

Towards practical, high-capacity, low-maintenance information storage in synthesized DNA

Nick Goldman¹, Paul Bertone¹, Siyuan Chen², Christophe Dessimoz¹, Emily M. LeProust², Botond Sipos¹ & Ewan Birney¹

Digital production, transmission and storage have revolutionized how we access and use information but have also made archiving an increasingly complex task that requires active, continuing maintenance of digital media. This challenge has focused some interest on DNA as an attractive target for information storage¹ because of its capacity for high-density information encoding, longevity under easily achieved conditions^{2–4} and proven track record as an information bearer. Previous DNA-based information storage approaches have encoded only trivial amounts of information^{5–7} or were not amenable to scaling-up⁸, and used no robust error-correction and lacked examination of their cost-efficiency for large-scale information archival⁹. Here we describe a scalable method that can reliably store more information than has been handled before. We encoded computer files totalling 739 kilobytes of hard-disk storage and with an estimated Shannon information¹⁰ of 5.2×10^6 bits into a DNA code, synthesized this DNA, sequenced it and reconstructed the original files with 100% accuracy. Theoretical analysis indicates that our DNA-based storage scheme could be scaled far beyond current global information volumes and offers a realistic technology for large-scale, long-term and infrequently accessed digital archiving. In fact, current trends in technological advances are reducing DNA synthesis costs at a pace that should make our scheme cost-effective for sub-50-year archiving within a decade.

Although techniques for manipulating, storing and copying large amounts of existing DNA have been established for many years^{11–13}, one of the main challenges for practical DNA-based information storage is the difficulty of synthesizing long sequences of DNA *de novo* to an exactly specified design. As in the approach of ref. 9, we represent the information being stored as a hypothetical long DNA molecule and encode this *in vitro* using shorter DNA fragments. This offers the benefits that isolated DNA fragments are easily manipulated *in vitro*^{11,13}, and that the routine recovery of intact fragments from samples that are tens of thousands of years old^{14,15} indicates that well-prepared synthetic DNA should have an exceptionally long lifespan in low-maintenance environments^{3,4}. In contrast, approaches using living vectors^{6–8} are not as reliable, scalable or cost-efficient owing to disadvantages such as constraints on the genomic elements and locations that can be manipulated without affecting viability, the fact that mutation will cause the fidelity of stored and decoded information to reduce over time, and possibly the requirement for storage conditions to be carefully regulated. Existing schemes used for DNA computing in principle permit large-scale memory^{1,16}, but data encoding in DNA computing is inextricably linked to the specific application or algorithm¹⁷ and no practical storage schemes have been realized.

As a proof of concept for practical DNA-based storage, we selected and encoded a range of common computer file formats to emphasize the ability to store arbitrary digital information. The five files comprised all 154 of Shakespeare's sonnets (ASCII text), a classic scientific paper¹⁸ (PDF format), a medium-resolution colour photograph of the European Bioinformatics Institute (JPEG 2000 format), a 26-s excerpt from Martin Luther King's 1963 'I have a dream' speech (MP3 format) and a Huffman code¹⁰ used in this study to convert bytes to base-3

digits (ASCII text), giving a total of 757,051 bytes or a Shannon information¹⁰ of 5.2×10^6 bits (see Supplementary Information and Supplementary Table 1 for full details).

The bytes comprising each file were represented as single DNA sequences with no homopolymers (runs of ≥ 2 identical bases, which are associated with higher error rates in existing high-throughput sequencing technologies¹⁹ and led to errors in a recent DNA-storage experiment⁹). Each DNA sequence was split into overlapping segments, generating fourfold redundancy, and alternate segments were converted to their reverse complement (see Fig. 1 and Supplementary Information). These measures reduce the probability of systematic failure for any particular string, which could lead to uncorrectable errors and data loss. Each segment was then augmented with indexing information that permitted determination of the file from which it originated and its location within that file, and simple parity-check error-detection¹⁰. In all, the five files were represented by a total of 153,335 strings of DNA, each comprising 117 nucleotides (nt). The perfectly uniform fragment lengths and absence of homopolymers make it obvious that the synthesized DNA does not have a natural (biological) origin, and so imply the presence of deliberate design and encoded information².

We synthesized oligonucleotides (oligos) corresponding to our designed DNA strings using an updated version of Agilent Technologies' OLS (oligo library synthesis) process²⁰, creating $\sim 1.2 \times 10^7$ copies of each DNA string. Errors occur only rarely (~ 1 error per 500 bases) and independently in the different copies of each string, again enhancing our method's error tolerance. We shipped the synthesized DNA in lyophilized form that is expected to have excellent long-term preservation characteristics^{3,4}, at ambient temperature and without specialized packaging, from the USA to Germany via the UK. After resuspension, amplification and purification, we sequenced a sample of the resulting library products at the EMBL Genomics Core Facility in paired-end mode on the Illumina HiSeq 2000. We transferred the remainder of the library to multiple aliquots and re-lyophilized these for long-term storage.

Our base calling using AYB²¹ yielded 79.6×10^6 read-pairs of 104 bases in length, from which we reconstructed full-length (117-nt) DNA strings *in silico*. Strings with uncertainties due to synthesis or sequencing errors were discarded and the remainder decoded using the reverse of the encoding procedure, with the error-detection bases and properties of the coding scheme allowing us to discard further strings containing errors. Although many discarded strings will have contained information that could have been recovered with more sophisticated decoding, the high level of redundancy and sequencing coverage rendered this unnecessary in our experiment. Full-length DNA sequences representing the original encoded files were then reconstructed *in silico*. The decoding process used no additional information derived from knowledge of the experimental design. Full details of the encoding, sequencing and decoding processes are given in Supplementary Information.

Four of the five resulting DNA sequences could be fully decoded without intervention. The fifth however contained two gaps, each a run

¹European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton CB10 1SD, UK. ²Agilent Technologies, Genomics-LSSU, 5301 Stevens Creek Boulevard, Santa Clara, California 95051, USA.

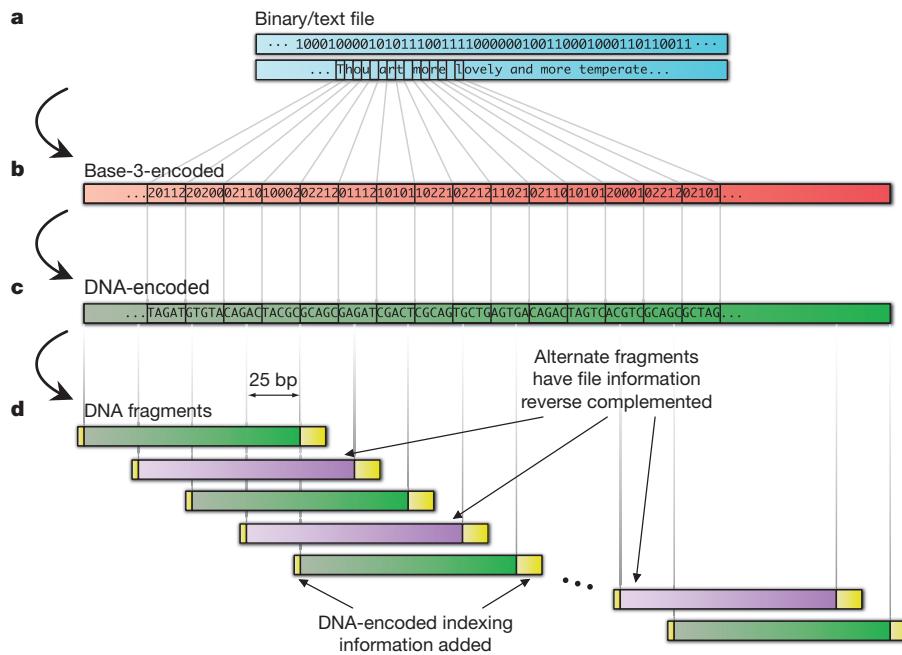


Figure 1 | Digital information encoding in DNA. Digital information (a, in blue), here binary digits holding the ASCII codes for part of Shakespeare’s sonnet 18, was converted to base-3 (b, red) using a Huffman code that replaces each byte with five or six base-3 digits (trits). This in turn was converted *in silico* to our DNA code (c, green) by replacement of each trit with one of the three nucleotides different from the previous one used, ensuring no homopolymers

of 25 bases, for which no segment was detected corresponding to the original DNA. Each of these gaps was caused by the failure to sequence any oligo representing any of four consecutive overlapping segments. Inspection of the neighbouring regions of the reconstructed sequence permitted us to hypothesize what the missing nucleotides should have been (see Supplementary Information) and we manually inserted those 50 bases accordingly. This sequence could also then be decoded. Inspection confirmed that our original computer files had been reconstructed with 100% accuracy.

An important issue for long-term digital archiving is how DNA-based storage scales to larger applications. The number of bases of synthesized DNA needed to encode information grows linearly with the amount of information to be stored, but we must also consider the indexing information required to reconstruct full-length files from short fragments. As indexing information grows only as the logarithm of the number of fragments to be indexed, the total amount of synthesized DNA required grows sub-linearly. Increasingly large parts of each fragment are needed for indexing however and, although it is reasonable to expect synthesis of longer strings to be possible in future, we modelled the behaviour of our scheme under the conservative constraint of a constant 114 nt available for both data and indexing information (see Supplementary Information). As the total amount of information increases, the encoding efficiency decreases only slowly (Fig. 2a). In our experiment (megabyte scale) the encoding scheme is 88% efficient; Fig. 2a indicates that efficiency remains >70% for data storage on petabyte (PB, 10^{15} bytes) scales and >65% on exabyte (EB, 10^{18} bytes) scales, and that DNA-based storage remains feasible on scales many orders of magnitude greater than current global data volumes²². Figure 2a also shows that costs (per unit information stored) rise only slowly as data volumes increase over many orders of magnitude. Efficiency and costs scale even more favourably if we consider the synthesized fragment lengths available using the latest technology (Supplementary Fig. 5).

As the amount of information stored increases, decoding requires more strings to be sequenced. A fixed decoding expenditure per byte of

were generated. This formed the basis for a large number of overlapping segments of length 100 bases with overlap of 75 bases, creating fourfold redundancy (d, green and, with alternate segments reverse complemented for added data security, violet). Indexing DNA codes were added (yellow), also encoded as non-repeating DNA nucleotides. See Supplementary Information for further details.

encoded information would mean that each base is read fewer times and so is more likely to suffer decoding error. But extension of our scaling analysis to model the influence of reduced sequencing coverage on the per-decoded-base error rate (see Supplementary Information) revealed that error rates increase only very slowly as the amount of information encoded increases to a global data scale and beyond (Supplementary Table 4). This also suggests that our mean sequencing coverage of 1,308 times was considerably in excess of that needed for reliable decoding. We confirmed this by subsampling from the 79.6×10^6 read-pairs to simulate experiments with lower coverage. Figure 2b indicates that reducing the coverage by a factor of 10 (or even more) would have led to unaltered decoding characteristics, which further illustrates the robustness of our DNA-storage method.

