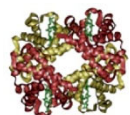# Proteogenomic workflows

**Tim Griffin**
**University of Minnesota**
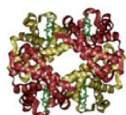
*tgriffin@umn.edu*

*Learn more at galaxyp.org*
z.umn.edu/itcrgalaxyvideo

# Acknowledgements

### Biochemistry, Molecular Biology & Biophysics

**Dr. Pratik Jagtap (Co-leader, Galaxy-P)**
Praveen Kumar
Subina Mehta
Caleb Easterly
Ray Sajulga
Andrew Rajczewski
Dr. Shane Hubler
Mark Esler
Dr. Art Eschenlauer
Dr. Candace Guerrero
Matt Chambers

### GalaxyP

*Collaborators*
David Largaespada
Frank Ondrey
Mo Heydarian/Karen Reddy
Brian Crooker/Wanda Weber
Bart Mesuere
Brook Nunn
Thilo Muth
Magnus Øverlie Arntzen

### Minnesota Supercomputing Institute

**James Johnson**
**Tom McGowan**
Dr. Getiria Onsongo
Dr. Michael Milligan

## COMMUNITY-BASED SOFTWARE DEVELOPMENT

**Harald Barsnes and Marc Vaudel**
*University of Bergen, Bergen, Norway*
**Bjoern Gruening (Galaxy community...)**
*University of Freiburg, Freiburg, Germany*
**Lennart Martens**
*VIB Department of Medical Protein Research, UGent, Belgium*
**Lloyd Smith/Michael Shortreed**
*University of Wisconsin-Madison*

