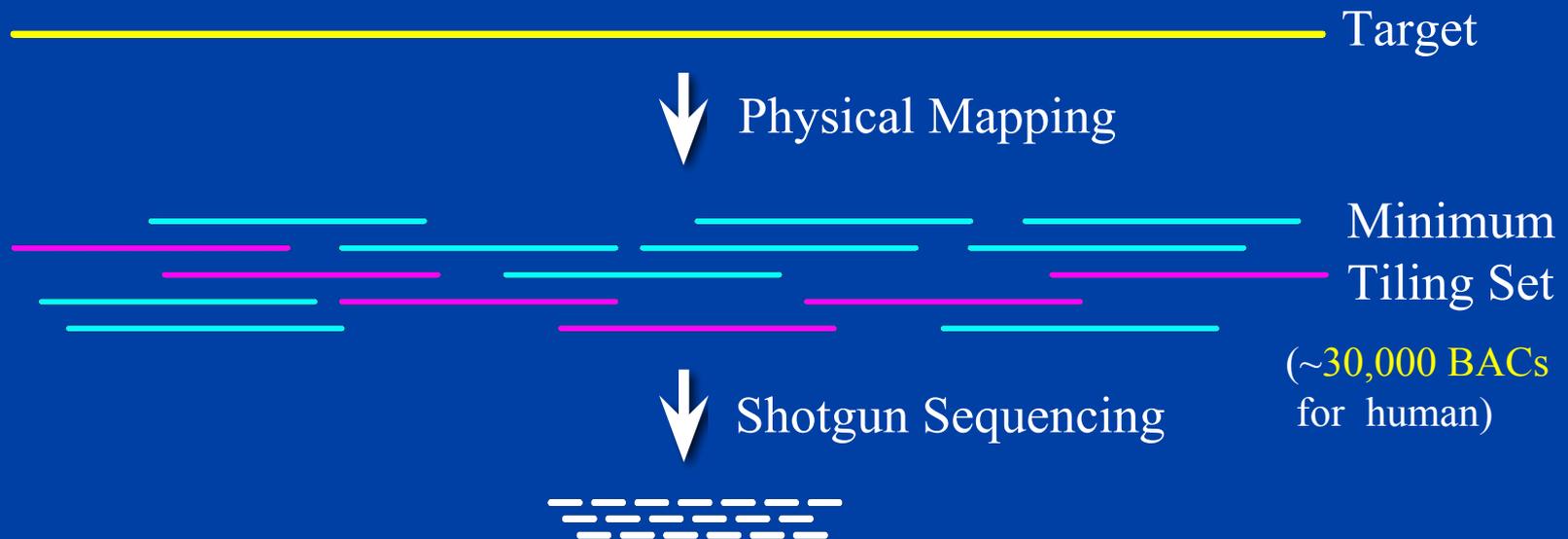




Whole Genome Sequencing Approaches

u Hierarchical HGP Approach:



- 2 separate processes
- maps very hard to complete, libraries unstable
- must make shotgun library of each BAC
- + infrastructure is already developed
- + quality of outcome is known



Sequencing Factory

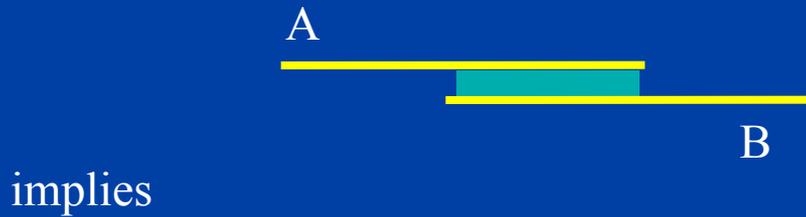
- **300 ABI 3700 DNA Sequencers installed**
- **50 Production Staff**
- **40 Support Staff (R&D, QC/QA, Service)**
- **20,000 sq. ft. of wet lab**
- **20,000 sq. ft. of sequencing space**
- **800 tons of A/C (160,000 cfm)**
- **4,000 amps electrical service**

The DNA is loaded into automated sequencers. Celera's automated sequencers run 24-7 and have the ability to decipher more than 100 million letters of genetic code per day - the equivalent of 3 percent of the entire human genetic code every day.

The sequencers create an image of the DNA samples being decoded. The four letters of the genetic code -- A, C, T, G -- each are assigned a color.



True vs. Repeat-Induced Overlaps



TRUE



OR

REPEAT-INDUCED



Assembly Pipeline

167:41 cpu hrs. for Dros

8:37
86:25
38:29
4:12
5:44+4:21+19:53
(~25)

8:37

86:25

38:29

4:12

5:44+4:21+19:53

(~25)

Screeners

Mask heterochromatin and ribo-DNA,
Tag known interspersed repeats.

