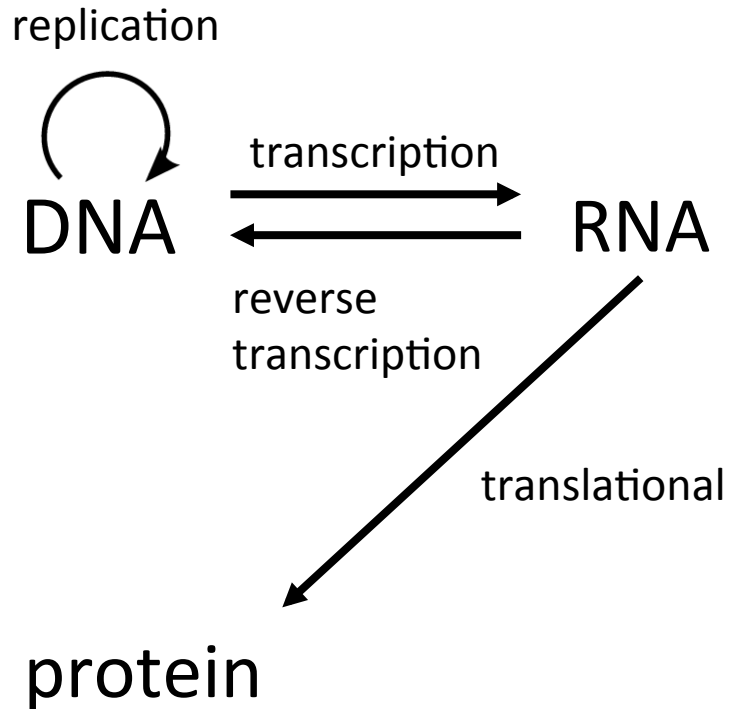


Building a low-budget public resource for large-scale proteomic analysis

Anoop Mayampurath
Computation Institute
April 17, 2014



The world of proteomics



Proteins carry out **most** cellular activity, including **control** (regulation) of transcription, translation, and replication of DNA.

Alzheimer's Disease

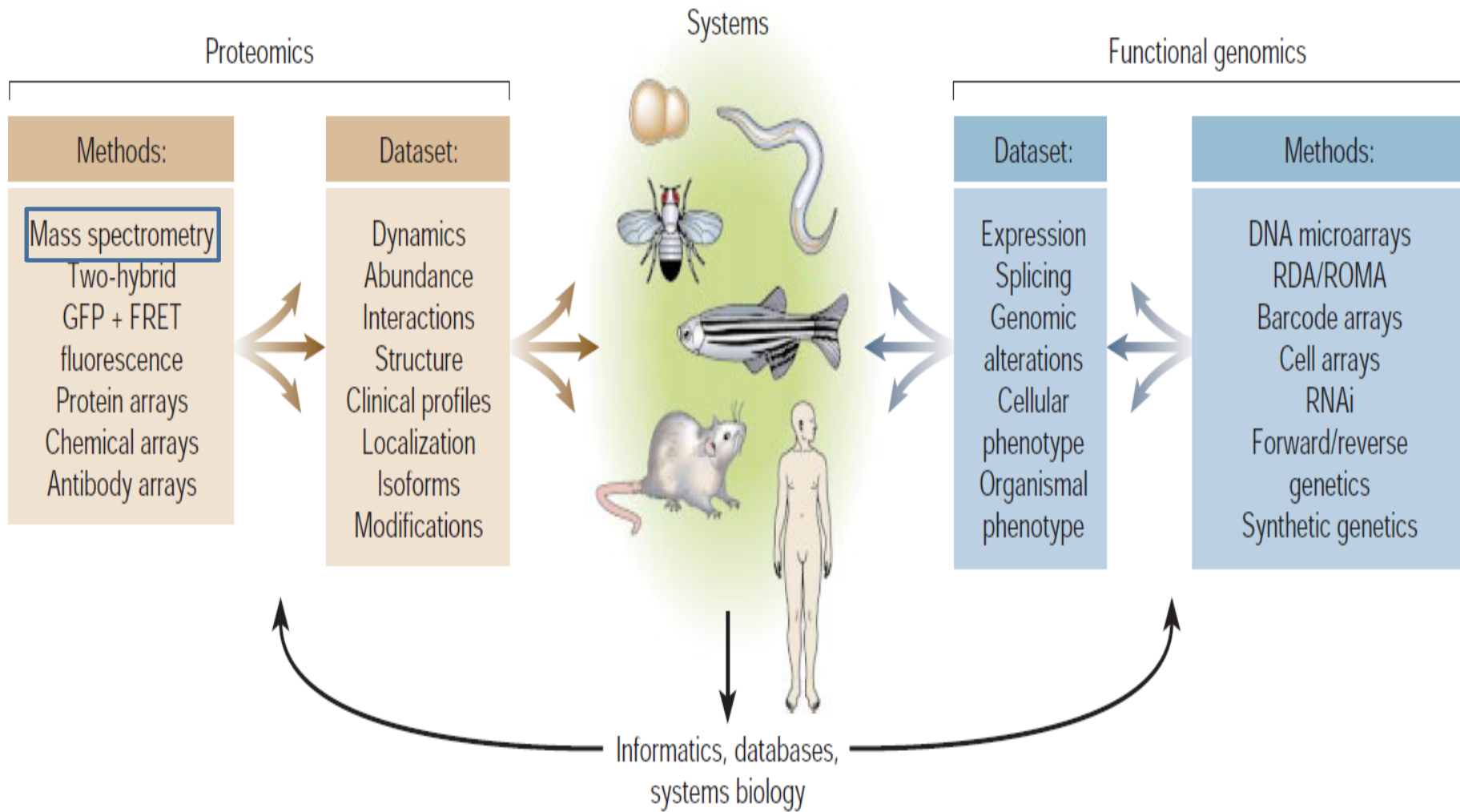
- Most common late-onset Alzheimer's gene : APOE (e2,e3, and e4)
- Plaque accumulation but do not develop Alzheimer's
- Repressor element-1 silencing transcription factor

