



# The CPTAC Python API

Caleb M. Lindgren, David W. Adams, Benjamin Kimball, Hannah Boekweg, Sadie Tayler, Samuel Pugh, Samuel H. Payne  
Department of Biology, Brigham Young University, Provo, UT, USA

## Abstract

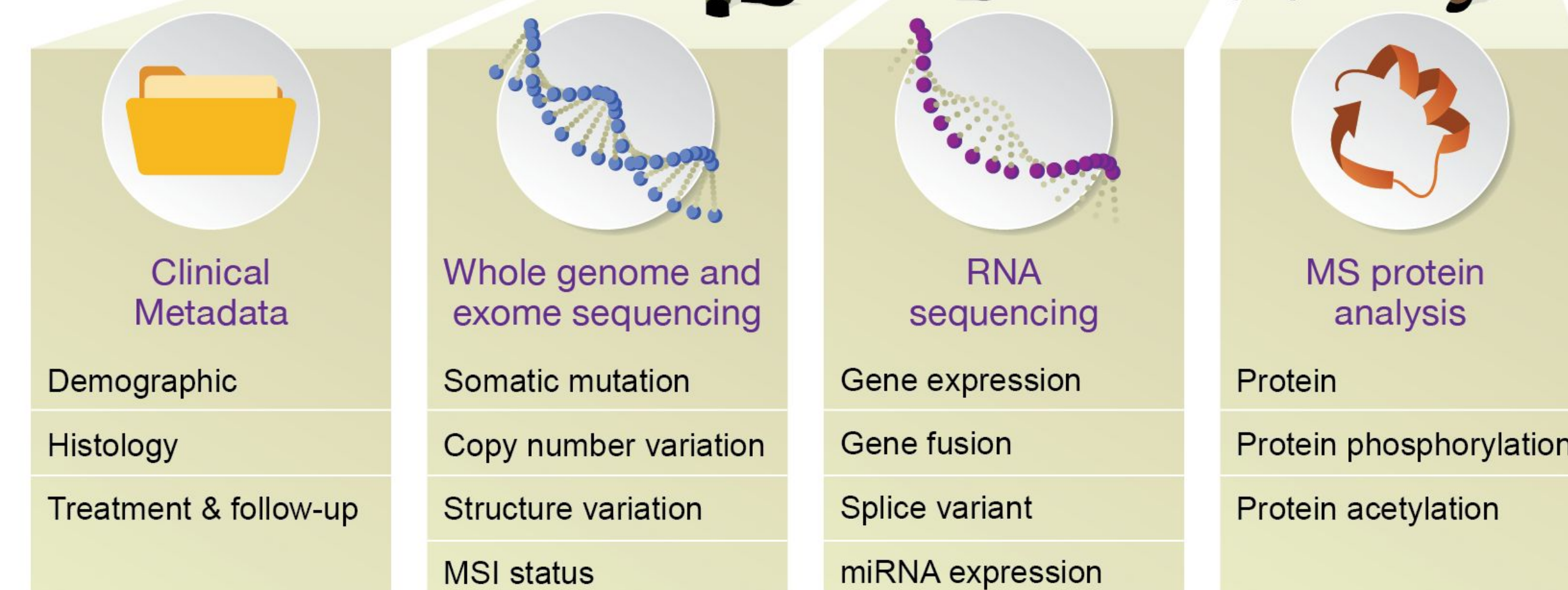
We present a new method for data sharing across large collaborations to improve reproducibility and transparency, by creating a Python package that serves as an interface (API) to the multi-omics characterization of tumors from NCI's CPTAC program.

## Introduction

Cancer data has many audiences, including clinicians, biologists, data scientists and patients. Sharing data and analyses across these diverse audiences is challenging. In particular, we want to simplify the link between data and analysis scripts to enable easier data exploration. We embed NCI's CPTAC data into a software API. Each tumor type describes samples with clinical, omics, and imaging data.

