# Entering the era of mega-genomics

Michael Schatz

March 20, 2012 JGI Users Meeting

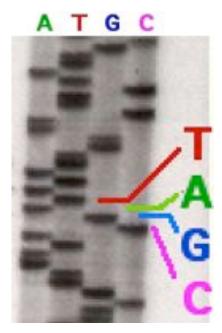




### **Outline**

- I. Milestones in genomics
  - I. Sanger to nanopore
  - 2. 21st Century Mega-Genomics
- 2. Applications of mega-genomics
  - I. Cloud-scale genomics for bioenergy
  - 2. Single molecule sequencing & assembly
  - 3. De novo mutations in autism

# Advances in Sequencing: Zeroth, First, Second Generation



1970s: 0th Gen

Radioactive Chain Termination

5000bp / week



1980s-1990s: 1st Gen

Automated Capillary Sequencing

384kbp / day



2000s: 2<sup>nd</sup> Gen

Pyrosequencing, SOLiD Sequencing-by-Synthesis

IGbp+ / day

# Advances in Sequencing: Now Generation Sequencing



Illumina HiSeq 2000 Sequencing by Synthesis

>60Gbp / day



**Ion Proton**Postlight Sequencing

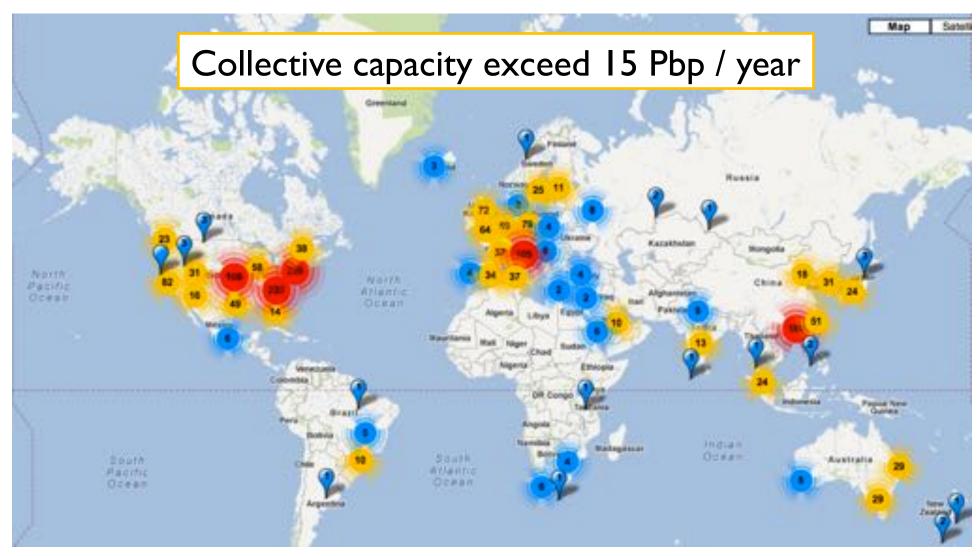
>100Gbp / day



Oxford Nanopore
Nanopore sensing

Many GB / day

# Sequencing Centers



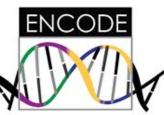
Next Generation Genomics: World Map of High-throughput Sequencers http://pathogenomics.bham.ac.uk/hts/