ARTICLE

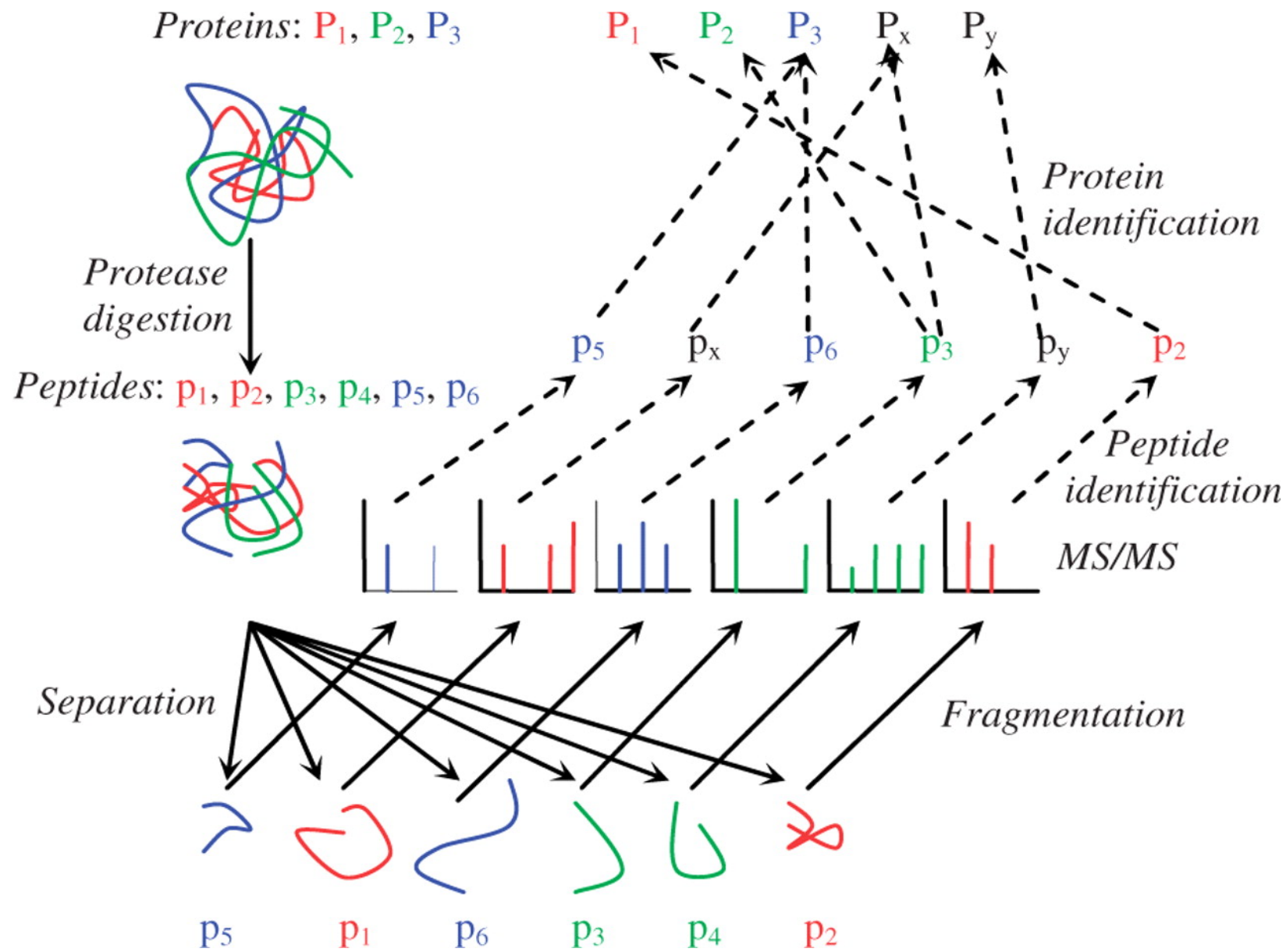
doi:10.1038/nature13163

REST and stress resistance in ageing and Alzheimer's disease

Tao Lu¹, Liviu Aron¹, Joseph Zullo¹, Ying Pan¹, Haeyoung Kim¹, Yiwen Chen², Tun-Hsiang Yang¹, Hyun-Min Kim¹, Derek Drake¹, X. Shirley Liu², David A. Bennett³, Monica P. Colaiácovo¹ & Bruce A. Yankner¹



