

modENCODE Data Coordination Center

www.modencode.org



modENCODE Scope

- Transcriptome
- Regulatory elements
- Chromosome and chromatin associated elements
- Small and microRNAs
- Origins of replication

modENCODE is a Resource



- DCC: Delivers data to you, the community
- We ask you to tell us what you want
- DCC is listening

~100 Year History of Community Cooperation

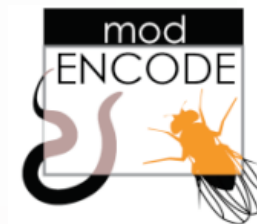


- Resource/ Stock Centres (BDGP, CGC, DGRC...)
- Shared reagent resources
- Genome projects (single, multiple)
- Annotation jamborees
- Online resources
- *Drosophila/ C. elegans* coordination (modENCODE)

modENCODE in Relation to ENCODE

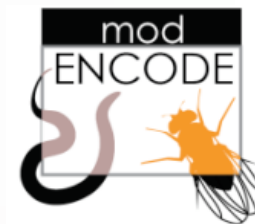


- Both working on functional elements of DNA
- Exchange of technologies with ENCODE (UCSC)
- modENCODE advantages:
 - can do experimental validation
 - identify functional elements in repetitive sequences
 - Comparisons between worm and fly (and human...)
- This community must set an example for ENCODE



Who are the DCC?

- Principal Investigators
 - Lincoln Stein (PI)
 - Suzi Lewis (coPI)
 - Gos Micklem (coPI)
 - Jim Kent (coPI)
- Wiki and Web Site
 - Sergio Contrino
 - François Guillier
- ENCODE Contacts
 - Kate Rosenbloom
 - Galt Barber
- Data Managers
 - E.O. Stinson (fly)
 - Nicole Washington (fly)
 - Sheldon McKay (worm)
 - Zheng Zha (worm)
- Infrastructure
 - Richard Smith
 - Kim Rutherford
 - Chris Mungall



The DCC's Charge

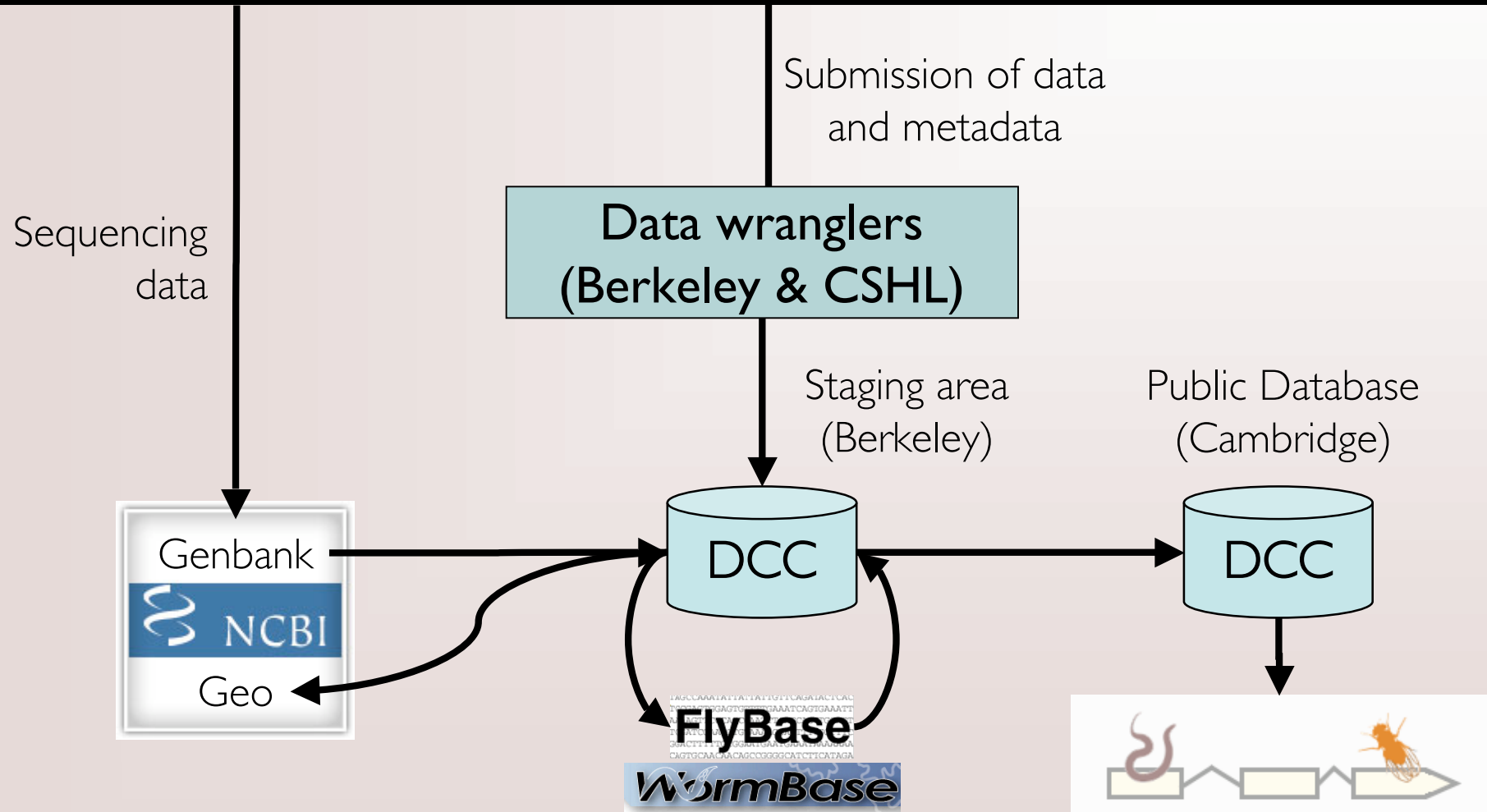
- Provide support for data submission and tracking
- Ensure data integrity and consistency
- Integrate with other information sources
- Make available to the research community