DNA-based storage might already be economically viable for long-horizon archives with a low expectation of extensive access, such as government and historical records^{23,24}. An example in a scientific context is CERN’s CASTOR system²⁵, which stores a total of 80 PB of Large Hadron Collider data and grows at 15 PB yr^{-1} . Only 10% is maintained on disk, and CASTOR migrates regularly between magnetic tape formats. Archives of older data are needed for potential future verification of events, but access rates decrease considerably 2–3 years after collection. Further examples are found in astronomy, medicine and interplanetary exploration²⁶. With negligible computational costs and optimized use of the technologies we employed, we estimate current costs to be $\$12,400 \text{ MB}^{-1}$ for information storage in DNA and $\$220 \text{ MB}^{-1}$ for information decoding. Modelling relative long-term costs of archiving using DNA-based storage or magnetic tape shows that the key parameters are the ratio of the one-time cost of synthesizing the DNA to the recurrent fixed cost of transferring data between tape technologies or media, which we estimate to be 125–500 currently, and the frequency of tape transition events (Supplementary Information and Supplementary Fig. 7). We find that with current technology and our encoding scheme, DNA-based storage may be cost-effective for archives of several megabytes with a ~ 600 –5,000-yr horizon (Fig. 2c). One order of magnitude reduction in synthesis costs reduces this to ~ 50 –500 yr; with two orders

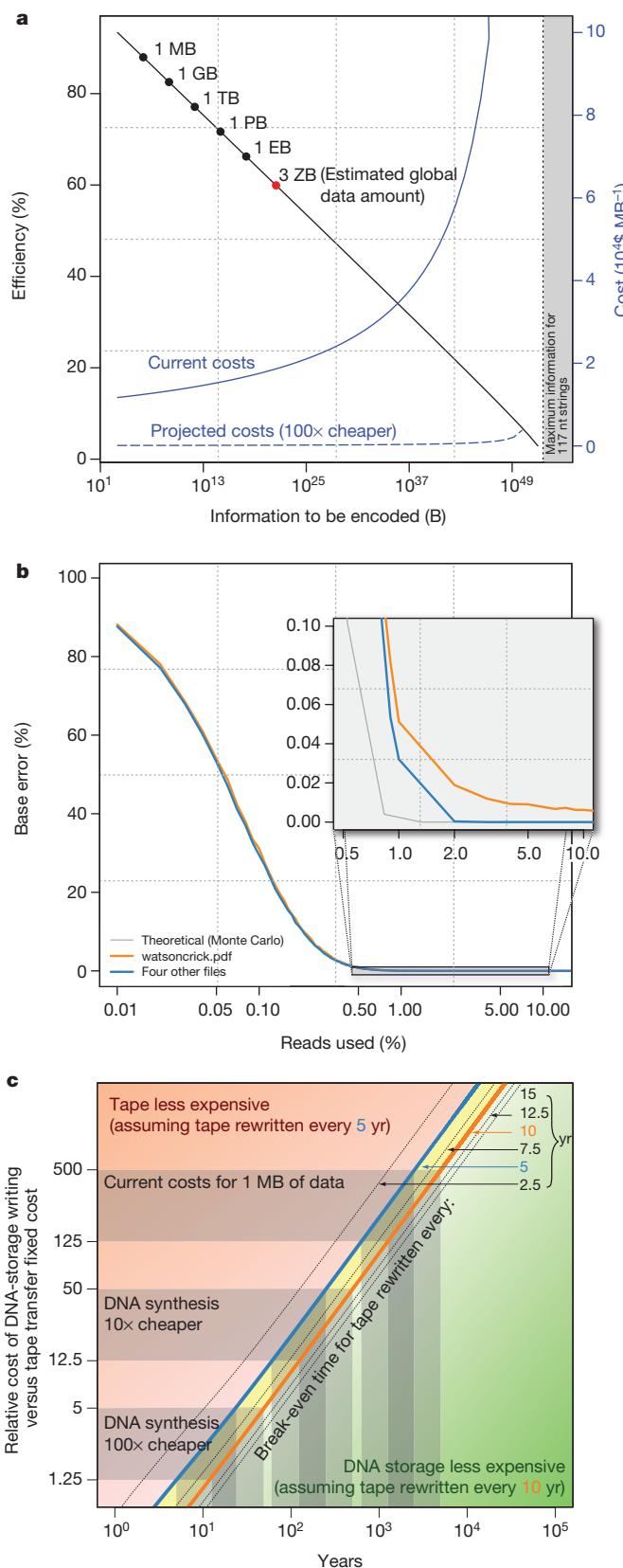


Figure 2 | Scaling properties and robustness of DNA-based storage. **a**, Encoding efficiency and costs change as the amount of stored information increases. The x axis (logarithmic scale) represents the total amount of information to be encoded. Common data scales are indicated, including the three zettabyte (3 ZB, 3×10^{21} bytes) global data estimate, shown red. The black line (y-axis scale to left) indicates encoding efficiency, measured as the proportion of synthesized bases available for data encoding. The blue curves (y-axis scale to right) indicate the corresponding effect on encoding costs, both at current synthesis cost levels (solid line) and in the case of a two-order-of-magnitude reduction (dashed line). **b**, Per-recovered-base error rate (y axis) as a function of sequencing coverage, represented by the percentage of the original 79.6×10^6 read-pairs sampled (x axis; logarithmic scale). The blue curve represents the four files recovered without human intervention: the error is zero when $\geq 2\%$ of the original reads are used. The grey curve is obtained by Monte Carlo simulation from our theoretical error rate model. The orange curve represents the file (watsoncrick.pdf) that required manual correction: the minimum possible error rate is 0.0036%. The boxed area is shown magnified in the inset. **c**, Timescales for which DNA-based storage is cost-effective. The blue curve indicates the relationship between break-even time beyond which DNA storage is less expensive than magnetic tape (x axis) and relative cost of DNA-storage synthesis and tape transfer fixed costs (y axis), assuming the tape archive has to be read and rewritten every 5 yr. The orange curve corresponds to tape transfers every 10 yr; broken curves correspond to other transfer periods as indicated. In the green-shaded region, DNA storage is cost-effective when transfers occur more frequently than every 10 yr; in the yellow-shaded region, DNA storage is cost-effective when transfers occur every 5–10 yr; in the red-shaded region tape is less expensive when transfers occur less frequently than every 5 yr. Grey-shaded ranges of relative costs of DNA synthesis to tape transfer are 125–500 (current costs for 1 MB of data), 12.5–50 (achieved if DNA synthesis costs are reduced by one order of magnitude) and 1.25–5 (costs reduced by two orders of magnitude). Note the logarithmic scales on both axes. See Supplementary Information for further details.

both processes can be accelerated through parallelization (Supplementary Information).

The DNA-based storage medium has different properties from traditional tape- or disk-based storage. As DNA is the basis of life on Earth, methods for manipulating, storing and reading it will remain the subject of continual technological innovation. As with any storage system, a large-scale DNA archive would need stable DNA management²⁷ and physical indexing of depositions. But whereas current digital schemes for archiving require active and continuing maintenance and regular transferring between storage media, the DNA-based storage medium requires no active maintenance other than a cold, dry and dark environment^{3,4} (such as the Global Crop Diversity Trust's Svalbard Global Seed Vault, which has no permanent on-site staff²⁸) yet remains viable for thousands of years even by conservative estimates. We achieved an information storage density of ~ 2.2 PB g $^{-1}$ (Supplementary Information). Our sequencing protocol consumed just 10% of the library produced from the synthesized DNA (Supplementary Table 2), already leaving enough for multiple equivalent copies. Existing technologies for copying DNA are highly efficient^{11,13}, meaning that DNA is an excellent medium for the creation of copies of any archive for transportation, sharing or security. Overall, DNA-based storage has potential as a practical solution to the digital archiving problem and may become a cost-effective solution for rarely accessed archives.

Received 15 May; accepted 12 December 2012.

Published online 23 January 2013.

1. Baum, E. B. Building an associative memory vastly larger than the brain. *Science* **268**, 583–585 (1995).
2. Cox, J. P. L. Long-term data storage in DNA. *Trends Biotechnol.* **19**, 247–250 (2001).
3. Anchordoquy, T. J. & Molina, M. C. Preservation of DNA. *Cell Preserv. Technol.* **5**, 180–188 (2007).
4. Bonnet, J. et al. Chain and conformation stability of solid-state DNA: implications for room temperature storage. *Nucleic Acids Res.* **38**, 1531–1546 (2010).
5. Clelland, C. T., Risca, V. & Bancroft, C. Hiding messages in DNA microdots. *Nature* **399**, 533–534 (1999).
6. Kac, E. *Genesis* (1999); available at <http://www.ekac.org/geninfo.html> (accessed 10 May 2012).

of magnitude reduction, as can be expected in less than a decade if current trends continue (ref. 13, and <http://www.synthesis.cc/2011/06/new-cost-curves.html>), DNA-based storage becomes practical for archives with a horizon of less than 50 yr. The speed of DNA-storage writing and reading are not competitive with current technology, but

7. Ailenberg, M. & Rotstein, O. D. An improved Huffman coding method for archiving text, images, and music characters in DNA. *Biotechniques* **47**, 747–754 (2009).
8. Gibson, D. G. *et al.* Creation of a bacterial cell controlled by a chemically synthesized genome. *Science* **329**, 52–56 (2010).
9. Church, G. M., Gao, Y. & Kosuri, S. Next-generation digital information storage in DNA. *Science* **337**, 1628 (2012).
10. MacKay, D. J. C. *Information Theory, Inference, and Learning Algorithms* (Cambridge Univ. Press, 2003).
11. Erlich, H. A., Gelfand, D. & Sninsky, J. J. Recent advances in the polymerase chain reaction. *Science* **252**, 1643–1651 (1991).
12. Monaco, A. P. & Larin, Z. YACs, BACs, PACs and MACs: artificial chromosomes as research tools. *Trends Biotechnol.* **12**, 280–286 (1994).
13. Carr, P. A. & Church, G. M. Genome engineering. *Nature Biotechnol.* **27**, 1151–1162 (2009).
14. Willerslev, E. *et al.* Ancient biomolecules from deep ice cores reveal a forested southern Greenland. *Science* **317**, 111–114 (2007).
15. Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010).
16. Kari, L. & Mahalingam, K. In *Algorithms and Theory of Computation Handbook* Vol. 2, 2nd edn (eds Atallah, M. J. & Blanton, M.) 31-1–31-24 (Chapman & Hall, 2009).
17. Päun, G., Rozenberg, G. & Salomaa, A. *DNA Computing: New Computing Paradigms* (Springer, 1998).
18. Watson, J. D. & Crick, F. H. C. Molecular structure of nucleic acids. *Nature* **171**, 737–738 (1953).
19. Niedringhaus, T. P., Milanova, D., Kerby, M. B., Snyder, M. P. & Barron, A. E. Landscape of next-generation sequencing technologies. *Anal. Chem.* **83**, 4327–4341 (2011).
20. LeProust, E. M. *et al.* Synthesis of high-quality libraries of long (150mer) oligonucleotides by a novel depurination controlled process. *Nucleic Acids Res.* **38**, 2522–2540 (2010).
21. Massingham, T. & Goldman, N. All Your Base: a fast and accurate probabilistic approach to base calling. *Genome Biol.* **13**, R13 (2012).
22. Gantz, J. & Reinsel, D. *Extracting Value from Chaos* (IDC, 2011).
23. Brand, S. *The Clock of the Long Now* (Basic Books, 1999).
24. Digital archiving. History flushed. *Economist* **403**, 56–57 (28 April 2012); available at <http://www.economist.com/node/21553410> (2012).
25. Bessone, N., Cancio, G., Murray, S. & Tarelli, G. Increasing the efficiency of tape-based storage backends. *J. Phys. Conf. Ser.* **219**, 062038 (2010).
26. Baker, M. *et al.* in *Proc. 1st ACM SIGOPS/EuroSys European Conf. on Computer Systems* (eds Berbers, Y. & Zwaenepoel, W.) 221–234 (ACM, 2006).
27. Yuille, M. *et al.* The UK DNA banking network: a “fair access” biobank. *Cell Tissue Bank.* **11**, 241–251 (2010).
28. Global Crop Diversity Trust. Svalbard Global Seed Vault. (2012); available at <http://www.croptrust.org/main/content/svalbard-global-seed-vault> (accessed 10 May 2012).