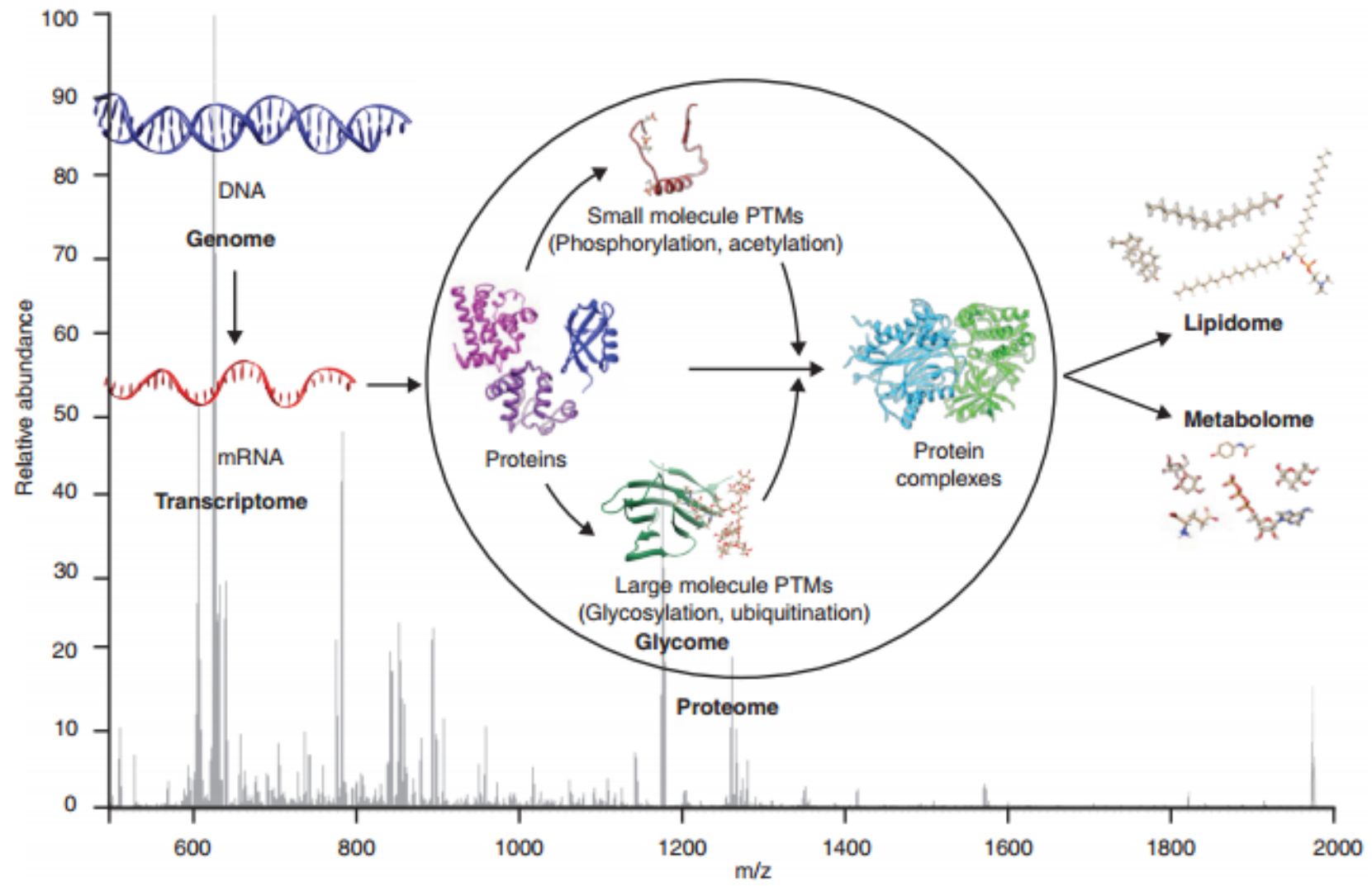
A little
background
on
mass spectrometry-proteomics



