
SRP Workshop: Bioinformatics Tools for Data Integration

SUSAN TILTON

Oregon State University
Corvallis, OR
September 13-14, 2017



Overview of omics data integration

Rationale

Data integration: use of multiple sources of information (or data) collectively within a system

Goal: to provide a better understanding of a system/situation/association to improve knowledge discovery

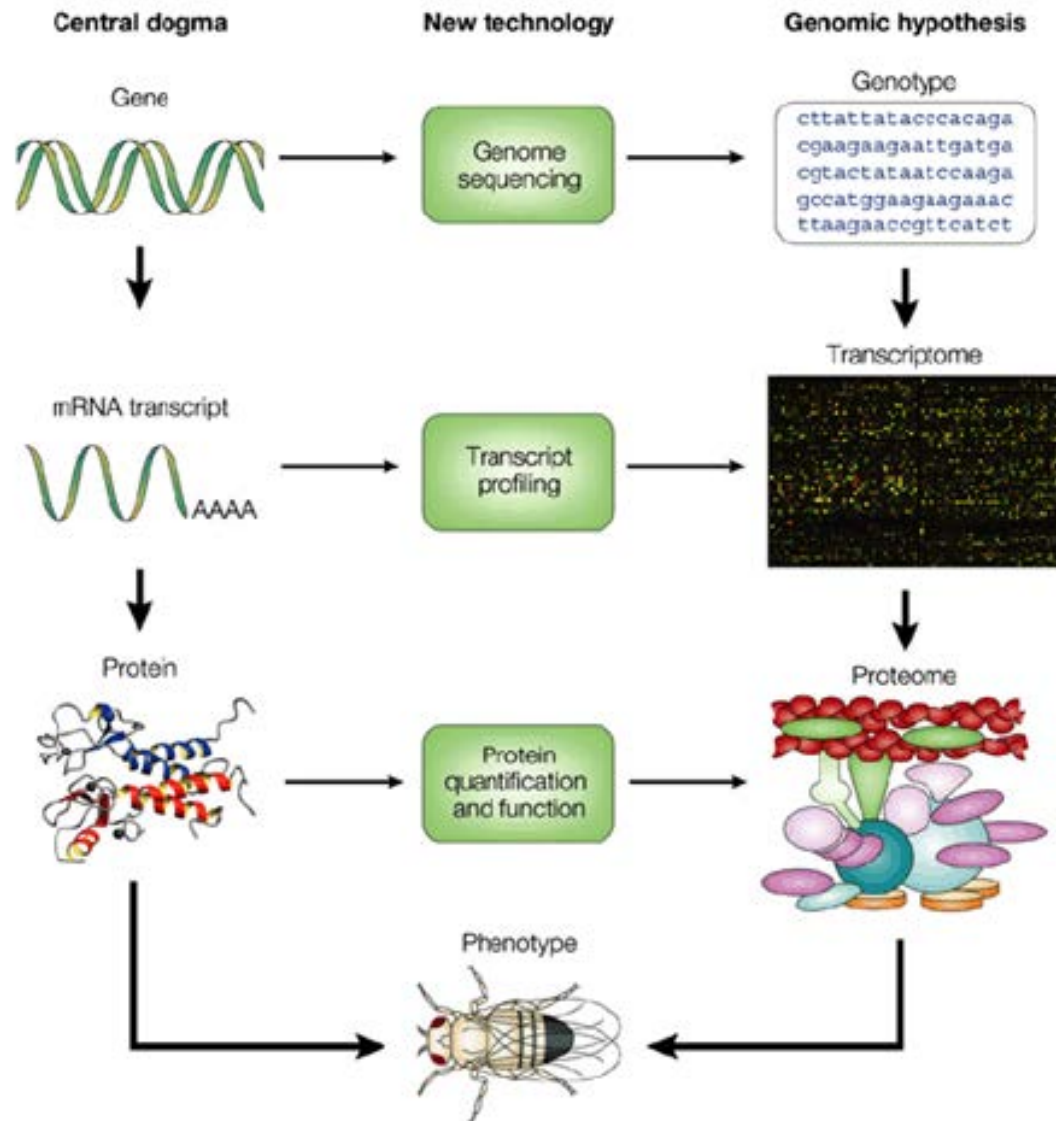
Challenges: data explosion associated with high content technologies

Rationale

Make comparisons among datasets:

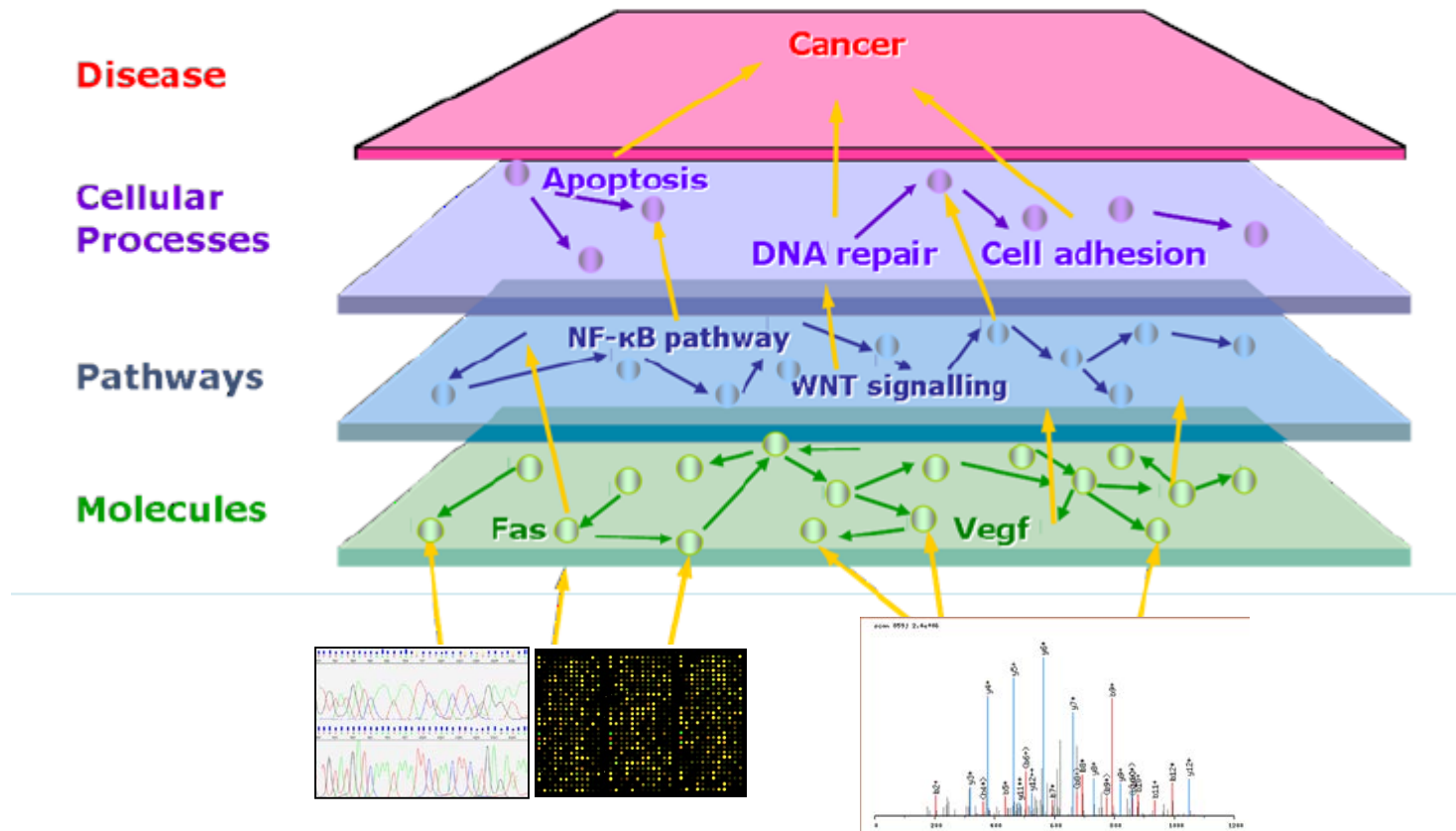
- Technologies (gene/proteins/metabolites)
- Species (mouse/human/zebrafish)
- Platforms (microarray/RNAseq)
- Phenotype (toxicity/disease)
- Experimental models (in vitro/in vivo)
- Dose or time
- Chemical class