### The rise of mega-genomics



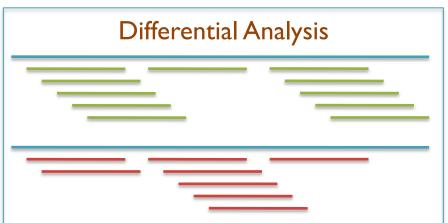


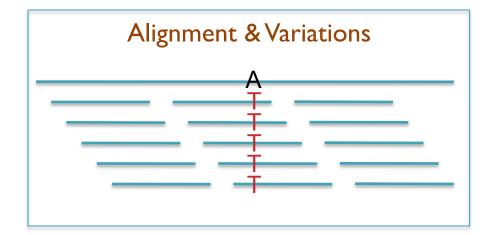


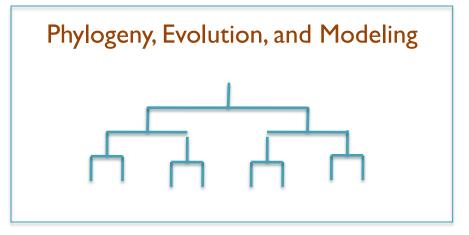










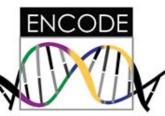


### The rise of mega-genomics



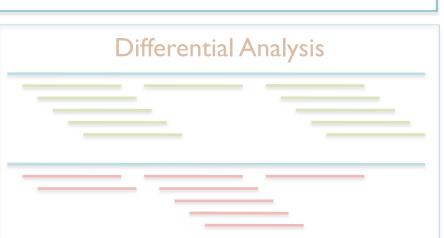


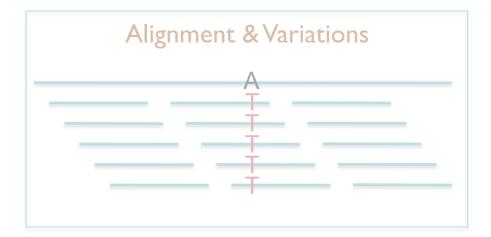


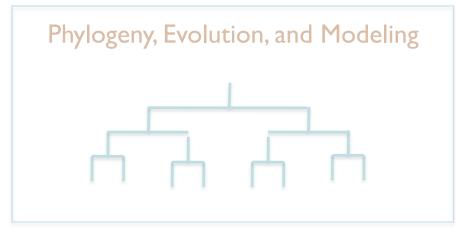




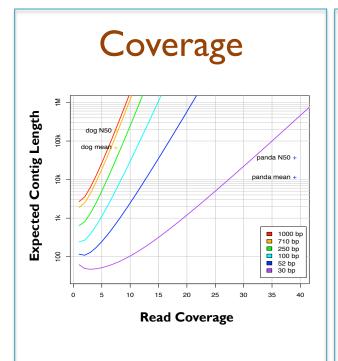






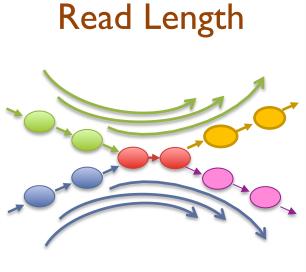


# Ingredients for a good assembly



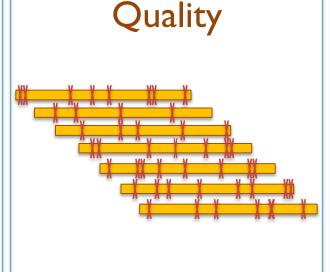
#### High coverage is required

- Oversample the genome to ensure every base is sequenced with long overlaps between reads
- Biased coverage will also fragment assembly



### Reads & mates must be longer than the repeats

- Short reads will have false overlaps forming hairball assembly graphs
- With long enough reads, assemble entire chromosomes into contigs



#### Errors obscure overlaps

- Reads are assembled by finding kmers shared in pair of reads
- High error rate requires very short seeds, increasing complexity and forming assembly hairballs

Current challenges in de novo plant genome sequencing and assembly Schatz MC, Witkowski, McCombie, WR (2012) In Press.

### Hybrid Sequencing



**Illumina**Sequencing by Synthesis

High throughput (60Gbp/day)
High accuracy (~99%)
Short reads (~100bp)



**Pacific Biosciences**SMRT Sequencing

Lower throughput (600Mbp/day)

Lower accuracy (~85%)

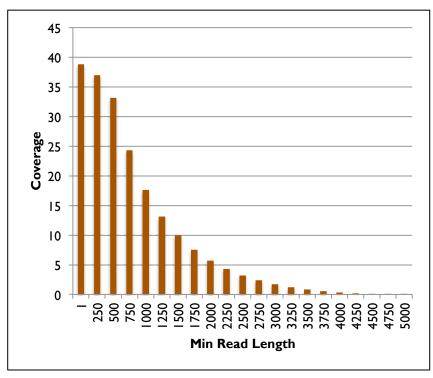
Long reads (10kbp+)