Supplementary Information is available in the online version of the paper.

Acknowledgements At the University of Cambridge: D. MacKay and G. Mitchison for advice on codes for run-length-limited channels. At CERN: B. Jones for discussions on data archival. At EBI: A. Löytynoja for custom multiple sequence alignment software, H. Marsden for computing base calls and for detecting an error in the original parity-check encoding, T. Massingham for computing base calls and advice on code theory and K. Gori, D. Henk, R. Loos, S. Parks and R. Schwarz for assistance with revisions to the manuscript. In the Genomics Core Facility at EMBL Heidelberg: V. Benes for advice on Next-Generation Sequencing protocols, D. Pavlinić for sequencing and J. Blake for data handling. C.D. is supported by a fellowship from the Swiss National Science Foundation (grant 136461). B.S. is supported by an EMBL Interdisciplinary Postdoctoral Fellowship under Marie Curie Actions (COFUND).

Author Contributions N.G. and E.B. conceived and planned the project and devised the information-encoding methods. P.B. advised on oligo design and Next-Generation Sequencing protocols, prepared the DNA library and managed the sequencing process. S.C. and E.M.L. provided custom oligonucleotides. N.G. wrote the software for encoding and decoding information into/from DNA and analysed the data. N.G., E.B., C.D. and B.S. modelled the scaling properties of DNA storage. N.G. wrote the paper with discussions and contributions from all other authors. N.G. and C.D. produced the figures.

Author Information Data are available at <http://www.ebi.ac.uk/goldman-srv/> DNA-storage and in the Sequence Read Archive (SRA) with accession number ERP002040. Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: details are available in the online version of the paper. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to N.G. (goldman@ebi.ac.uk).

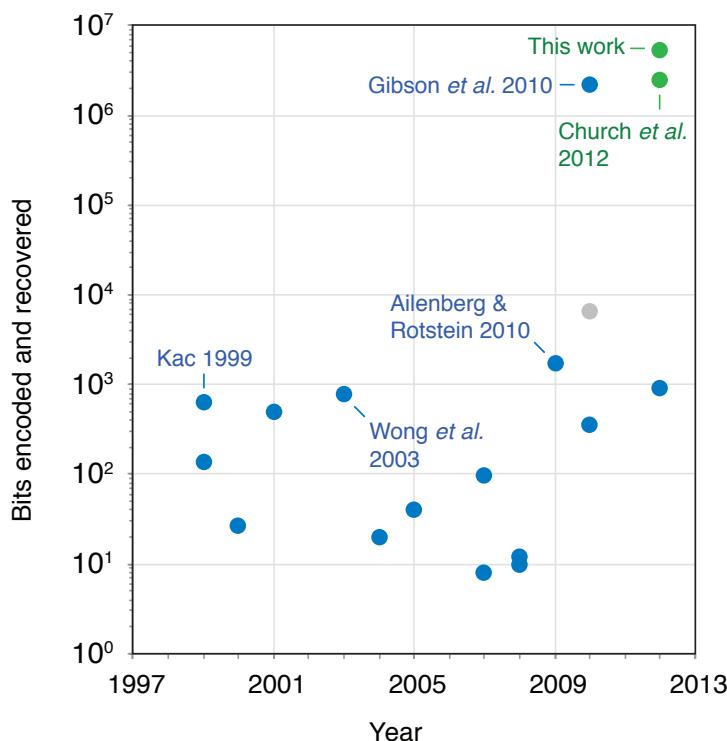
HUMAN-DESIGNED INFORMATION STORED IN DNA

Supplementary Table 1 and Supplementary Fig. 1 show the amounts of human-designed information stored in DNA and successfully recovered in this letter and 16 previous studies. Supplementary Fig. 2 illustrates some of the information encoded in this study. The Shannon information content¹⁰ of the designed messages was approximated by the minimum number of bits required to encode the message using any of the following methods:

- compress ASCII file containing the message in natural form, using Unix command `gzip --best`
- compress ASCII file containing the message in natural form, using Unix command `bzip2 --best`
- for DNA sequence, 2 bits per base
- for simple English text, 5 bits per character (permits use of $2^5 = 32$ characters, e.g. 26 letters of the alphabet plus space and simple punctuation)
- for English/Latin text using reduced or extended alphabets, the number of bits per character is calculated similarly (e.g. 3 bits per character for an alphabet of 8 = 2^3 characters, 6 bits/char for a 64-character alphabet)

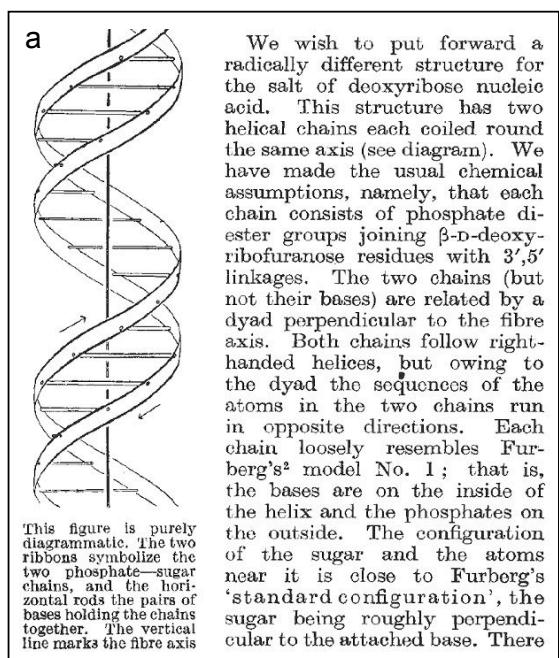
Ref.	Authors	Year	Message type	Message length	Bases used	Shannon information (bits)	Notes
5 Clelland et al.	1999	English text	23 characters	69	138		
6 Kac	1999	English text	129 characters	360	645	Biblical quotation encoded in mutating <i>E. coli</i> genome as a work of art; decoded with 3 character errors, attributed to mutation	
29 Leier et al.	2000	three 9-bit numbers	27 bits	810	27		
30 Bancroft et al.	2001	English text	106 characters	318	504		
31 Wong et al.	2003	English text (64 character alphabet)	185 characters*	560	800 *estimated		
32 Arita & Ohashi	2004	English text	4 characters	24	20		
33 Kashimawura et al.	2005	DNA string	20 bases	65	40		
34 Skinner et al.	2007	four 2-bit numbers	8 bits	231	8		
35 Yachie et al.	2007	mathematical equation and date (256 character alphabet)	12 characters	250	96		
36 Heider & Barnekow	2008	English text	2 characters	5	10		
37 Portney et al.	2008	Latin text (using 8 character alphabet)	5 characters	80	12	final character of 5 (i.e. 3 bits) lost in decoding	
7 Ailenberg & Rotstein	2009	English text, simple musical notation, simple line-drawing notation	349 characters	844	1715		
38 CUHK-iGEM	2010	English text	70 characters	438	350		
8 Gibson et al.	2010	bacterial genome with additional "watermark" sequences (see below)	1077947 bases	1077947	2155894	decoded with 8 base errors and two insertions of 768 and 85 bases, respectively	
of which:		"watermarks": English text plus programming symbols (64 character alphabet)	1280 characters	4658	6504		
9 Church et al.	2012	English text, JPEG images, computer code, all within HTML encoding	658776 characters (bytes)	6313270	2495760	decoded with 10 bit errors	
39 Jarvis & NPC	2012	English text	180 characters	540	900	Article 1 of the Universal Declaration of Human Rights encoded in <i>E. coli</i> genome as work of art	
Goldman et al.	this letter	Total	757051 bytes	17940195	5165800		
of which:		English text (all 154 Shakespeare sonnets)	107738 characters (bytes)	2533635	297856	file wssnt10.txt (from Project Gutenberg, http://www.gutenberg.org/ebooks/1041)	
		PDF document (Watson and Crick, 1953)	280864 bytes	6659172	2119848	file watsoncrick.pdf (from the Nature website, http://www.nature.com/nature/dna50/archive.html , modified to achieve higher compression); see Supplementary Fig. 2a	
		MP3 audio file (extract from Martin Luther King "I Have a Dream" speech)	168539 bytes	3997773	1227176	file MLK_excerpt_VBR_45-85.mp3 (from http://www.americanrhetoric.com/speeches/mlkihaveadream.htm , modified to achieve higher compression: variable bit rate, typically 48–56 kbps; sampling frequency 44.1 kHz)	
		JPEG 2000 image file (image of EBI, 640 x 480 pixels, 16.7M colours)	184264 bytes	4379076	1474000	file EBI.jpg2 (authors' own picture); see Supplementary Fig. 2b	
		ASCII file (Huffman code used to convert bytes to base-3; human readable)	15646 bytes	370539	46920	file View_huff3.cd.new	

Supplementary Table 1 | Amounts of human-designed information stored in DNA and successfully recovered. Message length uses the natural measurement according to the Message type. Bases used indicates the number of DNA bases designed to contain a single copy of the encoded message and ignores the number of copies synthesized.



Supplementary Figure 1 | Amounts of human-designed information stored in DNA and successfully recovered.

Information content is measured in bits; note the logarithmic scale on the y-axis. Blue points indicate studies not adapted to high-throughput data storage; green indicates high-throughput methods. The grey point indicates that part of the Gibson *et al.* (2010) experiment⁸ that encoded information of non-biological origin.