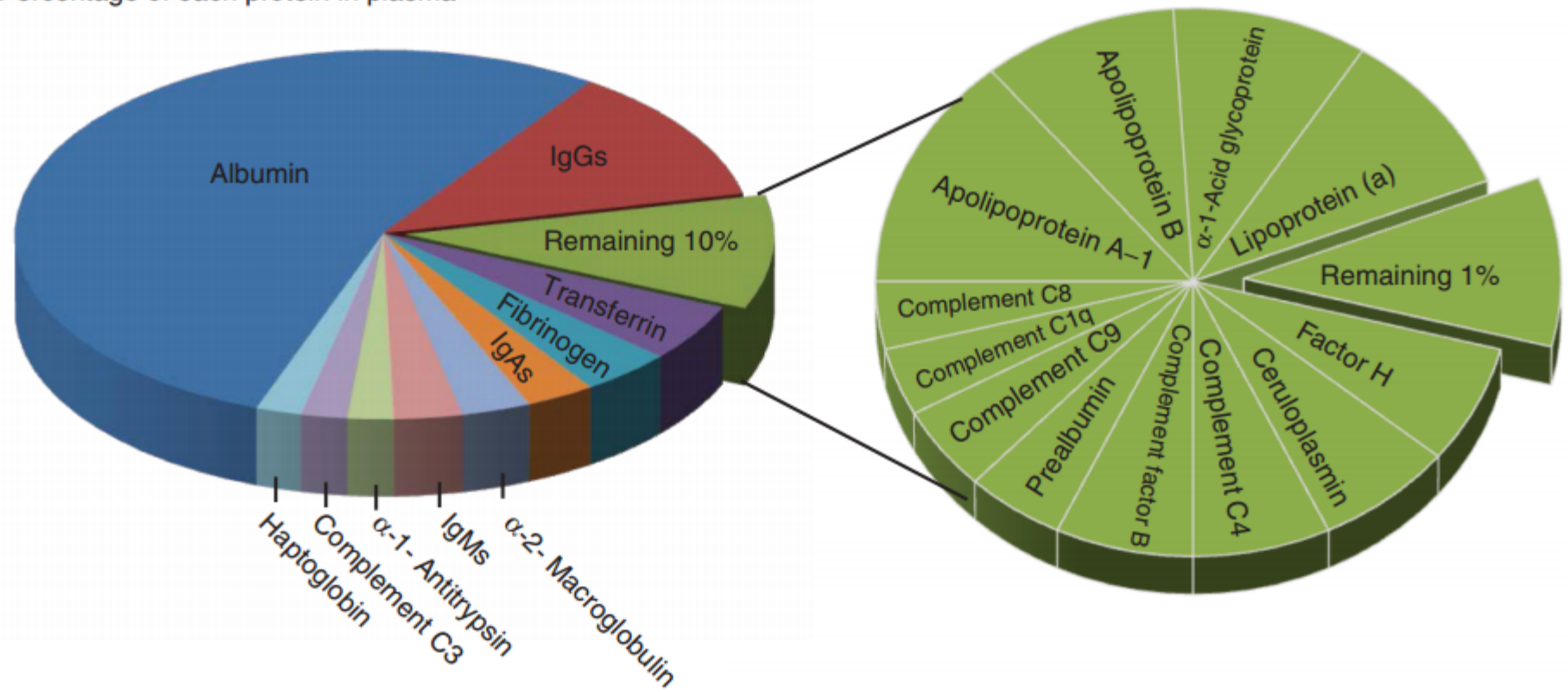
LC-MS/MS clinical proteomic challenges

Biology

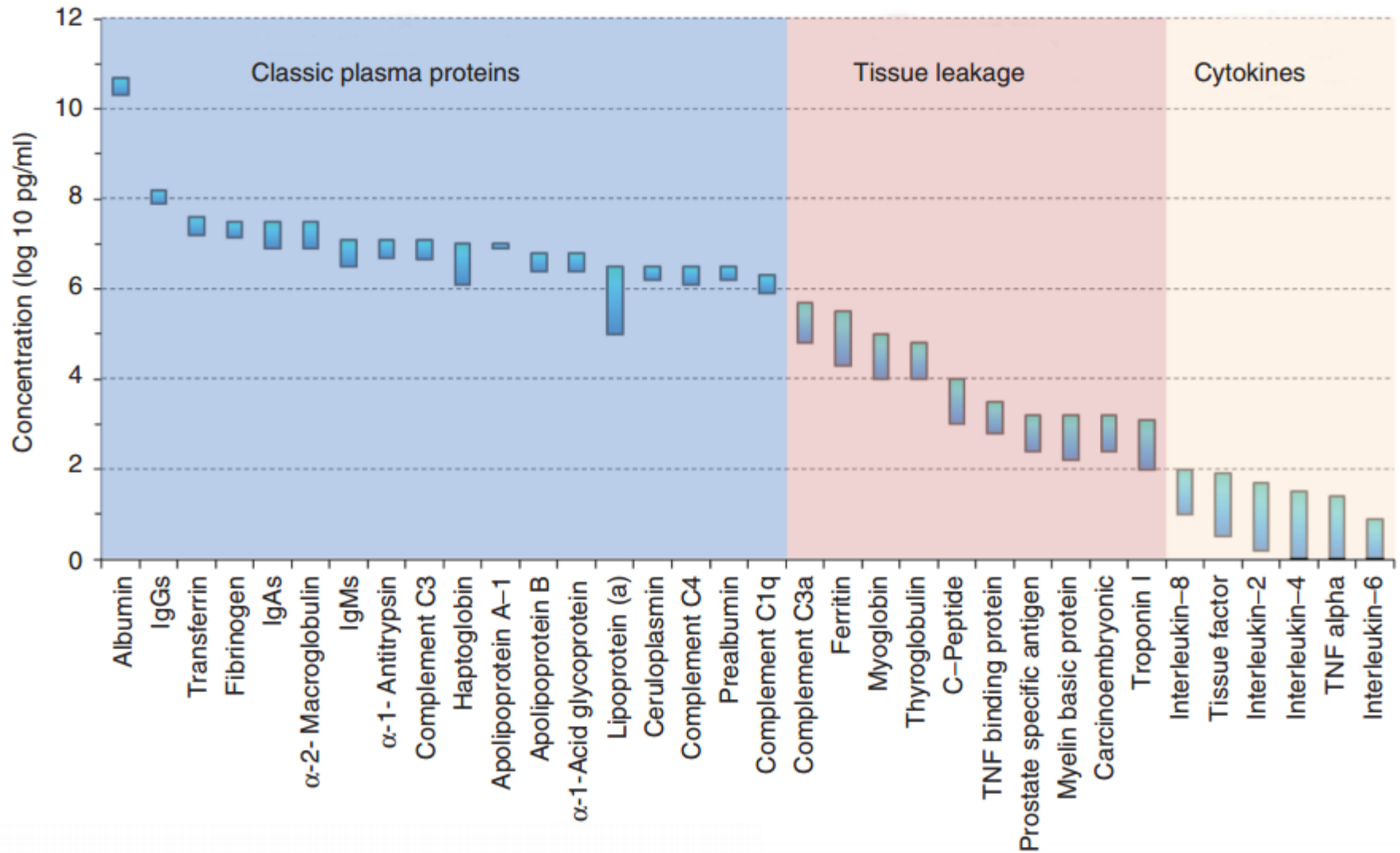
- Complexity of the proteome
- Dynamic range of plasma proteins



