool of plotting functions we will curate domain-specific name spaces, which can be augmented by downstream libraries to enable users to discover the functionality they need.

This is a huge change in the API that will require time to carefully consider the consequences.

1.2.3 Rich Composite Artists

Artists are the “middle layer” of Matplotlib that encode user-intent, style, and data. To update the style or data users interact with the **Artists** objects and **Artists** know how to turn the user’s input into colored pixels. Currently, Matplotlib has a mix of “primitive” **Artists** (e.g. lines, images, and text) and “composite” artists (e.g. the whole **Figure**).

The mapping between the user API and the **Artists** can be one-to-one, but often one user call will generate many decoupled **Artists**. For example `hist` displays a histogram, but returns a list of independent **Rectangle Artists** (one per bin). If the user wants to update the data or the style, they must adjust each **Rectangle** independently. Instead, we will have a single object that can have its data and style updated to simplify user code for interactive data exploration, consuming streaming data, and generating animations.

1.2.4 Data Model

“Structured data” combines multiple pieces of, possibly heterogeneous, data along with labels, metadata, and the the relationship between the components into single data structure is not natively supported by Matplotlib. Currently users must split the structures up and pass the components individually into plotting methods, destroying the structure.

Further, each plotting method and **Artists** handle sanitizing and storing data independently. This means that some common functionality, such as handling data with attached units (e.g., degrees Celsius, dates), is scattered throughout the code base. This leads to inconsistencies across library and makes it difficult to write code that updates the data or style for interactive exploration, streaming, and animation use cases.

We will re-organize the internal data representation in Matplotlib to a model appropriate for the base Matplotlib library and, more importantly, to be the technical underpinning to handle, exploit, and update structured data in a coherent fashion. By removing the direct data storage from the **Artists** and defining an API for data sources we will enable:

- native consumption of structured data;
- smart down sampling of plotted data based on view limits;
- seamless update the underlying data, either streaming or interactively;
- non-materialized data sources such as database queries or analytic functions

We will decouple the development of the data access from the **Artists**. Downstream libraries will be able to provide sophisticated data sources to the **Artists** in the core library or sophisticated **Artists** that use the data sources from core.

1.2.5 Additional Export Methods

Matplotlib **Figures** can be render to either raster or vector file formats. From there it is displayed to an interactive window or saved to disk and be used like any other image file. There is currently no good way to “reopen” a Matplotlib **Figure** or export it to another plotting library, such as

`bokeh`, `d3` or `QtCharts`. Due to the way Matplotlib internals are implemented, it is difficult to take advantage GPUs to accelerate drawing.

To address these problems we will investigate adding two additional export paths. One at a high-level, suitable for a Matplotlib-specific file format and interoperability with other high-level plotting libraries, and one at a low level scene-graph level, suitable to pass to a GPU.

1.3 Coordination with downstream projects

The most common visualizations in a domain need to be one or two simple lines of code for the end-practitioners with the “obvious” customization options surfaced. Our goal is to make these libraries as thin and easy to write as possible. This means there will always be a need for domain specific visualization libraries. Much of the domain-specific specialization is carried in the semantics of the structured data and the specific visualization needs of the domain.

To this end we will identify and engage with downstream libraries in the life sciences that are currently using Matplotlib for their visualization to identify their pain points and ensure that we are actually solving their problems. In particular we plan to engage with `scikit-learn`³, `CellProfiler`⁴, `scanpy`⁴, `starfish`⁴, `nipy`, and `scikit-image`⁴

2 Expected outcomes, success evaluation and metrics

2.1 Issue and PR curation

Quantitatively evaluating maintenance work can be tricky, some Issues or PRs can take minutes to review where as others can take days to months of effort, however we believe that there is value at looking at the net number of new Issues and Pull requests. We will reduce this number, ideally closing Issues and Pull Requests faster than they are opened. NumPy has had success in reversing the ever increasing trends in the number of Issues / Pull Requests with paid developers⁵.

We will evaluate and label every open Issue and Pull Request determining: assigning an action, a priority, and an estimated difficulty. Once that is done, we will aim to have all new Issues and Pull Requests labeled within 7 days of being opened.

2.2 Road-map and Architecture

We will write a white paper and road map documenting the proposed design, critical use-cases, and requirements for the data model and API overhaul.

To validate the design, we will develop end-to-end prototypes targeting one or more of the life-science libraries discussed above.

3 Work Plan

The funds will be paid to:

- Fund Thomas Caswell’s position at 40%. Caswell is currently the lead developer of Matplotlib and an Associate Computational Scientist at Brookhaven National Laboratory. His long-term experience, API design expertise, and project leadership are critical to the success of the work in this proposal. He will work on all aspects of the proposal.

