



# CRAVAT (Cancer-Related Analysis of Variants Toolkit)

Integration into Galaxy-P and extension towards proteogenomic studies



Ray Sajulga<sup>1</sup>, Subina Mehta<sup>1</sup>, Praveen Kumar<sup>1,2</sup>, James E. Johnson<sup>3</sup>, Candace Guerrero<sup>1</sup>, Michael Ryan<sup>4</sup>, Rachel Karchin<sup>5,6,7</sup>, Pratik D. Jagtap<sup>1</sup>, and Timothy J. Griffin<sup>1</sup>

<sup>1</sup>University of Minnesota, Biochemistry, Molecular Biology and Biophysics; <sup>2</sup>University of Minnesota-Rochester, Bioinformatics and Computational Biology Program; <sup>3</sup>University of Minnesota, Minnesota Supercomputing Institute; <sup>4</sup>In Silico Solutions, LLC; <sup>5</sup>Johns Hopkins University, Department of Biomedical Engineering; <sup>6</sup>Johns Hopkins University, The Institute for Computational Medicine; <sup>7</sup>Johns Hopkins University School of Medicine, Department of Oncology

## Introduction

- CRAVAT - an existing tool suite developed by the Karchin lab at Johns Hopkins University and *In Silico Solutions*
- Compiles variant information:
  - Literature/knowledgebase annotation:
    - PubMed, GeneCards, ClinVar, COSMIC, gnomAD, CGC, etc...
  - Harnesses machine learning algorithms for predicting cancer impact and pathogenicity
  - Visualizes variant location on 3-D protein models
  - Displays interaction gene networks (NDeX)
  - Offers extensibility for other -omics analyses:
    - Docker image
    - RESTful API

## Galaxy Tool

CRAVAT Submit, Check, and Retrieve Submits, checks for, and retrieves data for cancer annotation (Galaxy Version 0.1.0)

Source file: 1: FreeBayes.vcf

Include proteogenomic input?  Yes  No

Peptides with Genomic Coordinates (ProBED Format): 2: Peptide Genomic Coordinate.bed

Submit only intersected variants?  Yes  No

Output intersected genomic file?  Yes  No

Analysis Program: VEST

Genome Reference Consortium Human Build (GRCh): GRCh38/hg38

Execute

- Parameters:**
- Input file (CRAVAT or VCF format)
  - Proteogenomic input
    - Submit only intersected variants to CRAVAT
    - Output the intersected genomic file
  - Analysis Programs (machine learning predicting algorithms):
    - VEST (Variant Effect Scoring Tool) – evaluates variants based on pathogenicity
    - CHASM (Cancer-specific High-throughput Annotation of Somatic Mutations)
  - Genome Reference Consortium Human Build
- \* New parameters in red*