Types of Queries

- Retrieve data by:
 - Protocol type
 - Submitting lab
 - Platform (e.g. Affy Chip set 5.x)
 - Experimental attribute (e.g. Show me the expression array results that use a sliding window $>$ of 30bp)
- Ask useful questions:
 - What antibodies were produced in rabbit that work for ChIP-chip?
 - What is the EST evidence for a particular gene model?
 - What gene models lack cDNA evidence for their TSS?
 - What genomic regions show transfrags that don't have RNA Pol II binding sites within 500bp?

modENCODE & the DCC

Snyder Waterston Lai Henikoff Piano Karpen Celniker White Lieb MacAlpine



Data type summary



DATA TYPES	experimental variables	Data Files	Annotation Files
ChIP-chip	genotype	Array*	GFF3
ChIP-Seq	cell lines	GFF3	WIG
Expression Microarray	cell type	Images?	BED
CGH tiling array	dev stage	weight matrices	
CGH resolution array	tissue type	FASTA	
RACE sequence (454)	antibody	NCBI_XML	
RT-PCR sequence (Solexa)	depleted protein(s)		
cDNA sequence (454)	RNAi target		
Mass-Spec	growth conditions		
northern blots	sequence		
phenotypes	concentration		
	primers		
	gene target ID		
	label		
	threshold cutoffs		
	software version		
	genome version		

Data type summary



	DATA TYPES	experimental variables	Data Files	Annotation Files	
Cell Type Ontology	ChIP-chip	genotype	Array*	GFF3	Sequence Ontology
	ChIP-Seq	cell lines	GFF3	WIG	
	Expression Microarray	cell type	Images?		Drosophila Anatomy Ontology
	CGH tiling array	dev stage	weight matrices		
	CGH resolution array	tissue type	FASTA		
	RACE sequence (454)	antibody	NCBI_XML		
	RT-PCR sequence (Solexa)	depleted protein(s)			
	cDNA sequence (454)	RNAi target			
	Mass-Spec	growth conditions			
	northern blots	sequence			
phenotypes	concentration				
Units Ontology		primers			
		gene target ID			
		label			
ChEBI Ontology		threshold cutoffs			
		software version			
		genome version			

Data type summary



DATA TYPES	experimental variables	Data Files	Annotation Files
ChIP-chip	genotype	Array*	GFF3
ChIP-Seq	cell lines	GFF3	WIG
Expression Microarray	cell type	Images?	BED
CGH tiling array	dev stage	weight matrices	
CGH resolution array	tissue type	FASTA	
RACE sequence (454)	antibody	NCBI_XML	
sequence (Solexa)	depleted protein(s)		
sequence (454)	RNAi target		
Mass-Spec	growth conditions		
northern blots	sequence		
phenotypes	concentration		
	primers		
	gene target ID		
	label		
	threshold cutoffs		
	software version		
	genome version		

Permanent
public ID

sequence (454)

Mass-Spec

northern blots

phenotypes

depleted protein(s)

