# Advanced MS analysis: proteogenomics

**Tim Griffin**
**University of Minnesota**

**International Mass Spectrometry Conference**
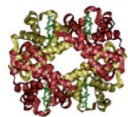**28 August, 2022**
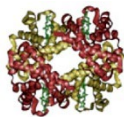
*tgriffin@umn.edu*

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

GalaxyP
*galaxyp.org*

# 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*

**Jetstream**

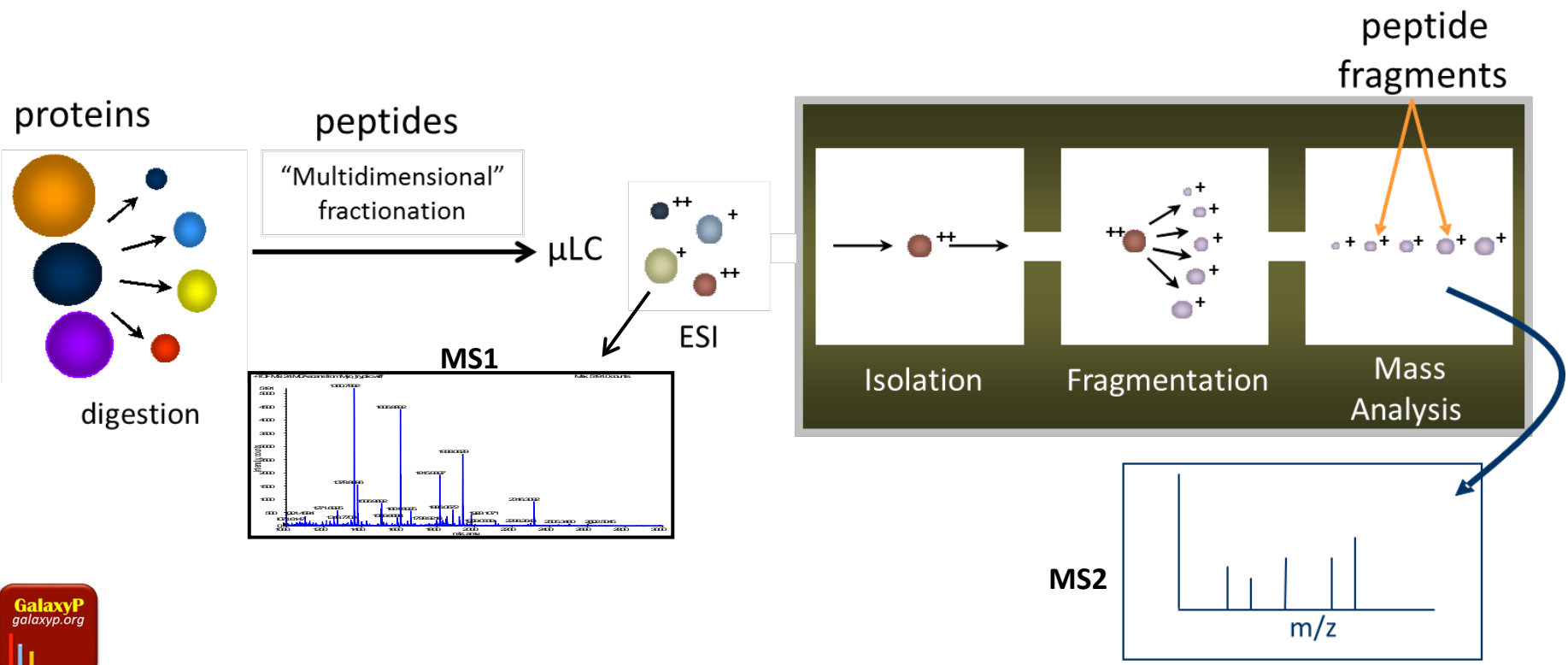**GalaxyP**
*galaxyp.org*

NSF

# Outline: Proteogenomics and bioinformatics

- **Background and informatics challenges**

- **Overview of components involved in proteogenomic workflows**
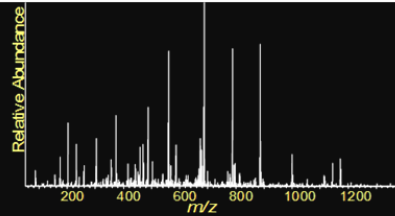
- **Hands-on workshop and tutorial**