**ITCR groups**
**Rachel Karchin/Michael Ryan**
*Johns Hopkins University/In-Silico Solutions*

**Tom Doake/Jeremy Fischer**
*Indiana University*

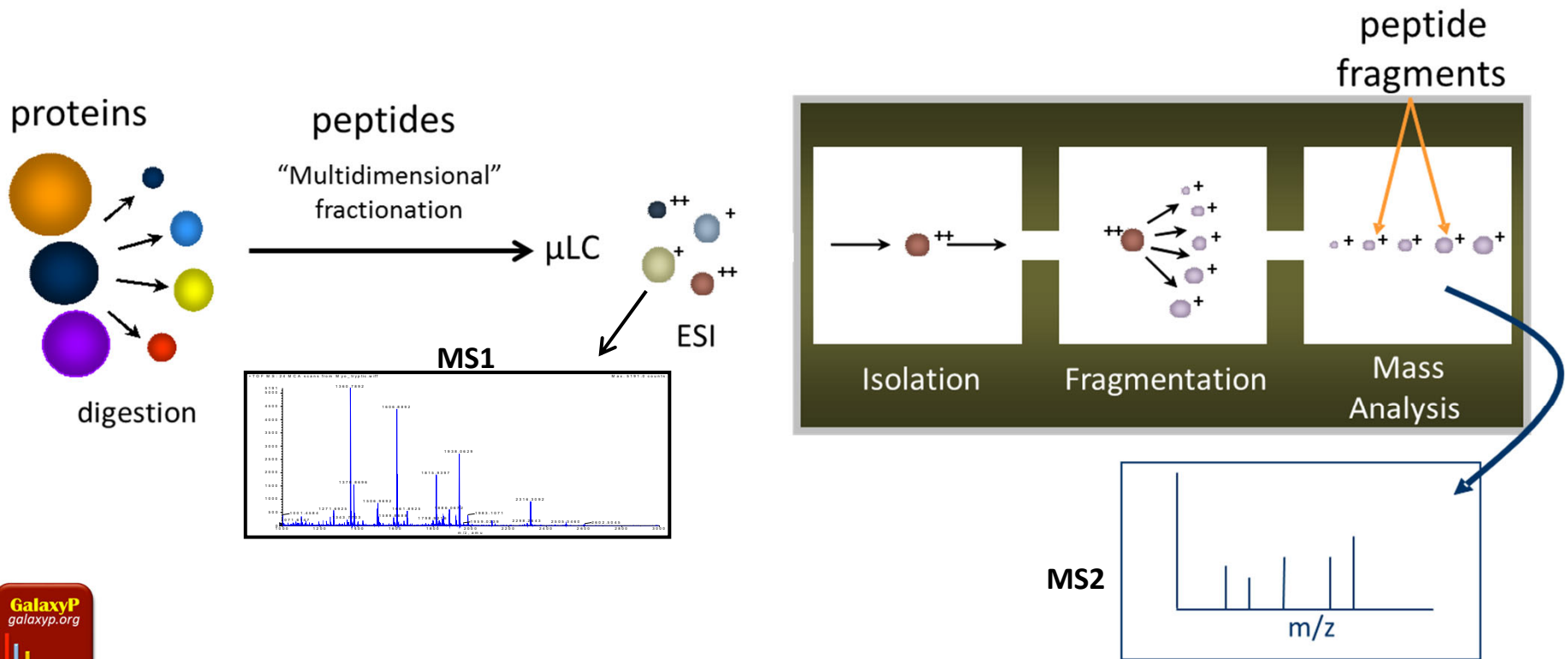GalaxyP
galaxyp.org

Jetstream

# Outline: Proteogenomic workflows

- **Background and informatics challenges**

- **Overview of existing software and workflows**

- **Access to the community**

GalaxyP
*galaxyp.org*
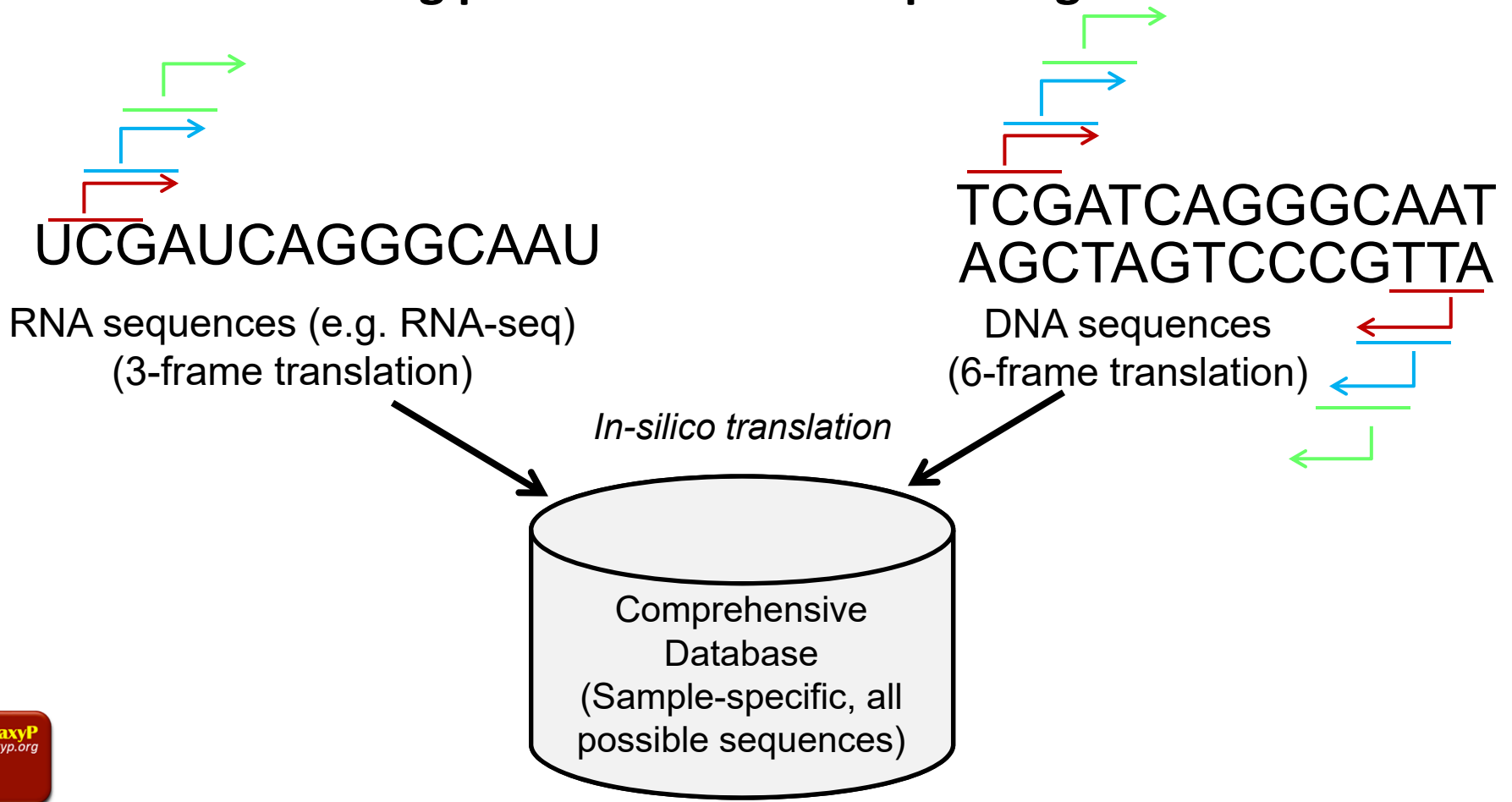
# Proteogenomics:  A primer



