

Proudly Operated by Baffelle Since 1965

# MultiAlign Tutorial 02 – Data Processing

**BRIAN LAMARCHE** 



#### **About this tutorial**

- This tutorial will describe the data that MultiAlign processes.
- You will learn:
  - Terminology
  - How data is linked in the resulting relational database



Proudly Operated by Baffelle Since 1965

#### **Data Sources**

THIS SECTION PROVIDES THE BASICS OF HOW DATA IS DERIVED DURING THE MULTIALIGN PROCESS

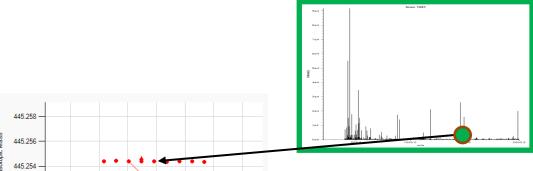
Pacific Northwest



#### **Terminology – MS Feature**

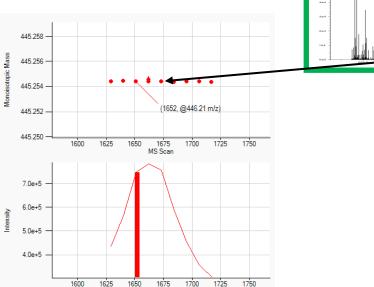
#### MS Feature

■ A feature (isotopic distribution) that was identified in a single MS¹ scan from a single dataset



Mono Mass vs. Scan

LC-Elution Profile (monoisotopic peak)





#### **Terminology – LC-MS Feature**

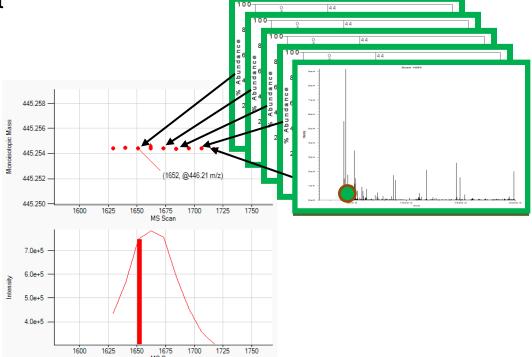
- MS Feature
  - A feature (isotopic distribution) that was identified in a single MS¹ scan from a single dataset
- LC-MS Feature

A collection of MS features that was grouped across multiple scans from a

single dataset



LC-Elution Profile (monoisotopic peak)



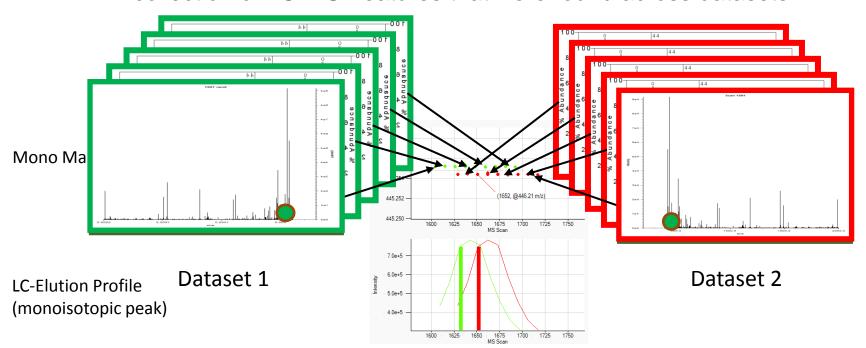


#### Terminology – LC-MS Feature Cluster



Proudly Operated by Battelle Since 1965

- MS Feature
  - A feature (isotopic distribution) that was identified in a single MS¹ scan
- LC-MS Feature
  - A collection of MS features that was grouped across multiple scans
- LC-MS Feature Cluster (consensus feature)
  - A collection of LC-MS Features that were found across datasets