(b) Percentage of each protein in plasma



(a) Dynamic range of proteins in plasma

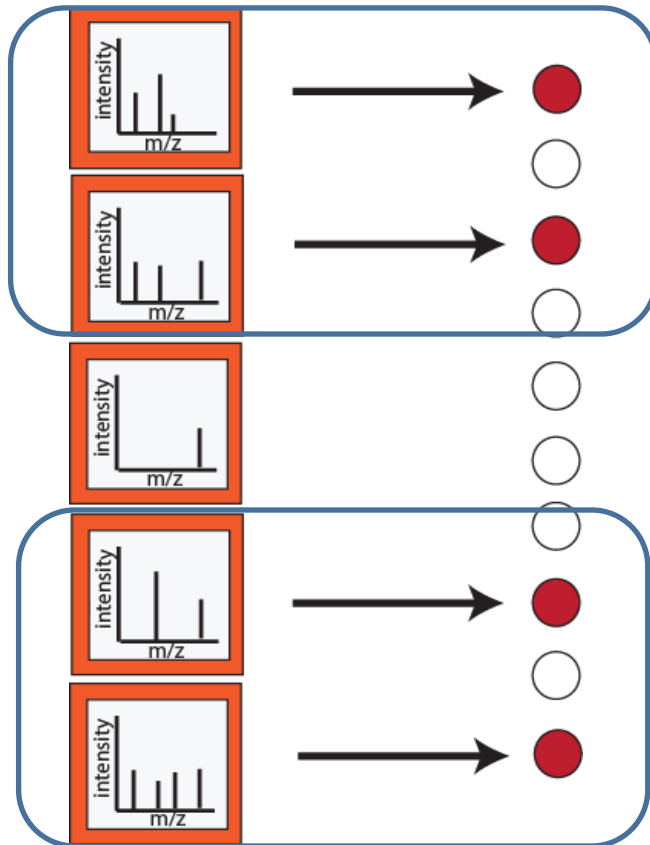
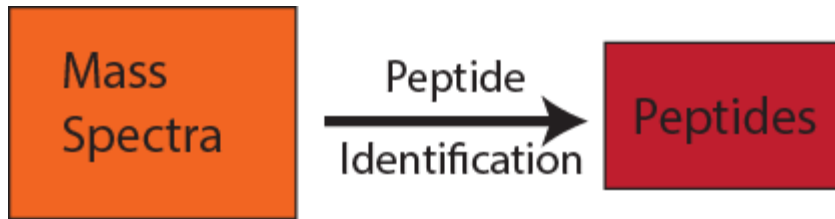