CPTAC Patient Cohort  
~100 tumors

- Representation of genomic & histologic subtypes
- Adjacent and normal tissue



## Results

A Python package, `cptac`, facilitates access to all CPTAC data. Using the package as the single point of data access unifies and simplifies analysis methods across diverse consortia. Publicly available data includes: colon, ovarian, endometrial, renal, lung adenocarcinoma. New datasets coming online soon.

### Features

- One step installation via pip
- Consistent data types and formats across cancers
- Data encapsulated in Pandas dataframes, meaning no need for writing parsers
- API handles table joining between data types
- Seamlessly works with numeric and graphing libraries (numpy, pandas, matplotlib, seaborn, etc.)
- Versioned data releases

```
In [6]: import cptac
en = cptac.Endometrial()
proteomics = en.get_proteomics()
proteomics.head(3)

Out[6]:
```

Sample_ID	A1BG	A2M	A2ML1	AGALT	AAAS	AACS	AADAT	AAED1	AAGAB	AAK1	...	ZSWIM8	ZSWIM9	ZW10	ZWLICH	ZWNT	ZXDC	ZYG11
S001	-1.180	-0.863	-0.802	0.222	0.256	0.665	1.2800	-0.3390	0.412	-0.664	...	-0.08770	NaN	0.0229	0.109	NaN	-0.332	-0.4330
S002	-0.685	-1.070	-0.684	0.884	0.135	0.334	1.3000	0.1390	1.330	-0.367	...	-0.03560	NaN	0.3650	1.070	0.737	-0.564	-0.0046
S003	-0.528	-1.320	0.435	NaN	-0.240	1.040	-0.0213	-0.0479	0.419	-0.500	...	0.00112	-0.145	0.0165	-0.116	NaN	0.151	-0.0740

```
In [7]: clinical = en.get_clinical()
clinical.head(3)

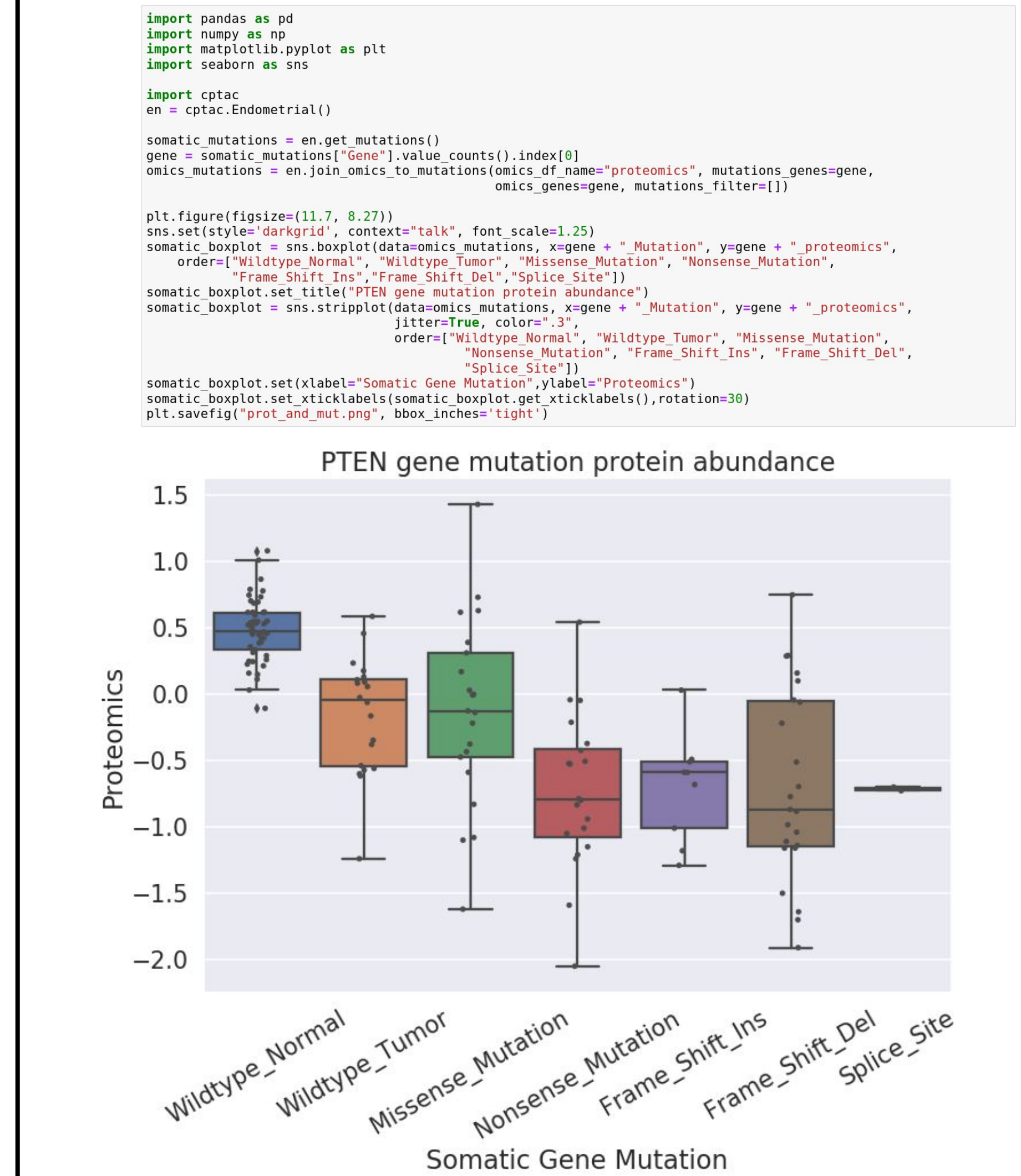
Out[7]:
```