Data and technologies



Data and technologies

Ultimate Goal:
Combine data to increase knowledge

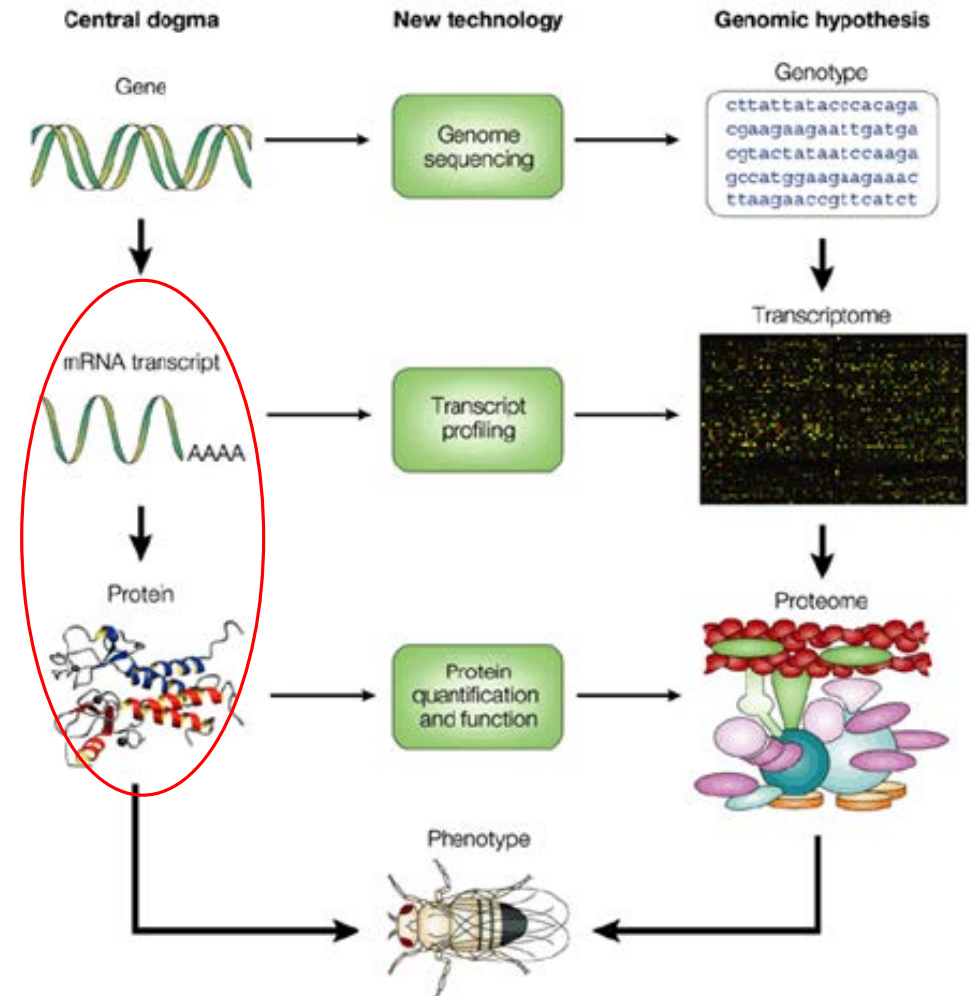


Approaches for data integration

- Direct integration
 - Genes/proteins
 - Cross-species
 - In vivo / in vitro
- Integration at functional or pathway level
- Integration based on common upstream regulators
 - Transcription factors
 - miRNAs
- Integration based on statistical and network-based approaches
 - Correlation
 - Clustering
 - Network structure

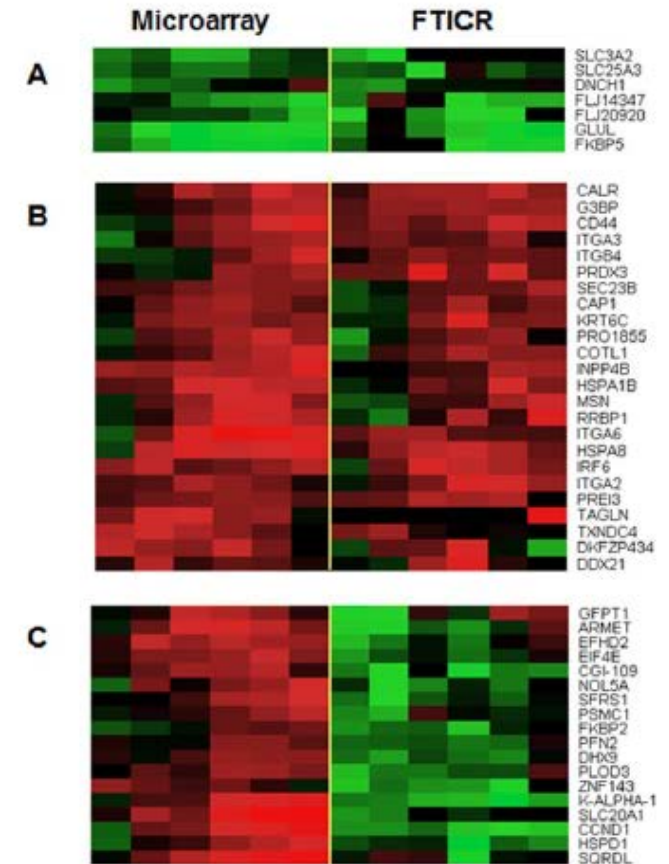
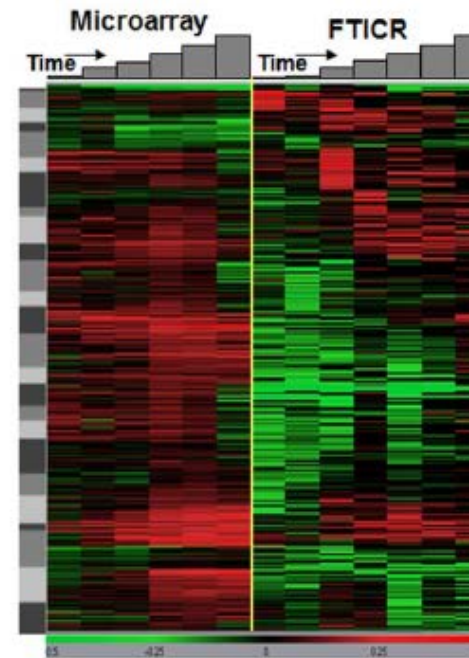
Direct integration of Transcriptomics and Proteomics

