

A Disease-Centric Draft Map of the Human Interactome

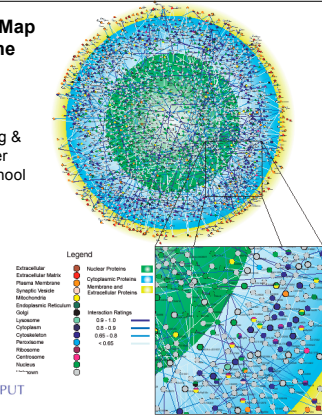
Joel S. Bader

Dept. of Biomedical Engineering & High-Throughput Biology Center
Johns Hopkins University & School of Medicine
joel.bader@jhu.edu
www.jhubiomed.org

AICHE
Nov. 1, 2005



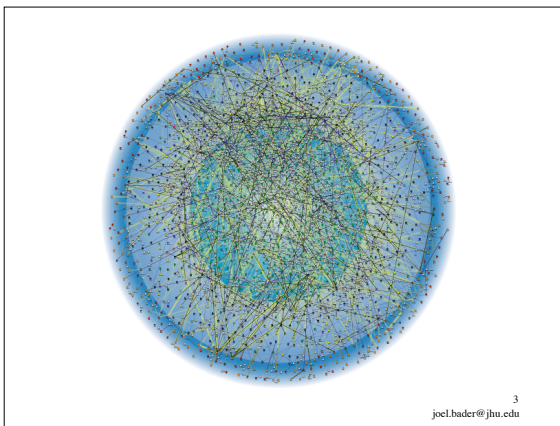
HIT HIGH THROUGHPUT BIOLOGY



Topics

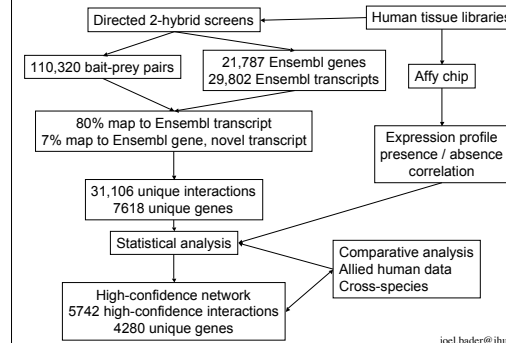
- Large-scale human protein interaction map
- Statistical framework for estimating coverage (false-negative rate) of sampled networks
- A new sequencing machine