Peptide fractionation coupled to tandem mass spectrometry (MS/MS)

# Matching amino acid sequences to MS/MS data
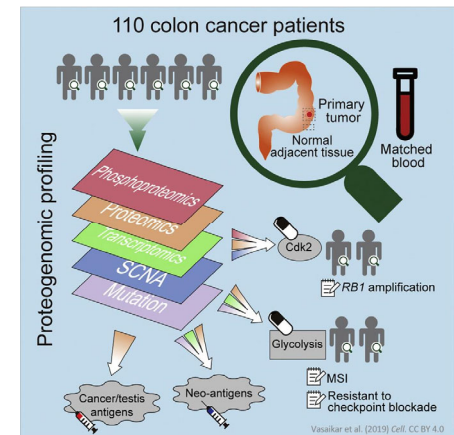
# Detecting protein variants via proteogenomics

UCGAUCAGGGCAAU

RNA sequences (e.g. RNA-seq)
(3-frame translation)

TCGATCAGGGCAAT
AGCTAGTCCCGTTA

DNA sequences
(6-frame translation)

*In-silico translation*

Comprehensive
Database
(Sample-specific, all
possible sequences)

GalaxyP
*galaxyp.org*

# Proteogenomic outcomes



- ✓ *Confirms translation of variants*
- ✓ *Direct evidence of potential functional variants*
- ✓ *Applications in neoantigen discovery (immuno-oncology)*

# Bringing proteogenomics to the masses: informatics challenges

- Many software tools, integration, automation….

# Proteogenomic informatics challenges

- *Assembly and variant calling from DNA/RNA sequencing data*

- *Customized protein sequence database generation*

- *Matching sequences to MS/MS data: best practices?*

- *Filtering, QC and verification of putative variant sequences*

- *Interpretation! Beyond a list....*

- *Access and usability by the research community*

# One workflow solution: Galaxy

- ✓ A web-based, community developed bioinformatics workbench for integrating disparate software -- flexible
- ✓ Geared towards use by bench scientists; many training resources available
- ✓ Already home to genomic/transcriptomic tools
- ✓ Provenance tracking, sharing and reproducibility
- ✓ Amenable to other 'omic tools (e.g. Galaxy for proteomics project, Galaxy-P)
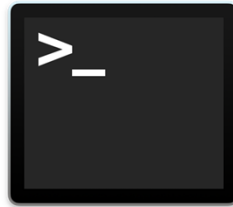
*Working philosophy:*

# Galaxy: an integrative workbench well-suited for multi-omics



*Courtesy Jeremy Goecks, OHSU*

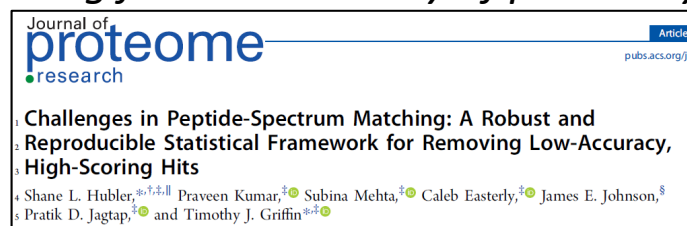# Integrative data processing: RNA-Seq + proteomics