# Proteogenomics: A primer

## Peptide fractionation coupled to tandem mass spectrometry (MS/MS)
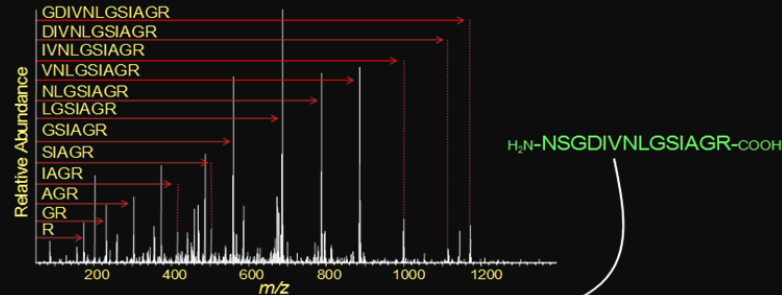
# Matching amino acid sequences to MS/MS data

# 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 and applications



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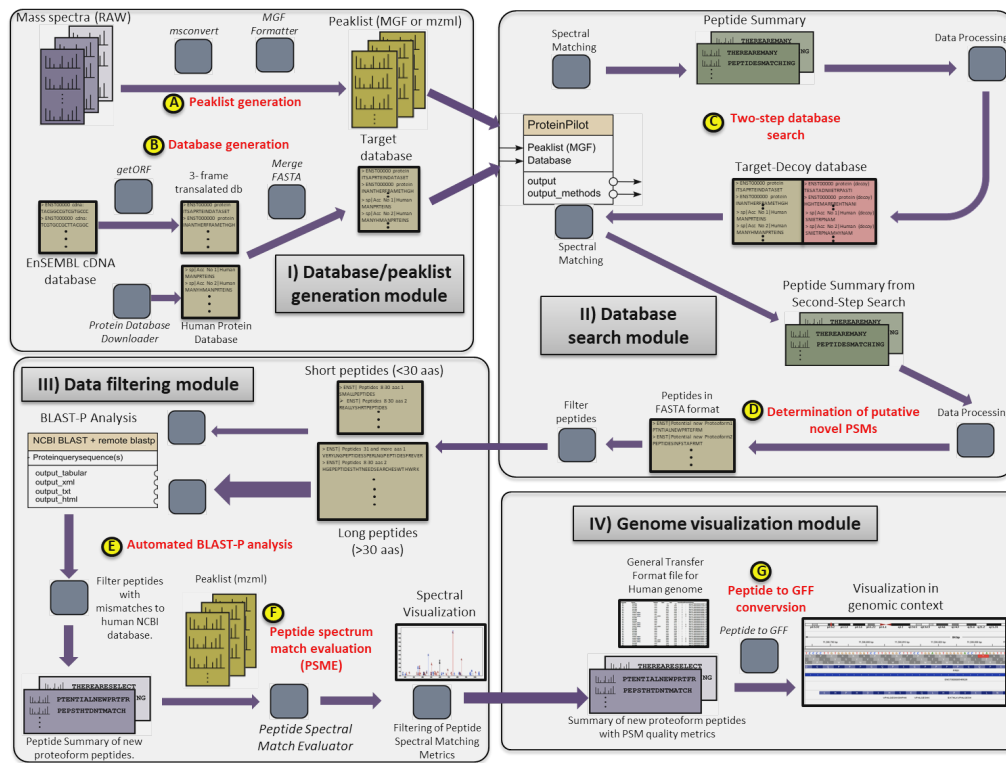- ✓ ***Applications in neoantigen discovery (immuno-oncology)***

# Bringing proteogenomics to the masses: informatics challenges

- Many software tools, integration, automation….