2
joel.bader@jhu.edu



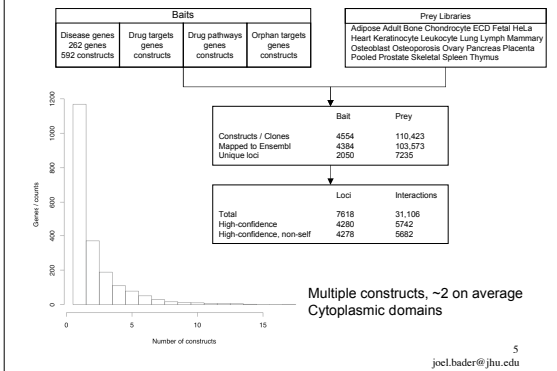
3
joel.bader@jhu.edu

Progress in the human network



4
joel.bader@jhu.edu

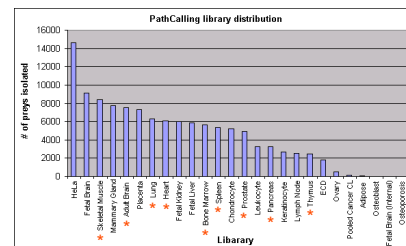
Bait design



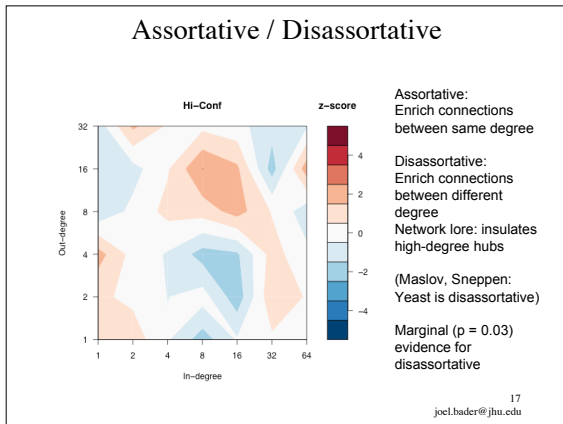
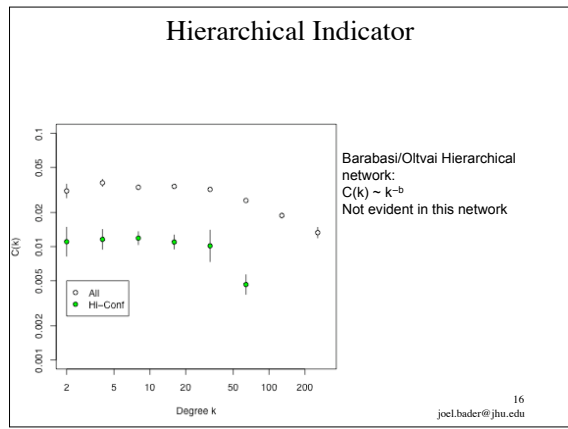
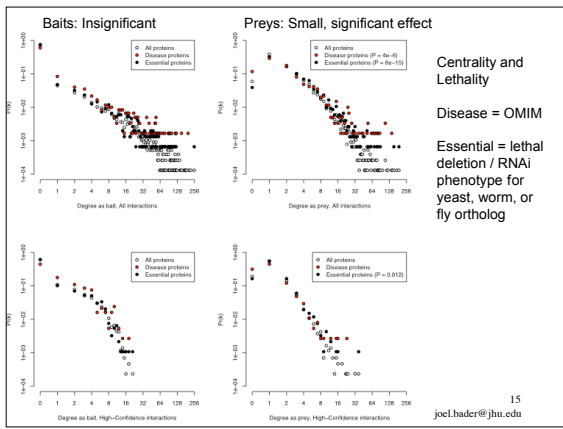
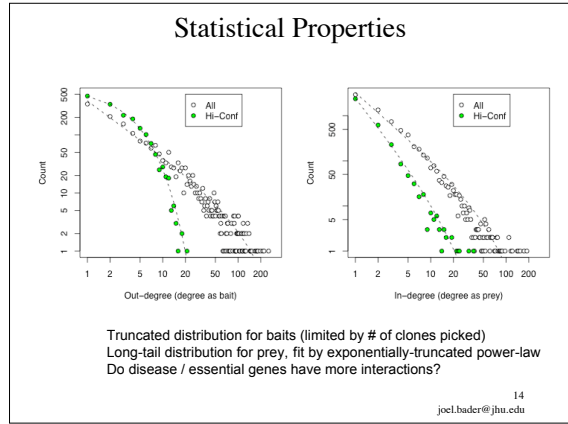
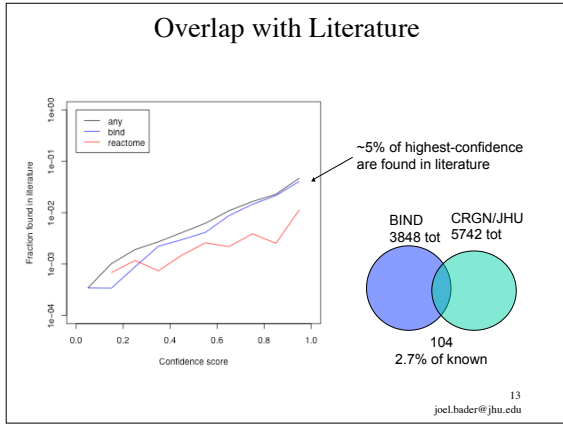
5
joel.bader@jhu.edu

Library distribution

- 26 PathCalling bait libraries used
- 9 of these have expression data available for same tissue
 - Affy U95A-E human chip hybridization
 - Data publicly available from NCBI GEO database
 - These 9 tissues cover 42.7% of PathCalling preys identified



* Normal tissue used for Affy U95 hybridization, public data set available from NCBI GEO database



Significant motifs

Motif	Observed	Expected	Z-Score	P-value	Motif	Observed	Expected	Z-Score	P-value
	104	20.9±4.6	18.2	3×10^{-24}		1103	440±130	5.1	1.6×10^{-7}
	557	76±1	44	0		344	2.2±1.8	190	0
	97	24.3±7.1	10.2	1×10^{-24}		56	2.6±3.1	17	4×10^{-55}
	17,333	3100±450	31.7	0		504	289±28	7.6	1.5×10^{-14}
	1521	680±230	3.6	0.00016		285	61±14	16	6.4×10^{-35}
						52	9.3±6.2	6.9	2.6×10^{-12}
						30	0.8±1.0	30.	0

