
This copy is for your personal, non-commercial use only.

If you wish to distribute this article to others, you can order high-quality copies for your colleagues, clients, or customers by [clicking here](#).

Permission to republish or repurpose articles or portions of articles can be obtained by following the guidelines [here](#).

The following resources related to this article are available online at www.sciencemag.org (this information is current as of August 8, 2011):

Updated information and services, including high-resolution figures, can be found in the online version of this article at:

<http://www.sciencemag.org/content/327/5961/78.full.html>

Supporting Online Material can be found at:

<http://www.sciencemag.org/content/suppl/2009/11/05/1181498.DC1.html>

This article **cites 32 articles**, 10 of which can be accessed free:

<http://www.sciencemag.org/content/327/5961/78.full.html#ref-list-1>

This article has been **cited by** 21 article(s) on the ISI Web of Science

This article has been **cited by** 31 articles hosted by HighWire Press; see:

<http://www.sciencemag.org/content/327/5961/78.full.html#related-urls>

This article appears in the following **subject collections**:

Genetics

<http://www.sciencemag.org/cgi/collection/genetics>

26. L. F. Robinson *et al.*, *Science* **310**, 1469 (2005).
27. L. F. Robinson, T. van de Fliet, *Geology* **37**, 195 (2009).
28. I. M. Gil, L. D. Keigwin, F. G. Abrantes, *Paleoceanography*, **24**, PA4101 (2009).
29. H. K. Evans, I. R. Hall, *Geochem. Geophys. Geosyst.* **9**, Q03023 (2008).
30. S. K. Praetorius, J. F. McManus, D. W. Oppo, W. B. Curry, *Nat. Geosci.* **1**, 449 (2008).
31. C. Wunsch, *Quat. Sci. Rev.* **22**, 371 (2003).
32. A. Ganopolski, S. Ramstorf, *Nature* **409**, 153 (2001).
33. C. Hillaire-Marcel, A. de Vernal, G. Bilodeau, A. J. Weaver, *Nature* **410**, 1073 (2001).
34. Supported by UK NERC grant NE/D002206/1, RG43765/LBZG021, "Assessing the role of ocean circulation in rapid climate change through the novel integration of high-resolution proxy records," and by NSF and the Comer Science and Education Fund (J.M.). We thank A. Galy, D. Hodell, and L. Skinner for discussion; E. Martin for advice on fish teeth cleaning; and J. Clegg and A. Scrivner for help in the lab.

Supporting Online Material

www.sciencemag.org/cgi/content/full/327/5961/75/DC1
SOM Text
Figs. S1 and S2
Tables S1 to S6
References

22 June 2009; accepted 26 October 2009
10.1126/science.1178068

Human Genome Sequencing Using Unchained Base Reads on Self-Assembling DNA Nanoarrays

Radoje Drmanac,^{1*} Andrew B. Sparks,^{1†} Matthew J. Callow,^{1†} Aaron L. Halpern,^{1†} Norman L. Burns,^{1†} Bahram G. Kermani,^{1†} Paolo Carnevali,^{1†} Igor Nazarenko,^{1†} Geoffrey B. Nilsen,^{1†} George Yeung,^{1†} Fredrik Dahl,^{1†‡} Andres Fernandez,^{1†} Bryan Staker,^{1†} Krishna P. Pant,^{1†} Jonathan Baccash,¹ Adam P. Borcharding,¹ Anushka Brownley,¹ Ryan Cedeno,¹ Linsu Chen,¹ Dan Chernikoff,¹ Alex Cheung,¹ Razvan Chirita,¹ Benjamin Curson,¹ Jessica C. Ebert,¹ Coleen R. Hacker,¹ Robert Hartlage,¹ Brian Hauser,¹ Steve Huang,¹ Yuan Jiang,¹ Vitali Karpinchy,¹ Mark Koenig,¹ Calvin Kong,¹ Tom Landers,¹ Catherine Le,¹ Jia Liu,¹ Celeste E. McBride,¹ Matt Morenzoni,¹ Robert E. Morey,^{1§} Karl Mutch,¹ Helena Perazich,¹ Kimberly Perry,¹ Brock A. Peters,¹ Joe Peterson,¹ Charit L. Pethiyagoda,¹ Kaliprasad Pothuraju,¹ Claudia Richter,¹ Abraham M. Rosenbaum,² Shaunak Roy,¹ Jay Shafto,¹ Uladzislau Sharanovich,¹ Karen W. Shannon,^{1||} Conrad G. Sheppy,¹ Michel Sun,¹ Joseph V. Thakuria,² Anne Tran,¹ Dylan Vu,¹ Alexander Wait Zaranek,² Xiaodi Wu,³ Snezana Drmanac,¹ Arnold R. Oliphant,¹ William C. Banyai,¹ Bruce Martin,¹ Dennis G. Ballinger,^{1*} George M. Church,² Clifford A. Reid¹