*J. Proteome Res.,* 2014, 13, pp 5898–5908

# 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*

# A solution: The Galaxy Ecosystem

✓ 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)
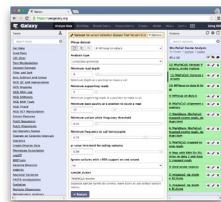
*Working philosophy:*

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



**Interfaces**

Web UI

Programmatic API

*Integrate datasets, analysis tools, visualizations, and computing resources for large-scale biomedical data science*

**Datasets**

**Computing Resources**

**Analysis Tools and Visualizations**

MS/MS sequence matching

Protein sequence database generation

ID of novel sequence variants

*Courtesy Jeremy Goecks, OHSU*

GalaxyP
*galaxyp.org*

# Integrative data processing:  RNA-Seq + proteomics

1. *Generate protein sequence database from RNA-Seq data*

2. *Match empirical MS/MS data to protein sequences*

# What's next?  Beyond a big list….

# Characterizing the nature of detected variants

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

# Training tutorials for main components of proteogenomics

**Proteogenomics 1: Database Creation**

**Proteogenomics 2: Database search**

**Proteogenomics 3: Novel peptide analysis**



*Please ask questions as they come up!*

# Training tutorials for main components of proteogenomics



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

# Accessing Galaxy via a public gateway

**①** **Login/Register: usegalaxy.eu**

**②** **Go to TIaaS link:** https://usegalaxy.eu/join-training/imsc_galaxy_training

**③** **Return to usegalaxy.eu site**

# Database generation from RNA-Seq data: SAV/Indels + unexpected novel translation events

# Working with Galaxy

1. Start with a new and blank History in usegalaxy.edu. Give it a name of your choosing. Go to the link below, which will open a Data Library in usegalaxy.eu. Select each of the three datasets <u>ending in the text shown below (a-c)</u> and hit "Export to History". Export each "As dataset" and select your newly named History as the destination.

https://z.umn.edu/imscdata

a) FASTQ_ProB_22LIST.fastqsanger
b) Mus_musculus.GRCm38.86.gtf
c) Trimmed_ref_5000_uniprot_cRAP.fasta (<u>On page 2 of the data library</u>)

2. In your active History, hit the "pen" icon and rename each History item as shown for a-c above, taking out the URL information to simplify the names

3. Listen to lecture on the details of the upcoming analysis
        (https://youtu.be/b_kZf8mXHdo)

# Running a workflow: Database creation

1. Go to Shared Data → Workflows and search for "IMSC Proteogenomics 1"

2. On the dropdown menu select "Run"

3. In the dialogue box select the correct three input files from your History, matching names to the input item in your History

4. Hit "Run Workflow"

# Part 2: Sequence database searching

After listening to the introductory lecture information, import the completed History for this module of the workshop:
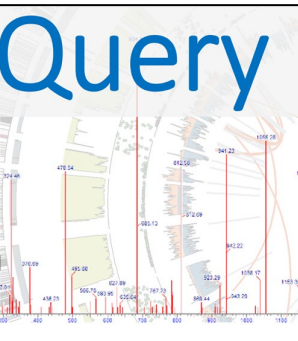
1.  Go to Shared Data → Histories and search for "IMSC2022"

2.  Select "IMSC2022 Proteogenomics 2: Database Search" and hit the "+Import" button in the upper right corner of the page

3.  Listen to the remaining lecture material and then we will explore some of the inputs, tools and outputs.
    (https://youtu.be/q1OjmTcbvBA)

# Ensuring confidence in MS/MS matches to novel peptide sequences: important!



# PepQuery

a peptide-centric search engine for novel peptide identification and validation

*Genome research* 29.3 (2019): 485-493.