RNAi target

growth conditions

sequence

concentration

primers

gene target ID

label

threshold cutoffs

software version

genome version

Data type summary



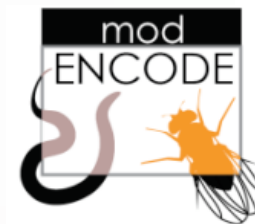
DATA TYPES	experimental variables	Data Files	Annotation Files
ChIP-chip	genotype	Array*	GFF3
ChIP-Seq	cell lines	GFF3	WIG
Expression Microarray	cell type	Images?	BED
CGH tiling array	dev stage	weight matrices	
CGH resolution array	tissue type	FASTA	
RACE sequence (454)	antibody	NCBI_XML	
RT-PCR sequence (Solexa)	depleted protein(s)		
cDNA sequence (454)	RNAi target		
Mass-Spec	growth conditions		
northern blots	sequence		
phenotypes	concentration		
	primers		
	gene target ID		
	label		
	threshold cutoffs		
	software version		
	genome version		

Previously provided
descriptions from
experimental project

Biological Investigation Report Tabular Format: BIR-TAB



- A superset of microarray submission formats (MAGE-TAB)
- Extensive support for controlled vocabulary
- Supports any protocol that can be abstracted as an operation with inputs and outputs



Steps in Validation

- Syntax checking
- Only terms from specified controlled vocabularies
- External identifiers exist
- Agrees with previously supplied protocol and reagents descriptions
- Check format of data files (currently WIG, GFF3, and BED supported)
- Load into database
- Generate downstream output files

Why Consistency



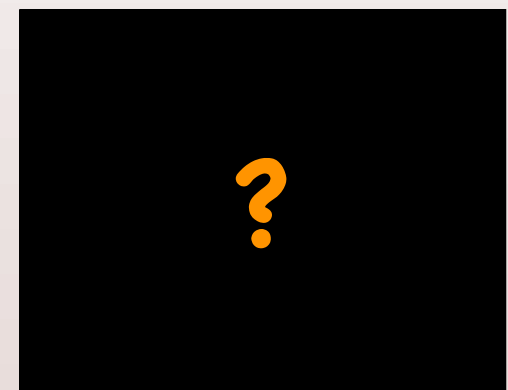
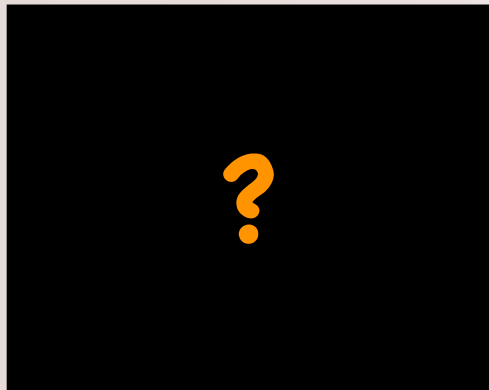
- Comparison, CompaRison, Comparison
- Without consistency there can be no comparison



Why Consistency



- Comparison, CompaRison, Comparison
- Without consistency there can be no comparison





Status of data delivery

Lai Piano Karpen Snyder Henikoff Lieb Waterston White MacAlpine Celniker

In
discussions

Protocols and
reagents being added

Sample data for
defining metadata

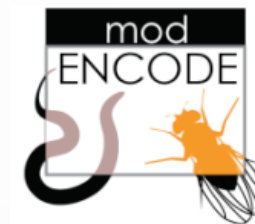
Raw ChIP data submission
being formalized

About to submit
first set of
microarray data

First submission arrived;
collecting metadata

Test submission;
soon to be
completed

First submission
in DCC



modENCODE Q & A

- Royal Palm Salon 4
- Personal question and answer sessions
- Ask for Richard, Nicole, Gos or EO





← Town and Country Guests
Have Preferred Tee Time Availability.
Dial Ext 1234 for information.



☒ FlyBase Genes



☒ BDGP 5 prime RACE

RM03535.5prime
☒ RM03935.5prime
☒ RM03635.5prime
☒ RM04035.5prime

RM02207.5prime
☒ RM02407.5prime
☒ RM02507.5prime
☒ RM02307.5prime
☒ RM01907.5prime
☒ RM02107.5prime
☒ RM01807.5prime
☒ RM02007.5prime

www.modencode.org/cgi-bin/gbrowse/fly/

☒ 1182-4H: Affy "Present" probes



☒ Transcribed Fragments by cell-line



☒ Transcription by cell line (bandwidth 50)

