# From Fragments to Finished Genome

#### Overview of Sequencing, Assembly, and Closure at TIGR

Luke Tallon





#### **Genome Sequencing Process**

rDNA

Molecules

**Library Construction** 

**Clone Picking** 

**Template Preparation** 

**Sequencing Reactions** 

Electrophoresis and Base Calling

**Genome Assembly** 







Genome Closure Order Contigs Close Gaps



#### **Identify Repeats Finish the Genome**

Annotation





# Library Construction

- Break the genomic DNA from the organism into many small fragments
- The fragments are cloned into plasmid vectors
- This collection of fragments is called a library



## Library Requirements

1. Free vector should be at low or undetectable level.

2. No chimeric clones. Chimeras occur two or more random fragments from separate parts of the genome recombine and end up next to each other.

3. The majority of the inserts should be of relatively uniform size.

4. Libraries need to be random and cover the whole genome.



#### **DNA Shearing with HydroShear**



Nebulized DNA vs HydroShear sheared DNA Size Selection Gels



Quant Gels Prior to Vector Ligation

THE INSTITUTE FOR GENOMIC RESEARCH

HydroShear Assembly Cleaning QC



PCR Products After 4 Wash Cycles 1 BAC DNA, 2 sheared BAC DNA, 3 sheared water, 4 sterile water

BACs assembled and checked for contaminant sequences XSCQB – 3 seqs (2- Lotus japonicus, 1- Homo sapiens) XSCQC – 0 XSCQD – 0 XSCQE – 0

#### Types of vectors used in genomic library construction at TIGR

- 1 pHOS2 plasmid is a modified pBR192. It is low copy number: 20 copies/ cell, used for small insert library (2-4 Kb) and medium insert library (10-12 Kb).
- 2 BACs (Bacteria Artificial Chromosome) are used for positional cloning, physical mapping and genomic sequencing of large DNAs. They can accomodate up to 350 Kb inserts. One copy/cell.
- 3 Fosmids are cosmids that have the F plasmid origin of replication and cos sites, the DNA should be cut into 40 Kb pieces and packaged into the phage head which is then transfected into E.coli to produce colonies. One copy/cell.
- 4 Linking libraries.



#### BstXI adaptor cloning system



#### Design Features of the BstXI Adaptor Cloning System





# pHOS Vector Features

• The sequencing primer sites immediately flank the cloning site to avoid excessive re-sequencing of vector DNA.

• PCR primer sites are located immediately outside of the Sequencing primer sites to allow PCR amplification for template preparation.

• The vector continues to have the regular, good vector features such as a gene for replication (origin of replication), antibiotic resistance and multiple cloning sites where we can cut the circular plasmid and insert desired features.



# **Template Production**

- Transform clones into *E. coli*
- Plate and pick colonies to isolate clones
- Replicate and isolate plasmid for sequencing
- 384 well HT plate processing







#### **Alkaline Lysis 384-well DNA Isolation Method**





Average concentration 30 - 40 ng/μL



#### **Template Production Laboratory**

Current Capacity: 22,000,000 plasmids/year Expansion Capacity: 54,000,000 plasmids/year



# Sequence Production

- Sequencing Reactions and Precipitations
- Samples are loaded into the capillaries on a sequencing machine
- Data Collection
- Data Analysis







#### **Sequence Production Laboratory**

Current Capacity: 40,000,000 sequences/year Expansion Capacity: 100,000,000 sequences/year



### Capillary array view





# Sample Electropherogram





#### **Basecalling & Quality Assignments**



Phred & TraceTuner

- •Read DNA sequencer traces
- •Call bases
- •Assign base quality values
- •Write basecalls and quality values to output files.

Lowering the error rate, averaging 40% - 50% fewer errors than ABI software independent of position in read, machine running conditions, or sequencing chemistry



THE INSTITUTE FOR GENOMIC RESEARCH

# What are *phred* quality values?

The quality value q assigned to a base call is defined as:

$$q = -10 \times \log_{10}(p)$$

where *p* is the estimated error probability for that base-call.