- Assumption that gene and protein regulation is directly correlated



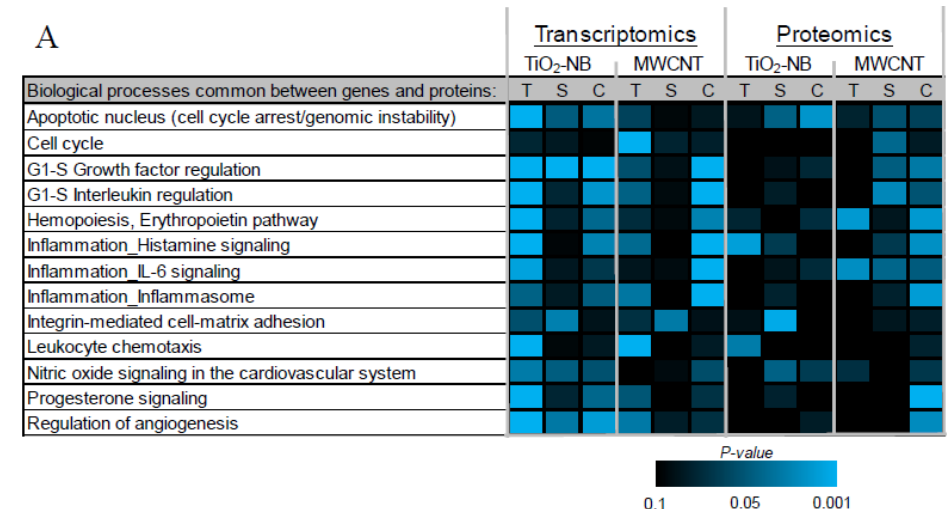
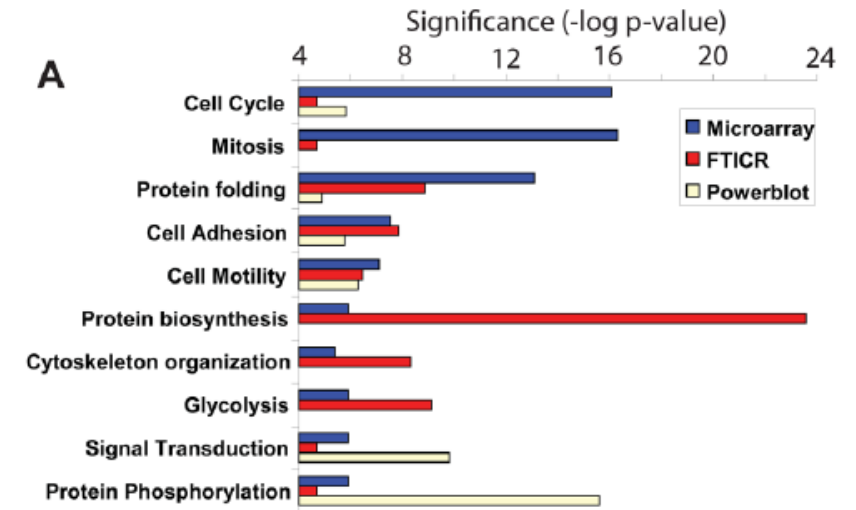
Direct integration of Transcriptomics and Proteomics

- Assumption that gene and protein regulation is directly correlated
- Lack of direct correlation between genes & proteins due to following challenges:
 - Post-translational modifications
 - Gene/protein turnover and degradation
 - Statistical filtering / platform differences and limitations → missing data



Integration at functional or pathway level

- Correlation among datasets is observed at the functional or pathway level
- Even though the same genes and proteins are not regulated simultaneously within the system, genes and proteins within the same pathways are co-expressed
- Measured through pathway or functional enrichment analysis

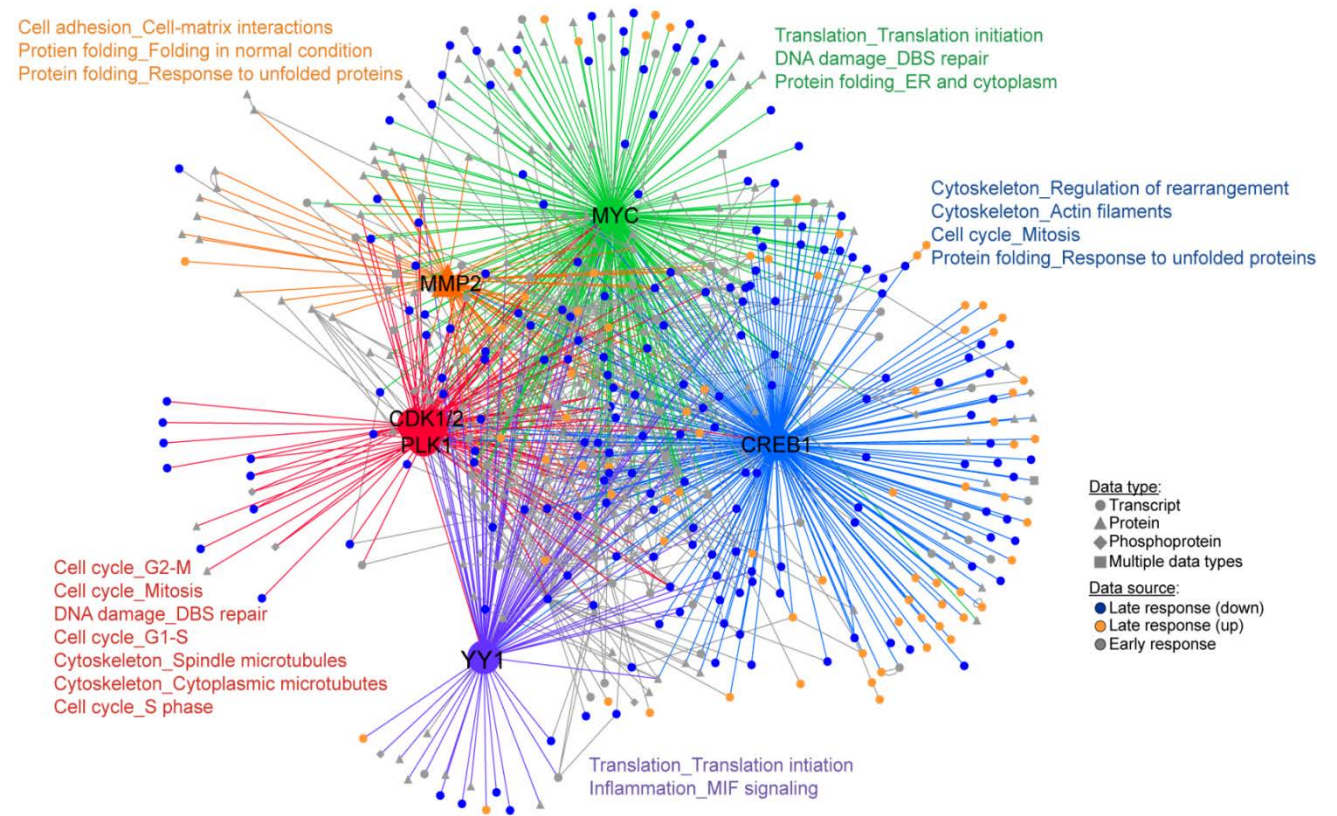


¹WATERS ET AL. 2012. PLOS ONE, 7(3):E34515

²TILTON ET AL. 2014. NANOTOXICOLOGY, 8(5):533-48.

Integration based on common upstream regulators

- Similarly, overlap among datasets is observed through common upstream regulators (e.g. transcription factors)
- Measured through transcription factor enrichment, hub analysis