Supplementary Figure 2 | Digital information encoded in DNA. a, An excerpt from the Watson and Crick (1953) paper¹⁸ (PDF format) and **b**, a digital photograph of the European Bioinformatics Institute (JPEG 2000 format) that were among the files encoded in DNA and successfully recovered in this study.

SUPPLEMENTARY METHODS

Digital information encoding. Five files of digital information stored on a hard disk drive were encoded using purpose-written computer software. Each byte of each file to be encoded was represented as a sequence of DNA bases via base-3 digits ('trits' 0, 1 and 2) using a purpose-designed Huffman code¹⁰. Each of the 256 possible bytes was represented by five or six trits; the Huffman code is given in Supplementary File huffman.pdf. Next, each trit was encoded as a DNA nucleotide selected from the three nucleotides different from the last one used, to exclude homopolymer runs. The resulting DNA sequence was converted to segments of length 100 bases, each overlapping the previous by 75 bases, to give strings of a length that was readily synthesized and to provide fourfold redundancy (each DNA base is included in four different segments). Alternate segments were reverse complemented. Indexing information, comprising two trits for file identification (permitting up to $3^2 = 9$ files to be distinguished, in this implementation), 12 trits for intra-file location information (permitting up to $3^{12} = 531,441$ locations per file, i.e. a total of up to $3^{14} = 4,782,969$ unique data locations) and one parity-check¹⁰ trit, again encoded as non-repeating DNA nucleotides, was appended to the 100 information storage bases. Each indexed DNA segment had one further base added to each end, consistent with the 'no homopolymers' rule, that would indicate whether the entire fragment was reverse complemented during the 'reading' stage of the experiment. A full formal specification of the digital information encoding scheme is given in Supplementary File file2features.pdf. In total, the five files were represented by a total of 153,335 strings of DNA, each comprising 117 nt ($= 1 + 100 + 2 + 12 + 1 + 1$) to encode original digital information plus indexing information. The fourfold redundancy provides simple but effective error correction: as each base is encoded in four of the DNA segments, two of which are reverse complemented, any systematic or chance errors in synthesis or sequencing may be corrected by majority vote or more complex decoding schemes. We used simple majority voting (see below).

The data-encoding component of each string can contain Shannon information at 5.07 DNA bases per byte (i.e. $8/5.07 = 1.58$ bits per base), close to the theoretical optimum capacity of 5.05 bits per DNA base (see huffman.pdf) for base-4 channels with runlength limited to 1 (i.e. no repeated nucleotides). After error-correction redundancy and addition of indexing and parity-check information, the data content of our encoding scheme was 4.94 bytes per string ($= 757,051/153,335$), or 0.0422 bytes per base ($= 4.94/117$) (i.e. 23.70 bases per byte). The 153,335 designed DNA strings are available online at <http://www.ebi.ac.uk/goldman-srv/DNA-storage>.

Our indexing scheme, with 14 nt per string available to record file identification and intra-file location, is easily extended by the addition of further indexing nucleotides. This is considered below, in our analysis of the scaling properties of our DNA-storage scheme. Increasing the number of indexing trits (and therefore bases) used to specify file and intra-file location by just two, to 16, gives $3^{16} = 43,046,721$ unique locations, in excess of the 16.8M that is the practical maximum for the Nested Primer Molecular Memory (NPMM) scheme^{40,16}. While these indexing schemes share the aim of encoding which part of a larger total message any one string contains, ours is simpler and more-readily-extensible as it does not incorporate any system by which the indexing information is used to physically extract a subset of the information-bearing strings prior to decoding, as the NPMM scheme does.

DNA synthesis. The synthesis process was also used to incorporate 33 nt paired-end adapter sequences at the 5' and 3' ends of each oligonucleotide (oligo) to facilitate PCR amplification and sequencing on the Illumina platform:

- 5' adapter: 5' -ACACTCTTCCCTACACGACGCTCTTCCGATCT-3'
- 3' adapter: 5' -AGATCGGAAGAGCGGTTCAGCAGGAATGCCGAG-3'

The 153,335 DNA oligo designs were synthesized in three distinct runs (with oligos randomly assigned to runs) using an updated version of Agilent Technologies' OLS (oligo library synthesis) process described previously^{41,20}. This adapts the phosphoramidite chemistry developed previously⁴² and employs inkjet printing and flowcell reactor technologies in the SurePrint *in situ* microarray synthesis platform. Inkjet printing within an anhydrous chamber allows the delivery of very small volumes of phosphoramidites to a confined coupling area on a 2D planar surface, resulting in the addition of hundreds of thousands of bases in parallel. Subsequent oxidation and detritylation are carried out in a flowcell reactor. Once DNA synthesis has been completed, the oligos are then cleaved from the surface and deprotected⁴³.

Up to ~99.8% coupling efficiency is achieved by using thousands-fold excess of phosphoramidite and activator solution. Similarly, millions-fold excess of detritylation agent drives the removal of the 5'-hydroxyl protecting group to near-completion. A novel controlled process in the flowcell reactor significantly reduces depurination, the most prevalent side reaction²⁰. With the latest platform, up to 244,000 unique sequences are synthesized in parallel and delivered as ~1–10 pmol pools of oligos. This is equivalent to $\sim 2.5\text{--}25 \times 10^6$ oligos for each designed sequence ($= 1\text{--}10 \times 10^{-12} \times 6.02 \times 10^{23}/244,000$). In our experiment, three runs were used to synthesize 153,335 designs, leading to the higher figure of $\sim 12\text{--}120 \times 10^6$ ($= 3\text{--}30 \times 10^{-12} \times 6.02 \times 10^{23}/153,335$).

Error rates in the Agilent OLS process are approximately 1 per 500 bases synthesized⁴¹ and synthesis errors are believed to occur independently in different oligos (SC and EML, unpublished data). Combined with our data encoding scheme, this gives further error tolerance. The probability that a given oligo is synthesized entirely correctly is ~0.79 ($= (1 - 1/500)^{117}$), giving a large pool of correct oligos; oligos with a small number of errors in their 100 nt data region may also contribute to correct decoding, with the majority of positions contributing correct information and a small number of errors being outweighed by contributions from other reads (see below).

Library preparation and sequencing. The three samples of lyophilised oligos were resuspended in Tris buffer to a concentration of 5 ng/ml. Samples were then purified from residual synthesis by-products on Ampure XP paramagnetic beads (Beckman Coulter). The reconstituted oligo library was amplified in a total of 22 cycles using thermocycler conditions selected for even A/T vs. G/C processing⁴⁴. PCR was performed with high-fidelity AccuPrime reagents (Invitrogen), a combination of Taq and *pyrococcus* polymerases with a thermostable accessory protein, and paired-end PCR primers (Illumina) complementary to the synthesized adapter sequences flanking each DNA-storage oligo to incorporate additional sequences necessary for flowcell attachment. PCR amplification enabled enrichment for full-length oligos with both 5' and 3' adapters correctly synthesized, and allowed us to achieve appropriate concentration for sequencing while simultaneously incorporating the additional sequences necessary for flowcell attachment and cluster formation. The amplified library products were bead-purified and quantified on the Agilent 2100 Bioanalyzer (concentration

determined to be 15.1 ng/μl, i.e. 86 nM given a peak construct size measured at 270 bp and approximating 650 pg/pmol per bp), diluted to a concentration of 16 pM for flowcell loading and sequenced in paired-end mode on the Illumina HiSeq 2000. The sequencing reaction consumed ~0.1% of the DNA in the initial library: 337 pg of DNA (120 μl at 16 pM) from a starting value of 302 ng (20 μl at 86 nM). Further details of the sample preparation are given in Supplementary Table 2.

Sample stage	Vol (μl)	Conc (ng/μl)	Conc (nM)	Amount of DNA (pg)	Description
A	20	15.1	86	302000	PCR products, quantified on the Agilent 2100 Bioanalyzer
B	2	15.1	86	30200	2 μl of (A) extracted for sequencing; remaining 18 μl reserved for future analysis
C	17.2	1.76	10	30200	(B) then diluted in Tris to 10 nM concentration
D	2	1.76	10	3512	2 μl of (C) retained (remainder discarded in this experiment)
E	20	0.176	1	3512	(D) then denatured in 100mM NaOH, giving dilution to 1 nM
F	16	0.176	1	2809	16 μl of (E) retained (remainder discarded in this experiment)
G	1000	0.00281	0.016	2809	(F) then diluted to 1000 μl
H	120	0.00281	0.016	337	120 μl of (G) sent for clustering (remainder discarded in this experiment)

Supplementary Table 2 | Sample preparation details. Products at stage A (blue) were measured; other values were computed from these.

Base calls were computed from observed intensities using the AYB software²¹, producing 79.6M read-pairs of 104 bases in length. (Illumina's base calling software Bustard produced 65.9M read-pairs, 17.2% fewer than AYB, but led to qualitatively identical decoding results.) Quality control of the reads was performed using FastQC (version 0.10.1; ref. 45). Overall the QC report (available as Supplementary File *FastQC.pdf*) indicated a high-quality sequencing run. Per-cycle quality scores were as expected for an Illumina HiSeq run. The mean quality (Q) score was 36.7, with 95% of quality scores \geq Q30. The GC content of the sequenced reads and the *k*-mer frequencies along the reads were consistent with the structure of the designed DNA strings. The read duplication levels were high, in concordance with the design of the library providing many reads covering any single string.

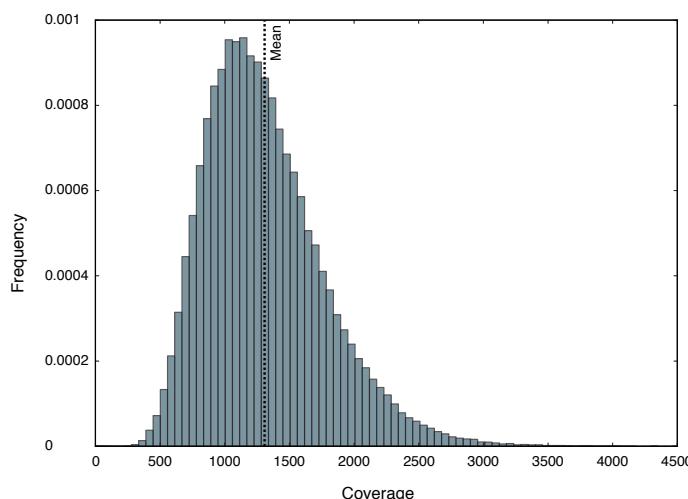
As further quality assessment, but not used for subsequent digital information decoding, the reads were aligned to the designed DNA strings using BWA version 0.6.1-r104 in paired-end mode⁴⁶. Per-cycle error rates were calculated from the resulting alignment using the ErrorRatePerCycle functionality of the GATK package (version 2.1-8-g5efb575; ref. 47). Per-cycle error rates were as expected for an Illumina HiSeq run (see Supplementary File *GATK.txt*) and the mean error rate of 0.001774 after 12.81% unmappable reads were discarded is in line with the combination of current estimates of synthesis error (1 base in 500, above) and sequencing error (1 base in 1,000; ref. 48).