9: CRAVAT Results: data 1 and data 2 using VEST

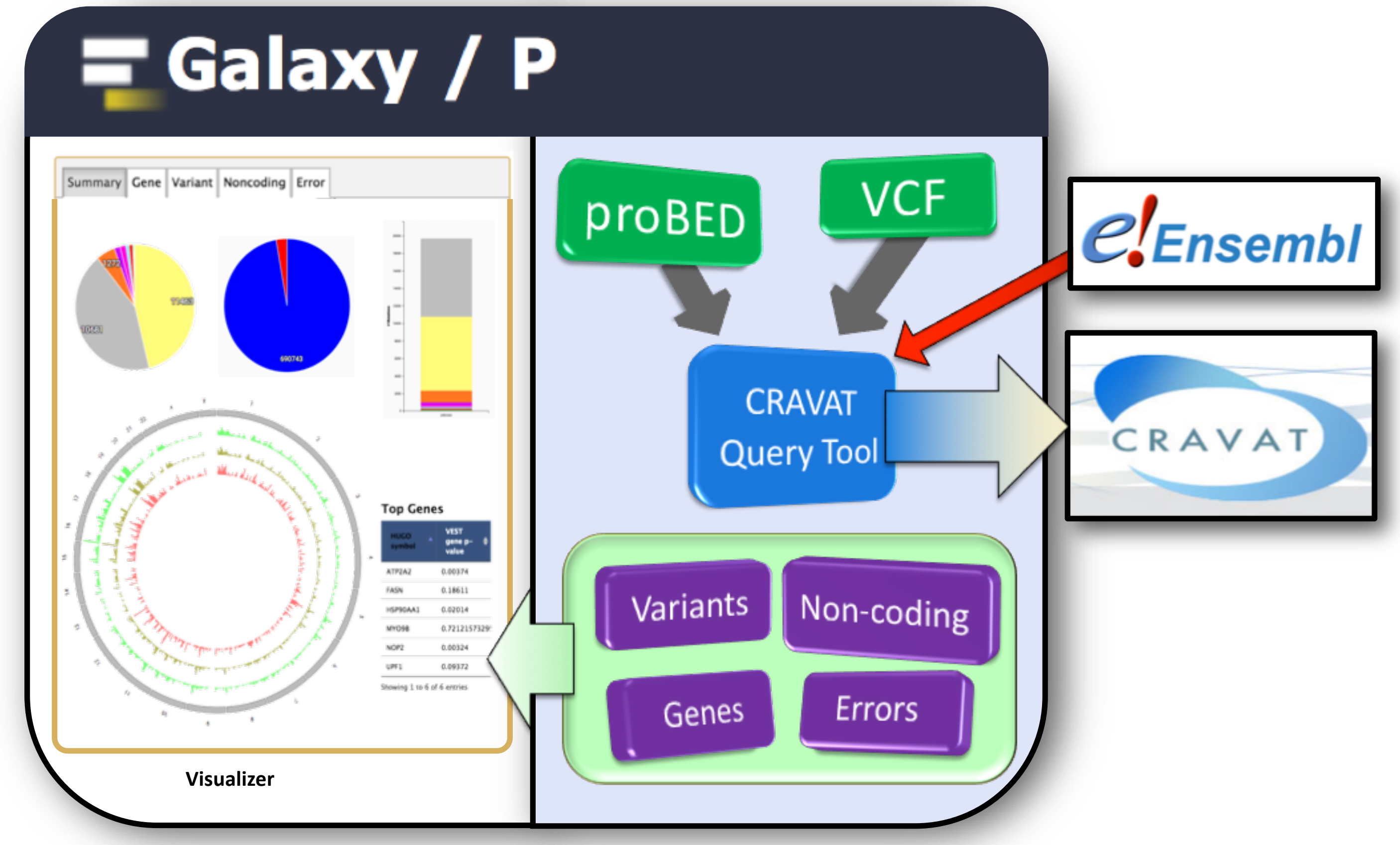
a list of 4 datasets

- Gene
- Variant
- Noncoding
- Error

## Methodology

ID	Chr.	Position	Strand	Ref. base	Alt. base	Chr.	Start	End	Peptide	Strand
TR1	chr14	94079142	-	G	T	chr14	94079127	94079178	ADVSAWKDLFVPGVLR	-
TR2	chr10	121520166	-	G	A	chr10	121520124	121520166	EADSPVSVFLVHQNQR	-
TR3	chr13	48459831	+	C	A	chr13	48459831		EWGSGSDILR	+
TR4	chr7	116777451	+	G	T	chr7	116777451		GVVDSENLPLNISR	+
TR5	chr7	140753336	-	T	A	chr7	140753336		IQSHCSYTYGRMGEPGAE	-
TR6	chr17	39724745	+	G	T				GHFVGVCDLSLTSK	-
Ins1	chr17	39724745	+	-	T					
Del1	chr17	39724745	+	A	-					
CSub1	chr2	39644095	+	ATGCT	GA					

- The updated Galaxy tool performs an intersection between the genomic (VCF) and proteogenomic (proBED) file to find which variants were shown to be expressed (via peptides).
- Utilizes the S.O. (sequence ontology) transcript IDs annotated by CRAVAT to obtain protein sequences from Ensembl for amino acid sequence confirmation (and alters the reference sequence based on the input mutations)



## Future Directions

- Integrate peptides into the circos plot as a "pircos" plot
- Support interpretation of more complex variants (e.g., alternative splicing)
- Include the user's e-mail within the CRAVAT API submission (may allow for notification of job completion and a unique job ID)
- Allow for exporting selected variants back into the Galaxy workspace
- Establish interactivity with our Multi-omics Visualization Plugin (MVP)
- Explore using a Galaxy Interactive Environment through CRAVAT's Docker image

## Visualization Plugin

- Sortable and filterable variants table with column visibility toggled in the sidebar
- Gene, positional, and disease-related information about the variant.
- Protein diagram that shows domains and known variants from TCGA (The Cancer Genome Atlas) displayed as a lollipop diagram.
- Bar meters displaying allele frequencies from 1000 Genomes, ESP6500, and gnomAD.
- Mapping of variant amino acid positions onto three-dimensional structure of protein using MuPIT (Mutational Position Imaging Toolbox)
- Visualization of network of known interactions for the variant gene and gene product using NDeX (Network Data Exchange)

## Conclusions

- We have enabled a Galaxy-based workflow for integrating peptide sequence variants (detected via proteogenomics) with the CRAVAT analysis suite.
- This extension demonstrates the value of the Galaxy platform for developing and disseminating sophisticated proteogenomic workflows.
- Also, this extension uses the recommended community standard output for proteogenomic, proBED, which can be generated from a Galaxy-based proteogenomics workflow.

## Acknowledgements

- This project is supported by the Informatics Technologies for Cancer Research (ITCR) program at the NIH/NCI, from grant U24CA204817 to R. Karchin and grant U24CA199347 to T. Griffin
- Demonstration proteogenomic data were generated with the assistance of the University of Minnesota Center for Mass Spectrometry and Proteomics and Genome Center
- We acknowledge use of the Jetstream cloud-based computing resource for scientific computing maintained at Indiana University for assistance in maintaining the publicly available Galaxy instance for demonstration purposes.