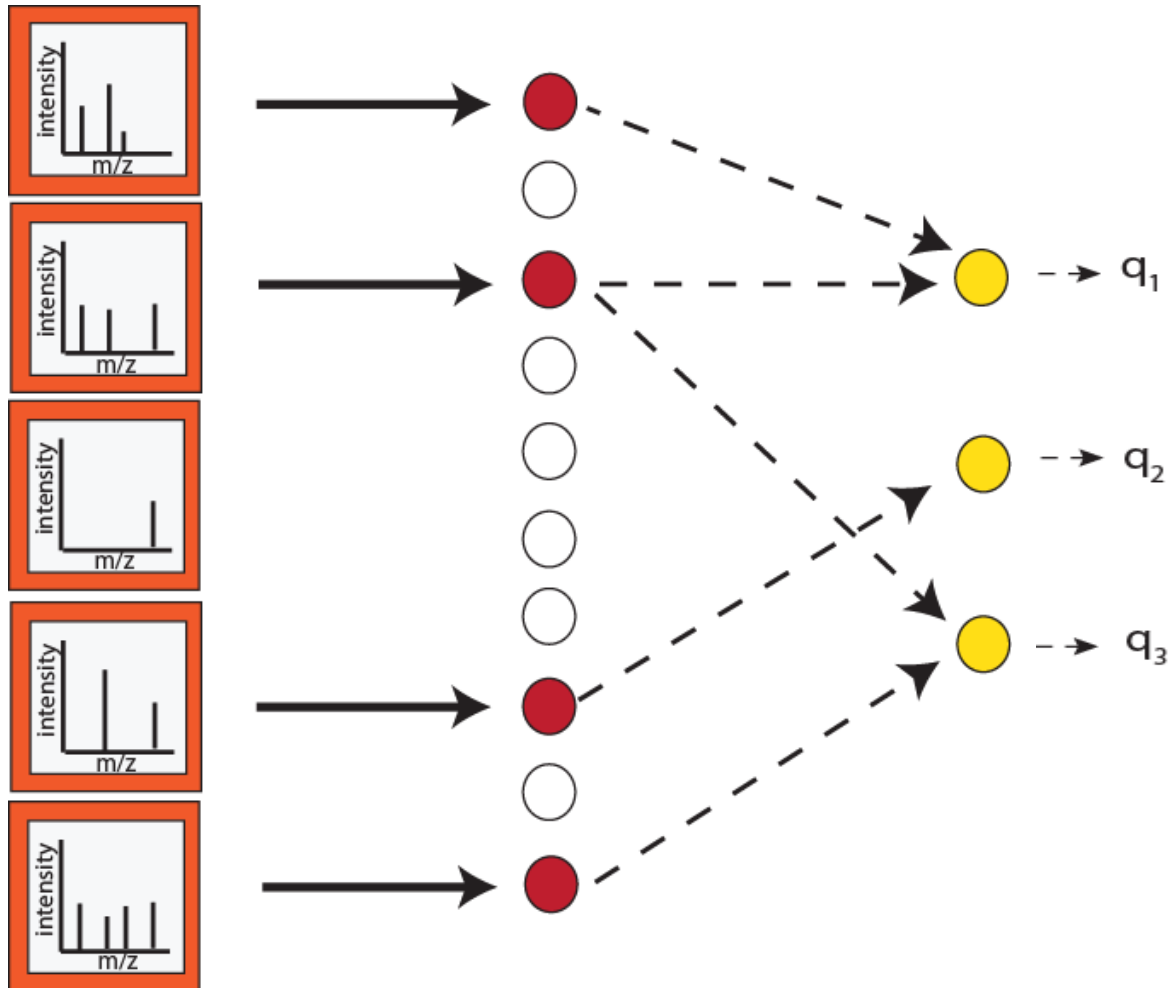
LC-MS/MS clinical proteomic challenges

Computational

- lack of a framework capable of performing large-scale proteomic analysis

The workflow







ms-utils.org

- [About](#)
- [Software List](#)
- [Editing Policies](#)
- [FAQ](#)
- [ChangeLog](#)

related sites

- [ExPASy tools](#)
- [NBIC BioAssist](#)
- [PNL Tools](#)
- [SPC Proteomic Tools](#)

Search:

Go

[View](#) [Edit](#) [History](#) [Print](#)

Software List

platforms, pipelines and libraries

CPAS	LIMS and analysis tools for proteomics data (includes msInspect)		
CPFP	Central Proteomics Facilities Pipeline [1] (demo here)		
GenePattern	platform for integrative genomics and proteomics (includes PEPPER [2] and other tools for proteomics)	Java	
InSilicoSpectro	open source proteomics library (of Perl functions) [3]	Perl	
libfbi	a fast implementation of box intersection for correspondence estimation in peak picking, alignment, etc.	C++	
Mass-up	utility with full GUI for proteomics data analysis, particularly MALDI-TOF	Java	
MASSyPup	a lightweight Linux live distribution prepackaged with X!Tandem, mMass, MZmine, PepNovo/UniNovo, PeptideShaker, mscouvert, XCMS etc.		
mspire	MS data processing in Ruby, including mzML reader/writer, <i>in-silico</i> digestion, isotopic pattern calculation etc. [4]	Ruby	
OpenMS	library for the analysis, reduction and visualization of LC-MS(/MS) data	C++	
PAPPSO	Plateforme d'Analyse Protéomique de Paris Sud-Ouest	Java	
Proteios	pipeline/LIMS for proteomics experiments and analysis	Java	
Proteomatic	platform for creating MS/MS data analysis workflows using scripts [5]	C++	
ProteoWizard	open source library for proteomics tools development (supports mzML) [6]	C++	
pymzML	Python module to parse mzML data based on cElementTree [7]	Python	
Pyteomics	framework for proteomics data analysis, supporting mzML, MGF, pepXML and more [8]	Python	
QuPE	integrated environment for storage, analysis and integration of proteomics data (requires login) [9]	Java	web
Rproteomics	set of routines for analyzing proteomics data, an XML database to store the results and a user interface	R	
TOPP	the OpenMS protein identification/quantitation pipeline	C++	
TPP	Institute for Systems Biology "Trans-Proteomic Pipeline"		
XCMS	software package (in R) for metabolite profiling from LC-MS data	R	



About

Team

Workflow integration

Downloads

Documentation

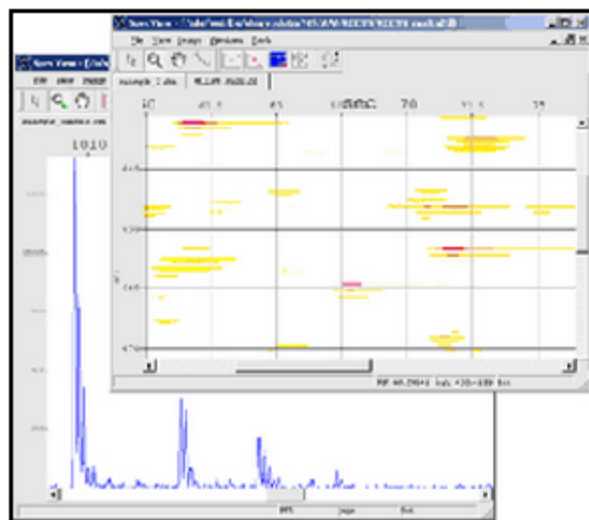
Publications

Support

About

OpenMS is an open-source software C++ library for **LC/MS data management and analyses**. It offers an infrastructure for the rapid development of mass spectrometry related software. OpenMS is free software available under the three clause BSD license and runs under Windows, MacOSX and Linux.

