### Unified access to cancer proteogenomics data



Caleb M. Lindgren, Hannah Boekweg, David W. Adams, Sadie Tayler, CPTAC Investigators, 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

### 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. Current tumor types include: colon, ovarian, endometrial, renal, breast, lung adeno, glioblastoma and head/neck.

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

#### nrogram

# 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.

#### Features

- One step installation via pip
- Consistent data and formats across cancers
- Data presented in Pandas dataframes, meaning no need for writing parsers
- API facilitates joining across data types
- Seamlessly works with numeric and graphing libraries (numpy, pandas, matplotlib, seaborn, etc.)
- Versioned data releases
- Package automatically checks that it's up to date, and that the latest version of the data

<pre>import cp en = cpta</pre>	ptac ac.End	ometri	.al()															
proteomi proteomi	cs = e cs.hea	n.get_ d(3)	proteo	mics()														
	A1BG	A2M	A2ML1	A4GALT	AAAS	AACS	AADAT	AAED1	AAGAB	AAK1		ZSWIM8	ZSWIM9	ZW10	ZWILCH	ZWINT	ZXDC	ZYC
Sample_ID																		
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	-0.4
S002	-0.685	-1.070	-0.684	0.984	0.135	0.334	1.3000	0.1390	1.330	-0.367		-0.03560	NaN	0.3630	1.070	0.737	-0.564	-0.0
clinical	= en. .head(	get_cl 3)	inical.	()														
ctinicat				Tumor N	lormal	Country	Histolo	gic_Grad	e_FIGO	Myomet	rial_	invasion_	Specify H	listologio	c_type T	reatment	_naive	Tum
CLINICAL	Patient	LID Pr	oteomics															
Sample_ID	Patient	t_ID Pr	oteomics															
Sample_ID	Patient C3L-00	t_ <b>ID Pr</b> 4	oteomics		Tumor	United States		FIGO	grade 1			und	er 50 %	Endom	netrioid		YES	
Sample_ID S001 S002	Patient C3L-00 C3L-00	2_ <b>ID</b> Pr 006 008	oteomics		Tumor Tumor	United States United States		FIGO	grade 1 grade 1			und	er 50 % er 50 %	Endom Endom	netrioid netrioid		YES YES	

Col	mparing clinical attributes
<pre>import pand</pre>	las <mark>as</mark> pd
import matp	>lotlib.pyplot <mark>as</mark> plt
import seat	porn <mark>as</mark> sns
import cpta	ac
en = cptac.	Endometrial()
clinical_da	ata = en.get_clinical()
simplified_	_clinical = clinical_data[[' <mark>FIGO_stage', 'BMI'</mark> ]]
pd.set_opti simplified_ simplified_ simplified_ simplified	<pre>lon('mode.chained_assignment', None) _clinical.loc[simplified_clinical.FIG0_stage == 'IA', 'FIG0_stage'] = 'I' _clinical.loc[simplified_clinical.FIG0_stage == 'IB', 'FIG0_stage'] = 'I' _clinical.loc[simplified_clinical.FIG0_stage == 'IIIA', 'FIG0_stage'] = 'III' _clinical.loc[simplified_clinical.FIG0_stage == 'IIIA', 'FIG0_stage'] = 'III'</pre>

**CPTAC Endometrial** Carcinoma Cohort 87 endometrioid tumors 13 serous tumors 49 normal uterine samples • 18 normal endometrium • 31 mixed endometriummyometrium ...... **RNA MS** protein Whole genome and analysis exome sequencing sequencing Gene expression Protein Somatic mutation Copy number variation Gene fusion Protein phosphorylation Splice variant Protein acetylation Structure variation MSI status miRNA expression

is installed

#### Mutation effects on the proteome





A1BG\_cross = en.join\_omics\_to\_omics(df1\_name="proteomics", df2\_name="transcriptomics", genes1="A1BG", genes2="A1BG")

title='Transcriptomics vs. Proteomics for the A1BG gene')
plt.savefig("prot and tran.png")



## 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; hannahboekweg@gmail.com; sam\_payne@byu.edu