Digital information decoding. As the central 91 bases of each oligo were sequenced from both ends, rapid computation of full-length (117 base) oligos and removal of reads

inconsistent with our designs was straightforward. Sequencing reads were decoded using purpose-written software that exactly reverses the encoding process. The numbers of reads used in different stages of the information decoding process are given in Supplementary Table 3. At the final stage of decoding, the five files were reconstructed from 50.1M strings, giving a mean sequencing depth of 1,308 \times coverage (standard deviation 459). Supplementary Fig. 3 shows the distribution of sequencing depths over encoded data locations (bases of the files' DNA representations). Virtually every location within each decoded file was detected in hundreds or thousands of different sequenced DNA oligos.

Analysis stage	Number of reads	% of total	% of previous stage	Notes	Possible reasons for losses relative to previous analysis stage
A	79564267	100		read-pairs from AYB base caller	
B	55047046	69.19	69.19	117 nt fragment reads recovered from combining 104 nt paired-end reads with 91 nt overlap as expected, with at most 6 mismatches within the overlap region	synthesis error, sequencing error, contamination
C	50145113	63.02	91.10	reads with indexing information indicating they belong to one of the five files encoded in this experiment	synthesis error or sequencing error leading to dinucleotide repeat, invalid file identification in indexing information or parity-check failure
<i>of which:</i>					
	18270252	22.96	33.19	file watsoncrick.pdf	
	8064484	10.14	14.65	file wssnt10.txt	
	11966357	15.04	21.74	file EBI.jp2	
	802908	1.01	1.46	file View_huff3.cd.new	
	11041112	13.88	20.06	file MLK_excerpt_VBR_45–85.mp3	
D	50141326	63.02	99.99	reads contributing 'votes' to final decoding of file	synthesis error or sequencing error in indexing information leading to invalid location in file
<i>of which:</i>					
	18269250	22.96	99.99	file watsoncrick.pdf	
	8063761	10.13	99.99	file wssnt10.txt	
	11965715	15.04	99.99	file EBI.jp2	
	802414	1.01	99.94	file View_huff3.cd.new	
	11040186	13.88	99.99	file MLK_excerpt_VBR_45–85.mp3	

Supplementary Table 3 | Numbers of reads used during decoding.



Supplementary Figure 3 | Distribution of sequencing depths over encoded locations. The mean is at a coverage of 1,308; the standard deviation is 459.

Majority voting was used to resolve any discrepancies caused by DNA synthesis or sequencing errors. The error rate amongst the ‘votes’ used to reconstruct the five files was 0.004004 (20.08M errors in 5,014M bases counted). This is higher than the combined synthesis and sequencing error rate reported above because of cases where one or a small number of errors in indexing information led to a read being misplaced in its correct file, or placed in an incorrect file, generating on average 75 incorrect votes (100 misplaced votes, each with probability ~0.75 of being incorrect).

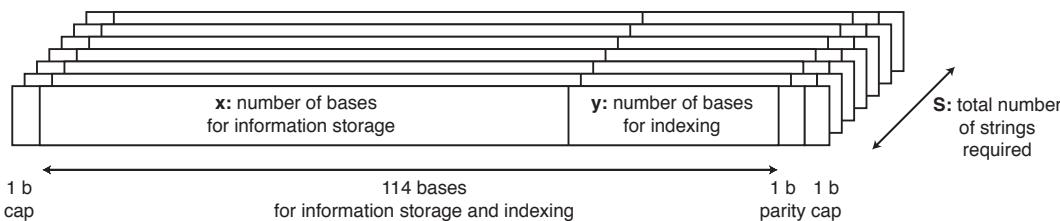
On completion of this procedure, four of the five original files were reconstructed perfectly. The fifth file required manual intervention to correct two regions each of 25 bases that were not recovered from any sequenced read, as described below.

Scaling properties of the DNA-storage scheme. In the following sections we demonstrate that, even constrained to today’s technology, the scheme presented here scales nearly linearly, well beyond any realistically needed range, i.e. beyond 20 orders of magnitude larger than the estimated global amount of digital data of 3 ZB (3×10^{21} bytes).

The global data volume was estimated by adding the 1.8 ZB estimate of data produced in 2011 (ref. 22) to the estimated production in 2010, calculated assuming a doubling time of two years²², to give 3 ZB ($= 1.8 + 1.8/\sqrt{2}$).

As each 117 base string can be synthesized and sequenced independently, it is reasonable to assume that the costs associated with these processes is linear in the number of strings. Therefore, to demonstrate the scaling behaviour of our scheme, we show that the number of strings required to reliably store the data increases nearly linearly with the amount of data. First, we focus on the relationship between the number of strings required and the amount of data to be stored. Second, we show that increase the amount of data to be stored does not lead to higher error rates. This is achieved by both theoretical and empirical estimates of the error rate as a function of amount of data and sequencing coverage.

Scaling of the total number of strings required. Let I be the information to be stored, in bytes. We now show that the number of 117 nt strings required to encode these data scales nearly linearly with respect to I . Recall that in each 117 nt string, 114 bases can be used to store data and indexing information. As the amount of data to be stored increases, we may need to use more than the current 14 bases for indexing. Supplementary Fig. 4 shows how, in general, the strings may be partitioned into x data bases and y indexing bases. Note that y depends on the total number of strings to be used, S . Optimally, $y = \lceil \log_3(S) \rceil$ and therefore $x = 114 - \lceil \log_3(S) \rceil$.



Supplementary Figure 4 | Schematic representation of information encoding in DNA. Multiple strings are used, each comprising 117 bases of which x may be used for storing information and y for indexing, with $x + y = 114$.

Intuitively, we can anticipate that the relationship between S and I is not linear: as the amount of information increases, the number of strings required also increases; this in turn requires more indexing information (and thus greater y), which leaves fewer bases to store information in each string (thus smaller x).

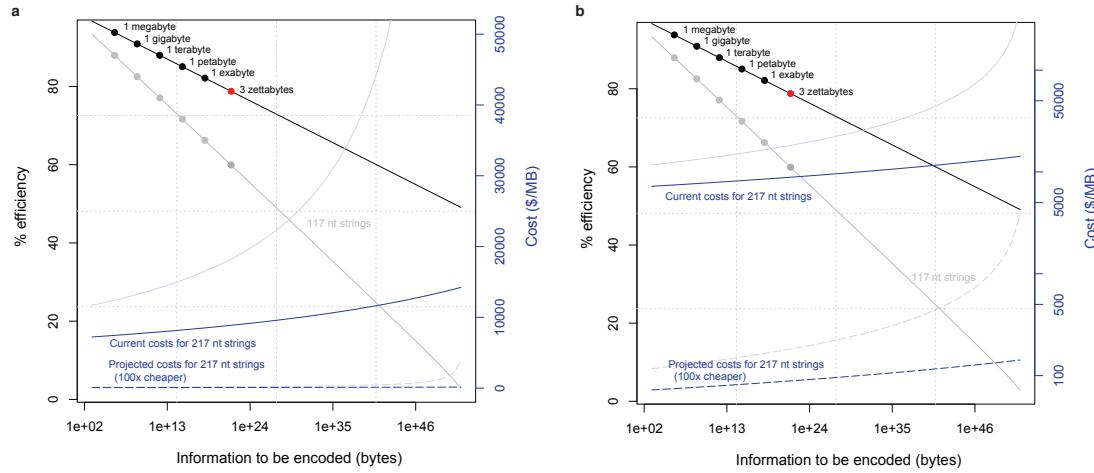
Expressed in bases, the information to be stored is $B = 5.07I$, as we can encode 5.07 bases/byte (see above). Since the encoding strings are 75% overlapping, the relationship between B and S is:

$$B = Sx / 4 = S(114 - \lceil \log_3(S) \rceil) / 4 \quad (1)$$

This can be solved for S numerically. Fig. 2a (Main Text) illustrates the relationship of information stored (I) and the efficiency of encoding, measured as the proportion of bases synthesized (data and indexing) that are used to hold data ($eff = B / 117S$). Two features become apparent. First, the proportion of the total DNA available for encoding data decreases slowly, and is reasonable across the entire relevant data size range. Second, even constrained to 117 base-long strings, the current encoding scheme makes it possible to encode > 20 orders of magnitude more data than is currently practically relevant. (The theoretical limit comes when every base is required for indexing, and none remains to store information.)

DNA-storage costs are affected by the efficiency achieved for different information volumes. Current costs are about \$12,400/MB (below), based on our efficiency of $B / 117S = 0.88$. Costs scale as the inverse of efficiency, and Fig. 2a (Main Text) also shows the cost function $12,400(0.88 / eff)$.

We repeated the above calculations based on longer synthesized strings. The Agilent OLS process can already produce 300-base oligos, with 244,000 designed strings costing approximately \$30,000. Assuming that this would provide 217 nt strings for DNA-storage, an increase of 100 nt, we get twice as many available bases for 30/25 times the price. This would give a current cost of about \$7,440/MB ($= 12,400 \times (30,000 / 25,000) / 2$) based on an efficiency of 0.94 (the encoding efficiency that would be achieved repeating our experiment with 217 nt strings) and consequently a cost function of $7,440(0.94 / eff)$. Supplementary Fig. S5 repeats the information of Fig. 2a (Main Text), and adds the corresponding results for these longer strings. This shows that with achievable improvement in DNA synthesis technology the scalability of DNA-storage is substantially improved, with higher efficiency, lower cost and slower decline in efficiency and increase in cost for larger data volumes.



Supplementary Figure 5 | Scaling properties of DNA-storage. The graphs show how encoding efficiency and costs change as the amount of stored information increases, for longer strings with 217 nt available for data and indexing. The x-axis (logarithmic scale) represents the total amount of information to be encoded. Common data scales are indicated, including the 3 ZB global data estimate. The black line (y-axis scale to left) indicates encoding efficiency, measured as the proportion of synthesized bases available for data encoding. The blue curves (y-axis scale to right) indicate the corresponding effect on encoding costs, both at current synthesis cost levels (solid line) and in the case of a two-order of magnitude reduction (dashed line). The pale grey and pale blue lines give the corresponding results for 117 nt strings, for ease of comparison with Fig. 2a (Main Text).

a, Linear cost scale; **b**, logarithmic cost scale.

Scaling of the decoded data error rate. Assuming that all synthesized strings have the same probability of being sequenced, the mean error rate per base of encoded data depends on three variables:

- ϵ : the mean error rate per base at the level of a sequencing read (due to synthesis and sequencing error), set to 0.004 as determined for our experiment (above)
- S : the total number of designed strings
- c_B : the base coverage (mean number of times each base of encoded information is sequenced)

Recall that the decoding scheme calls each base of encoded information based on a majority vote of all the read bases corresponding to its position. Because each base of encoded information is represented in four strings (due to the 75% overlap in encoded data between neighbouring strings), the mean string coverage is $c_S = c_B / 4$. Thus, in total, there are $S_{c_S} = S c_B / 4$ reads. The probability that any read covers base i of encoded information is $4 / S$. Thus, the number x_i of base reads for encoded base i follows a binomial distribution with number of trials $S c_B / 4$ and probability of success $4 / S$, which we write as $B(S c_B / 4, 4 / S)$.