Genome sequencing of large numbers of individuals promises to advance the understanding, treatment, and prevention of human diseases, among other applications. We describe a genome sequencing platform that achieves efficient imaging and low reagent consumption with combinatorial probe anchor ligation chemistry to independently assay each base from patterned nanoarrays of self-assembling DNA nanoballs. We sequenced three human genomes with this platform, generating an average of 45- to 87-fold coverage per genome and identifying 3.2 to 4.5 million sequence variants per genome. Validation of one genome data set demonstrates a sequence accuracy of about 1 false variant per 100 kilobases. The high accuracy, affordable cost of \$4400 for sequencing consumables, and scalability of this platform enable complete human genome sequencing for the detection of rare variants in large-scale genetic studies.

Genotyping technologies have enabled the routine assessment of common genetic variants at up to a million sites across the genome in thousands of individuals (1) and have increased our understanding of human genetic diversity and its biological and medical impact. Whole-genome sequencing costs have

dropped from the >\$100 million cost of the first human genomes (2, 3) to the point where individual labs have generated genome sequences in a matter of months for material costs of as low as \$48,000 (4–12) (table S5). Sequencing technologies, which use a variety of genomic microarray construction methodologies and sequencing chemistries (13–32), can determine human genetic diversity over an entire genome and identify common as well as rare single-nucleotide polymorphisms (SNPs), insertions, and deletions. Despite these advances, improvements are still needed to enable the cost-effective characterization of the many hundreds of genomes required for genetic studies of complex diseases and for personalized disease prevention, prognosis, and treatment.

We generated sequencing substrates [Fig. 1A and supporting online material (SOM)] by means of genomic DNA (gDNA) fragmentation and recursive cutting with type IIS restriction

enzymes and directional adapter insertion (Fig. 1B and fig. S1). The resulting circles were then replicated with *Phi*29 polymerase (RCR) (33). Using a controlled, synchronized synthesis, we obtained hundreds of tandem copies of the sequencing substrate in palindrome-promoted coils of single-stranded DNA, referred to as DNA nanoballs (DNBs) (Fig. 1C). DNBs were adsorbed onto photolithographically etched, surface-modified (SOM) 25- by 75-mm silicon substrates with grid-patterned arrays of ~300-nm spots for DNB binding (Fig. 1C). The use of patterned arrays increased DNA content per array and image information density relative to random genomic DNA arrays (6, 9, 11, 14, 28). High-accuracy combinatorial probe anchor ligation (cPAL) sequencing chemistry was then used to independently read up to 10 bases adjacent to each of eight anchor sites (Fig. 1D), resulting in a total of 31- to 35-base mate-paired reads (62 to 70 bases per DNB). cPAL is based on unchained hybridization and ligation technology (15, 27, 28, 31), previously used to read 6 to 7 bases from each of four adapter sites (26 total bases) (28), here extended using degenerate anchors to read up to 10 bases adjacent to each of the eight inserted adapter sites (Fig. 1D, right) with similar accuracy at all read positions (fig. S3). This increased read length is essential for human genome sequencing.