Overlapper

Find all overlaps ≥ 40 bp allowing 6% mismatch.
(1000X Blast)

Unitiger

ASSEMBLER CORE:

- Compute all consistent sub-assemblies = **unitigs**
- Identify those that cover unique DNA = **U-unitigs**
- Scaffold U-unitigs with confirmed shorts & longs
- Then with BAC ends
- Fill repeat gaps with:
 - I. Doubly anchored mates
 - II. O-path confirmed singly-anchored mates
 - III. Greedy path completion using QVs

Scaffolder

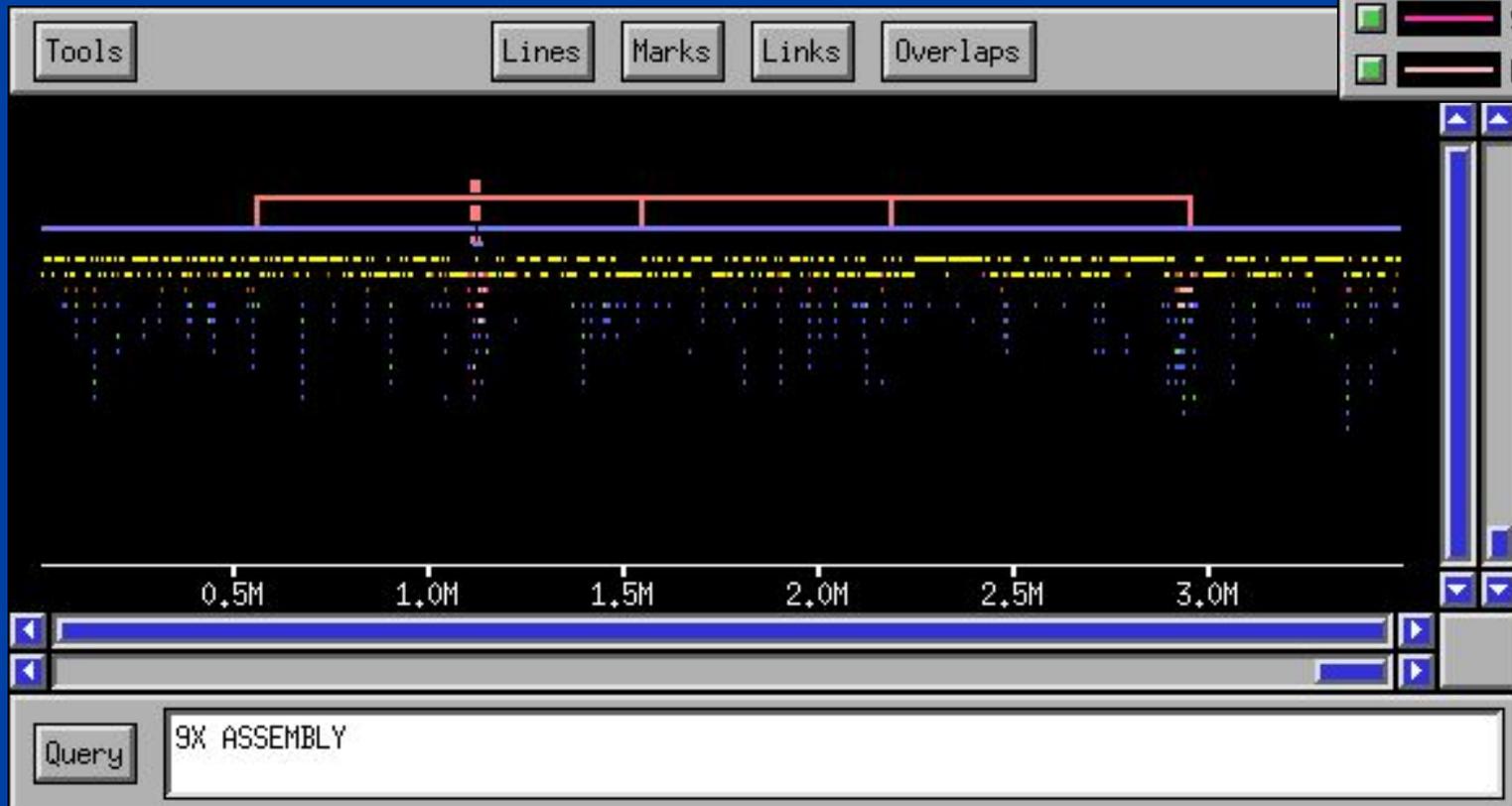
Repeat Rez I, II,
III

Consensus

Bayesian "SNP" consensus using quality values.
Occurs throughout assembler core.

Assembly Progression (Macro View)

- Read
- Valid
- Invalid
- U-Unitig
- Contig
- Rock
- Stone
- Pebble



- Contig 8677013
- Contig 8677012
- Contig 8677011
- Contig 8677010
- Contig 8677009
- Contig 8677008
- Contig 8677007
- Contig 8677006
- Contig 8677005
- Contig 8677004
- Contig 8677003
- Contig 8677002
- Contig 8677001
- Contig 8677000
- Contig 8676999
- Contig 8676998
- Contig 8676997
- Contig 8676996
- Contig 8679936
- Contig 8677295

Unique identifiers

start and stop site predictions

Splice site predictions

Homology based exon predictions

computational exon predictions

TRANSCRIPT Feature 8000

Property	Value
url	http://ohio.c...
celera a...	CT32501
gene ac...	null
transcript...	null
feature id	8000
feature ty...	TRAN
display p...	High
is comp...	true
assigned...	promote de
assigned...	11/30/1999
approved...	null
approved...	null
is child	false
feature s...	null
orientation	Reverse
contig id	500000609
contig b	64039

Tracking information

Consensus gene structure (both strands)