Next, consider that for encoded base i to be correctly called, the majority of the x_i read bases (votes) need to be correct. The distribution of correct bases $x_{i,\text{correct}}$ is $B(x_i, 1 - \epsilon)$. The majority vote regarding encoded base i is wrong if $x_{i,\text{correct}} < x_i / 2$. (This is the worst case scenario that all incorrect votes are for the same incorrect base. Our results are not significantly altered when other, more-favourable, scenarios are considered.) Because of this

dependency, the expected encoded base error rate is not straightforward to compute analytically, but it is easily estimated by Monte Carlo simulation. Supplementary Fig. 6 shows an R function that performs this estimation.

```
# R function to estimate the encoded base error rate:
estimateBaseErrorRateMC = function(eps,S,cS,nsamples) {
  x = rbinom(nsamples, cS*S ,4/S);
  s = sapply(x,function(t) rbinom(1,t,1-eps));
  e = 1 - sum (s > x/2) / nsamples;
  return(e);
}
```

Supplementary Figure 6 | R function to estimate encoded base error rate.

Supplementary Table 4 provides estimates of the encoded base error rate, as a function of data size and sequencing effort (per encoded byte) relative to our experiment, based on 10^5 Monte Carlo samples (`nsamples` in Supplementary Fig. 6) per cell. This shows that, keeping sequencing effort per encoded byte constant, the error rate increases only very slowly with the increase in amount of data.

		Information size in byte															
		1000	1MB	1GB	1TB	1PB	1EB	1ZB	1.00E+24	1.00E+27	1.00E+30	1.00E+33	1.00E+36	1.00E+39	1.00E+42	1.00E+45	1.00E+48
% sequencing effort (per encoded byte) relative to experiment	1000%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	316%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	100%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	31%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	10%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	3%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	1%	0	0	1.00E-05	4.00E-05	3.00E-05	7.00E-05	9.00E-05	0.00016	0.00054	0.00118	0.00317	0.00657	0.0145	0.03268	0.07506	0.17551
	0.3%	0.01221	0.01764	0.02209	0.02725	0.03568	0.04084	0.05129	0.06733	0.08702	0.11427	0.14752	0.19205	0.24854	0.32076	0.42418	0.5758
	0.1%	0.24756	0.2693	0.29461	0.31652	0.34522	0.34376	0.37356	0.40983	0.44767	0.50255	0.54485	0.59052	0.63929	0.69121	0.7514	0.81493
	0.03%	0.63294	0.66131	0.67665	0.69713	0.71112	0.70758	0.72623	0.74469	0.76418	0.78336	0.80058	0.8224	0.84375	0.86439	0.88685	0.90861
	0.01%	0.859	0.87701	0.88289	0.89178	0.90067	0.87071	0.8787	0.88407	0.89109	0.89681	0.90474	0.91015	0.91856	0.9251	0.93250	0.94003
	0.003%	0.9575	0.95954	0.96263	0.96321	0.96711	0.92772	0.9297	0.9313	0.93422	0.9357	0.93776	0.94094	0.9441	0.94557	0.94768	0.94916
	0.001%	0.97811	0.98709	0.98762	0.9881	0.98941	0.99047	0.94677	0.94583	0.9464	0.94967	0.94959	0.94987	0.95178	0.95045	0.95207	0.95358

Supplementary Table 4 | Error rate as a function of data size and sequencing effort.

Percentage sequencing effort is measured relative to our experiment; the highlighted row corresponds to the same sequencing effort per encoded byte as realised in our experiment.

In contrast, the error rate depends strongly on the coverage. Supplementary Table 4 suggests that the effective coverage of our actual experiment ($1,308\times$; see above) could be substantially lowered without impacting on the error rate. To confirm this theoretical analysis using our empirical data, we subsampled the 79.6M read-pairs at varying fractions and attempted to reconstruct the five encoded files using our original protocol. Fig. 2b (Main Text) presents results on the per-encoded-base error rate for recovery of the file `watsoncrick.pdf`, which always has at least 50 base errors due to encoded bases that were not recovered from any sequenced read (see below), for the recovery of the other four encoded files combined and for our theoretical predictions based on the analysis above. The plot shows the error rate (y-axis, as a percentage) as a function of subsampling percentage (x-axis, logarithmic scale). This indicates good agreement of our theoretical and empirical

results. The difference between the `watsoncrick.pdf` results and the other four files is explained by the unrecoverable 50 bases described above. In this case, the minimum possible error rate is 0.0036% (10 bytes not recovered, out of 280,864). The discrepancy between the theoretical and empirical curves for relatively high subsampling fractions is probably due to model violation. For example, unlike in our model, the true sampling probability of strings is almost certainly not uniform due to unequal DNA synthesis and sequencing efficiencies. Note however that the discrepancy corresponds to a difference of only a few per cent in terms of reads used, and the model remains a good approximation. The plots confirm that we could reduce sequencing coverage by a factor of 10 or even 100 without significantly impacting on our ability to recover the encoded information.

Modelling cost-effectiveness of DNA-storage. We modelled the costs of DNA-storage over time according to:

$$C_D(t) = D_0 + Ft \quad (2)$$

where $C_D(t)$ is the cost to archive 1 MB of information for a period of t years. D_0 is the initial cost to write this information to DNA-storage — this will decrease over time as DNA synthesis technology improves — and F is the cost per year to maintain a DNA-storage facility (per MB of information stored). Storage on magnetic tape was modelled according to:

$$C_T(t) = T_0 + Ft + \sum_{i=1}^{ft} \left[R_{\text{fix}} + \frac{T_0 + R_{\text{dim}}}{\left(2^{f^{-1}/2.5}\right)^i} \right] \quad (3)$$

where $C_T(t)$ is the cost to archive 1 MB. T_0 is the initial cost to write this information to tape and F is the cost per year to maintain a tape storage facility (e.g. data centre), assumed equal to the corresponding cost for a DNA-storage archive. This assumption is likely to strongly favour tape over DNA due to the costs of power and of recurrent replacement of computing and tape hardware. The parameter f is the frequency of ‘tape transfer events’, i.e. how often it is necessary to read and re-write the information using the current technology as the previous one becomes obsolete. Industry standards suggest f is likely to be approximately $1/5-1/10 \text{ yr}^{-1}$, i.e. data must be read and re-written to new technology every 5–10 years^{49,50}. The summation represents the costs (per MB) of the ft transfer events occurring in t years, each comprising fixed cost R_{fix} (e.g. finite labour cost of retrieval of existing tape archive, set-up of copying process, storage of new archive material) and diminishing costs proportional to T_0 (for new storage media) and R_{dim} (other expenses, e.g. costs proportional to time spent reading and re-writing information), both of which are assumed to halve every 2.5 years due to technological improvements⁵¹.

Tape costs are already very low, and so we set $T_0 = 0$. The break-even point when DNA-storage achieves the same cost as tape comes when $C_D(t) = C_T(t)$; equating equations (2) and (3) allows us to write:

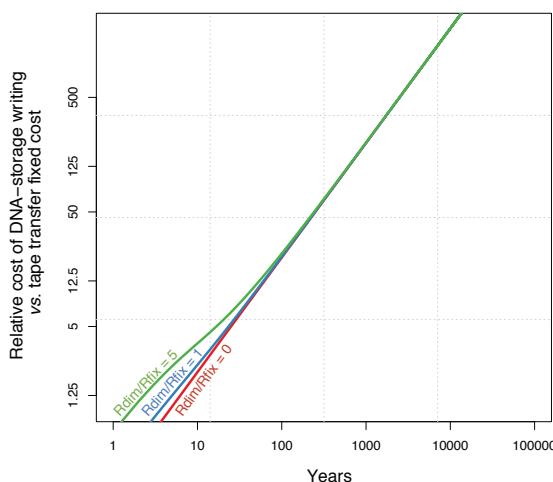
$$\frac{D_0}{R_{\text{fix}}} = ft + \frac{R_{\text{dim}}}{R_{\text{fix}}} \left(\frac{1 - 2^{-t/2.5}}{2^{f^{-1}/2.5} - 1} \right) \quad (4)$$

D_0/R_{fix} is the relative cost of writing DNA-storage compared to the fixed costs of tape transfer events. Equation (4) indicates the balance between values of $D_0, f, t, R_{\text{fix}}$ and R_{dim} that leads to

break-even point of DNA-storage; smaller values of D_0/R_{fix} or greater values of t correspond to conditions where DNA-storage is more cost-effective than tape; conversely, larger D_0/R_{fix} or smaller t make tape favourable.

The current commercial cost of the Agilent OLS process is approximately \$25,000 for 244,000 designed oligos of length 200 bases (approximately \$0.05/100 bases). We encoded 739 kB in 153,335 DNA strings of length 117 bases, leading to a value for D_0 of approximately \$12,400/MB ($= 25,000/(0.739 \times (244,000/153,335) \times (200/117))$). For archives of a few megabytes, we estimate that the cost in personnel, labour and management of a corresponding tape technology transition might be of the order of \$25–100, leading to a current estimate of D_0/R_{fix} in the range 125–500. Other current DNA synthesizing methods, e.g. maskless photolithography, can be used to produce shorter oligos with potentially higher error rates⁵², but are less expensive per base synthesized. It is possible that these could be used to reduce the cost per MB (D_0).

For realistic values of f (above), the second term on the right-hand side of equation (4) rapidly becomes small. Supplementary Fig. 7 shows the relationship of D_0/R_{fix} and break-even points (timescale on which DNA-storage and tape storage costs equate) when $R_{\text{dim}}/R_{\text{fix}} = 0, 1$ and 5 and $f = 1/5$. It is clear that even for $R_{\text{dim}}/R_{\text{fix}} = 5$, which is unrealistically large, the effect of $R_{\text{dim}}/R_{\text{fix}}$ is negligible for values of D_0/R_{fix} that are likely to be achieved in the near future (see Main Text). The same is true for $f = 1/10$ (not shown). Consequently, we have assumed $R_{\text{dim}}/R_{\text{fix}} = 1$ for illustrative purposes. Fig. 2c (Main Text) plots D_0/R_{fix} against time t in this case, highlighting the break-even points for $f = 1/5$ and $f = 1/10$.



Supplementary Figure 7 | Effect of $R_{\text{dim}}/R_{\text{fix}}$ on DNA-storage break-even timescale. The x -axis is the break-even time beyond which DNA-storage is less expensive than magnetic tape, assuming the tape archive has to be read and re-written every 5 years ($f = 1/5$); the y -axis is the relative cost of DNA-storage synthesis and tape transfer fixed costs. Lines plotted are for $R_{\text{dim}}/R_{\text{fix}} = 0$ (red), 1 (blue) and 5 (green). Note the logarithmic scales on both axes.