Sample_ID	Patient_ID	Proteomics_Tumor_Normal	Country	Histologic_Grade_FIGO	Myometrial_Invasion_Specify	Histologic_Type	Treatment_naive	Tumor_J
S001	C3L-00006	Tumor	United States	FIGO grade 1	under 50 %	Endometrioid	YES	N
S002	C3L-00008	Tumor	United States	FIGO grade 1	under 50 %	Endometrioid	YES	N
S003	C3L-00032	Tumor	United States	FIGO grade 2	under 50 %	Endometrioid	YES	N

Table 1: User documentation

Tutorial 1: Data intro	Basics of installation, data access
Tutorial 2: Pandas	Exploring the data using Pandas
Tutorial 3: Join dataframes	Built-in functions for joining different data types
Tutorial 4: MultiIndex	Unique aspects of multi-level column indexes.
Tutorial 5: Updates	Working with data and package version updates.
Use case 1: Multi-omic integration	Data access and integration for multiple omics data types
Use case 2: Clinical covariates	Explores meta-data for correlation between clinical attributes
Use case 3: Clinical and acetylation	Compares acetylation levels between tumor subtypes
Use case 4: Mutations and omics	Studies the effects of DNA mutations on protein abundance
Use case 5: Enrichment analysis	Uses GSEAPy to find enriched pathways
Use case 6: Derived molecular	Work with attributes derived from molecular data, e.g. MSI status
Use case 7: Trans genetic effect	Effect of DNA mutations on the expression of a different protein
Use case 8: Outliers	Uses Blacksheep to study outliers in expression
Use case 9: Clinical outcomes	Explores patient follow up data, molecular data and patient survival