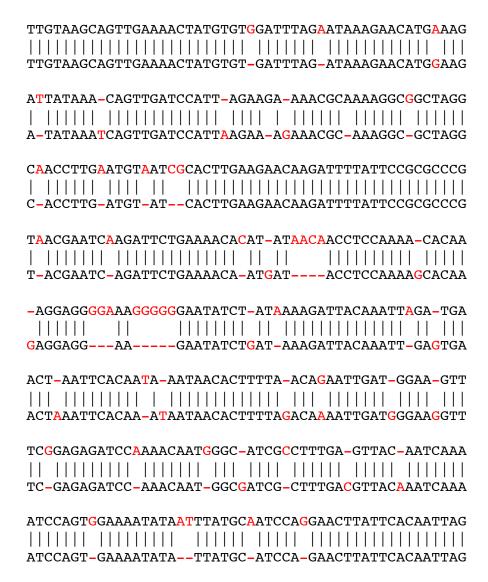
### **SMRT** Sequencing Data

## Yeast (12 Mbp genome)

65 SMRT cells 734,151 reads after filtering Mean: 642.3 +/- 587.3

Median: 553 Max: 8,495





Sample of 100k reads aligned with BLASR requiring > 100bp alignment Average overall accuracy: 83.7%, 11.5% insertions, 3.4% deletions, 1.4% mismatch

### PacBio Error Correction

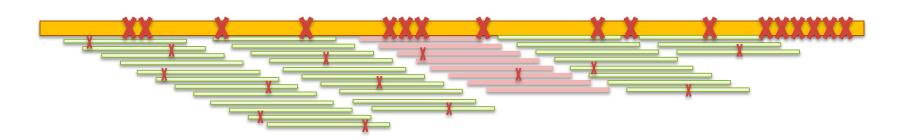
http://wgs-assembler.sf.net

### I. Correction Pipeline

- I. Map short reads (SR) to long reads (LR)
- 2. Trim LRs at coverage gaps
- 3. Compute consensus for each LR

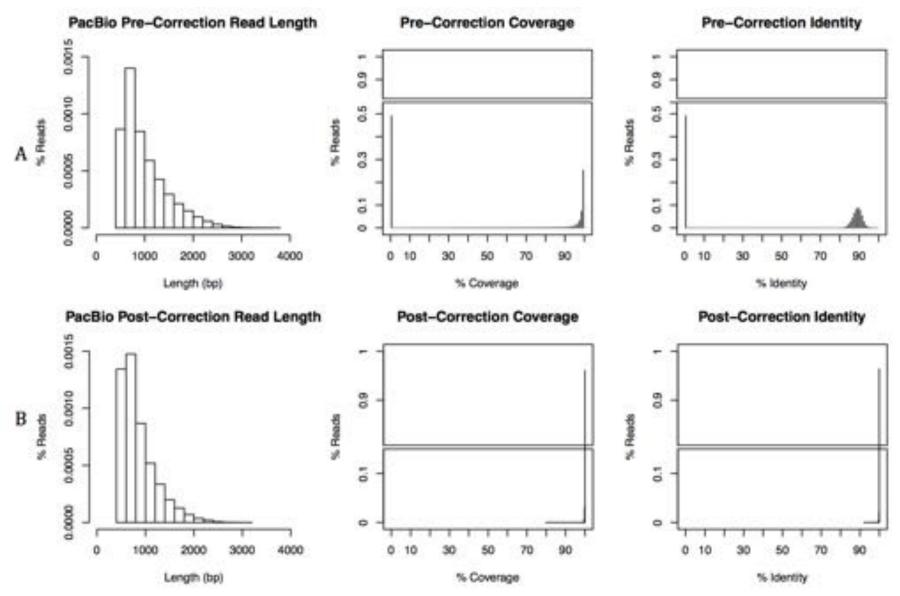


2. Error corrected reads can be easily assembled, aligned



Hybrid error correction and de novo assembly of single-molecule sequencing reads. Koren, S, Schatz, MC, Walenz, BP, Martin, J, Howard, J, Ganapathy, G, Wang, Z, Rasko, DA, McCombie, WR, Jarvis, ED, Phillippy, AM. (2012) *Under Review* 