18
joel.bader@jhu.edu

Model Variables

Variable	Cards	Y2H
k	Number in deck	# of interaction partners
k_i	Number seen i times	# preys seen i times
$\tilde{k} = \sum_{i>0} k_i$	Number seen at least once	# unique preys
$n = \sum_i i k_i$	Number sampled	# clones picked
\hat{k}	Estimate of k	Estimate of k
\tilde{k} / \hat{k}	Fraction of deck seen	Coverage of network

25
joel.bader@jhu.edu

Bayesian Formula

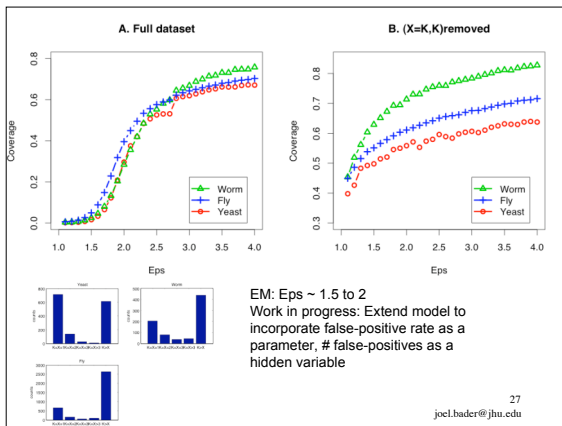
$$\Pr(\{k_i\} | k, n) = \frac{k!}{\prod_i k_i!} \cdot \frac{n!}{\prod_{i>0} (k_i!)^i} \cdot k^{-n}$$

$$\Pr(k | \{k_i\}, n) = \frac{[k! / (k - \tilde{k})!] k^{-n} \Pr(k)}{\sum_{k'} [k'! / (k' - \tilde{k})!] k'^{-n} \Pr(k')}$$

Pr(k) = Prior estimate for k, power-law for scale-free network
 Possible estimators:
 MAP = k that maximizes Pr(k|{k_i}, n) + 95% Confidence Interval
 PME = $\sum_i k \Pr(k|{k_i}, n)$
 Problems when n = k_i (no partner observed twice): k = infinity
 Choice for prior: Pr(k) = k^{-b} / zeta(b) regularizes when b > 1

- (1) Can use EM to find best exponent b considering k as hidden variable
- (2) Can look at subset of baits that detect at least 1 prey twice

26
joel.bader@jhu.edu



DNA Sequencing Breakthrough

nature

ARTICLES

Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies¹, Michael Egholm¹, William E. Altman¹, Said Attiya¹, Joel S. Bader², Lisa A. Bemben¹, Jan Berka¹, Michael S. Braverman¹, Yi-Ju Chen¹, Zhoutao Chen¹, Scott B. Dewell¹, Lei Du¹, Joseph M. Fierro¹, Xavier V. Gomes¹, Brian C. Godwin¹, Wen He¹, Scott Helgeson¹, Chun He Ho¹, Gerard P. Irzyk¹, Shivshanker C. Jando¹, Maria L. L. Akenquist¹, Thomas P. Jarvis¹, Kshama B. Jirage¹, Jong-Bum Kim¹, James R. Knight¹, Janna R. Lanza¹, John H. Leamon¹, Steven M. Lefkowitz¹, Ming Lei¹, Jing Li¹, Kenton L. Lohman¹, Hong Lu¹, Vinod B. Makhijani¹, Keith E. McDade¹, Michael P. McKenna¹, Eugene W. Myers¹, Elizabeth Nickerson¹, John R. Noble¹, Ramona Plant¹, Bernard P. Puc¹, Michael T. Ronan¹, George T. Roth¹, Gary J. Sarkis¹, Jan Fredrik Simons¹, John W. Simpson¹, Muthuvyan Srinivasan¹, Karrie R. Tartaro¹, Alexander Tomas², Kari A. Vogt¹, Greg A. Volkmer¹, Shally H. Wang¹, Yong Wang¹, Michael P. Weiner¹, Pengguang Yu¹, Richard F. Begley¹ & Jonathan M. Rothberg¹

¹MSL Life Sciences Corp., 20 Commercial Street, Branford, Connecticut 06405, USA. ²University of California, Berkeley, California 94720, USA. ³Laboratory of Microbiology, The Rockefeller University, New York, New York 10021, USA. ⁴The Rothberg Institute for Childhood Diseases, 530 Whitefield Street, Guilford, Connecticut 06437, USA. *These authors contributed equally to this work.