# OR

# A base-call having a probability of 1/1000 of being incorrect is assigned a quality value of 30. Probability Quality Value 1/100 20 1/10 10





TIGR THE INSTITUTE FOR GENOMIC RESEARCH

# QC Tools

- 1. Web based document control system including SOP's and Process Control Forms
- 2. Receipt testing and lot control of sequencing reagents
- 3. QC audits of reagents and protocols used by production teams
- 4. Record of operators performing each step of the process
- 5. Equipment calibration and preventive maintenance program
- 6. QC/R&D group and QC representative on each team



# Assembling the fragments





### Celera Assembler

- creates high confidence "uniquely assembleable contigs" = unitigs
- marks those that appear repetitive w.r.t. arrival-rate statistics (surrogates and degenerates)
- uses insert (clone-mate) information to build contigs & mark ambiguous unitigs (surrogates)





overlap - region of similarity between regionsoverhang - un-aligned ends of the sequences

The assembler screens merges based on:

- length of overlap
- % identity in overlap region
- maximum overhang size.



### Arrival rate statistics

8x coverage, 800bp reads - a read every 100bp

200 300 100 400 ()

800

A-stat = log (p(single copy) / p(two-copy))  $p(single copy) = ((pF/G)^k/k!) exp(-pF/G)$  $p(two-copy) = ((2pF/G)^k/k!) exp(-2pF/G)$  p/k - average arrival rate

F - # fragments G - genome size

ITUTE FOR GENOMIC RESEARCH

#### Forward-reverse constraints

- The sequenced ends are facing towards each other
- The distance between the two fragments is known (within certain experimental error)







THE INSTITUTE FOR GENOMIC RESEARCH

#### Degenerate contigs

• Too deep

• Too few reads

#### (NOT IN SCAFFOLDS)



# TIGR Assembler Greedy

- Build a rough map of fragment overlaps
- Pick the largest scoring overlap
- Merge the two fragments
- Repeat until no more merges can be done





# Scaffolding

• Given a set of <u>non-overlapping</u> contigs <u>order and orient</u> them along a chromosome





#### Clone-mates



# Linking information



THE INSTITUTE FOR GENOMIC RESEARCH

## Grouping the contigs









**sequencing gap** - we know the order and orientation of the contigs and have at least one clone spanning the gap

**physical gap** - no information known about the adjacent contigs, nor about the DNA spanning the gap



### Scaffolder output




# Problems with the data

- Incorrect sizing of inserts
  - cut from gel sizing is subjective
  - error increases with size
- Chimeras (ends belong to different inserts)
  - biological reasons (esp. for large sized inserts)
  - sample tracking (human error)
- Software must handle a certain error rate.



## Ambiguous scaffold







### BAC-by-BAC Sequencing and Closure

#### Why?

- Aids complete closure by avoiding genomic repeat problems
- More streamlined collaboration with other sequencing and closure centers
- Allows for data release and annotation of completed sequence throughout the project

Why Not?

- Complete closure not required
- Small genome size
- Time/cost concerns



### BAC-by-BAC Sequencing and Closure

- Generate BAC-end sequence data
- Choose and sequence "seed" BACs
- Select minimally overlapping "tiling" BACs for contig extension



Select BACs based on BAC-end overlaps, physical mapping data, etc.



## Life of a BAC at TIGR



## BAC Shotgun Sequencing

BAC DNA





## Genome Closure/Finishing



# Why Completeness is Important

- Improves characterization of genome features
  - Gene order, replication origins
- Better comparative genomics
  - Genome duplications, inversions
- Determination of presence and absence of particular genes and features is less subjective
- Missing sequence might be important (e.g., centromere)
- Allows researchers to focus on biology not sequencing
- Facilitates large scale correlation studies
- Controls for contamination



# What Is Closure?

- Obtaining sequence that was not obtained during random sequencing which resulted in:
  - Sequencing Gaps
  - Physical Ends
- Confirming the integrity of assemblies
  - Repeats and misassemblies
  - Verification of Clone Coverage
- Confirming the base sequence of the consensus
  - Editing
  - Verification of Sequence Coverage