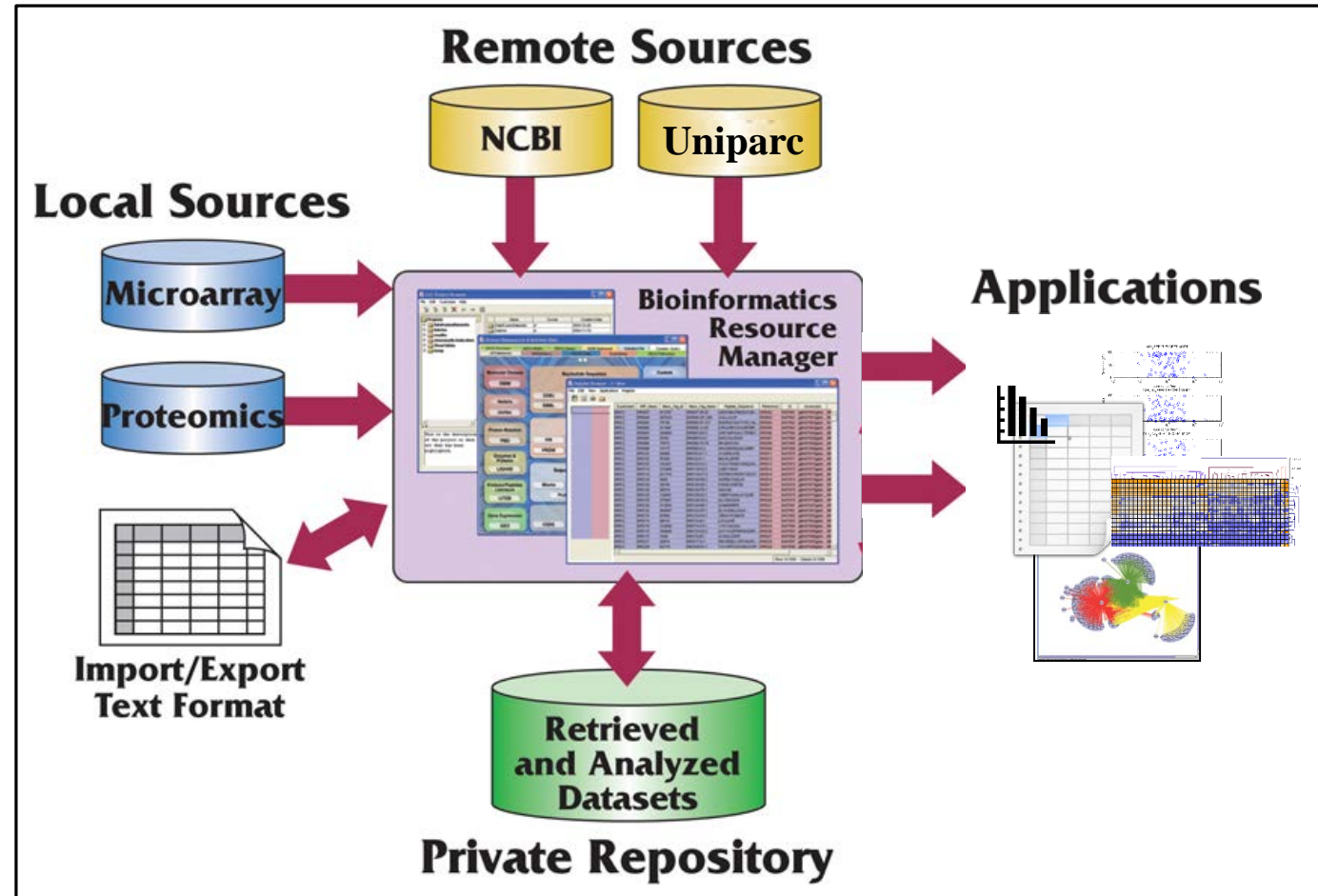
Practical challenges for data integration

- Choosing appropriate approach to answer biological questions
 - Handle missing data
 - Applying statistical/filters to data
- Common bioinformatics problem of translating identifiers across data tables
 - Genes → proteins
 - Mouse → Human
 - Entrez Gene ID → Ensembl Gene ID

Bioinformatics Resource Manager (BRM): A
systems biology web tool for miRNA and omics
with data integration

BRM Overview

A web tool developed to provide biological scientists with computation tools for annotation retrieval, cross-referencing and integration of high-throughput (HTP) data.

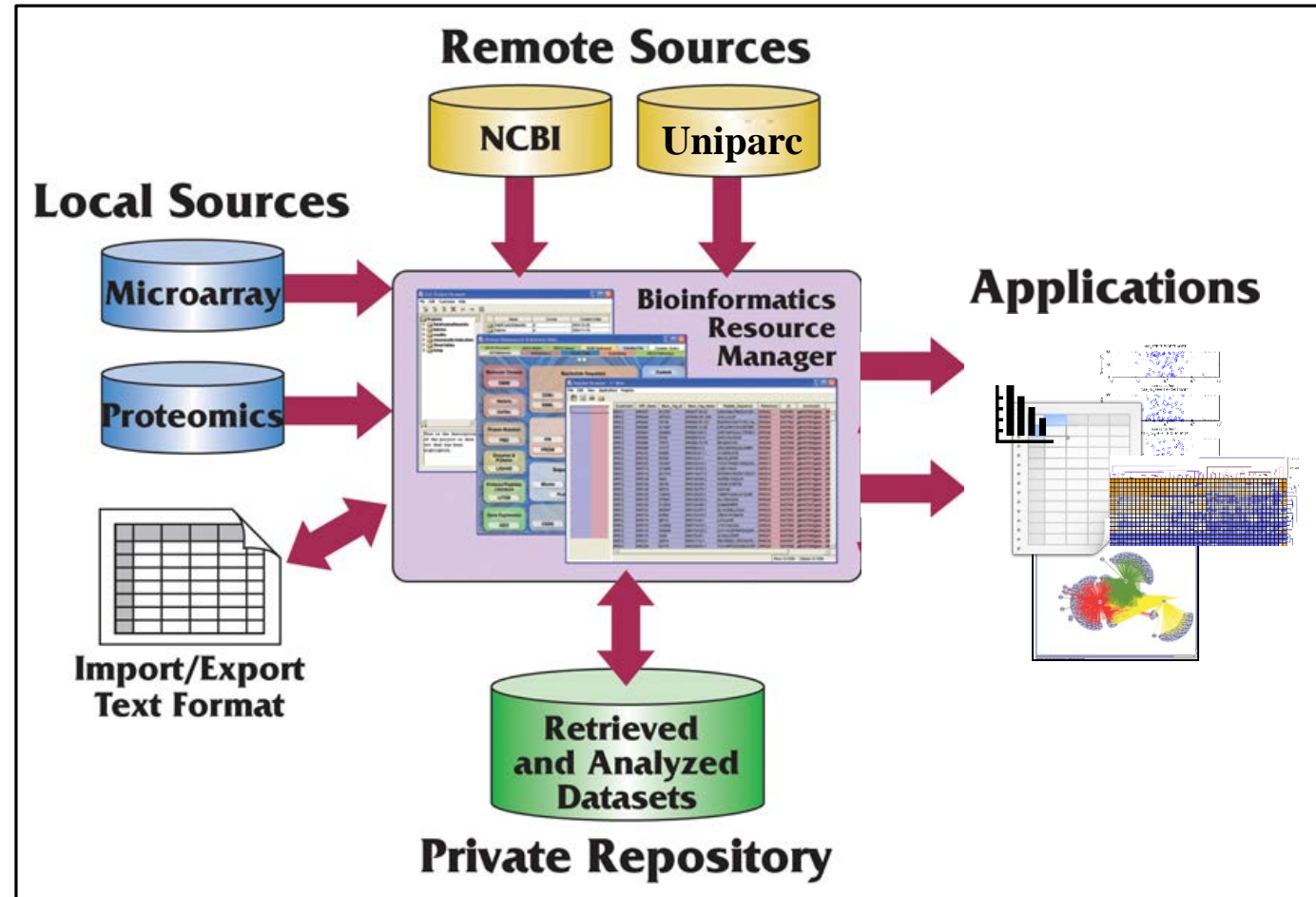


BRM Overview

Provide a platform for integration across HTP datasets

Examples include:

- miRNA/mRNA
- Transcriptomics/Proteomics
- Mouse/Human/Rat/Zebrafish/Macaque
- Cross-Identifier



BRM Overview

Brown et al. *BMC Bioinformatics* (2019) 20:255
<https://doi.org/10.1186/s12859-019-2805-6>

BMC Bioinformatics

DATABASE

Open Access

Bioinformatics Resource Manager: a systems biology web tool for microRNA and omics data integration



Joseph Brown^{1,5}, Aaron R. Phillips², David A. Lewis², Michael-Andres Mans³, Yvonne Chang³, Robert L. Tanguay^{3,4}, Elena S. Peterson², Katrina M. Waters^{1,3,4*} and Susan C. Tilton^{3,4*} 

BRM Capabilities

Bath Annotation Retrieval

- Gene Annotation (NCBI)
- Protein Annotation (Uniparc)
- miRNA Annotation (MiRBase/microCosm)

XREF Identifier Retrieval

- Gene Identifiers (NCBI)
- Protein Identifiers (Uniparc)
- Cross Species Orthologs (Ensembl)

Integration of Data Tables

- Includes all data in output

BRM Capabilities

miRNA Datasource Retrieval

- miRNA Predicted Targets (TargetScan, microcosm, microRNA)
- miRNA Identifiers
- miRNA Metadata

Data Workflows *(performs bioinformatics tasks in background)*

- miRNA Predicted Target Retrieval
- miRNA Target Prediction and Integration with Co-expressed Data
- Data Integration Workflow
 - Transcriptomic and Proteomic Data
 - Cross-Species Data

BRM Web Access⁵

(<http://cbb.pnnl.gov/brm/>)

Bioinformatics Resource Manager

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM Software

Tilton SC, Tal TL, Scroggins SM, Franzosa JA, Peterson ES, Tanguay RL and Waters KM. 2012. Bioinformatics Resource Manager v2.3: An integrated software environment for systems biology with microRNA and cross-species analysis tools. BMC Bioinformatics. 13:311. doi: 10.1186/1471-2105-13-311. PMCID: PMC3534564. PMID: 23174015.