### **Error Correction Results**

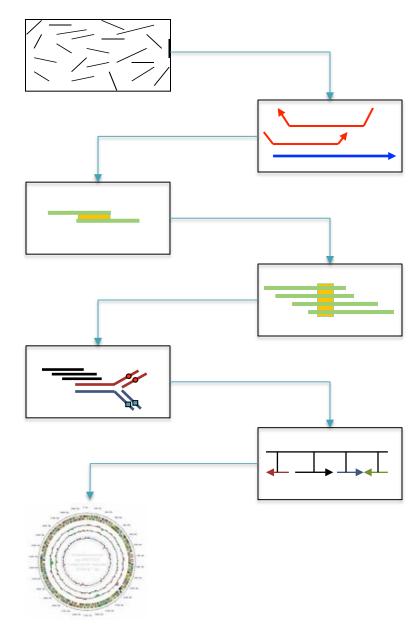


Correction results of 20x PacBio coverage of E. coli K12 corrected using 50x Illumina

### Celera Assembler

#### http://wgs-assembler.sf.net

- I. Pre-overlap
  - Consistency checks
- 2. Trimming
  - Quality trimming & partial overlaps
- 3. Compute Overlaps
  - Find high quality overlaps
- 4. Error Correction
  - Evaluate difference in context of overlapping reads
- 5. Unitigging
  - Merge consistent reads
- 6. Scaffolding
  - Bundle mates, Order & Orient
- 7. Finalize Data
  - Build final consensus sequences



# **SMRT-Assembly Results**











Organism	Technology	Reference bp	Assembly bp	# Contigs	Max Contig Length	NS	
Lambda NEB3011	Illumina 100X 200bp	48 502	48 492	3	48 492 / 48 492	48 492 / 48 492 (100%) *	
(median: 727 max: 3 280)	PacBio PBcR 25X		48 440	.1	48 444 / 48 444	48 444 / 48 440 (100%)	
E.col/ K12	Illumina 100X 500bp	4 639 675	4 462 836	61	221 615 / 221 553	100 338 / 83 037 (82.76%)	
(median; 747 max; 3 068)	PacBio PBcR 18X		4 465 533	77	239 058 / 238 224	71 479 / 68 309 (95.57%) *	
	Both 18X PacBio PBcR + Illumina 50X 500bp		4 576 046	65	238 272 / 238 224	93 048 / 89 431 (96.11%) *	
E. coli C227-11	PacBio CCS 50X	5 504 407	4917717	76	249 515	100 322	
(median: 1 217 max: 14 901)	PacBio 25X PBcR (corrected by 25X CCS)		5 207 946	80	357.234	98 774	
	Both PacBio PBcR 25X + CCS 25X		5 269 158	39	647 362	227 302	
	PacBio 50X PBcR (corrected by 50X CCS)		5 445 466	35	1 076 027	376 443	
	Both PacBio PBcR 50X + CCS 25X		5 453 458	33	1 167 060	527 198	
	Manually Corrected ALLORA Assembly <sup>8</sup>		5 452 251	23	653 382	402 041	
S. cereviniae S228c	Illumina 100X 300bp	12 157 105	11 034 156	192	266 528 / 227 714	73 871 / 49 254 (66.68%) *	
(median: 674 max: 5 994)	PacBio PBcR 13X		11 110 420	224	224 478 / 217 704	62 898 / 54 633 (86.86%) *	
	Both PacBio PBcR 13X + Illumina 50X 300bp		11 286 932	177	262 846 / 260 794	82 543 / 59 792 (72.44%) *	
Melopsittacus undulatus	Illumina 194X (220/500/800 paired-end 2/5/10Kb mate-pairs)	1.23 Gbp	1 023 532 850	24 181	1 050 202	47 383	
	454 15.4X (FLX + FLX Plus + 3/8/20Kbp paired-ends)		999 168 029	16.574	751 729	75 178	
(median 997, max 13 079)	454 15.4X + PacBio PBcR 3.75X		1 071 356 415	15 081	1 238 843	99 573	

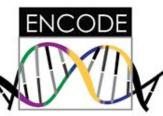
Hybrid assembly results using error corrected PacBio reads
Meets or beats Illumina-only or 454-only assembly in every case
\*\*\* Also useful for transcriptome, repeat, and other analysis \*\*\*