### **Levels of Genome Closure**



## **TIGR Finishing Criteria**

A genome is considered finished and ready if it satisfies the following criteria:

- 1. A continuous consensus of DNA sequence
- 2. No ambiguous consensus basepairs
- 3. At least 2X sequence coverage over the entire genome.
  - a) both strands of one clone are sequenced
  - b) two different clones sequenced
  - c) Same clone sequenced with dual chemistry
- 4. At least 2X clone coverage over the entire genome
- 5. Complete confidence in all repetitive areas



# Causes for gaps

- Non-random shotgun library
  - Toxicity of genes or promoters in E. coli
  - Genomic DNA difficult to clone (capsular polysaccharides)
  - Unstable regions (low complexity)
- Sequencing problems
  - Hard stops
    - Secondary structures
    - Very high or low GC content
    - Small unit tandem repeats
  - Loss of signal
    - Homopolymeric tracts
    - Very high or low GC content

THE INSTITUTE FOR GENOMIC RESEARCH



# Sequencing Gaps



To close Sequencing gaps:

- A. Resequence short reads at contig ends
- B. Perform sequencing reactions:
  - 1. Design primers at contig ends
  - 2. Do sequencing reactions with:
    - Insert 1 + primer 1
    - Insert 2 + primer 1

Insert 1 + primer 2 Insert 2 + primer 2



# **Physical Ends**

No known links exist; there is no template available for sequencing.



To link and close Physical ends use PCR:

- 1. Design primers at contig ends
- 2. If PCR products are obtained, sequence PCR products completely



# **Repetitive Areas**

- Repetitive areas are regions of high similarity within the genome/BAC.
- Sequences in these areas may be misassembled by the Assembler.
- Verification of the sequence of repetitive areas:
  - A. Identify potential repetitive areas, using *repeatFinder* and other tools.
  - B. Classify repeats based on length, copy number, % similarity, structure and complexity.
  - C. If repeats are misassembled, transpose spanning clones or obtain PCR products and sequence to verify assembly.
- Misassemblies can occur in two ways:



The base sequence of the repeat may be misassembled:



• Repeat is spanned by linking clones

• Underlying sequences are unlinked clones, with mates in other assemblies (branched clones) or mates that did not sequence (dead end clones).

• Results in incorrect consensus sequence

ITUTE FOR GENOMIC RESEARCH







# Identifying repetitive areas

Repetitive areas are identified by the program *repeatFinder* <u>ftp://ftp.tigr.org/pub/software/repeatFinder/</u>

- Analyzes sets of contigs for repeats greater than 50 bp and then groups these repeats into classes
- *repeatFinder* calls *REPuter*

[Copyright© University of Bielefeld, Germany <u>www.genomes.de</u>; Used with permission]



## Sequence editor

In this example there is an obvious discrepancy between the base calls of several of the underlying clones in this region.





# Assembly Viewer

- Several clones appear to be misassembled
- There are both obvious size violations and orientation issues