Cell lines derived from two individuals previously characterized by the HapMap Project (34), a Caucasian male of European descent (NA07022) and a Yoruban female (NA19240), were sequenced. NA19240 was selected to allow for a comparison of our sequence to the sequence of the same genome currently being assembled by the 1000 Genome Project. In addition, lymphoblast DNA from a Personal Genome Project Caucasian male sample, PGP1 (NA20431) was sequenced because substantial data are available for biological comparisons (35–37). Automated cluster analysis of the four-dimensional intensity data produced raw base reads and associated raw base scores (SOM).

We mapped these sequence reads to the human genome reference assembly with a custom alignment algorithm that accommodates our read structure (fig. S4), resulting in between 124 and 241 Gb mapped and an overall genome coverage of 45- to 87-fold per genome.

To assess representational biases during circle construction, we assayed genomic DNA and intermediate steps in the library construction process by quantitative polymerase chain reaction

¹Complete Genomics, Inc., 2071 Stierlin Court, Mountain View, CA 94043, USA. ²Department of Genetics, Harvard Medical School, Cambridge, MA 02115, USA. ³School of Medicine, Washington University, St. Louis, St. Louis, MO 63110, USA.

*To whom correspondence should be addressed. E-mail: rdrmanac@completegenomics.com (R.D.); dballinger@completegenomics.com (D.G.B.)

†These authors contributed equally to this work.

‡Present address: Ion Torrent Systems, San Francisco, CA 94158, USA.

§Present address: San Diego State University, San Diego, CA 92115, USA.

||Present address: Life Technologies, Carlsbad, CA 92008, USA.

(QPCR) (fig. S2). This and mapped coverage showed a substantial deviation from Poisson expectation with excesses of both high and low coverage regions (fig. S5), but only a few percent of bases have coverage insufficient for assembly (Table 1). Much of this coverage bias is accounted for by local GC content in NA07022, a bias that was significantly reduced by improved adapter ligation and PCR conditions in NA19240 (fig. S5); the fraction of the genome with less

than 15-fold coverage was accordingly reduced from 11% in NA07022 to 6.4% in NA19240, despite the latter having 25% less total coverage (Table 1).
Discordance with respect to the reference genome in uniquely mapping reads from NA07022 was 2.1% (range 1.4% to 3.3% per slide). However, considering only the highest scoring 85% of base calls reduced the raw read discordance to 0.47%, including about 0.1% of true variant positions.

Mapped reads were assembled into a best-fit, diploid sequence with a custom software suite employing both Bayesian and de Bruijn graph techniques (SOM). This process yielded diploid reference, variant, or no-calls at each genomic location with associated variant quality scores. Confident diploid calls were made for 86% to 95% of the reference genome (Table 1), approaching the 98% that can be reconstructed in simulations. The 2% that is not reconstructed

Fig. 1. Amplified DNA nanoarray platform. (A) Schematic flow diagram of the process used. (B) Library construction schematic (fig. S1). r1 to r8 are gDNA regions adjacent to distinct adapter ends; Ad1 to Ad4 indicate adapters 1 to 4. (C) DNB production using *Phi*29 DNA polymerase (fig. S11) and nanoarray formation (SOM) schematics. (D) Schematic of cPAL products (SOM).