Information decoding costs. In our experiment, we decoded 739 kB of information using one lane of the Illumina HiSeq 2000, at a sequencing cost of approximately \$1,600. This gives a decoding cost of ~\$2,200/MB ($= 1,600/0.739$). As shown above, we could have sequenced 10 times as much encoded information in the same run and still recovered our data. This suggests that ~\$220/MB ($= 2,200/10$) is a reasonable approximation for the decoding costs for optimised use of existing technologies.

SUPPLEMENTARY DISCUSSION

Repair of file with missing reads. During decoding, one file (ultimately determined to be *watsoncrick.pdf*) reconstructed *in silico* at the level of DNA (prior to decoding, via base-3, to bytes) contained two regions, each of 25 bases in length, that were not recovered from any sequenced read. Given the overlapping segment structure of our encoding, each such region indicates the failure of any oligo representing any of four consecutive segments to be synthesized or sequenced successfully, as any one of four consecutive overlapping segments would have contained the bases corresponding to this location. Inspection of the two regions indicated that the non-detected bases fell within long repeats of the following 20-base motif:

5' -GAGCATCTGCAGATGCTCAT-3'

(colours used to highlight motif repeats; see below). We noticed that repeats of this motif have a self-reverse complementary pattern (Supplementary Fig. 8) and we hypothesised that long, self-reverse complementary DNA fragments might not be readily sequenced using the Illumina process. In terms of DNA synthesis, we know no reason to expect self-reverse complementary fragments to be problematic in the Agilent OLS system. The PCR conditions used for library construction involved denaturing the template oligos at 98 °C, initially for a period of 3 min and for 80 s per cycle thereafter. These conditions are more than sufficient to denature and amplify any self-annealing oligos. A further denaturing step separates the double-stranded products prior to dilution and flowcell loading, and bridge amplification should also proceed nominally to form clusters of clonally amplified library constructs.

5' -... GAGCATCTGCAGATGCTCAT GAGCATCTGCAGATGCTCAT GAGCATCTGCAGATGCTCAT ...-3'
.....|||||.....|||||.....|||||.....|||||.....|||||.....|||||.....|||||.....|||||.....|||||.....
3' -... TACTCGTAGACGTCTACGAG TACTCGTAGACGTCTACGAG TACTCGTAGACGTCTACGAG ...-5'

Supplementary Figure 8 | Self-reverse complementary nature of the 20-base motif.

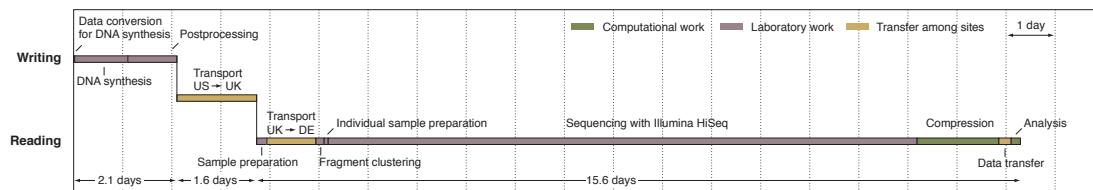
However, the subsequent sequencing conditions are intended to promote DNA hybridisation/annealing, and so in principle internal secondary structures might form and remain stable throughout the process. We therefore reason that each of the two self-complementary regions produced by our encoding scheme led to the formation of stem-loop structures within the target sequences, inhibiting the sequencing-by-synthesis reaction over these nonlinear stretches. The associated clusters would be omitted from the imaging readout, thereby resulting in a gap in the sequencing data when the files were reconstructed.

Examination of our coding methods indicated that such long repeats of the 20-base motif would arise when the original computer file contained repeats of byte value 255 (hexadecimal FF). Consequently, we modified the *in silico* reconstructed DNA sequence to repair the repeating motif pattern and subjected this to subsequent decoding steps. No further problems were encountered, and the final decoded file matched perfectly the file *watsoncrick.pdf*.

With hindsight, we should have devised a code that ensured that no long self-complementary regions existed in any of our designed DNA segments. One way to achieve this would be to pre-process the files to be encoded using a one-time pad or other stream

cipher with a standard or known key stream, leading to the DNA segments having random properties⁵³.

Timescale of DNA-storage experiment. Supplementary Fig. 9 shows the time taken for each stage of our DNA-storage experiment. The experiment was not optimised for speed. All encoding and decoding computations were performed on one core of an Intel i5-2540M processor running at 2.60 GHz, except for the reconstruction of full-length (117 base) oligos from paired-end (104 base) reads which was performed using one core of an Intel Xeon X5650 at 2.67GHz. In a large-scale DNA-storage archive, transfer periods could be eliminated by having encoding, storage and decoding taking place at one site. Both computer software and laboratory procedures could readily be optimised and parallelised; laboratory procedures could also be automated with liquid-handling robotics for high-throughput applications⁵⁴. Both computational and laboratory equipment are subject to continual innovation, improving their speed.



Supplementary Figure 9 | Timeline of DNA-storage experiment. We report only periods of active work on the experiment. We have omitted time taken to devise repairs for the file with two information gaps (above).

Information storage density. We recovered 757,051 bytes of information from 337 pg of DNA (above), giving an information storage density of $\sim 2.2 \text{ PB/g} (= 757,051/337 \times 10^{-12})$. We note that this information density is enough to store the US National Archives and Records Administration's Electronic Records Archives' 2011 total of $\sim 100 \text{ TB}$ (ref. 55) in $< 0.05 \text{ g}$ of DNA, the Internet Archive Wayback Machines's 2 PB archive of web sites⁵⁶ in $\sim 1 \text{ g}$ of DNA, and CERN's 80 PB CASTOR system for LHC data²⁵ in $\sim 35 \text{ g}$ of DNA.

REFERENCES FOR SUPPLEMENTARY INFORMATION

29. Leier, A., Richter, C., Banzhaf, W. & Rauhe, R. Cryptography with DNA binary strands. *Biosystems* **57**, 13–22 (2000)
30. Bancroft, C., Bowler, T., Bloom, B. & Clelland, C. T. Long-term storage of information in DNA. *Science* **293**, 1763–1765 (2001)
31. Wong, P. C., Wong, K.-K. & Foote, H. Organic data memory. Using the DNA approach. *Comm. ACM* **46**, 95–98 (2003)
32. Arita, M. & Ohashi, Y. Secret signatures inside genomic DNA. *Biotechnol. Prog.* **20**, 1605–1607 (2004)
33. Kashiwamura, S., Yamamoto, M., Kameda, A., Shiba, T. & Ohuchi, A. Potential for enlarging DNA memory: the validity of experimental operations of scaled-up nested primer molecular memory. *Biosystems* **80**, 99–112 (2005)

34. Skinner, G. M., Visscher, K. & Mansuripur, M. Biocompatible writing of data into DNA. *J. Bionanoscience* **1**, 17–21 (2007)
35. Yachie, N., Sekiyama, K., Sugahara, J., Ohashi, Y. & Tomita, M. Alignment-based approach for durable data storage into living organisms. *Biotechnol. Prog.* **23**, 501–505 (2007)
36. Heider, D. & Barnekow, A. DNA watermarks: a proof of concept. *BMC Mol. Biol.* **9**, 40 (2008)
37. Portney, N. G., Wu, Y., Quezada, L. K., Lonardi, S. & Ozkan, M. Length-based encoding of binary data in DNA. *Langmuir* **24**, 1613–1616 (2008)
38. CUHK iGEM 2010. Bacterial-based storage and encryption device (2010) http://2010.igem.org/Team:Hong_Kong-CUHK accessed online, 10 May 2012
39. Jarvis, C. & Netherlands Proteomics Centre. Blighted by Kenning (2012) <http://www.artforeating.co.uk/restaurant/index.php?/blighted-by-ken/project-overview/> accessed online, 8 November 2012
40. Yamamoto, M., Kashiwamura, S., Ohuchi, A. & Furukawa, M. Large-scale DNA memory based on the nested PCR. *Natural Computing* **7**, 335–346 (2008)
41. Kosuri, S. *et al.* A scalable gene synthesis platform using high-fidelity DNA microchips. *Nature Biotech.* **28**, 1295–1299 (2010)
42. Beaucage, S. L. & Caruthers, M. H. Deoxynucleoside phosphoramidites — a new class of key intermediates for deoxypolynucleotide synthesis. *Tetrahedron Lett.* **22**, 1859–1862 (1981)
43. Cleary, M. A. *et al.* Production of complex nucleic acid libraries using highly parallel *in situ* oligonucleotide synthesis. *Nature Methods* **1**, 241–248 (2004)
44. Aird, D. *et al.* Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biol.* **12**, R18 (2011)
45. Andrews, S. FastQC. A quality control tool for high throughput sequence data (2012) <http://www.bioinformatics.babraham.ac.uk/projects/fastqc> accessed online, 5 September 2012
46. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009)
47. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010)
48. Minoche, A. E., Dohm, J. C. & Himmelbauer, H. Evaluation of genomic high-throughput sequencing data generated on Illumina HiSeq and Genome Analyzer systems. *Genome Biol.* **12**, R112 (2011)
49. Vican, I. & Stančić, H. Long-term inactive data retention through tape storage technology. In Stančić, H. *et al.* (eds.) *INFUTURE2009: the Future of Information Sciences — Digital Resources and Knowledge Sharing*. pp. 105–114. (Department of Information Sciences, University of Zagreb, 2009)
50. Klingler, S. L. Information storage technologies. In Bates, M. J. (ed.) *Understanding Information Retrieval Systems: Management, Types and Standards*. pp. 245–258. (CRC Press, 2012)
51. Fujifilm. Fujifilm to manufacture 5TB tape cartridge for Oracle's StorageTek T10000C drive (2011) <http://www.fujifilm.com/news/n110201.html> accessed online, 10 May 2012
52. Agbavwe, C. *et al.* Efficiency, error and yield in light-directed maskless synthesis of DNA microarrays. *J. Nanobiotechnology* **9**, 57 (2011)
53. Paar, C. & Pelzl, J. *Understanding Cryptography. A Textbook for Students and Practitioners*. (Springer, 2010)

54. Quail, M. A., Swerdlow, H. & Turner, D. J. Improved protocols for the Illumina genome analyzer sequencing system. *Curr. Protoc. Hum. Genet.* **62**, 18.2.1–18.2.27 (2009)
55. US National Archives & Records Administration (2012) <http://www.archives.gov/era/status.html> accessed online, 10 May 2012
56. Internet Archive. The Wayback Machine (2012) <http://archive.org/about/faqs.php> accessed online, 7 September 2012

Method for Encoding and Decoding Arbitrary Computer Files in DNA Fragments

1 Encoding

1.1: An arbitrary computer file is represented as a string S_\emptyset of bytes (often interpreted as a number between \emptyset and $2^8 - 1$, i.e. a value in the set $\{\emptyset \dots 255\}\}^*$.