I HOUR DOORTAINS		
	RPT2D RPT2	C RPT2A
218	6 4371	6557
L34 (4): 15324	TBEFH17 (1) TBE	EFL20 (4): 15324 TBEFQ18 (4): 15324
<mark>651 (4</mark> ): 15324	TBEFP95 (1)	TBEFQ76 (1)
Q10 (4): 15324	TBEFN10 (1) T	➡ BEFK59_(4): 15324 TBEFD10 (2)
E89 (4): 15324 TBEFA89 (4): 1	5324 TBEFJ11 (1)	TBEFA56 (4): 15324
B86 (4): 15324 TBFFD38 (4): 15	324 TRFFD®à (1)	TBEEE75 (2)
IDEFUIZ (4): 15324		IBERP77 (3)
TBEF025 (4): 15324	TBEFD73 (1)	TBEF023 (4): 15324 TBEFN71 (4): 1532
TBEFN34 (4): 15324	TBEF083 (1)	TBEFK79 (4): 15324 TBEFD01 (
TBEFP12 (4): 15324	TBEF091 (1)	TBEFQ46 (4): 15324 TBEFB07 (4): 15
TBEFE	34 (1)	TBEFQ07 (1)
TBEFL 09 (4): 15324	TBEFN18 (1)	TBEFP47 (4): 15324
TBEFE49 (4): 15324TBE	FJ91 (4): 15324	TBEF017 (1)
TBEFJ54 (4): 15324	TBEFG94 (1)	TBEFB53 (4): 15324
TBEF003 (4): 1	5324 TBEFC87 (4): 15324 TBEFH26 (4):	15324 TBEFL53 (4): 15324 TBEFP52 (4):
TBEFH65 (4):	15324 TBEFB37 (4): 15324 TBEFG4	19 (4); 15324 TBEFH30 (4); 15324
TBEFD18	(4): 15324 TB	EFG57 (1) TBEFI76 (4):
	TBEFE15 (4): 15324	TBEF650 (4): 15324 TBEFH38 (4): 15324
	TBEFD29 (4): 15324	4 TBEFB31 (4): 15324 TBEFL86 (4): 15324
		TBEFE64 (4): 15324
		TBEFE69 (4): 15324
		TRFF.127 (4): 15334



# Resulting contig after repeat resolution





THE INSTITUTE FOR GENOMIC RESEARCH

## AMOS Assembly Investigator



THE INSTITUTE FOR GENOMIC RESEARCH

**Current Contig Position** 

## **Collapsed Repeat**



THE INSTITUTE FOR GENOMIC RESEARCH

# **Resolved Repeat**

- Unique flank order is correct
  - Use linking information across the repeat (large insert clones or PCR)
- Consensus sequence is correct
  - Use linked clones that have one mate in the repeat and the other anchored in unique sequence
  - Transposon mediated libraries



# Sequence Validation: Editing

- Finishers look at all areas where the consensus sequence is below a certain quality threshold.
- Electropherograms are used to determine which base calls may need to be changed.
- Only supporting sequences are edited, not the consensus.
- Editing provides an integrity check, may aid in assembly, identifying repeat areas and other unique problems.





# Sequence Validation: Sequence coverage



#### Sequence coverage rule:

Every base in an assembly must be covered by at least two sequences of high quality.

#### Why?

Validating sequence coverage provides a high degree of confidence in the consensus base calls.

# Sequence Coverage in Cloe

<ul> <li>Cloe</li> </ul>	- mtg2: as	m 1,153 [4	15,000 : 50	0,000]												
<u>File</u> <u>E</u> c	dit <u>V</u> iew	Windows	QA													
Assem	bly Seque	nce List	Screen	Screenshot Ctrl-Space												
Dirty?	Se	q name	Rev	/ersed?	Begin	End	Edi	Edit length	Phr	ed date						
	MIDVASATE	>					75	65.0	2004 04 21 2							
	······															
Bases	500	1000	1500	2000	2500	2000	2500	4000	45.00	500						
	500	1000	1500	2000	2300	5000	5500	4000	4300	5000						
18																
17																
16																
15																
100								ſ								
14									4 4							
13									U							
12																
11																
10								1								
0								Γ								
								ſ								
6																
7	П		1				Г	JU								
6																
5000			, ] [,													
4																
3																
-																
2																

TIGR THE INSTITUTE FOR GENOMIC RESEARCH

# Sequence Coverage in Cloe

- C	loe - mtg2: a	sm 1,153 [	45,000 : 50	,000]												
File	Edit View	Window	s QA													
Ass	embly Seque	ence List	Screen	Screenshot Ctrl-Space												
Dir	tv?l s	ea name	Rev	/ersed?	Begin	End	- 1	Edit length	Phred date							
		.D			500		75	650	2004 04	21.20%						
Bas	es Arrow	s Cover	age													
	500	1000	1500	2000	2500	3000	3500	4000	4500	5000						
<u> </u>	•		_	·		i —										
<u> </u>			←	<u> </u>		•	-			←						
		•		270		<u>→</u>	<u>+</u>									
								÷ —	•							
+		-	+					,	·							
	·							, <b></b>		-						
								+								
									<del>`</del>							
										<u> </u>						
								•		<b></b>						
									•							
									-							
										•—						
										1.2						



# Assembly Verification: Clone Coverage

#### **<u>Clone Coverage Rule:</u>**

Every base must be spanned by two valid subclones.