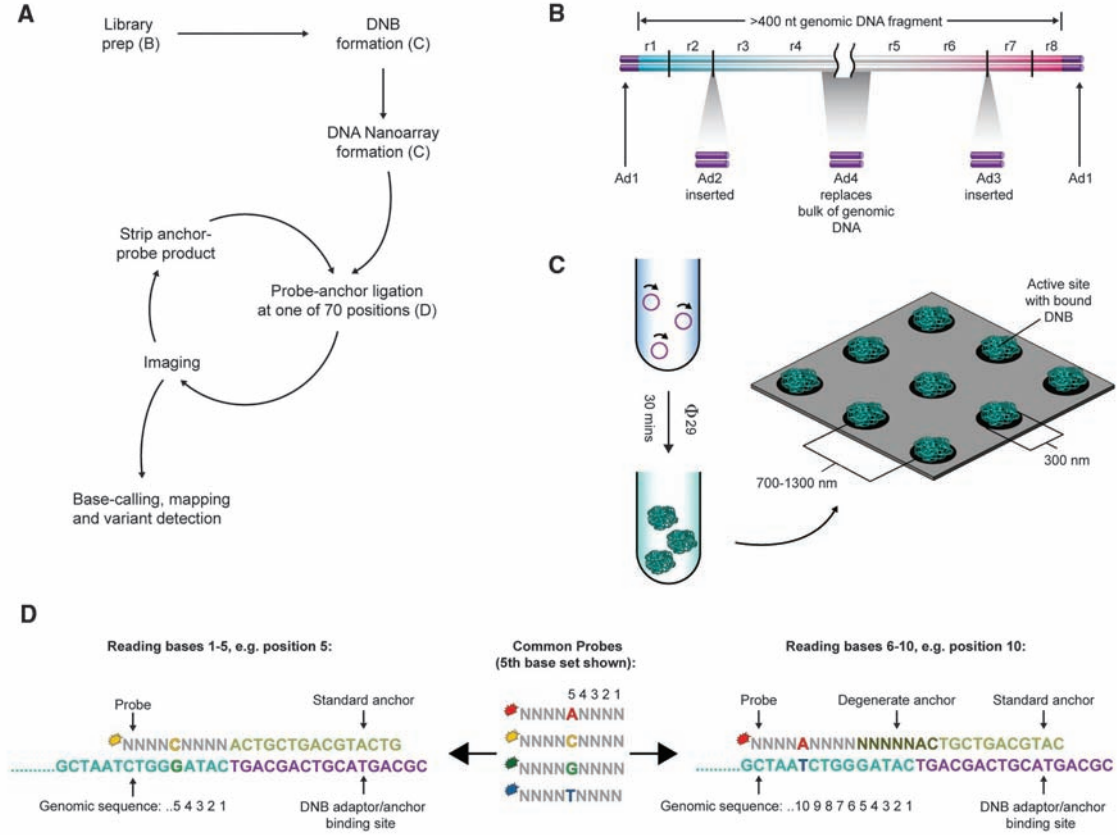


Table 1. Summary information from mapping and assembly of three genomes. All variations are with respect to the National Center for Biotechnology Information (NCBI) version 36 human genome reference assembly. Novel variations were ascertained by comparison to dbSNP [JDW, release 126; NA18507 (6), release 128; all other genomes, re-

lease 129]. NA18507 and NA19240 are Yoruban HapMap samples, which may explain the number of SNPs and novelty rates. In partially called regions of the genome, one allele could be called confidently but not the other. The high call rate in NA19240 reflects reduced library bias (fig. S5).

Sample	Mapped sequence (Gb)	Average coverage depth (fold)	Percent of genome called		SNPs		Indels		Insertion: deletion ratio
			Fully	Partially	Total	Novel	Total	Novel	
Genomes sequenced by Complete Genomics									
NA07022 (35)	241	87	91%	2%	3,076,869	10%	337,635	37%	1.0
NA19240 (36)	178	63	95%	1%	4,042,801	19%	496,194	42%	0.96
NA20431 (37)	124	45	86%	3%	2,905,517	10%	269,794	37%	1.0
Genomes previously published									
NA18507 (6)	131	46	—	—	4,139,196	26%	404,416	50%	0.77
NA18507 (9)	87	31*	—	—	3,866,085	19%	226,529	33%	0.72
JCV (3)	21	7	—	—	3,213,401	15%	851,575	—	—
JDW (4)	21	7	—	—	3,322,093	18%	222,718	51%	0.4

*This is 18x when constrained to nonduplicated and properly mated reads, which were those used for variant calling.