### The rise of mega-genomics



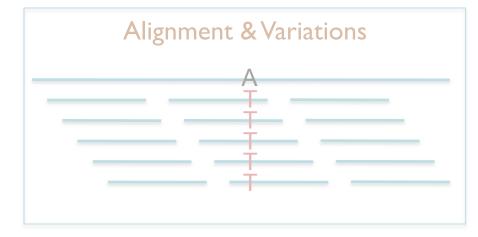




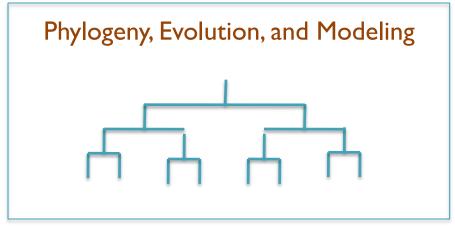






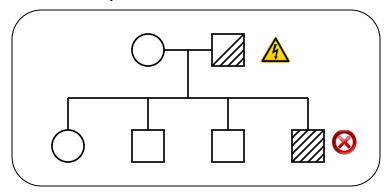






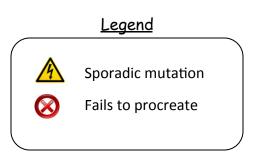
### Unified Model of Autism

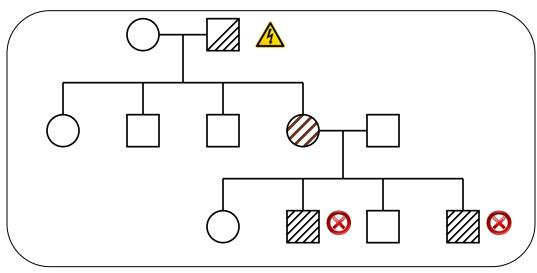
#### Sporadic Autism



De novo mutations of high penetrance contributes to autism, especially in low risk families with no history of autism.

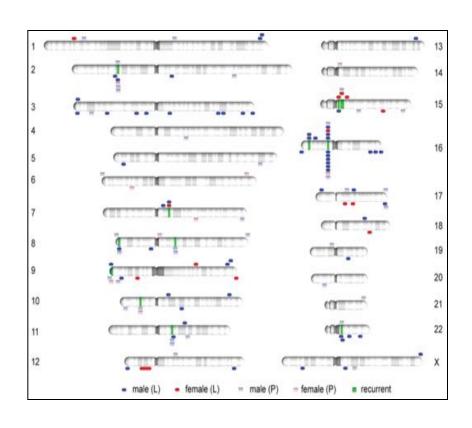
#### Familial Autism





A unified genetic theory for sporadic and inherited autism Zhao et al. (2007) PNAS. 104(31)12831-12836.

### Autism and de novo CNVs



Analysis of Simons Simplex Collection

- CGH arrays of 510 family quads
- 94 total de novo CNVs discovered

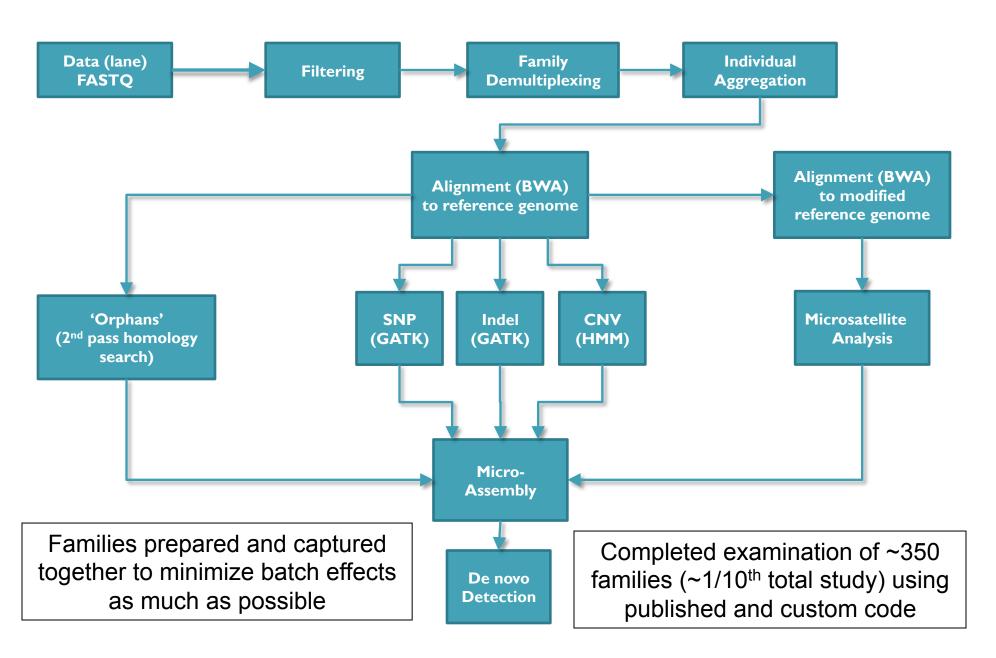
De novo CNVs are more common in autistic children