28
joel.bader@jhu.edu

The founding patent

US006274320B1

(12) **United States Patent** (10) Patent No.: **US 6,274,320 B1**
 Rothberg et al. (45) Date of Patent: **Aug. 14, 2001**

(54) **METHOD OF SEQUENCING A NUCLEIC ACID** OTHER PUBLICATIONS

(75) Inventors: **Jonathan M. Rothberg, Guilford, Joel S. Bader, New Haven, both of CT (US)**

(73) Assignee: **CuratGen Corporation, New Haven, CT (US)**

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.

(21) Appl. No.: **09/398,833**

(22) Filed: **Sep. 16, 1999**

(51) Int. Cl.: **C12Q 1/68; C12P 19/34; C12M 1/34**

(52) U.S. Cl.: **435/6; 435/91.2; 435/287.2**

(58) Field of Search: **435/6; 287.2; 91.2**

Baner et al., "Signal amplification of padlock probes by rolling circle replication." *Nucleic Acids Research* 29(22): 5073-5078 (1998).

Barshop et al., "Luminescent immobilized enzyme test systems for inorganic pyrophosphate: Assays using firefly luciferase and acetylcholinesterase-monomucleotide adenylyl transferase or adenosine-5'-triphosphate sulfurylase." *Analytical Biochemistry* 197: 266-272 (1991).

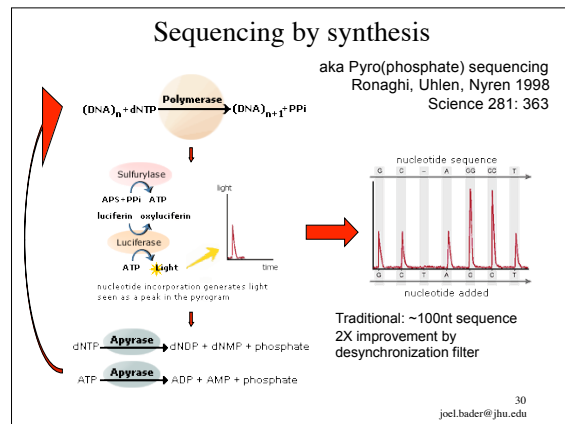
Brands et al., "Slow rate of phosphodiester bond formation accounts for the strong bias that Taq DNA polymerase shows against 2',3'-dideoxynucleotide terminators." *Biochemistry* 55: 2189-2200 (1990).