Table 2. Concordance with genotypes for NA07022 generated by the HapMap Project (release 24) and the highest quality Infinium assay subset of those genotypes, as well as genotyping on illumina Infinium 1M assay. Discordances with reported HapMap Infinium genotypes were verified by Sanger sequencing (SOM).

		Infinium 1M	HapMap phase I and II SNPs	HapMap Infinium subset	HapMap Infinium SNPs tested for accuracy by Sanger sequencing			
Published concordance*		—	99.03%	99.94%				
NA07022	Number reported	1 M	3.9 M	143 K				
	Called (%)	95.98%	94.39%	96.00%	These	These	Affirmed	
	Locus concordance (%)	99.89%	99.15%	99.88%	data correct	data incorrect	(%)	
	HapMap	Homozygous reference	99.96%	99.34%	99.96%	18	2	90%
	genotype	Heterozygous	99.78%	99.39%	99.80%	28	46	38%†
	calls	Homozygous alternate	99.81%	98.14%	99.84%	28	12	70%†

*Published concordance between Sanger sequencing and genotyping (34). †Average false negative rate is <0.19.

in simulations is composed of repeats that are longer than the ~400 base inserts used here and of high enough identity to prevent attribution of mappings to specific repeat copies. Longer mate-pair inserts minimize this limitation (6, 9). Similar limitations affect other short-read technologies.

We identified a range of 2.91 to 4.04 million SNPs with respect to the reference genome, 81% to 90% of which are reported in the database dbSNP, as well as short indels and block substitutions (Table 1 and table S6). Because of the use of local de novo assembly, indels were detected in sizes ranging up to 50 bp. As expected, indels in coding regions tend to occur in multiples of length 3, indicating the possible selection of minimally impacting variants in coding regions (fig. S6).

As an initial test of sequence accuracy, we compared our called SNPs with the HapMap phase I/II SNP genotypes reported for NA07022 (1). We fully called 94% of these positions with an overall concordance of 99.15% (Table 2) (the remaining 6% of positions were either half-called or not called). Furthermore, we fully called 96% of the Infinium (Illumina, San Diego, CA) subset of the HapMap SNPs with an overall concordance rate of 99.88%, reflecting the higher reported accuracy of these genotypes (34). Similar concordance rates with available SNP genotypes were observed in NA19240 (with a call rate of over 98%) and NA20431 (table S7).

We further characterized 134 of the 168 calls that were discordant with Infinium loci and Sanger sequencing of PCR products in NA07022, demonstrating that 55% of these discordances are errors in the reported HapMap genotypes (Table 2). The relationship between detection rate and read depth for about 1M Infinium HD SNPs that we subsequently genotyped in NA07022 shows that coverage of 30-fold at a position is sufficient to detect 90% of SNPs at heterozygous loci and 97% of SNPs at homozygous loci (fig. S5).

Because the whole-genome false-positive rate cannot be accurately estimated from known SNP loci, we tested a random subset of novel nonsynonymous variants in NA07022, a category that is enriched for errors (10). We extrapolated error rates from the targeted sequencing of 291

Table 3. False-positive rates and false discovery rates (FDRs) were calculated for the entire set of variations called in NA07022 by extrapolating the heterozygous (Het) FDRs calculated from comparative Sanger sequencing of 291 selected novel variants (table S8) to all variants. This is a conservative approach (detailed in SOM). The total number of all types of false-positive variants is estimated at 7.5 to 16.1 per Mb.

Variation type	Total detected	Novel	Het novel FDR (table S8)	Estimated false positives on genome	Estimated false positives / Mbp	Estimated FDR
SNP	3,076,869	310,690	2–6%	7k–17k	2.3–6.1	0.2–0.6%
Deletion	168,726	61,960	8–14%	5k–8k	1.8–3.0	3.0–5.0%
Insertion	168,909	61,933	11–18%	7k–11k	2.3–3.9	3.9–6.5%
Block substitution	62,783	30,445	11–29%	3k–9k	1.1–3.1	5.2–13.9%

such loci and estimated the false-positive rate at about one variant per 100 kb, including <6.1 substitution-, <3.0 short deletion-, <3.9 short insertion- and <3.1 block variants per Mb (Table 3 and table S8).