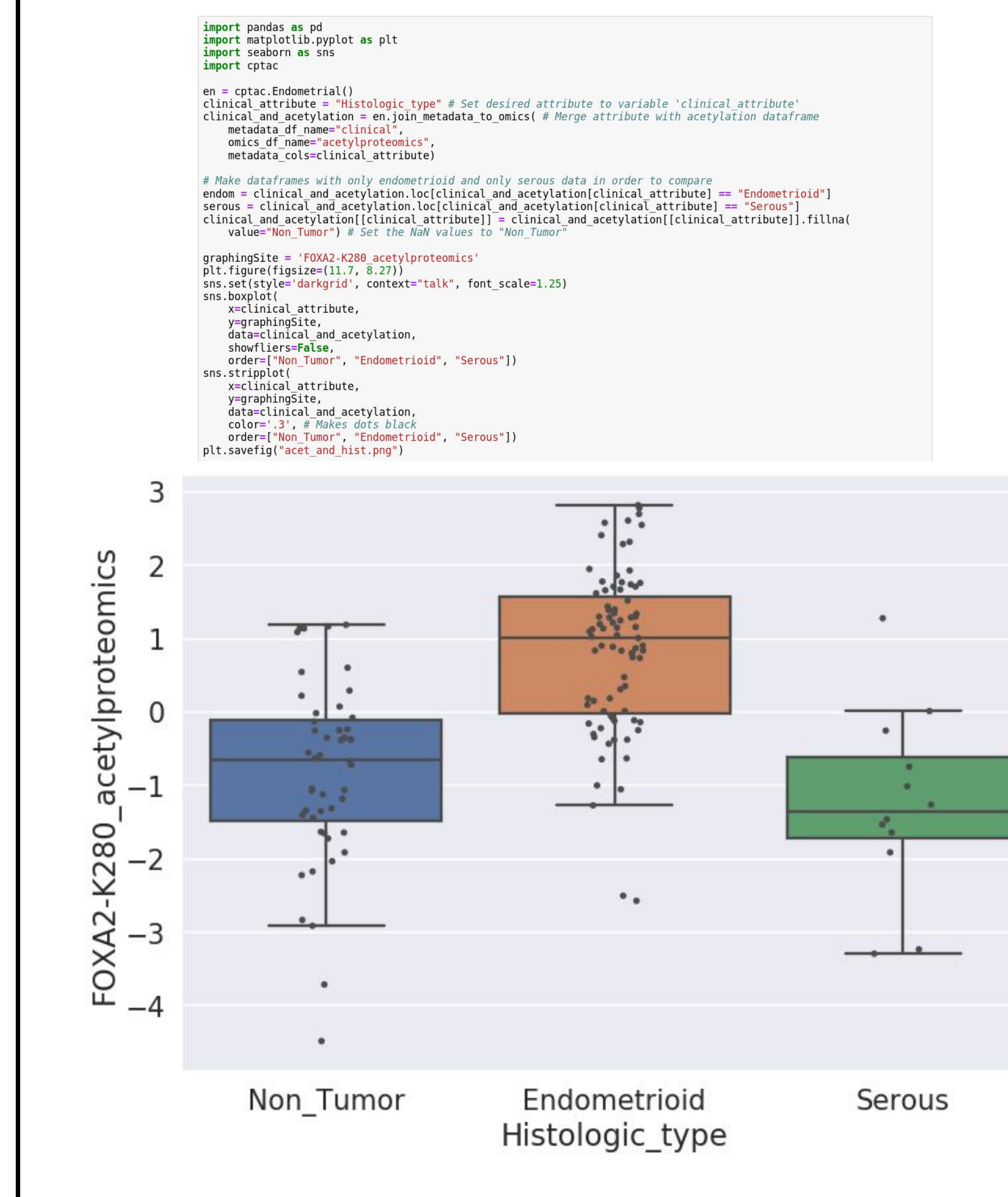
Mutation effects on the proteome



Comparing clinical attributes



Comparing histology with acetylation



## Conclusions

The `cptac` Python package brings cancer data to dispersed collaborative groups. Our package incorporates multiple data sets and lowers the entry barrier, expanding our audience while improving reproducibility and transparency.

**Acknowledgments:** National Cancer Institute (NCI) CPTAC award U24 CA210972.

**Contact:** calebmlindgren@gmail.com; sam\_payne@byu.edu

<https://payne.byu.edu>

[github.com/PayneLab](https://github.com/PayneLab)

@byu\_sam