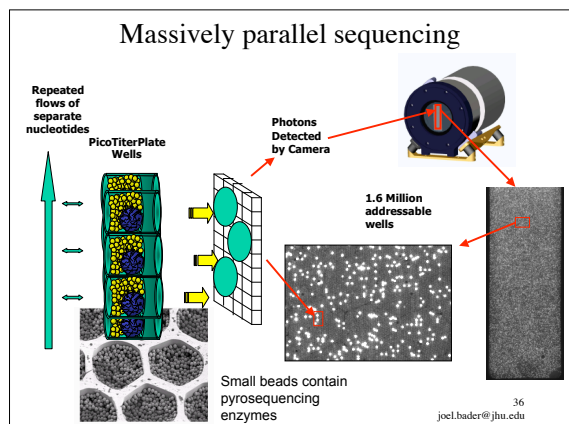
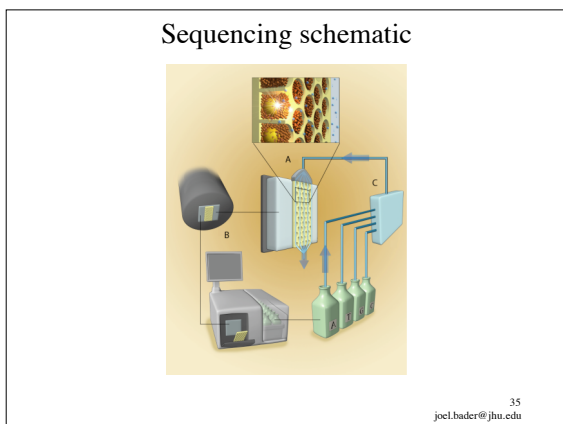
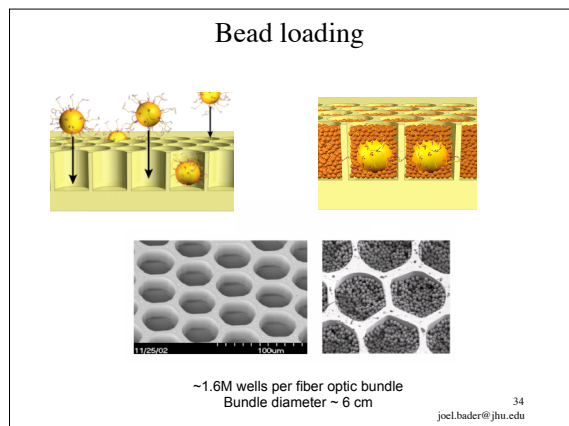
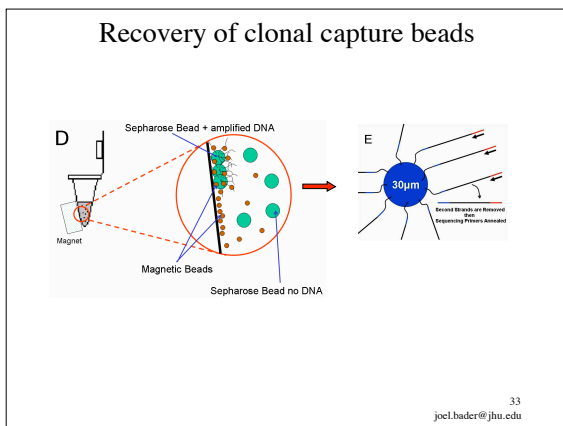
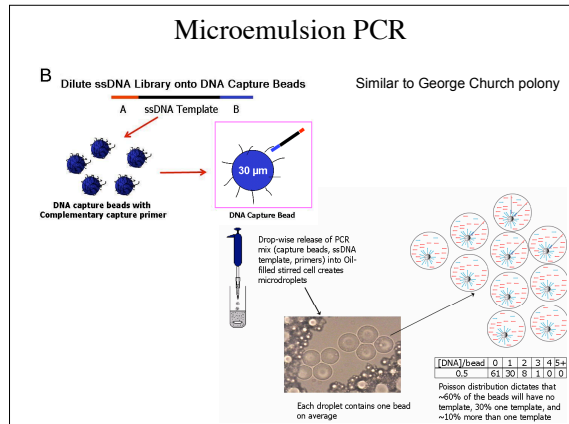
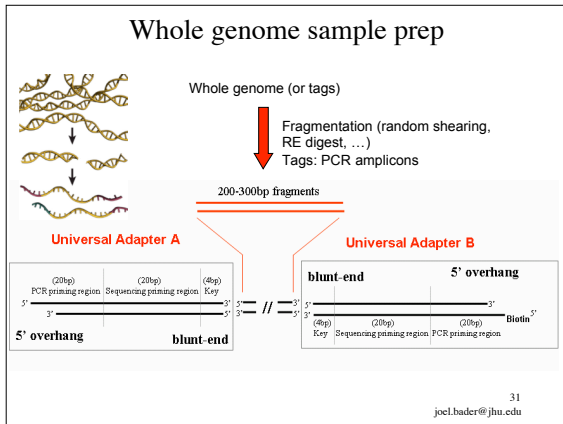
Bronk et al., "Combined imaging and chemical sensing using a single optical imaging fiber." *Anal. Chem.* 67: 2750-2757 (1995).

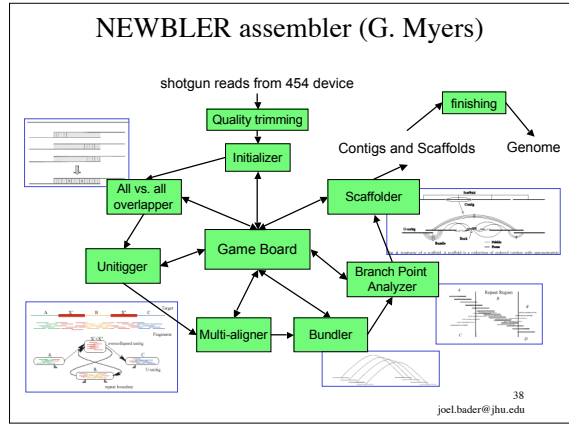
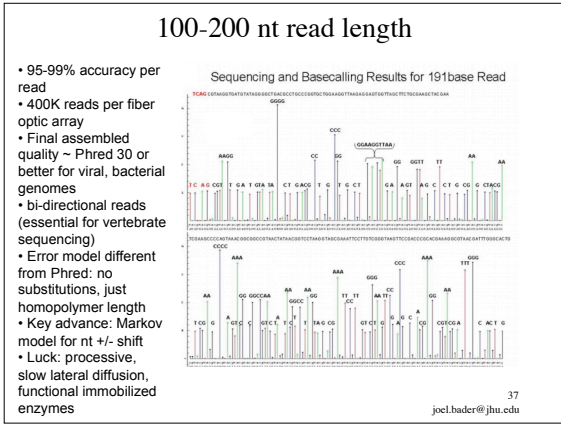
Burns et al., "An Integrated Nanoliter DNA Analysis Device." *Science* 282: 484-487 (1998).