Aberrant mate-pair gaps may indicate the presence of length-altering structural variants and rearrangements with respect to the reference genome. A total of 2,126 clusters of such anomalous mate-pairs were identified in NA07022. We performed PCR-based confirmation of one such heterozygous 1500-base deletion (fig. S7). More than half of the clusters are consistent in size with the addition or deletion of a single *Alu* repeat element.

Some applications of complete genome sequencing may benefit from maximal discovery rates, even at the cost of additional false positives, whereas for others, a lower discovery rate and lower false-positive rate may be preferable. We used the variant quality score to tune call rate and accuracy (fig. S8). Additionally, novelty rate (relative to the dbSNP) is also a function of variant quality score (fig. S9).

We processed the NA07022 data with Trait-o-matic (Scalable Computer Experts, Somerville, Massachusetts, and Church Laboratory, Harvard Medical School, Cambridge, Massachusetts) automated annotation software [as in (12)], yielding 1159 annotated variants, 14 of which may have disease implications (table S10).

Because the DNB sequencing substrates are produced by rolling-circle replication (33) in a uniform-temperature, solution-phase reaction with high template concentrations (>20 billion

per ml), this system avoids substantial selection bottlenecks and nonclonal DNBs. This circumvents the stochastic inefficiencies of approaches that require precise titration of template concentrations for in situ clonal amplification in emulsion (9, 14, 29) or bridge PCR (6, 19).

Our patterned arrays include high-occupancy and high-density nanoarrays self-assembled on photolithography-patterned, solid-phase substrates through electrostatic adsorption of solution-phase DNBs and yield a high proportion of informative pixels (site occupancies >95%) (fig. S12A) compared with random-position DNA arrays. This results in several hundred reaction sites in the compact (~300-nm diameter) DNB that produce bright signals useful for rapid imaging of the sequences (SOM). Such small DNBs also allow for high-density arrays. The data set reported herein was generated with arrays with ~350 million spots at a pitch of 1.29 μm. Such a spot density and higher ones achieved in proof-of-concept experiments (fig. S12B) result in high image efficiency and reduced reagent consumption that enable high sequencing throughput per instrument critical for high-scale human genome sequencing for research and clinical applications.

Both sequencing by synthesis (SBS) and sequencing by ligation (SBL) use chained reads, wherein the substrate for cycle *N* + 1 is dependent on the product of cycle *N*; consequently, errors may accumulate over multiple cycles and data quality may be affected by errors (especially incomplete extensions) occurring in previous cycles. Thus, reactions need to be driven to

near completion with high concentrations of expensive high-purity labeled substrate molecules and enzymes. The independent, unchained nature of cPAL avoids error accumulation and tolerates low-quality bases in otherwise high-quality reads, thereby decreasing reagent costs. The average sequencing consumables cost for these three genomes was under \$4400 (table S5). The raw base and variant call accuracy achieved compares favorably with other reported human genome sequences (2–12).

Because the sequencing substrates are produced by a DNA engineering process based on modified nick-translation for directional adapter insertion (SOM), we obtained over 90% yield in adapter ligation and low chimeric rates of about 1% (SOM). DNA molecules with an inserted adapter are further enriched with PCR (SOM). This recursive process can be implemented in batches of 96 samples and extended by inserting additional adapters to read 120 bases or more per DNB (fig. S10). The current read length is comparable to other massively parallel sequencing technologies (6–12).

The sequence data reported here achieve sufficient quality and accuracy for complete genome association studies, the identification of potentially rare variants associated with disease or therapeutic treatments, and the identification of somatic mutations. The low cost of consumables and efficient imaging may enable studies of hundreds of individuals. The higher

accuracy and completeness required for clinical diagnostic applications provides an incentive for continued improvement of this and other technologies.