# Part 3: Novel peptide analysis

After listening to the introductory lecture information, import the completed History for this module of the workshop:

1.  Go to Shared Data → Histories and search for "IMSC2022"

2.  Select "IMSC2022 Proteogenomics 3: Novel Peptide Analysis" and hit the "+Import" button in the upper right corner of the page

3.  Listen to the remaining lecture material and guidance on exploring results (https://youtu.be/Ku274KwFh1Y)
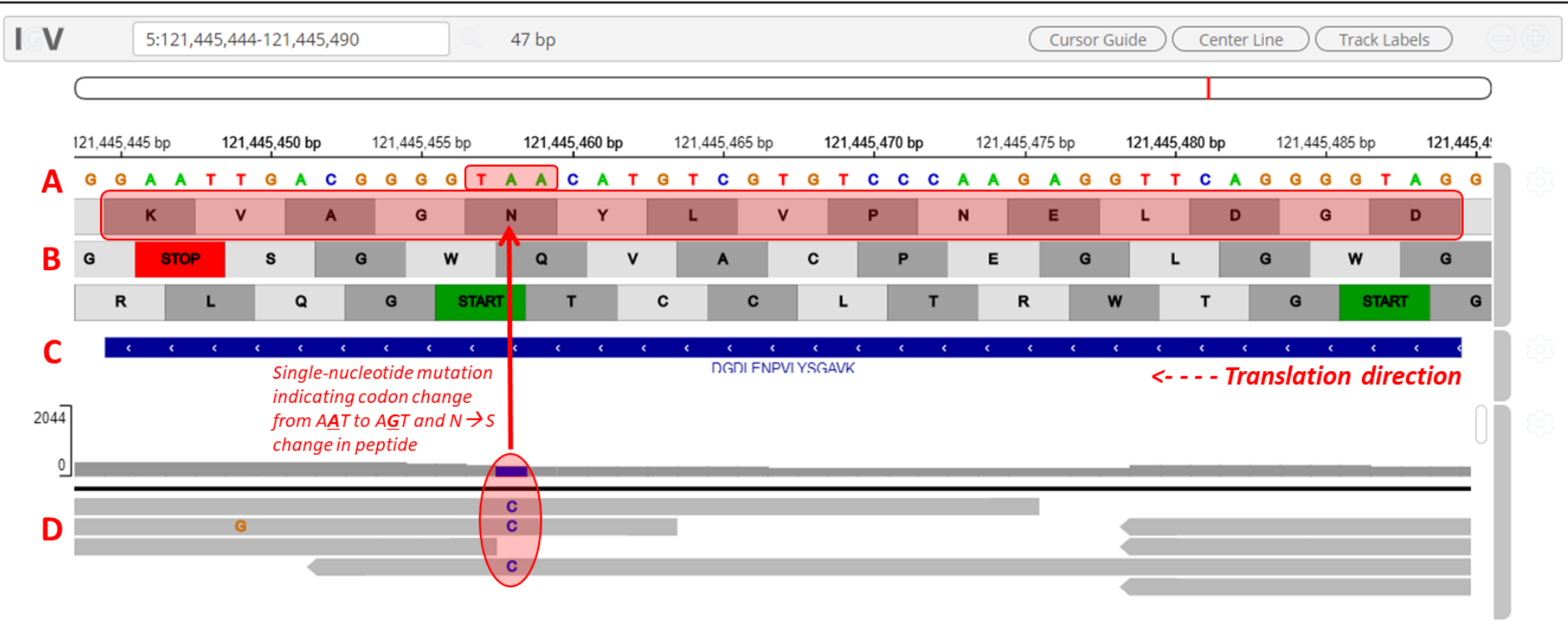
# Viewing and characterizing novel peptide sequences



*Gigascience* (2020), **9**:giaa025

# IGV viewer for visualizing peptide, RNA and DNA sequences

# Splice isoform example