A valid subclone is a pair of forward/reverse mates that are:

- Oriented correctly (pointing toward each other)
- Separated by an approximate length of sequence determined by the expected template size for the library.

### <u>Why?</u>

Validating clone coverage provides a high degree of confidence that the contig was assembled correctly.





• Blue: 2x clone coverage

 $\bullet$ 



### Verify Closure by Fingerprint Comparison

adder	32115	3598	38K24	45m8	46d18	52e11	-53d3	-53f22	-55111	adder	56114	62016	62p11	-64d17	71a20	74a13	93b1	-101c1	adder	
1 kb	mth2-	mth2-	mth2-	mth2-	mth2-	mth2	mth2	mth2	mth2	1 Kb	mth2	mth2	mth2-	mth2	mth2	mth2	mth2	mth2	1 Kb	
				Π	Π	T					Π	Π	Π	T		T				
		H				1			П		I	1	П					П		
		1				11	11					Ш	1	Ш						
12		-	_	-	-		-	=	-	=		-	#	=	=		۲		=	
2	-	-	-	=	=		=	=	=	-	-	=	-	I	-	-	=	=	-	
2	=	=	=	=			-		Π	-	=	Ξ	=		-	-	-	=	-	
-	-	-		-	=	=			1	_	_	Ξ		_	H		-		_	
	-		=	-	=	=		-	-		=	-		-	=	H	Ξ	-		
-	_			-	=	-	*	-		-		=		E	-	-	Π	=	-	
	1	-		=	=		1.1	Ľ			Ξ	-	Ц	-	H	H	Ξ	=		
-	-			-	-	-	-			_	=	=	=	-	H		-		_	
		-					E	1				H	_		1		-	-		
-	-	H		-	1	-	-	-	-		-	=	-	Еŝ	5.3			***	_	
		600				-	-	1						-	818		-	****		
						-	-							-	2:1		1	1		
		-					-	11	1			7	_						1	
						-		-	11				-							
																	100			

Italion@italion:~	X
<u>File E</u> dit <u>V</u> iew <u>T</u> erminal <u>G</u> o <u>H</u> elp	
With 1 enzymes: HINDIII	
October 5, 2004 14:59	
HindIII A'AGCT_T	
Cuts at: 0 1007 6197 9864 13017 19240 20329 20656 22722 Size: 1007 5190 3667 3153 6223 1089 327 2066	
Cuts at: 22722 26311 28424 33210 36482 43233 55712 57131 63240 Size: 3589 2113 4786 3272 6751 12479 1419 6109	
Cuts at: 63240 63756 63958 64175 65281 65817 71575 74772 78980 Size: 516 202 217 1106 536 5758 3197 4208	
Cuts at: 78980 81215 81242 85588 90065 90833 92037 103047 104202 Size: 2235 27 4346 4477 768 1204 11010 1155	
Cuts at: 104202 105088 107394 110136 111828 115437 115681 124021 124074 Size: 886 2306 2742 1692 3609 244 8340 53	
Fragments arranged by size:	
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	
Enzymes that do cut:	
HindIII	
Enzymes that do not cut:	
NONE [/usr/local/projects/MTG2/closure.dir/CLOSED.dir/13026_H222C10/FINAL_CHECK/pseud omol/DIGEST_CHECK] a15:~>]	-

THE INSTITUTE FOR GENOMIC RESEARCH

## Finished BAC Ready for Annotation!

**BAC DNA** 





# EXTRA SLIDES START HERE...


#### Graph/top navigation

machine catalyst db gel\_run\_pn rxn\_run\_pn rxn\_lot primer\_lot



#### Data navigation

machine catalyst db gel\_run\_pn rxn\_run\_pn rxn\_lot primer\_lot bad gels group gels