References and Notes

1. T. A. Manolio, L. D. Brooks, F. S. Collins, *J. Clin. Invest.* **118**, 1590 (2008).
2. International Human Genome Sequencing Consortium, *Nature* **409**, 860 (2001).
3. J. C. Venter *et al.*, *Science* **291**, 1304 (2001).
4. S. Levy *et al.*, *PLoS Biol.* **5**, e254 (2007).
5. D. A. Wheeler *et al.*, *Nature* **452**, 872 (2008).
6. D. R. Bentley *et al.*, *Nature* **456**, 53 (2008).
7. J. Wang *et al.*, *Nature* **456**, 60 (2008).
8. S. M. Ahn *et al.*, *Genome Res.* **19**, 1622 (2009).
9. K. J. McKernan *et al.*, *Genome Res.* **19**, 1527 (2009).
10. T. J. Ley *et al.*, *Nature* **456**, 66 (2008).
11. D. Pushkarev, N. F. Neff, S. R. Quake, *Nat. Biotechnol.* **27**, 847 (2009).
12. J. I. Kim *et al.*, *Nature* **460**, 1011 (2009).
13. R. Drmanac *et al.*, *Genomics* **4**, 114 (1989).
14. R. Drmanac, R. Crkvenjakov, *Sci. Yugosl.* **16**, 97 (1990); pdf accessible at www.rdrmanac.com.
15. R. Drmanac *et al.*, *Science* **260**, 1649 (1993).
16. P. C. Cheesman, U.S. patent 5,302,509 (1994).
17. R. Drmanac, World Intellectual Property Organization WO/1995/009248 (1995).
18. M. Ronaghi, S. Karamohamed, B. Pettersson, M. Uhlen, P. Nyren, *Anal. Biochem.* **242**, 84 (1996).
19. C. P. Adams, S. J. Kron, U.S. Patent 5,641,658 (1997).
20. P. M. Lizardi *et al.*, *Nat. Genet.* **19**, 225 (1998).
21. S. C. Macevitz, U.S. Patent 5,750,341 (1998).
22. S. Drmanac *et al.*, *Nat. Biotechnol.* **16**, 54 (1998).
23. R. D. Mitra, G. M. Church, *Nucleic Acids Res.* **27**, e34 (1999).
24. S. Brenner *et al.*, *Nat. Biotechnol.* **18**, 630 (2000).
25. I. Braslavsky, B. Hebert, E. Kartalov, S. R. Quake, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 3960 (2003).
26. R. D. Mitra, J. Shendure, J. Olejnik, O. E. Krzyminska, G. M. Church, *Anal. Biochem.* **320**, 55 (2003).
27. R. Drmanac *et al.*, World Intellectual Property Organization WO/2004/076683 (2004).
28. J. Shendure *et al.*, *Science* **309**, 1728 (2005).
29. M. Margulies *et al.*, *Nature* **437**, 376 (2005).
30. T. D. Harris *et al.*, *Science* **320**, 106 (2008).
31. A. Pihlak *et al.*, *Nat. Biotechnol.* **26**, 676 (2008).
32. J. Shendure, H. Ji, *Nat. Biotechnol.* **26**, 1135 (2008).
33. L. Blanco *et al.*, *J. Biol. Chem.* **264**, 8935 (1989).
34. The International HapMap Consortium, *Nature* **449**, 851 (2007).
35. K. Zhang *et al.*, *Nat. Methods* **6**, 613 (2009).
36. M. P. Ball *et al.*, *Nat. Biotechnol.* **27**, 361 (2009).
37. J. B. Li *et al.*, *Genome Res.* **19**, 1606 (2009).
38. We acknowledge the ongoing contributions and support of all Complete Genomics employees and R. Mercado for manuscript preparation. Some of this work was supported by PersonalGenomes.org and the National Heart, Lung, and Blood Institute. Data has been deposited at the National Center for Biotechnology Information: reads in the Short Read Archive (SRA), accession SRA008092, and SNPs in dbSNP, accessions ss161884913 to ss175323894. Employees of Complete Genomics have stock options in the company, and D.G.B. has stock in Perlegen Sciences. Complete Genomics has filed several patents on this work.