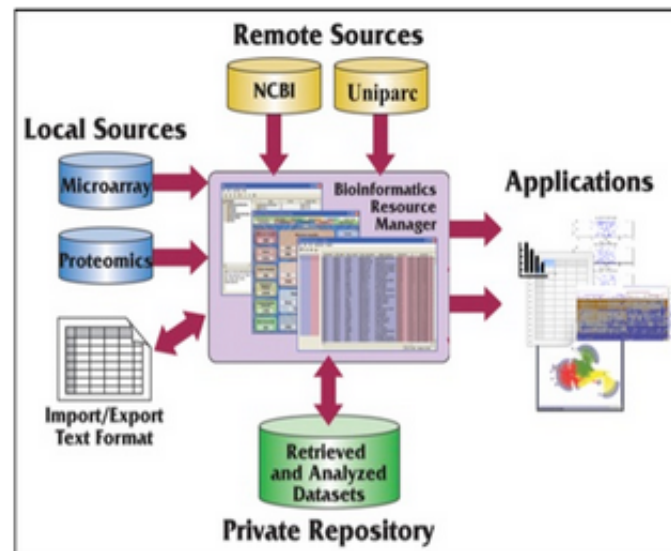
Shah AR, M Singhal, KR Klicker, EG Stephan, HS Wiley, KM Waters. 2007. "Enabling high-throughput data management for systems biology: The Bioinformatics Resource Manager." Bioinformatics 2007 23(7):906-909. doi:10.1093/bioinformatics/btm031. PMID: 17324940.

Bioinformatics Resource Manager Info

The Bioinformatics Resource Manager (BRM) is a software environment developed for data retrieval, integration and analysis of high-throughput (HTP; transcriptomic, proteomic or sequencing) biological data. BRM provides computational tools for biologists to merge datasets, cross-reference gene or protein identifiers, map identifiers across species and add functional annotation from NCBI, UniProt, Ensembl or microRNA databases, including predicted miRNA targets from multiple sources. BRM utilizes easy to navigate workflows for identification of predicted miRNA gene targets and integration with experimental mRNA or protein datasets. BRM further provides generic workflows for integrating cross-platform, cross-technology or cross-species. These tools provide biological researchers with a platform for straightforward integration and analysis of heterogeneous HTP datasets critical for biomedical research.

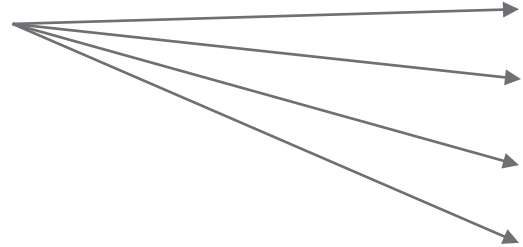
Use BRM for integrating across HTP experimental datasets. Examples include:

- o miRNA/mRNA
- o Transcriptomics/Proteomics
- o RNAseq/Microarray
- o Mouse/Human/Rat/Zebrafish/Macaque
- o Cross-Identifier (e.g. Entrez gene ID/Ensembl gene ID)



BRM Organization: Homepage

Workflows:



Bioinformatics Resource Manager ⓘ

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM

Version Info

Upload Data ▾

- Upload your data as a `.txt/.tab` tab-separated file and include a header ([Load Example](#))
- File uploads larger than 200MB are not recommended

⊞ Upload File

Conversion Options ▲

Output Options ▲

BRM Organization: Homepage

Workflows:

The screenshot shows the Bioinformatics Resource Manager (BRM) homepage. The header is orange with the text "Bioinformatics Resource Manager" and a refresh icon. A sidebar menu on the left contains the following items: "Add Identifiers" (highlighted with a red circle), "Merge Tables", "miRNA Targets", "miRNA Convert", "Cite BRM", and "Version Info". Four arrows point from the word "Workflows:" to the "Add Identifiers", "Merge Tables", "miRNA Targets", and "miRNA Convert" items. The main content area on the right includes a section titled "Upload Data" with a dropdown arrow, containing two bullet points: "Upload your data as a .txt/.tab tab-separated file and include a header (Load Example)" and "File uploads larger than 200MB are not recommended". Below this is an "Upload File" button. Further down are sections for "Conversion Options" and "Output Options", both with upward-pointing arrows. The footer consists of an orange bar on the left and a grey bar on the right.

Bioinformatics Resource Manager

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM

Version Info

Upload Data ▾

- Upload your data as a .txt/.tab tab-separated file and include a header ([Load Example](#))
- File uploads larger than 200MB are not recommended

Upload File

Conversion Options ▲

Output Options ▲

BRM: Batch Annotation Retrieval

Cross-identifier/platform

The screenshot displays the Bioinformatics Resource Manager interface. On the left is a vertical sidebar with navigation links: "Add Identifiers", "Merge Tables", "miRNA Targets", "miRNA Convert", "Cite BRM", and "Version Info". The main area is titled "Dataset Preview" and contains a table with 7 columns: Tracking_ID, gene_id, gene, locus, BAP10_log2(fold_change), BAP10_q_value, and BAP10_sig. Below the table are two sections: "Conversion Options" and "Select identifier types and their respective header values:". The "Conversion Options" section includes dropdown menus for "Identifier Type:" and "Input Column:", along with minus and plus buttons. A red arrow points from the "Version Info" link in the sidebar to the "Select identifier types..." text. Below this, there are labels for "Select a specific input and output to limit search space:", "Input Restriction:", "Output Restriction:", "Output Columns:", and "Select identifiers to include/exclude from output table:". A dropdown menu is open under "Identifier Type:", listing various database identifiers like Ensembl Gene ID, Refseq Protein ID, etc.

Bioinformatics Resource Manager

- Add Identifiers
- Merge Tables
- miRNA Targets
- miRNA Convert
- Cite BRM
- Version Info

Dataset Preview

Tracking_ID	gene_id	gene	locus	BAP10_log2(fold_change)	BAP10_q_value	BAP10_sig
ENSDARG00000028039	ENSDARG00000028039	cyp1a	18:4974675-4985566	7.45376	0.00461398	yes
ENSDARG00000058980	ENSDARG00000058980	cyp1c1	Zv9_NA892:103569-106640	4.06631	0.00461398	yes
ENSDARG00000068934	ENSDARG00000068934	cyp1b1	13:42664263-42671929	3.77738	0.00461398	yes
ENSDARG00000018298	ENSDARG00000018298	cyp1c2	Zv9_scaffold3548:296362-300627	3.49876	0.00461398	yes
ENSDARG00000005039	ENSDARG00000005039	gstp1	4:29632401-29639074	1.97748	0.00461398	yes
ENSDARG00000086826	ENSDARG00000086826	sult5b1	11:44488332-44500653	2.47401	0.00461398	yes

Conversion Options

Select identifier types and their respective header values:

Identifier Type: Select... Input Column: Select... - +

Select a specific input and output to limit search space:

Input Restriction: Select... Output Restriction: Select...

Output Columns: Select...