³Also applying for Essential Open Source Software for Science

⁴Currently funded by CZI

⁵https://github.com/seberg/numpy-talk-plots/blob/master/plots-used.in.talk/issues-prs_delta.pdf

- Fund Hannah Aizenman’s position for 12 months. Aizenman has been a core-contributor Matplotlib for three years and has previously contributed support for string-categorical values. She is a PhD candidate in computer science studying visualization at The City College of New York. Her work on the architecture of Matplotlib will be the basis of her PhD thesis. Aizenman will take the lead on the data model design and new-contributor on-boarding.
- Fund 12 months of a yet-to-be identified software engineer to support all aspects of the proposal but focusing on maintenance, prototyping, and engaging down-stream libraries.
- Travel to key Scientific and Python conferences (such as SciPy or PyCon) and for in-person meetings if required.

We want to use this dedicated effort to leverage and empower the Matplotlib developer community. In terms of direct work on the code base an equal amount of time will be spent mentoring and reviewing code from community members rather than directly implementing features or fixing bugs. All of the design work will be done in public with input from the community.

Part of this work is to develop the project road-map.

4 Existing Support

Thomas Caswell has 4hrs/wk from Brookhaven National Lab to work on Matplotlib.

References

- [APE⁺14] Alexandre Abraham, Fabian Pedregosa, Michael Eickenberg, Philippe Gervais, Andreas Mueller, Jean Kossaifi, Alexandre Gramfort, Bertrand Thirion, and Gael Varoquaux. Machine learning for neuroimaging with scikit-learn. *Frontiers in Neuroinformatics*, 8:14, 2014.
- [CJL⁺06] Anne E. Carpenter, Thouis R. Jones, Michael R. Lamprecht, Colin Clarke, In Han Kang, Ola Friman, David A. Guertin, Joo Han Chang, Robert A. Lindquist, Jason Moffat, Polina Golland, and David M. Sabatini. Cellprofiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biology*, 7(10):R100, Oct 2006.
- [CRDZ⁺14] Thomas M. Carlile, Maria F. Rojas-Duran, Boris Zinshteyn, Hakyung Shin, Kristen M. Bartoli, and Wendy V. Gilbert. Pseudouridine profiling reveals regulated mrna pseudouridylation in yeast and human cells. *Nature*, 515:143 EP –, Sep 2014.
- [GHWB09] Ryan N. Gutenkunst, Ryan D. Hernandez, Scott H. Williamson, and Carlos D. Bustamante. Inferring the joint demographic history of multiple populations from multidimensional snp frequency data. *PLOS Genetics*, 5(10):1–11, 10 2009.
- [Hun07] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.
- [HWSY12] Tamar Hashimshony, Florian Wagner, Noa Sher, and Itai Yanai. Cel-seq: Single-cell rna-seq by multiplexed linear amplification. *Cell Reports*, 2(3):666 – 673, 2012.
- [JSC⁺15] Xiaolong Jiang, Shan Shen, Cathryn R. Cadwell, Philipp Berens, Fabian Sinz, Alexander S. Ecker, Saumil Patel, and Andreas S. Tolias. Principles of connectivity among morphologically defined cell types in adult neocortex. *Science*, 350(6264), 2015.
- [KR12] Johannes Kster and Sven Rahmann. Snakemakea scalable bioinformatics workflow engine. *Bioinformatics*, 28(19):2520–2522, 08 2012.
- [LRA⁺14] Arthur Laganowsky, Eamonn Reading, Timothy M. Allison, Martin B. Ulmschneider, Matteo T. Degiacomi, Andrew J. Baldwin, and Carol V. Robinson. Membrane proteins bind lipids selectively to modulate their structure and function. *Nature*, 510:172 EP –, Jun 2014.
- [SIW⁺11] Nicola Segata, Jacques Izard, Levi Waldron, Dirk Gevers, Larisa Miropolsky, Wendy S. Garrett, and Curtis Huttenhower. Metagenomic biomarker discovery and explanation. *Genome Biology*, 12(6):R60, Jun 2011.
- [vdWSNI⁺14] Stfan van der Walt, Johannes L. Schnberger, Juan Nunez-Iglesias, Franois Boulogne, Joshua D. Warner, Neil Yager, Emmanuelle Gouillart, and Tony and Yu. scikit-image: image processing in python. *PeerJ*, 2:e453, June 2014.
- [WAT18] F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. Scanpy: large-scale single-cell gene expression data analysis. *Genome Biology*, 19(1):15, Feb 2018.