Supporting Online Material

www.sciencemag.org/cgi/content/full/1181498/DC1
Materials and Methods
Figs. S1 to S12
Tables S1 to S9
References

3 September 2009; accepted 23 October 2009
Published online 5 November 2009;
10.1126/science.1181498
Include this information when citing this paper.

Structure and Mechanisms of a Protein-Based Organelle in *Escherichia coli*

Shiho Tanaka,¹ Michael R. Sawaya,^{2,3} Todd O. Yeates^{1,3,4*}

Many bacterial cells contain proteinaceous microcompartments that act as simple organelles by sequestering specific metabolic processes involving volatile or toxic metabolites. Here we report the three-dimensional (3D) crystal structures, with resolutions between 1.65 and 2.5 angstroms, of the four homologous proteins (EutS, EutL, EutK, and EutM) that are thought to be the major shell constituents of a functionally complex ethanolamine utilization (Eut) microcompartment. The Eut microcompartment is used to sequester the metabolism of ethanolamine in bacteria such as *Escherichia coli* and *Salmonella enterica*. The four Eut shell proteins share an overall similar 3D fold, but they have distinguishing structural features that help explain the specific roles they play in the microcompartment. For example, EutL undergoes a conformational change that is probably involved in gating molecular transport through shell protein pores, whereas structural evidence suggests that EutK might bind a nucleic acid component. Together these structures give mechanistic insight into bacterial microcompartments.

Bacterial microcompartments are present in diverse bacteria, where they function as protein-based organelles (1–6). They range in size from just under 1000 Å to around 1500 Å, and they typically have a polyhedral shape (Fig. 1A). Each type of microcompartment contains a few different enzymes that catalyze sequential metabolic reactions. The enzymes are encapsulated by a shell formed from a few thousand shell protein

subunits. The simplest microcompartment is the carboxysome, which encapsulates the two enzymes carbonic anhydrase and ribulose-1,5-bisphosphate carboxylase-oxygenase (RuBisCO) in order to enhance cellular CO₂ fixation (1, 7, 8). Recent structural studies on the carboxysome and its shell proteins have provided a basic understanding of how that type of microcompartment is assembled and how it operates (9–14).

Other bacterial microcompartments with more complex metabolic functions have also been discovered (5, 15). One of these, which is dedicated to ethanolamine utilization (Eut), is present in several bacteria, including *Salmonella enterica* and *Escherichia coli* (4, 16). The Eut microcompartment shares a number of homologous enzymes with a propanediol utilization (Pdu) microcompartment; both metabolic pathways proceed via aldehyde intermediates, propionaldehyde in the case of Pdu and acetaldehyde in the case of Eut (3). Experiments in *Salmonella* have shown that the cellular function of the Eut microcompartment is to metabolize ethanolamine without allowing the release of acetaldehyde into the cytosol (Fig. 1B), thus mitigating the potentially toxic effects of excess aldehyde in the bacterial cytosol (17–19) and also preventing the volatile acetaldehyde from diffusing across the cell membrane and leading to a loss of carbon (20).

¹Department of Chemistry and Biochemistry, University of California Los Angeles, 611 Charles Young Drive East, Los Angeles, CA 90095, USA. ²Howard Hughes Medical Institute, University of California Los Angeles, 611 Charles Young Drive East, Los Angeles, CA 90095, USA. ³University of California Los Angeles, Department of Energy Institute for Genomics and Proteomics, 611 Charles Young Drive East, Los Angeles, CA 90095, USA. ⁴Molecular Biology Institute, University of California Los Angeles 611 Charles Young Drive East, Los Angeles, CA 90095, USA.

*To whom correspondence should be addressed. E-mail: yeates@mbi.ucla.edu