Select identifiers to include/exclude from output table:

- ☐ Ensembl Gene ID
- ☐ Ensembl Protein ID
- ☐ Ensembl Transcript ID
- ☐ Entrez Gene ID
- ☐ Gene Symbol
- ☐ PIR
- ☐ Refseq Gene ID
- ☐ Refseq Protein ID
- ☐ Refseq Transcript ID

BRM: Batch Annotation Retrieval

Cross-species

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM

Version Info

Dataset Preview

Tracking_ID	gene_id	gene	locus	BAP10_log2(fold_change)	BAP10_q_value	BAP10_sig
ENSDARG00000028039	ENSDARG00000028039	cyp1a	18:4974675-4985586	7.45376	0.00461398	yes
ENSDARG00000058980	ENSDARG00000058980	cyp1c1	Zv9_NA892:103569-106640	4.06631	0.00461398	yes
ENSDARG00000068934	ENSDARG00000068934	cyp1b1	13:42664263-42671929	3.77738	0.00461398	yes
ENSDARG00000018298	ENSDARG00000018298	cyp1c2	Zv9_scaffold3548:298352-300627	3.49876	0.00461398	yes
ENSDARG00000005039	ENSDARG00000005039	gstp1	4:29832401-29839074	1.97748	0.00461398	yes
ENSDARG00000086826	ENSDARG00000086826	sult5b1	11:44488332-44500653	2.47401	0.00461398	yes

Conversion Options ^

Select identifier types and their respective header values:

Identifier Type:

Ensembl Gene ID

 Input Column:

gene_id

 - +

Select a species restriction on the input and output to limit search space:

Input Restriction:

None

Output Restriction:

None

Rat

Human

Macaque

Mouse

Zebrafish

Output Options

Select identifier types and their respective header values:

☐ Ensembl Gene ID

☐ Ensembl Protein ID

☐ Ensembl Transcript ID

☐ Entrez Gene ID

☐ Gene Symbol

☐ PIR

☐ Refseq Gene ID

☐ Refseq Protein ID

BRM: Batch Annotation Retrieval

Bioinformatics Resource Manager

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM


Version Info

Identifier Type: Ensembl Gene ID Input Column: gene_id – +

Select a species restriction on the input and output to limit search space:

Input Restriction: Zebrafish

Output Restriction: Zebrafish

Output Options 

Select identifiers to add onto input table:

☐ Ensembl Gene ID

☐ Ensembl Protein ID

☐ Ensembl Transcript ID

☒ Entrez Gene ID

☐ Gene Symbol

☐ PIR

☐ Refseq Gene ID

☐ Refseq Protein ID

☐ Refseq Transcript ID

☐ Swissprot Accession

☐ Swissprot ID

☐ Tax ID

Select how to handle hits with multiple results:

☒ Return 1 result

☒ The first result

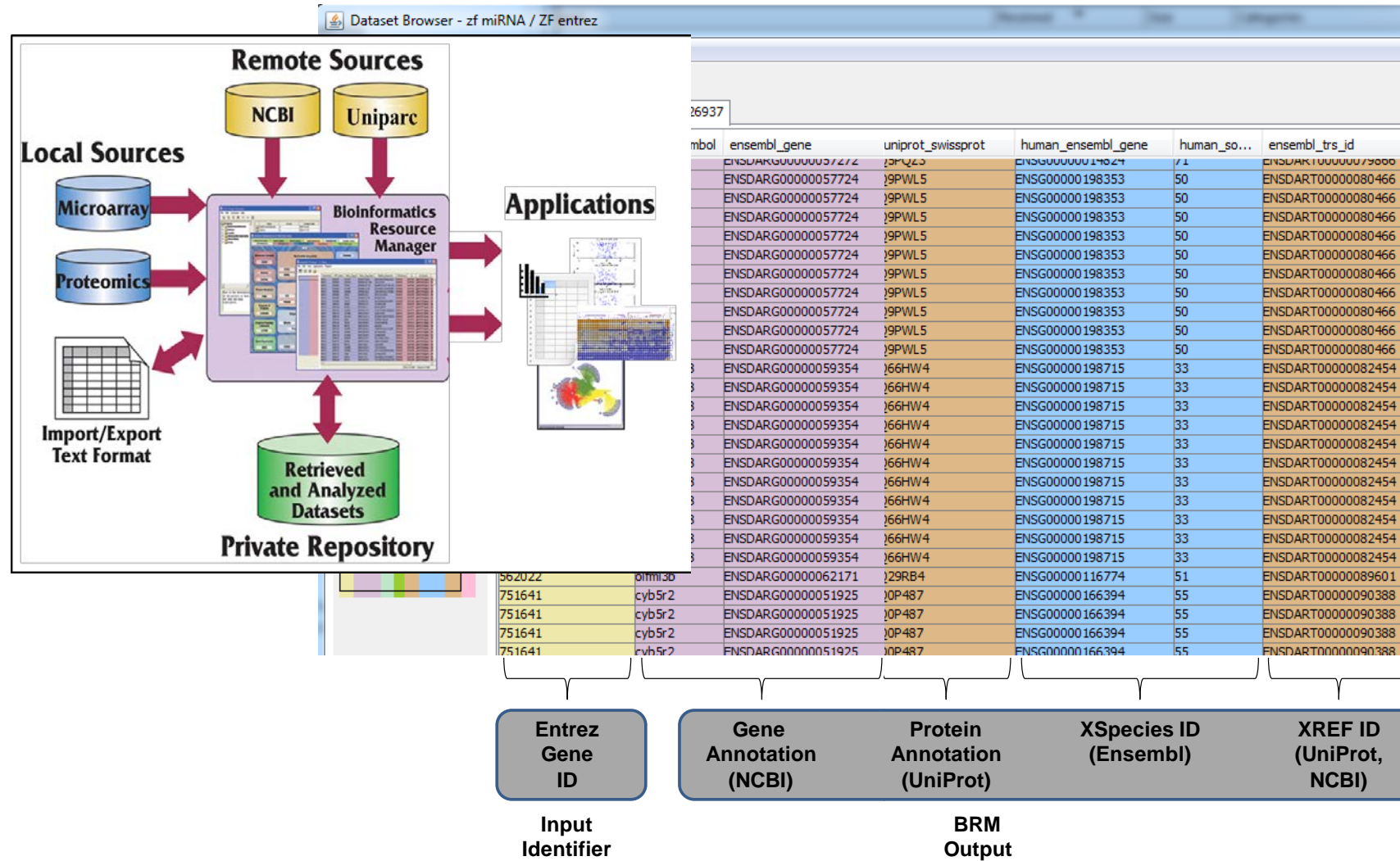
☐ The alphabetically last result

☐ Multiple entries per row (separated by semi-colons)

☐ Generate extra rows for multiple results

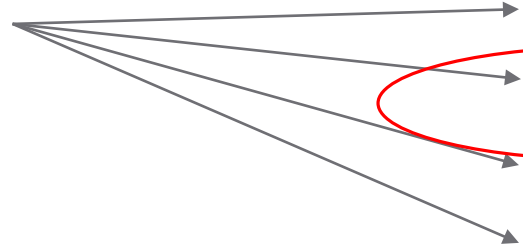
Run Query

BRM: Batch Annotation Retrieval



BRM Organization: Homepage

Workflows:



Bioinformatics Resource Manager ⓘ

Add Identifiers

Merge Tables

miRNA Targets

miRNA Convert

Cite BRM

Version Info

Upload Table 1 ▾

- Upload your data as a .txt/.tab tab-separated file and include a header ([Load Example](#))

Upload File

Upload Table 2 ▲

Upload File

BRM Organization: Homepage

Workflows:

The screenshot shows the Bioinformatics Resource Manager (BRM) homepage. On the left is a vertical navigation menu with the following items: 'Add Identifiers', 'Merge Tables', 'miRNA Targets', 'miRNA Convert', 'Cite BRM', and 'Version Info'. The 'miRNA Targets' item is highlighted with a red oval. Four arrows originate from the word 'Workflows:' and point to each of the four workflow-related menu items: 'Add Identifiers', 'Merge Tables', 'miRNA Targets', and 'miRNA Convert'. The main content area on the right features an orange header bar with the text 'Bioinformatics Resource Manager' and a refresh icon. Below the header, the 'Upload and Configure Dataset' section is expanded, showing instructions for identifying microRNA to gene target interactions and an 'Upload File' button. The 'Choose miRNA Target Databases' section is also expanded, displaying checkboxes for 'miRTarBase', 'Microcosm', 'MicroRNA', and 'TargetScan'. Below these, a 'Required Hits From' dropdown is set to '1 Of 0 Databases'. The 'Merge miRNA results with Gene ID Table (Optional)' section is collapsed, showing an 'Upload File' button and a 'Find Targets' button.