#### data for Jul 12 2003 to Jul 12 2003

db	Total Good	Total Lanes	% Good	VEIR	% VEIR	Average Edit Length	Average Quality
<u>gssea</u>	12375	14782	83.7 -4.5	698	4.7	827 <mark>-3</mark>	40.2 +0.6
mag	403	479	84.1 -4.1	30	6.3	853 +23	39.4 <mark>-0.2</mark>
pcg	572	672	85.1 <mark>-3.1</mark>	20	3.0	850 +20	38.7 <mark>-0.9</mark>
<u>smag</u>	2877	3157	91.1 +2.9	4	0.1	846 +16	40.1 +0.5
<u>tbg</u>	1112	1248	89.1 +0.9	68	5.4	874 +44	39.9 +0.3
<u>tcrg</u>	7021	7488	93.8 +5.6	28	0.4	829 -1	39.0 <mark>-0.6</mark>
<u>tseest</u>	2986	3552	84.1 -4.1	90	2.5	781 -49	35.7 <mark>-3.9</mark>
ttg	1905	2112	90.2 +2.0	2	0.1	839 +9	40.4 +0.8
tvg	2342	2592	90.4 +2.2	2	0.1	833 +3	40.5 +0.9
<u>zmg</u>	6602	7200	91.7 +3.5	203	2.8	829 <b>- 1</b>	39.9 +0.3
Totals 53409 (avg: 60551 (avg: 2967) 3364)		88.2	1467	2.4	830	39.6	

THE INSTITUTE FOR GENOMIC RESEARCH



TIGR THE INSTITUTE FOR GENOMIC RESEARCH







Assembly table: consensus w/ and w/o gaps Asmbl\_link table: asmbl\_id seq\_name seq\_lend, seq\_rend asm\_lend, asm\_rend lsequence Sequence table: sequence end5, end3



### Schematic View of 3 Typical Scaffolds



# Types of repeats

- Short Tandem Repeats
- Large Multi Unit Tandem Repeats
- Multi Class Repeats
- Inverted Repeats





## Short Tandem Repeats





# Large Multi Unit Tandem Repeats



## **Inverted Repeats**



### Construction of Transposon-Mediated Libraries

- Transposons are genetic elements that provide mobile primer sites within the target DNA for sequencing
- The GPS-1 Genome Priming System (New England Biolabs E7100S) enables transposon insertions by using TnsABC\* transposase
- Only one random insertion occurs per target DNA





### Use of Transposon-mediated libraries:

•Make a library using a transposon kit on a spanning clone.

- •Sequence the clones using the transposon primers
- •Assemble the sequences

Physical map of the transposon inserts:





### **Assembly of Transposon-Mediated Clone Walks**



- forward and reverse clone mates

- tandem repeat region
- transposon insertion site



## Closure Challenges: Sequencing Through Secondary Structures



TAGT6C66CC66CC6CCC66GCAAT66CC6TATC6CTAC66A66666AA6C66CCCCTCAA666C6CC6CTTCC6CC6C6666CA666TCA6CC6ACCTT6C6TT6

AAGGGCGCCGCTTCCGCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG

#### 

TTAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGG

**FTAGTGAGGCCGG** 

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGGAAGCGGCGCCCTCAAGGGCgCcgcTTCcGCCGCGCGGGGCAGGGTCAGCCGACCTTGCGTTG

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAG

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGG

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGGA

FTAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATC

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTA

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGG

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGGAAG

ITAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGGAAGCGGCGCCCT

TTAGTGCGGCCGGCGCGCCTGGGCAATGGCCGTATCGCTACGGAGGGGAAGCgGCGCCCTCAAGGGCgccGCTTCcGCCGCGGGGCAGGGTCAGCCGACCTTGCGTTG

CGCTTCCGCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG GCTTCCGCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG

CTTC-GCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG

TCCGCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG

CCGCCGCGCGGGCAGGGTCAGCCGACCTTGCGTTG

CGGGCAGGGTCAGCCGACCTTGCGTTG





## Homopolymeric tracts

THE INSTITUTE FOR GENOMIC RESEARCH

#### Cloe - stg: asm 1,995 [104,000 : 105,000]

<u>File Edit View Windows QA</u>

Accomply Coquence List

Assembly sequence List											
Dirty?	Seq name	Reversed?	Begin	End	Edit length	Phred date	Edit person	modify date			
	SDCR987TF	<b></b>	-531	252	782	2004-03-13 12:26:16	jessicah	2004-06-10 14:22:44 📤			
	SDCR205TR		-462	263	724	2004-02-29 07:31:37					
	SDCRE55T1914SBC	<b>~</b>	-152	417	564	2004-07-14 09:10:48	jessicah	2004-07-23 10:08:07.1			
	SDCRE55T1911SEC		92	295	203	2004-07-14 09:10:48	jessicah	2004-07-23 10:08:07			
	SDCRB67T1911SEC		106	296	191	2004-07-14 09:10:47	jessicah	2004-07-23 10:08:07			
	SDCR394TF		174	300	127	2004-02-29 06:24:07	jessicah	2004-07-23 10:08:07 🧾			
	SDCRB67T1914SBC	<b></b>	222	416	190	2004-07-14 09:10:47	jessicah	2004-07-23 10:08:07			
	SDCRE62TF	<b>~</b>	284	445	162	2004-03-13 12:39:31	jessicah	2004-07-23 10:08:07			
	SDCR614TF	<b>~</b>	287	703	417	2004-03-01 04:13:19	jessicah	2004-07-23 10:08:07 💌			



- - -

# Solutions

- Apply different sequencing chemistries
  - Big-Dye terminator (default)
  - Dye-primer
  - dGTP mix (GC rich regions)
- Denature structures Additives
  - Betaine
  - DMSO
- Break structure
  - Restriction digest
  - Transposon insertion
  - Micro-libraries



# Automated Finishing at TIGR

- Closure is a **feature driven** activity
- Feature Classification System
  - Identify/Classify Features
    - SEQ gaps
    - PHY gaps / scaffold ends
    - Low coverage
    - Repeat features
- Feature Reaction Design Strategy
  - Bin by Reaction Type and Feature Size
    - Primer walks
    - PCR
    - Transposon Bombing

THE INSTITUTE FOR GENOMIC RESEARCH

# Automated Finishing (cont'd)

- Laboratory Process Integration
  - Order reactions from the lab
  - Group reactions by lab process
  - Integrate changing laboratory automation
  - Monitor results
- Reaction Product Resolution Strategy
  - Resolve reaction products back to features & contigs: conduct targeted assembly.
  - Evaluate resolution criteria for feature type.

THE INSTITUTE FOR GENOMIC RESEARCH

### System Architecture



# **BAC** Project Challenges

- Dead Ends
- Filling small gaps between adjacent supercontigs
- bp mismatches in BAC sequence overlaps
- Project management/organization



### Dead Ends

A Dead End occurs when a BAC at the end of a sequence contig has no BAC-end hits or FPC BACs to extend the contig further

- Design primers and PCR small unique region near end of BAC
- Hybridize to BAC filters to identify candidate BACs
- End sequence and fingerprint candidate BACs to select minimal tiling BAC
- Sequence and Close new tiling BAC

### Filling small gaps in a chromosome



THE INSTITUTE FOR GENOMIC RESEARCH

# bp Mismatches In BAC Overlaps

### Why?

- Basecalling error in sequencing
- Point mutation in BAC
- Misassembly of repeats near BAC end

### What to do?

- Examine sequence quality in each BAC
- PCR region using genomic DNA template
- Sequence PCR to verify correct basecall

TIGR THE INSTITUTE FOR GENOMIC RESEARCH

### **Project Organization**

### • Database to store various BAC data:

- Tiling, overlap info
- Mapping, FPC info
- Sequencing and Closure data
- Collaborator's BAC data
- Web-based BAC status tracking and closure task tracking tools