Proudly Operated by Baffelle Since 1965

#### **Data Processing Workflow**

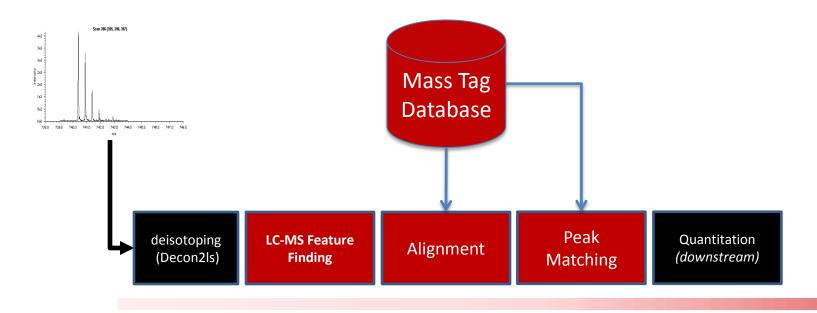
THIS SECTION DESCRIBES HOW MULTIALIGN ANALYZES DATA IN COMPARISON TO VIPER



#### AMT tag approach – single dataset



Proudly Operated by Battelle Since 1965



#### Traditional workflows would be to:

- 1. Deisotope an LC-MS instrument run into MS Features
- Detect LC-MS Features
- 3. Align these features to a reference, i.e. mass tag database
- 4. Identify Peptides, i.e. Peak Match
- 5. Perform Quantitation using downstream tools

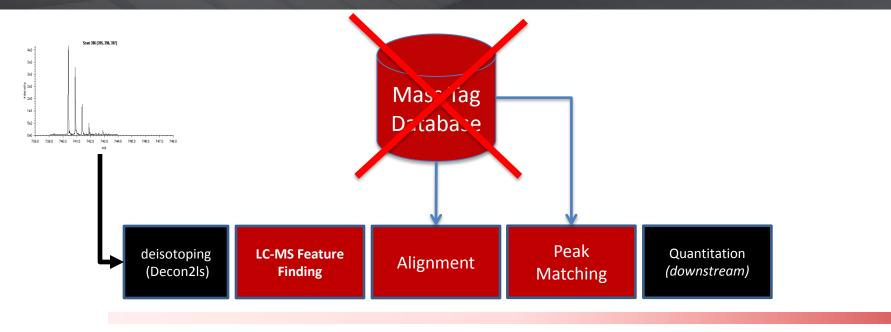
December 18, 2012 8



#### AMT tag approach variants



Proudly Operated by Baffelle Since 1965



But sometimes AMT tag databases are not available because genomic information is unknown, as with some microbial communities or if performing metabolomic based analysis.