It comes with a vast variety of pre-built and ready-to-use tools for proteomics and metabolomics data analysis (TOPPTools) and powerful 2D and 3D visualization(TOPPView).



OpenMS offers analyses for various quantitation protocols, including **label-free quantitation, SILAC, iTRAQ, SRM, SWATH,**

It provides built-in algorithms for **de-novo identification and database search**, as well as adapters to other state-of-the-art tools like **XTandem, Mascot, OMSSA, etc.**

It supports easy integration of OpenMS built tools into workflow engines like **Knime, Galaxy, WS-Pgrade, and TOPPAS** via the TOPPtools concept and a unified parameter handling (CTD).

News

- ▶ [OpenMS 1.11.1 released](#)
- ▶ [7th OpenMS User Meeting – High-performance software for high-throughput proteomics and metabolomics](#)
- ▶ [OpenMS 1.11 released](#)
- ▶ [OpenMS 1.10 released](#)
- ▶ [6th OpenMS User Meeting – High-performance software for high-throughput proteomics and metabolomics](#)

Get Galaxy-P

The [Minnesota Supercomputing Institute](#) provides a [public Galaxy-P server](#) capable of limited analyses, testing, and demo. For heavy use you will likely need to install your own instance of Galaxy-P or spin up a Galaxy-P cluster on the cloud.

Public Server

Advantages Immediately accessible. Easiest option for publically sharing data and pages.

Limitations Limited computational and disk resources. Potential problems associated with uploading protected or sensitive data to any public resource.

usegalaxyp.org

Install Your Own

Advantages Full control of computational resources. Easy to modify existing tools or add your own. Use our open source Galaxy-P tools targeting commercial applications that are not available on the public server. Right now these include [ProteinPilot](#) and [Scaffold](#).

Limitations Because of its flexibility, Galaxy can be time-consuming to install and maintain. Galaxy-P adds more tools and can be configured to utilize remote Windows resources adding additional complexity.

[Install Instructions](#)

Take to the Cloud

Advantages BioCloudCentral provides interface for creating a Galaxy-P cluster. CloudMan a likewise easy-to-use interface into the disk and computational resources of the cloud.

Limitations Currently no access to Windows such as MaxQuant or vendor software.

[Launch Now](#)

What is Galaxy-P?

Galaxy-P is a multiple 'omics' data analysis platform with particular emphasis on mass

News

For the latest Galaxy-P news, please follow us on Twitter.

proteomics.globusgenomics.org

The screenshot displays the Galaxy proteomics workflow editor interface. The browser address bar shows the URL: `proteomics.globusgenomics.org/workflow/editor?id=5715598cb1ceae2c#`. The top navigation bar includes links for `Analyze Data`, `Workflow`, `Shared Data`, `Visualization`, `Admin`, `Help`, and `User`. The main workspace is titled `Workflow Canvas | MM_Test_4` and features a grid background. On the left, a sidebar lists tool categories: **DATA TRANSFER** (Globus Data Transfer, Get Data, Send Data), **PROTEOMICS APPLICATIONS** (Conversion Tools, Search Tools, TPP Tools, Processing Tools, Volchemboum tools), and **NGS APPLICATIONS** (QC and manipulation, Mapping, RNA Analysis, Peak Calling, SAM Tools, BAM Tools, Picard, Indel Analysis, GATK Tools, Variant Detection, Interval Tools, VCF Tools, CGA Tools, Simulation, SNP/WGA: Data; Filters, SNP/WGA: QC; LD; Plots, SNP/WGA: Statistical Models, Phenotype Association). The workflow canvas contains five tools connected in a linear sequence from left to right: 1. `Get Data via Globus Online` (input: `out_file1 (txt)`), 2. `X!Tandem MSMS Search` (input: `MSMS File`, output: `output (raw_pepxml)`), 3. `Peptide Prophet` (input: `Raw Search Results`, output: `output (peptideprophet_pepxml)`), 4. `Protein Prophet` (input: `Peptide Prophet Results`, output: `output (protxml)`), and 5. `Send Data via Globus Online` (input: `Send this dataset`, output: `out_file1 (txt)`).

An example : multiple myeloma

Metric	One fraction (~800 MB)	All 22 fractions (~ 16GB)
Wall Time	19 minutes	24 minutes
Total CPU Time	19 minutes	418 minutes
AWS Nodes Used (m2.4xlarge)	1	3
On-demand Cost (\$1.64)	\$1.64	\$4.92
Spot instance Cost (\$0.14 per hour)	\$0.14	\$0.42*

Current pricing for off-campus LC-MS/MS is as follows:

University of Illinois at Chicago (UIC): Orbitrap-Velos (\$100/sample run), (\$75/hourly)

Mayo Clinic Proteomics Resource Center (MPRC): Orbitrap-LTQ (\$100/sample run)