# Best practices for peptide spectrum matches (PSMs) in proteogenomics

- *Utilize multiple database searching programs (e.g. SearchGUI)*

- *Multi-stage database search to obtain variant specific FDR estimates*



*J Proteome Res*. 2015 14:3555-67

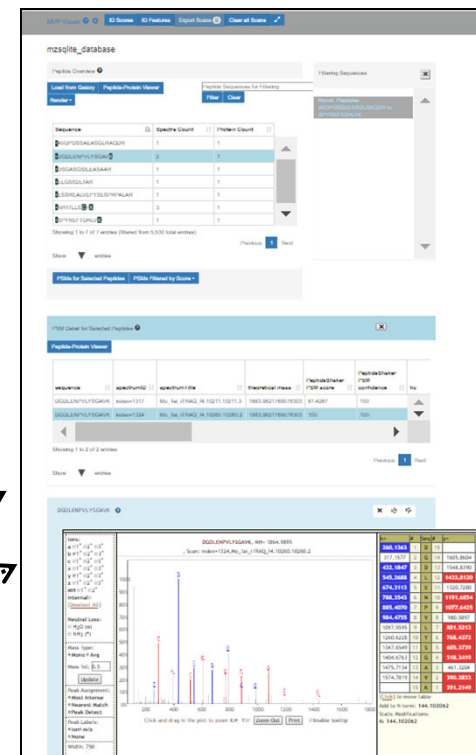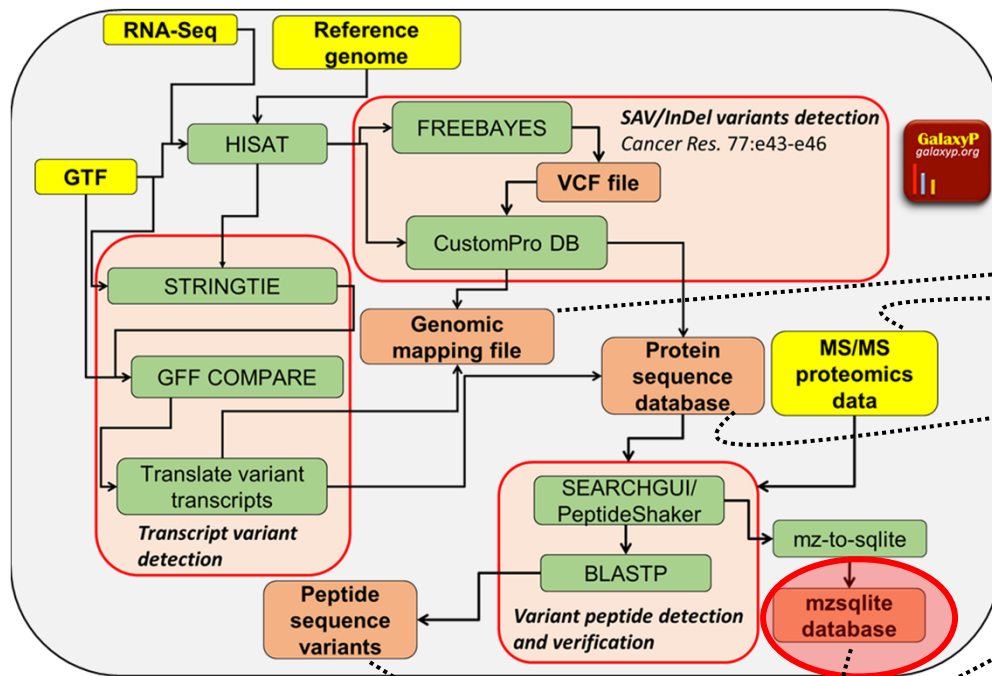- *Post-processing filtering for extra scrutiny of putatively confident PSMs*



Challenges in Peptide-Spectrum Matching: A Robust and Reproducible Statistical Framework for Removing Low-Accuracy, High-Scoring Hits

Shane L. Hubler,[*,†,‡,||] Praveen Kumar,[‡] Subina Mehta,[‡] Caleb Easterly,[‡] James E. Johnson,[§] Pratik D. Jagtap,[‡] and Timothy J. Griffin[*,‡]