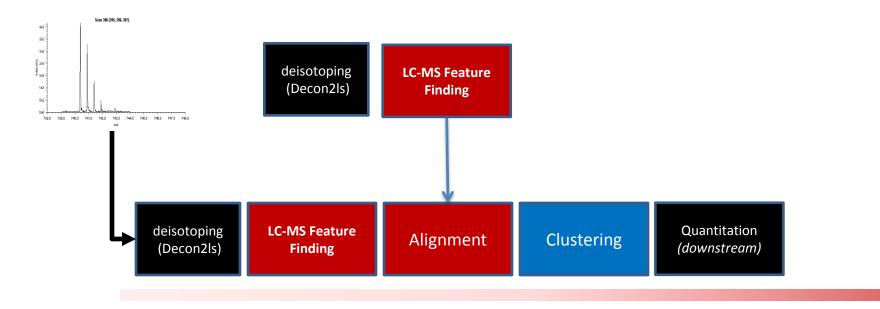
December 18, 2012



#### **AMT** tag approach variants



Proudly Operated by Battelle Since 1965



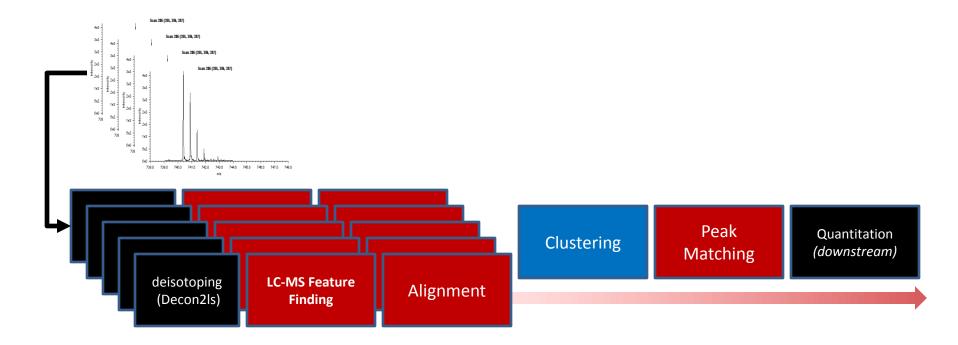
Instead one may be interested in purely comparative dataset analysis. In this case, a dataset can be used as the reference for alignment.

December 18, 2012



## MultiAlign – Tool for Multiple dataset analysis





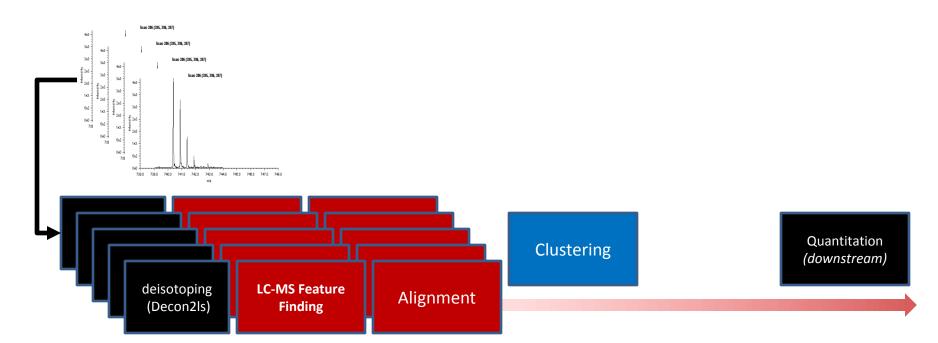
MultiAlign is a multiple dataset analysis tool (LC-MS, LC-IMS-MS) whose intention is to find common features across multiple datasets. It also can perform feature identification, e.g. for peptides) via the STAC algorithm.

MultiAlign includes an additional processing step clustering LC-MS features across datasets after alignment



# MultiAlign – Tool for Multiple dataset analysis





MultiAlign can skip the identification (peak matching) step if no AMT tag database is available.



Proudly Operated by Baffelle Since 1965

#### Traceback

THIS SECTION DESCRIBES TRACEBACK AND HOW DATA IS LINKED TOGETHER IN THE MULTIALIGN RESULTS DATABASE

Proudly Operated by Battelle Since 1965





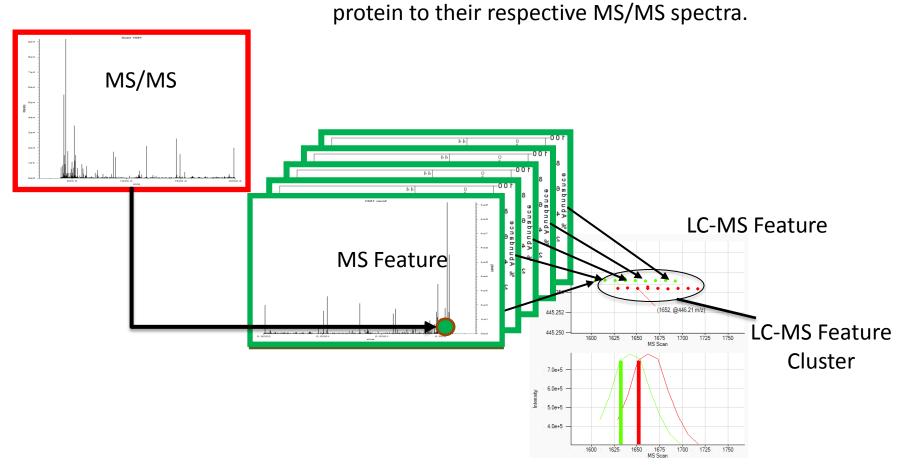
► Traceback simply implies the ability to start with a feature identification (e.g. peptide identification) or a cluster of features and trace all of the way to the raw spectra (MS¹) and MS/MS spectra if available.