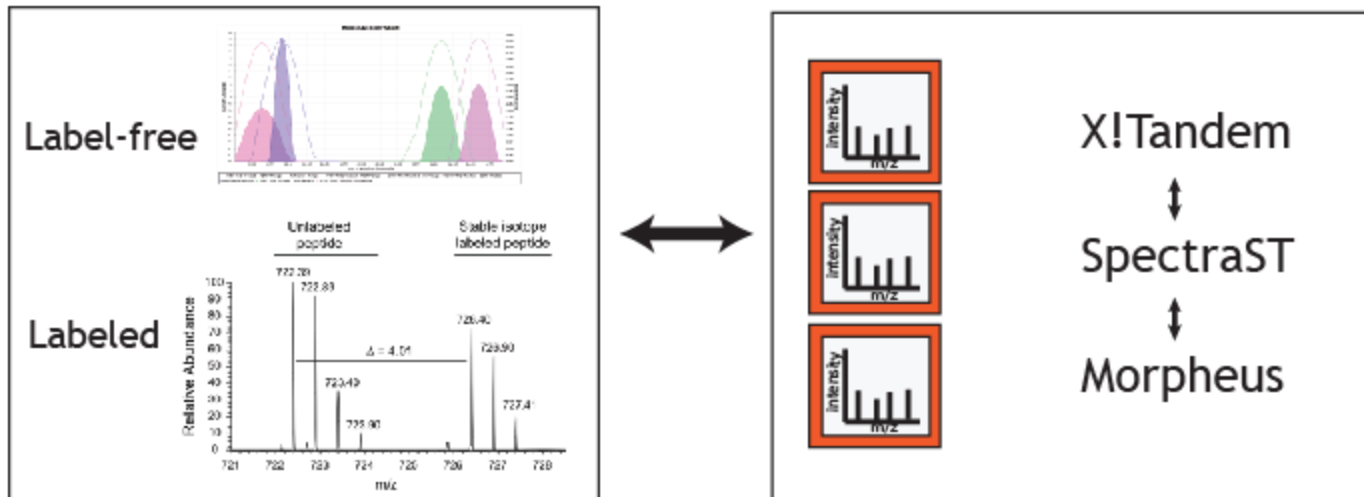
Northwestern University (NU) Proteomics Core: Orbitrap-Velos (\$140/sample run)

Northwestern University (NU) Proteomics Center of Excellence (PCE): Velos Orbitrap Elite with ETD, Velos-FT12T Ultra.


Future work

Functional category	Tools currently available	Tools to be incorporated in future
Dataset Tools	uf-mzML, , MGF-DTA file converters	msconvert Decon2LS (console version), DeconMSn, MultiAlign(command line), SuperHirn
Search Tools	X! Tandem, OMSSA, Mascot, MSGF+	SpectraST, Morpheus (command line), MyriMatch
Identification Tools	PeptideProphet, ProteinProphet , iProphet, Validator-MAX	MaxQuant, IDPicker3
Quantification Tools	Quantifier, XPRESS, Libra, ASAPRatio	MaxQuant, IDPQuantify
Other	PepXML/ProtXML to Table	Tools for: processing proteomics catalogs, combining outputs from different workflows, intelligent inclusion predictors.

Workflow integration



Globus Proteomics catalog

 globus online Manage Data | Groups | Support | kyle

[manage datasets](#) | [start transfer](#) | [view transfer activity](#) | [manage endpoints](#) | [dashboard](#)

Catalog
Proteomics

[Create Catalog](#)

Filter by Annotation

- Response
 - Response not present
 - Response present
 - CR
 - nCR
 - PR
 - VGPR
 - <VGPR
 - >=VGPR
- Instrument
- Experiment
- Treatment
 - Treatment not present
 - Treatment present
 - RVD
 - VDD
- owner

2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	6301DA_VGPR Owner: u:kyle label:	<input type="button" value="v"/>												
2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	VDD2_nVGPR Owner: u:kyle label:	<input type="button" value="v"/>												
2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	VDD15_nVGPR Owner: u:kyle label:	<input type="button" value="v"/>												
<div>Overview Tags Sharing Select Files Members X</div> <div>Edit Tags Add Tags</div> <table><tbody><tr><td>Experiment</td><td><input type="text" value="Label-Free"/></td><td><input type="button" value="edit"/></td></tr><tr><td>Instrument</td><td><input type="text" value="LTQ"/></td><td><input type="button" value="edit"/></td></tr><tr><td>Response</td><td><input type="text" value="<VGPR"/></td><td><input type="button" value="edit"/></td></tr><tr><td>Treatment</td><td><input type="text" value="VDD"/></td><td><input type="button" value="edit"/></td></tr></tbody></table>						Experiment	<input type="text" value="Label-Free"/>	<input type="button" value="edit"/>	Instrument	<input type="text" value="LTQ"/>	<input type="button" value="edit"/>	Response	<input type="text" value="<VGPR"/>	<input type="button" value="edit"/>	Treatment	<input type="text" value="VDD"/>	<input type="button" value="edit"/>
Experiment	<input type="text" value="Label-Free"/>	<input type="button" value="edit"/>															
Instrument	<input type="text" value="LTQ"/>	<input type="button" value="edit"/>															
Response	<input type="text" value="<VGPR"/>	<input type="button" value="edit"/>															
Treatment	<input type="text" value="VDD"/>	<input type="button" value="edit"/>															
2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	0311HI_nVGPR Owner: u:kyle label:	<input type="button" value="v"/>												
2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	4731BM_VGPR Owner: u:kyle label:	<input type="button" value="v"/>												
2013-11-05	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	8599KA_nVGPR	<input type="button" value="v"/>												

Acknowledgements

- Sam Volchenboum
- Kolbrun Kristjansdottir
- Don Wolfgeher
- Alex Rodriguez
- Kyle Chard
- Ravi Madduri
- Paul Dave