# What's next?  Beyond a big list....

# Multi-Omics Visualization Platform:
# Characterizing the nature of detected variants

- *HTML-based Galaxy plugin*
- *Interactive reading of mzsqlite dB*

# Multi-Omics Visualization Platform:  Characterizing the nature of detected variants

# Other visualization and interpretation tools for protegenomics

- *Leveraging other tools and knowledgebases to assses impact of protein sequence variants (depends on good API development)*
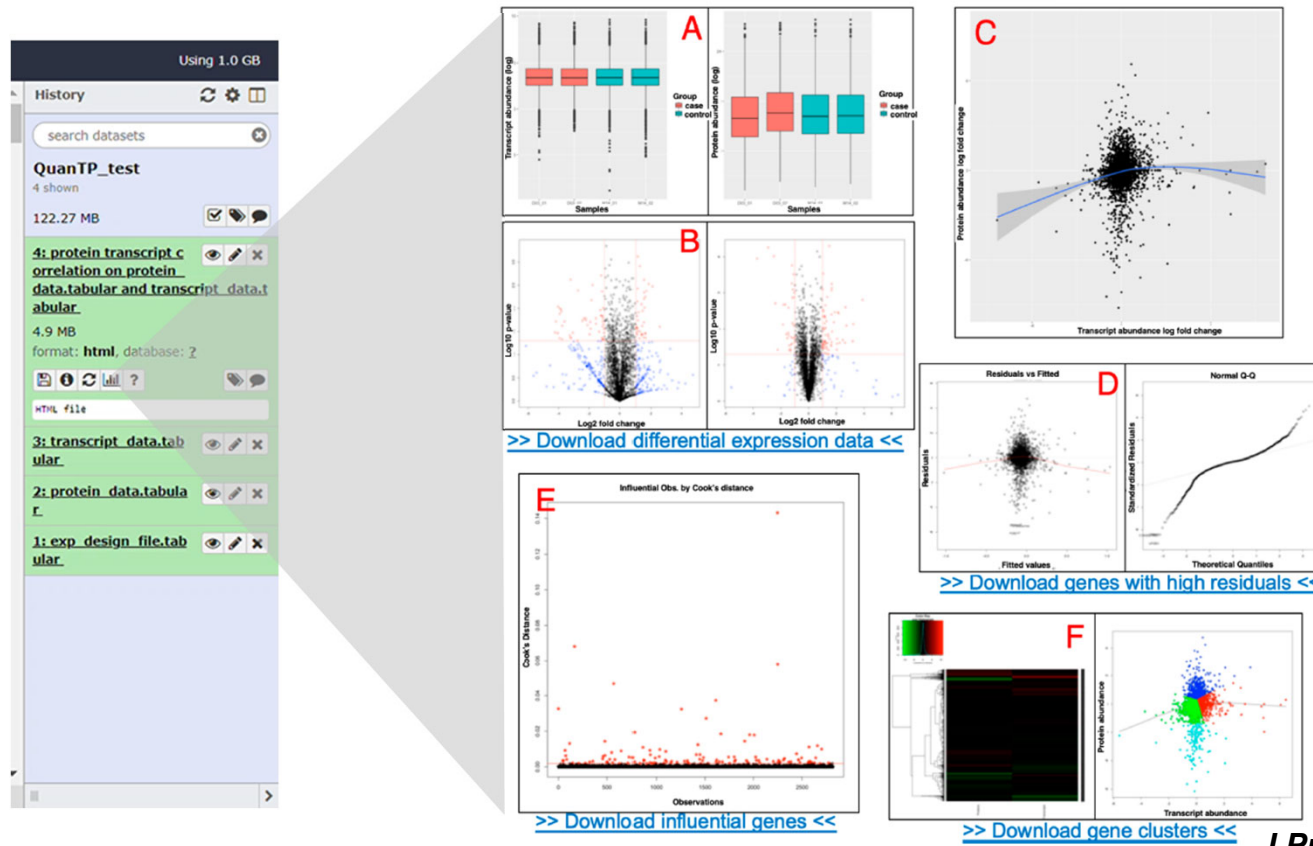


"CRAVAT-P"
***J Proteome Res***. 2018 ,17:4329-4336

# Other visualization and interpretation tools for protegenomics

- *Quantitative proteo-transcriptomics: comparing RNA and protein response*

# Providing access to the research community

**Tools and Workflows available on** : **https://proteomics.usegalaxy.eu/**



**Proteogenomics Gateway**: **z.umn.edu/proteogenomicsgateway**



**Training!**

**https://training.galaxyproject.org/training-material/topics/proteomics/**