# Traceback Visually with spectra



With traceback we are trying to link a cluster, mass tag, or

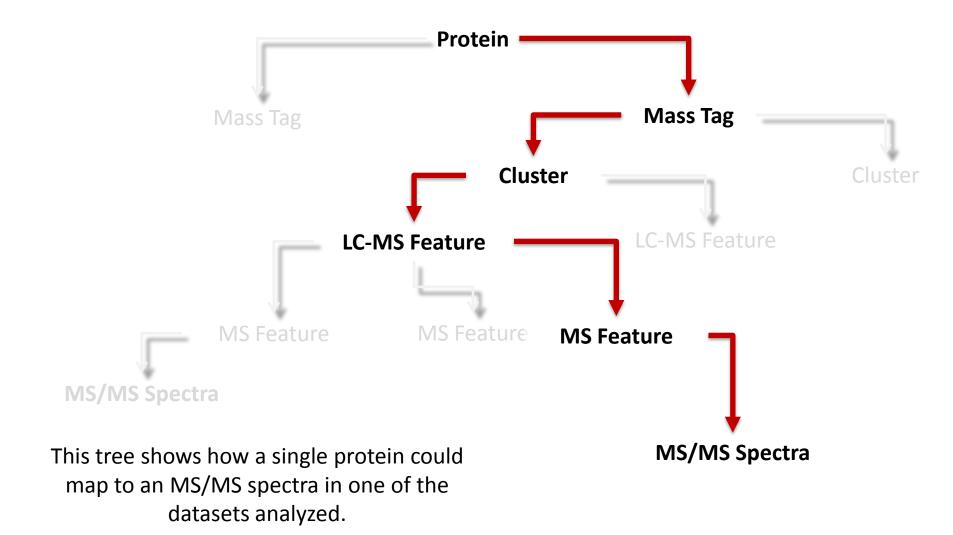




## Mapping a protein to MS/MS spectra



Proudly Operated by Battelle Since 1965



December 18, 2012

16

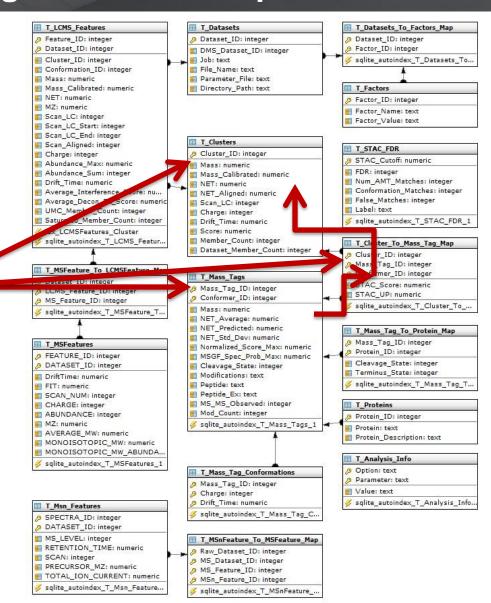




Proudly Operated by Battelle Since 1965

Here we show how this data is linked in the SQLite result database.