Chan and Nie, "Quantum dot bioconjugates for ultrasensitive nonisotopic detection." *Science* 281: 2016-2018 (1998).

Chee et al., "Accessing Genetic Information with High-Density DNA Arrays." *Science* 274(5287).



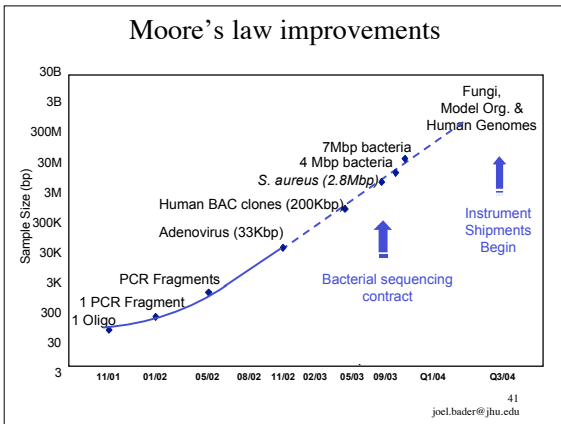
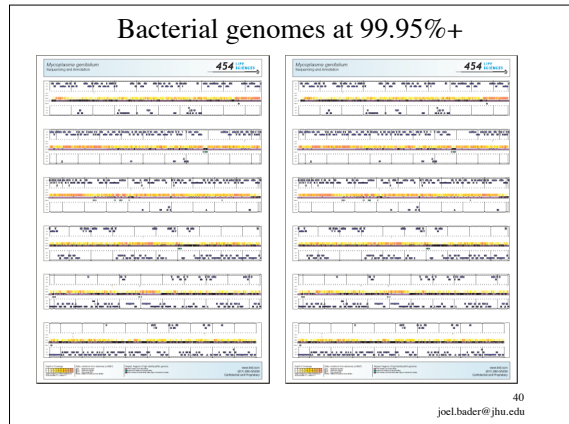




High-quality consensus sequence

	Adenovirus	M. Gon.	S. aureus
Sequencing Run Summary			
Size of fiber optic slide	30x60 mm	60x60 mm	60x60 mm
Run Time / Number of cycles	244min/42	244 min/42	244 min/42
High Quality reads	21,308	20,306	272,602
Average read length	105 bp	115 bp	102 bp
Total number of bases sequenced	2,300,957	33,659,471	28,508,698
Individual Reads			
Reads mapped to single locations in the reference genome	20,645	221,401	255,240
Insertion error rate (Total number of overcalls/Total number of bases aligned)	2.20%	2.42%	1.31%
Deletion error rate (Total number of undercalls/Total number of bases aligned)	1.72%	2.08%	3.08%
Substitution error rate	0.09%	0.70%	0.13%
Consensus Sequence			
Number of bases aligned	2,166,283	25,743,641	27,456,492
Average overcalling	48x	37x	7x
Genome coverage	33,367	568,836	2,506,618
Consensus accuracy	99.97%	99.1%	99.2%
Bases overcalled	0	0.001%	0.003%
Bases undercalled	0.006%	(6 bp)	(86 bp)
	(2 bp)	(5 bp)	(124 bp)

40
joel.bader@jhu.edu



- ### Future of 454 device
- Performance
 - ~100x higher throughput than ABI
 - 10 runs x 400K wells x 100 nt = 400 Mbp / day
 - ABI: 10 runs x 100 wells x 1000 nt = 1 Mbp / day
 - 100x improvements through feature size reductions, bi-directional reads
 - With another 10x, \$1000 human
 - Limitations
 - 100 to 200 nt read length
 - Errors at homonucleotide run, not at ends of sequence
 - Buy one today!
- 42
joel.bader@jhu.edu