Bioinformatics Resource Manager

Workflows:

- Add Identifiers
- Merge Tables
- miRNA Targets
- miRNA Convert
- Cite BRM
- Version Info

Upload and Configure Dataset

- Identify microRNA to gene target interactions from mature miRNA names ([Load Example](#))
- Identify gene target to microRNA interactions from gene names, e.g. FOXP2 ([Load Example](#))
- Upload your data as a .txt/ tab tab-separated file and include a header

Upload File

Choose miRNA Target Databases

- miRTarBase
- Microcosm
- MicroRNA
- TargetScan

Required Hits From: 1 Of 0 Databases

Merge miRNA results with Gene ID Table (Optional)

Upload File

Find Targets

BRM: microRNA tools

- Query predicted and experimentally validated miRNA targets from multiple data sources
 - TargetScan (<http://www.targetscan.org/>)
 - MicroCosm/miRBase (<http://mirbase.org/>)
 - miRNA (<http://www.microrna.org/>)
 - miRTarBase (<http://mirtarbase.mbc.nctu.edu.tw/>)
- Retrieve potential regulatory miRNAs for known genes
- Retrieve updated miRNA annotation
- Integrate experimentally derived miRNA and mRNA datasets for
 - *Homo sapiens* (human)
 - *Mus musculus* (mouse)
 - *Danio rerio* (zebrafish)
- Use output with other bioinformatics tools:
 - Visualize miRNA regulated gene signatures as heatmap or network
 - Determine functional consequences of miRNA regulation

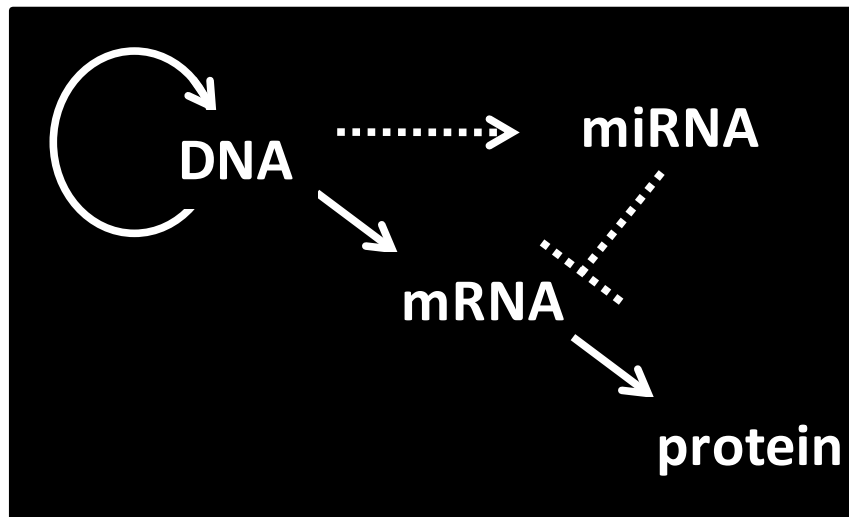
Overview

miRNAs – small non-coding RNAs that regulate gene expression in a sequence specific manner

Importance in biology:

- Developmental timing
- Cell differentiation, proliferation and apoptosis
- Energy metabolism
- Antiviral defense

→ Cancer biology, neurobiology, infectious disease, etc.

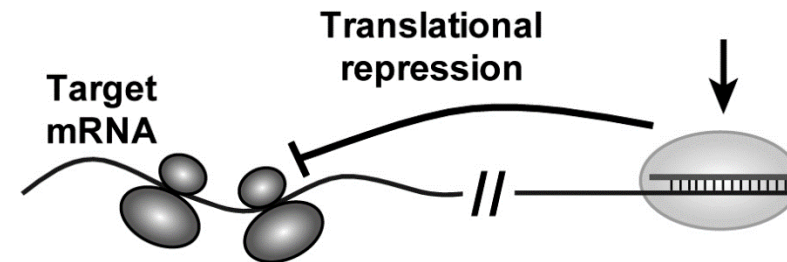
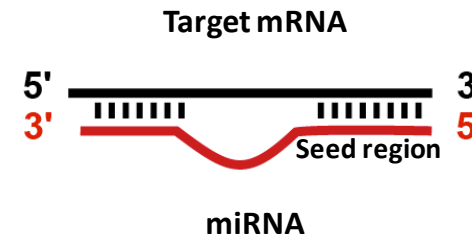
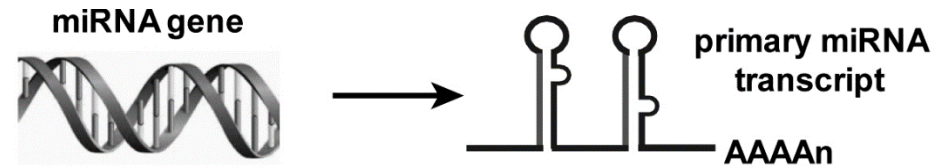


→ Important to identify miRNAs, targets and functional consequences

Important Computational Considerations

For miRNA identification and target predictions:

1. Sequence-based
2. Reliance on secondary structure/interactions
3. Evolutionary species conservation



BRM: microRNA tools

➤ Query predicted miRNA targets from multiple data sources

- TargetScan (<http://www.targetscan.org/>)
- MicroCosm/miRBase (<http://mirbase.org/>)
- miRNA (<http://www.microrna.org/>)
- miRTarBase (<http://mirtarbase.mbc.nctu.edu.tw/>)



➤ Retrieve potential regulatory miRNAs for known genes

➤ Retrieve updated miRNA annotation

➤ Integrate experimentally derived miRNA and mRNA datasets for

- *Homo sapiens* (human)
- *Mus musculus* (mouse)
- *Danio rerio* (zebrafish)

➤ Use output with other bioinformatics tools:

- Visualize miRNA regulated gene signatures as heatmap or network
- Determine functional consequences of miRNA regulation

miRNA Target Prediction Resources⁷⁻⁹

Resource	Type	Website
DIANA-microT-CDS	miRNA target prediction	http://www.microrna.gr/microT-CDS
EvoFold	miRNA target prediction	http://users.soe.ucsc.edu/~jsp/EvoFold/
MicroCosm	miRNA target prediction	http://www.ebi.ac.uk/enright-srv/microcosm/htdocs/targets/v5/
microRNA.org	miRNA target prediction	http://microrna.org
miRDB	miRNA target prediction	http://www.mirdb.org
miRiam	miRNA binding	http://ferrolab.dmi.unict.it/miriam.html
MiRscan	miRNA target prediction	http://genes.mit.edu/mirscan/
PicTar	miRNA target prediction	http://pictar.mdc-berlin.de/
PITA	miRNA target prediction	http://genie.weizmann.ac.il/pubs/mir07/
RNA22	miRNA target prediction	http://cbcsrv.watson.ibm.com/rna22.html
RNAz	miRNA target prediction	http://www.tbi.univie.ac.at/~wash/RNAz/
TargetMiner	miRNA target prediction	http://www.isical.ac.in/~bioinfo_miu/targetminer20.htm
Targetscan	miRNA target prediction	http://www.targetscan.org/
DIANA-TarBase	Manually curated validated mirna target database	http://www.microrna.gr/tarbase
miRecords	Manually curated validated mirna target database	http://mirecords.biolead.org
miRTarBase	Manually curated validated mirna target database	http://mirtarbase.mbc.nctu.edu.tw