Starting with a mass tag, query data joining on mass tag id and cluster id through the mapping table



December 18, 2012

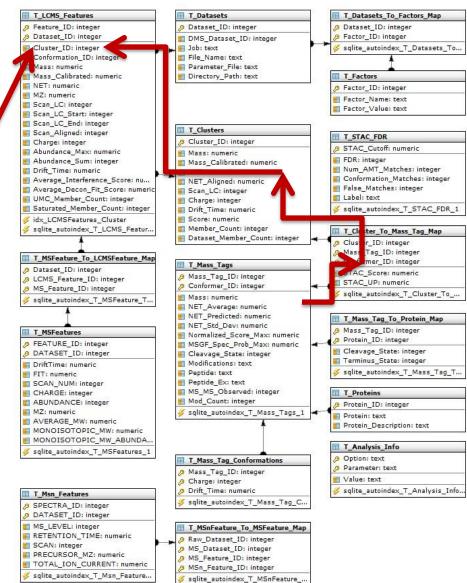




Proudly Operated by Battelle Since 1965

18

Then map the cluster to all possible LC-MS features based on cluster id in the T\_LCMS\_Features Table

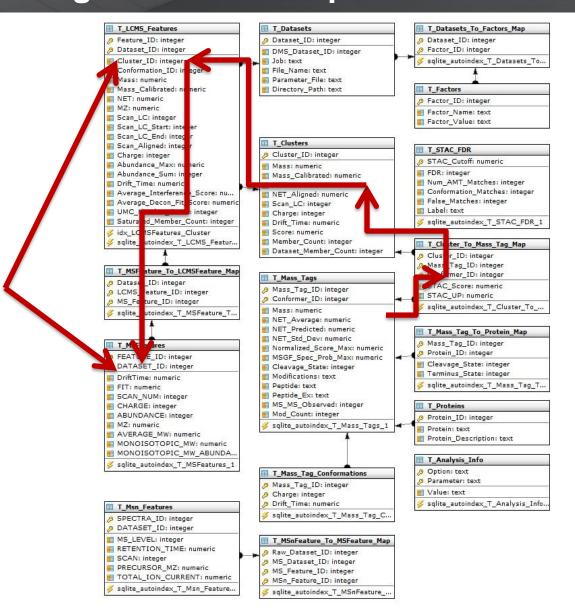






Proudly Operated by Battelle Since 1965

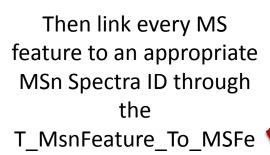
Then map each LC-MS
features to its
appropriate set of MS
Features based on
Dataset ID and LC-MS
feature ID