- 4:1 ratio in autistic kids relative to their non-autistic siblings
- Some recurrence at genes related to other psychiatric conditions

Counts of De Novo Events			Children with De Novo Events			Frequency in Children		
Combined	Del	Dup	Combined	Del	Dup	Combined	Del	Dup
75	46	29	68	44	27	7.9%	5.1%	3.1%
19	9	10	17	8	9	2.0%	0.9%	1.0%
	75	75 46	75 46 29	75 46 29 68	75 46 29 68 44	75 46 29 68 44 27	75 46 29 68 44 27 7.9%	75 46 29 68 44 27 7.9% 5.1%

Rare de novo and transmitted copy-number variation in autism spectrum disorders. Levy et al. (2011) Neuron. 70:886-897.

# Exome Sequencing Pipeline



### Scalpel: Haplotype Microassembly

G. Narzisi, D. Levy, I. Iossifov, J. Kendall, M. Wigler, M. Schatz

- Use assembly techniques to identify complex variations from short reads
  - Improved power to find indels
  - Trace candidate haplotypes sequences as paths through assembly graphs





```
Father: ...TCAGAACAGCTGGATGAGATCTTAGCCAACTACCAGGAGATTGTCTTTGCCCGGA...

Mother: ...TCAGAACAGCTGGATGAGATCTTAGCCAACTACCAGGAGATTGTCTTTGCCCGGA...

Sib: ...TCAGAACAGCTGGATGAGATCTTAGCCAACTACCAGGAGATTGTCTTTGCCCGGA...

Aut(1): ...TCAGAACAGCTGGATGAGATCTTAGCCAACTACCAGGAGATTGTCTTTGCCCGGA...

Aut(2): ...TCAGAACAGCTGGATGAGATCTTAGCCAACTACCAGGAGATTGTCTTTTGCCCGGA...
```

6bp heterozygous indel at chr13:25280526 ATP12A

### De novo Genetics of Autism

- In 343 family quads so far, we see significant enrichment in de novo *likely gene killers* in the autistic kids
  - Overall rate basically 1:1 (432:396)
  - 2:1 enrichment in nonsense mutations
  - 2:1 enrichment in frameshift indels
  - 4:1 enrichment in splice-site mutations
- Observe strong overlap with the 842 genes known to be associated with fragile X mental retardation.
  - These genes relate to neuron and brain development
  - Suggest these genes are under strong purifying selection and we hypothesize particularly dosage sensitive

Exome sequence analysis of simplex families with children on the autism spectrum lossifov et al. (2012) Under review

# Mega-Genomics Challenges



# The foundations of genomics will continue to be observation, experimentation, and interpretation

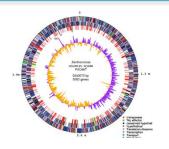
- Technology will continue to push the frontier
- Measurements will be made digitally over large populations,
   at extremely high resolution, and for diverse applications

#### Rise in Quantitative and Computational Demands

- 1. Experimental design: selection, collection & metadata
- 2. Observation: measurement, storage, transfer, computation
- 3. Integration: multiple samples, assays, analyses
- 4. Discovery: visualizing, interpreting, modeling

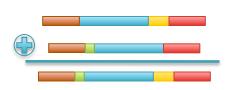
Ultimately limited by the human capacity to execute extremely complex experiments and interpret results

### Acknowledgements



James Gurtowski Matthew Titmus

Ware Lab KBase Members



Paul Baranay (CSHL/ND) Scott Emrich (ND) Steven Salzberg (JHU)

Adam Phillippy (NBACC) Sergey Koren (NBACC)



Giuseppe Narzisi Mitch Bekritsky

Ivan Iossifov Wigler Lab







# Thank You!

http://schatzlab.cshl.edu/apply@mike\_schatz / #JGIUM7