Balancing specificity and sensitivity in target prediction

TargetScan

- TargetScanS algorithm
- higher specificity, lower sensitivity
- Considers many parameters (seed matching, 3' complementarity, local AU content, cross-species conservation)
- Omits sites with poor seed pairing

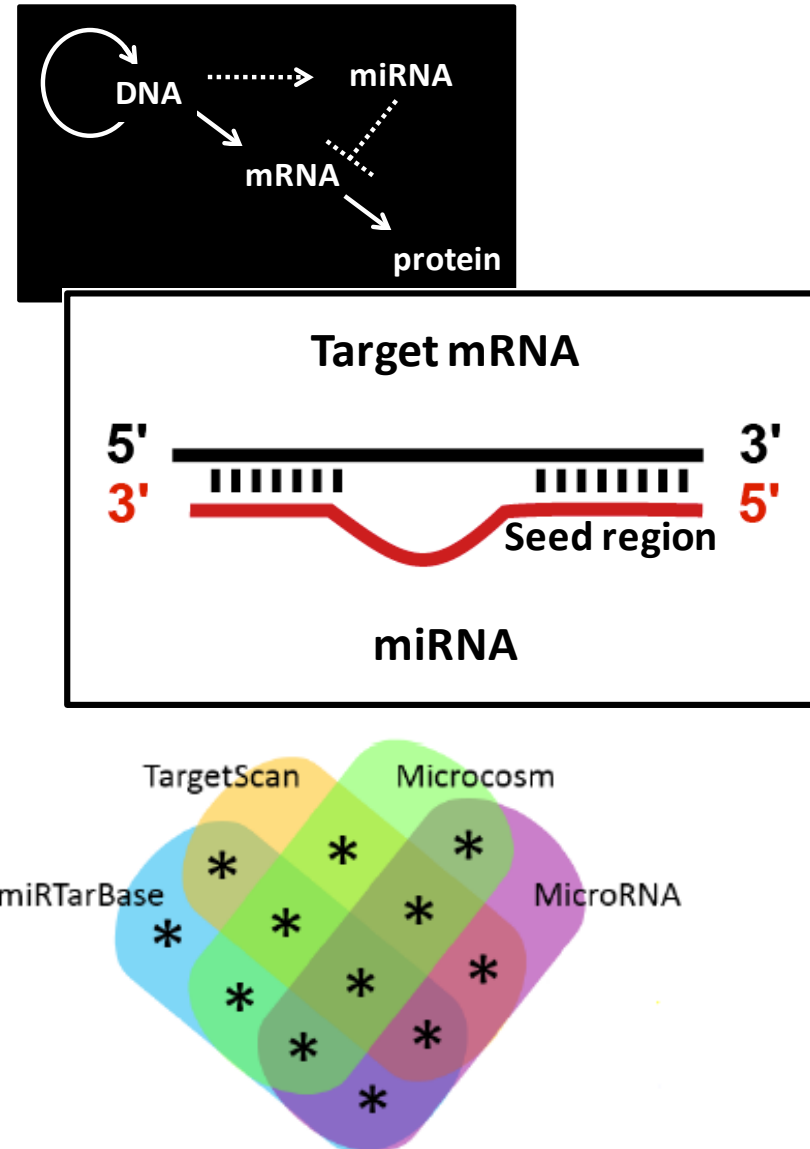
microRNA.org

- miRanda algorithm
- lower specificity, higher sensitivity
- Considers complementarity and free energy binding, species conservation
- More false positives

Challenges and Considerations

False positives in miRNA prediction:

- Some targets with perfect seed pairing are not repressed experimentally
- Continued discussion about which attributes are most important (sequence v. structure)
 - Low seed pairing stability
 - High target site abundance
- Methods to reduce false positives:
 - Species conservation
 - Combinatorial approaches (multiple tools for prediction)
 - mRNA coexpression



BRM: microRNA tools

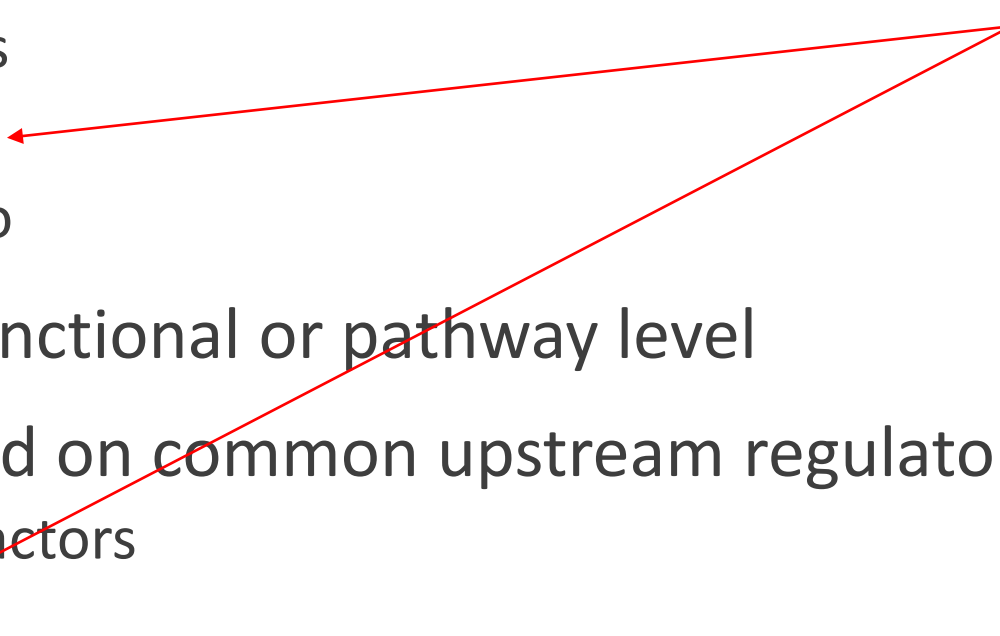
- Query predicted and experimentally validated miRNA targets from multiple data sources
 - TargetScan (<http://www.targetscan.org/>)
 - MicroCosm/miRBase (<http://mirbase.org/>)
 - miRNA (<http://www.microrna.org/>)
 - miRTarBase (<http://mirtarbase.mbc.nctu.edu.tw/>)
- Retrieve potential regulatory miRNAs for known genes
- Retrieve updated miRNA annotation
- Integrate experimentally derived miRNA and mRNA datasets for
 - *Homo sapiens* (human)
 - *Mus musculus* (mouse)
 - *Danio rerio* (zebrafish)
- Use output with other bioinformatics tools:
 - Visualize miRNA regulated gene signatures as heatmap or network
 - Determine functional consequences of miRNA regulation

BRM as a tool for data integration

BRM is a web-based software platform that allows biological researchers the necessary computational and bioinformatics tools for integration and comparison of multiple HTP omics datasets through easy-to-navigate workflows.

- Incorporates cross-reference IDs and annotation directly into your dataset
- Integrates data tables that have no common identifiers using NCBI, Uniparc and Ensembl
- Integrates data tables using multiple IDs and data sources for optimal matching
- Integrates data tables with any common information (string match)
- Works independent of data type and format
- Output includes user's full dataset(s) (all rows and columns) after data integration or annotation retrieval

Approaches for data integration

- Direct integration
 - Genes/proteins
 - Cross-species
 - In vivo / in vitro
 - Integration at functional or pathway level
 - Integration based on common upstream regulators
 - Transcription factors
 - miRNAs
 - Integration based on statistical and network-based approaches
 - Correlation
 - Clustering
 - Network structure
- Hands-on session
- 

BRM: Acknowledgements, publications and funding

BRM Acknowledgements and Funding:

- NIEHS Superfund Research Program P42_ES016465
- BRM was developed through funding from the Pacific Northwest National Laboratory, a multiprogram national laboratory operated by Battelle for the U.S. Department of Energy under Contract DE- AC05-76RL01830.

BRM Web Tool Development Team

Katrina Waters
Elena Peterson
Aaron Phillips
Joe Brown
Susan Tilton