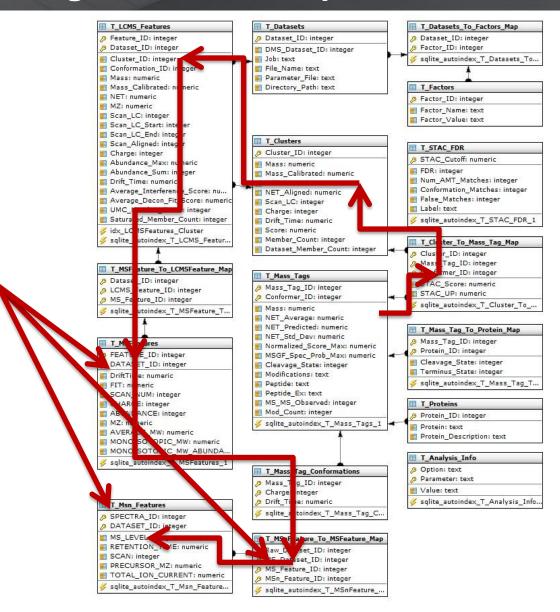




Proudly Operated by Battelle Since 1965



ature Map



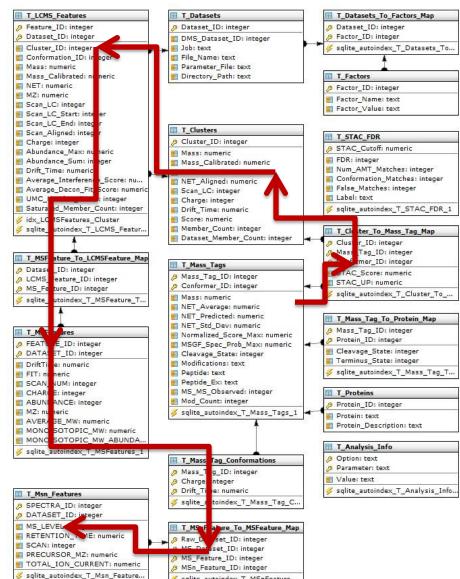




Proudly Operated by Battelle Since 1965

Now every identification is mapped back through to a possible MS/MS spectra.

If AMT is not performed, you can start from the T\_Clusters step and trace backwards



December 18, 2012

y sqlite\_autoindex\_T\_Msn\_Feature....
y sqlite\_autoindex\_T\_MsnFeature....
21

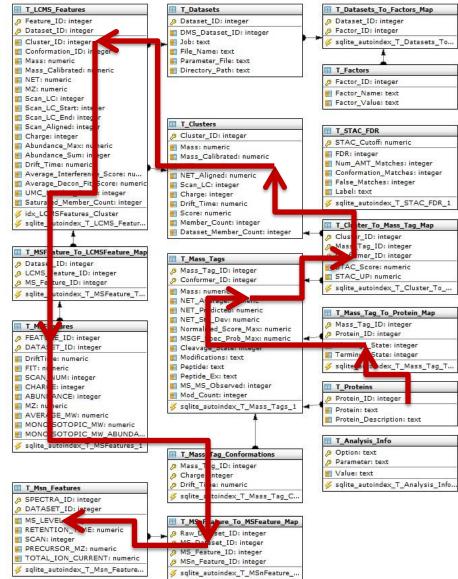


## Database Schema – Mapping a protein to MS/MS spectra



Proudly Operated by Battelle Since 1965

Or you could start from a significantly changing protein ID and move backwards, this additional step is shown with the additional link through the T\_Mass\_Tag To Protein Map table



December 18, 2012

y sqlite\_autoindex\_T\_MSnFeature\_...

g sqlite\_autoindex\_T\_MSnFeature\_...

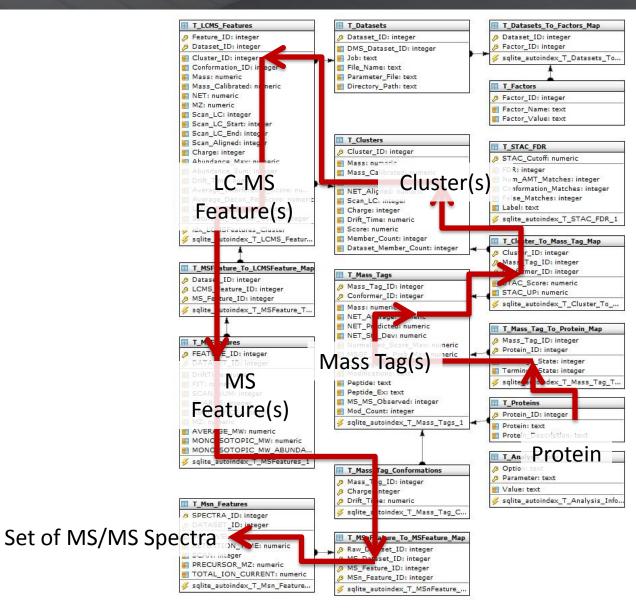
22



## Database Schema – Mapping a protein to MS/MS spectra



Proudly Operated by Battelle Since 1965





For more information see the MultiAlign website:

http://omics.pnl.gov/software/MultiAlign.php