1.2: S_\emptyset is encoded using a given Huffman code, converting it to base-3. This generates the string S_1 of characters in $\{\emptyset, 1, 2\}$, each such character called a ‘trit’.

1.3: Write $\text{len}()$ for the function that computes the length (in characters) of a string, and define $n = \text{len}(S_1)$. Represent n in base-3 and prepend \emptyset s to generate a string S_2 of trits such that $\text{len}(S_2) = 2\emptyset$. Form the string concatenation $S_4 = S_1 . S_3 . S_2$, where S_3 is a string of at most 24 \emptyset s chosen so that $\text{len}(S_4)$ is an integer multiple of 25.

1.4: S_4 is converted to a DNA string S_5 of characters in $\{A, C, G, T\}$ with no repeated nucleotides (nt) using the scheme illustrated in Figure 1. (The first trit of S_4 is coded using the ‘A’ row of the table. For each subsequent trit, characters are taken from the row defined by the previous character conversion.)

previous nt written	next trit to encode		
	\emptyset	1	2
A	C	G	T
C	G	T	A
G	T	A	C
T	A	C	G

Figure 1: Base-3 to DNA encoding ensuring no repeated nucleotides. For each trit t to be encoded, select the row labelled with the previous nt used and the column labelled t and encode using the nt in the corresponding table cell.

1.5: Define $N = \text{len}(S_5)$, and let ID be a 2-trit string identifying the original file and unique within a given experiment (permitting mixing of DNA from different

* \emptyset is used throughout to represent the number zero, to avoid confusion with letters o and O.

files S_\emptyset in one experiment). Split S_5 into overlapping segments of length 100 nt, each offset from the previous by 25 nt. This means there will be $\frac{N}{25} - 3$ segments, conveniently indexed $i = \emptyset \dots \frac{N}{25} - 4$; segment i is denoted F_i and contains (DNA) characters $25i \dots 25i + 99$ of S_5 .

Each segment F_i is further processed as follows:

1.6: If i is odd, reverse complement F_i .

1.7: Let $i3$ be the base-3 representation of i , appending enough leading \emptyset s so that $\text{len}(i3) = 12$. Compute P as the sum (mod 3) of the odd-positioned trits in ID and $i3$, i.e. $ID_1 + i3_1 + i3_3 + i3_5 + i3_7 + i3_9 + i3_{11}$. (P acts as a ‘parity trit’—analogous to a parity bit—to check for errors in the encoded information about ID and i .)

1.8: Form the indexing information string $IX = ID . i3 . P$ (comprising $2+12+1 = 15$ trits). Append the DNA-encoded version of IX to F_i using the same strategy as at step (1.4) above, starting with the code table row defined by the last character of F_i , to give indexed segment F'_i .

1.9: Form F''_i by prepending A or T and appending C or G to F'_i —choosing between A and T, and between C and G, randomly if possible but always such that there are no repeated nt. (This ensures that we can distinguish a DNA segment that has been reverse complemented during DNA sequencing from one that has not—the former will start with G|C and end with T|A; the latter will start A|T and end C|G.)

See Figure 2 for a schematic representation of the DNA-encoding of computer files.

1.10: The segments F''_i are synthesized as actual DNA oligonucleotides and stored, and may be supplied for sequencing.

2 Example

2.1: As it is difficult to represent all possible bytes in this document, we use a simple example of a file comprising just 18 bytes that happen to be easily

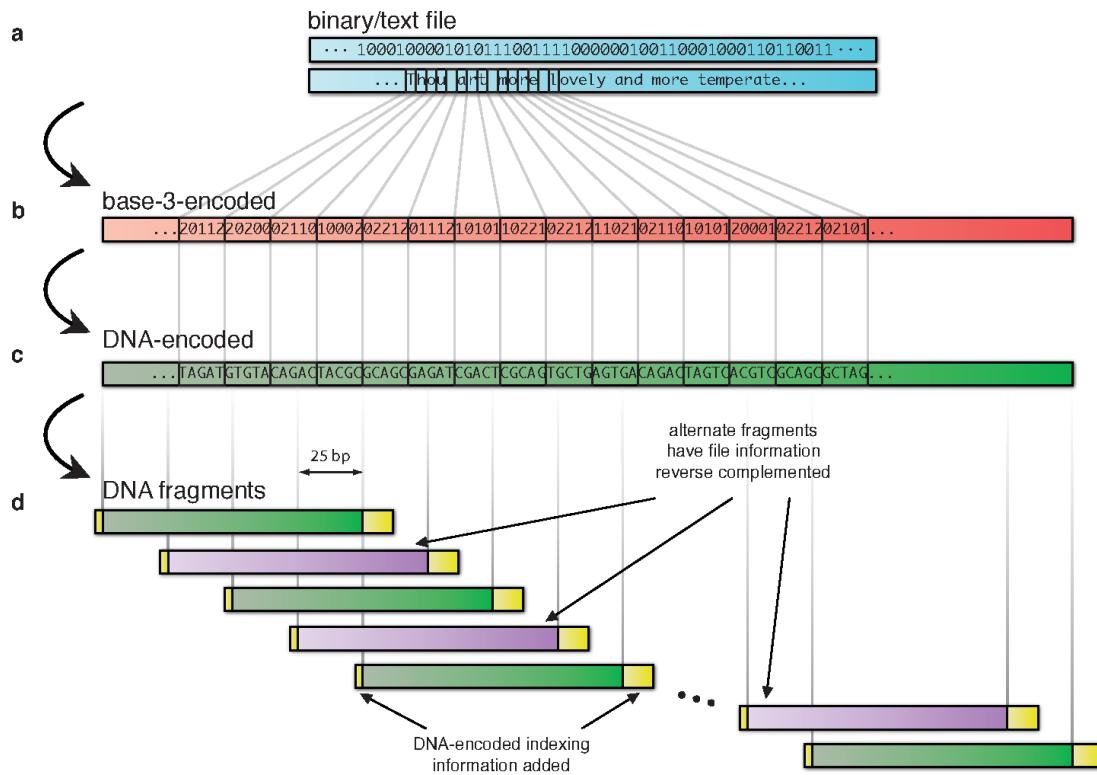


Figure 2: Schematic of DNA-fragment-encoding of computer files. The computer file (in any format, e.g. binary or text) shown in blue in (a) (step 1.1 above) is encoded in base-3 (red in b; 1.2–1.3) and then as DNA with no repeated nt (green in c; 1.4). This representation of the complete file is then split into overlapping segments (1.5), alternate segments reverse complemented (mauve; 1.6), and segment indexing and direction-determining information added (yellow; see d; 1.7–1.9).

represented via ASCII codes (for clarity, we show spaces as \sqcup):

$$S_\emptyset = \text{Birney} \sqcup \text{and} \sqcup \text{Goldman}$$

2.2: Using the Huffman code defined in the example file `View_huff3.cd`, we convert the bytes of S_\emptyset (shown above/below S_1 for illustrative purposes) into base-3[†]:

[†]Throughout what follows, spaces are not part of base-3 or DNA strings but are included to assist in ‘reading’ how strings have been derived in each step.

	B	i	r	n	e	y	u	a	n
$S_1 =$	$2\theta 10\emptyset$	$2\theta 21\emptyset$	$1\theta 1\theta 1$	$\emptyset\theta\theta 21$	$2\theta\theta\theta 1$	222111	$\emptyset 2212$	$\emptyset 1112$	$\emptyset\theta\theta 21$
	$221\emptyset\emptyset$	$\emptyset 2212$	222212	$\emptyset 211\emptyset$	$\emptyset 21\theta 1$	$221\emptyset\emptyset$	$11\emptyset 21$	$\emptyset 1112$	$\emptyset\theta\theta 21$
	d	u	G	o	l	d	m	a	n

2.3: $n = \text{len}(S_1) = 92$, which is 10102 in base-3. So:

2.4: Using the DNA coding strategy and table shown at step (1.4) above, convert S_4 to DNA:

$S_5 = \text{TAGTATATCGACTAGTACAGCGTAGCATCTCGCAGCGAGATAACGCTGCTACGCAGCATGCTGTGAGTATCGATGACGAGTGA} \dots$

2.5: $N = \text{len}(S_5) = 125$, and we choose (e.g.) $ID = 12$. S_5 will be split into overlapping segments F_i of length 100 nt for $i \in \{\emptyset, \dots, \frac{125}{25} - 4\}$, i.e. $i \in \{\emptyset, 1\}$. With overlapping parts underlined for illustration, F_\emptyset and F_1 are:

F_\emptyset = TAGTATATCGACTAGTACAGCGTAG CATCTCGCAGCGAGATAACGCTGCTA
 CGCAGCATGCTGTGAGTATCGATGA CGAGTGA~~T~~TGTACAGTACGTACG

$F_1 = \underline{\text{CATCTCGCAGCGAGATA CGCTGCTA}}$ $\underline{\text{CCGAGCATGCTGTGAGTATCGATGA}}$
 $\underline{\text{CGAGTCACTCTGTACAGTACGTACG}}$ $\text{TACGTACGTACGTACGTACGACTAT}$

2.6: Only $i = 1$ is odd, so F_1 is reverse complemented:

$F_1 = \text{ATAGTCGTACGTACGTACGTACGTACGTACTGTACAGAGTCAC TCG}$
 $\text{TCATCGATACTCACAGCATGCTGCGTAGCAGCGTATCTCGCTGCGAGATG}$

2.7: For $i = \emptyset$, $i3 = \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset$ (length 12) and the sum (mod 3) of the odd-positioned trits of ID and $i3$ is $P = 1 + \emptyset + \emptyset$ (mod 3) = 1. For $i = 1$, $i3 = \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$ and $P = 1 + \emptyset + 1$ (mod 3) = 1.

2.8: For $i = \emptyset$, $IX = ID \cdot i3 \cdot P = 12 \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$; for $i = 1$, $IX = 12 \emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset\emptyset 1$.

So:

$F'_\emptyset = \text{TAGTATATCGACTAGTACAGCGTAGCATCTCGCAGCGAGATA CGCTGCTA}$
 $\text{CGCAGCATGCTGTGAGTATCGATGACGAGTGACTCTGTACAGTACGTACG}$
 AT ACGTACGTACGT C (length $100 + 15 = 115$)

$F'_1 = \text{ATAGTCGTACGTACGTACGTACGTACGTACTGTACAGAGTCAC TCG}$
 $\text{TCATCGATACTCACAGCATGCTGCGTAGCAGCGTATCTCGCTGCGAGATG}$
 AT ACGTACGTACGA G (length 115)

2.9: Prepend A|T and append C|G (note that in this example we have only one random choice, at the end of F''_\emptyset):

$F''_\emptyset = \text{A TAGTATATCGACTAGTACAGCGTAGCATCTCGCAGCGAGATA CGCTGCTA}$
 $\text{CGCAGCATGCTGTGAGTATCGATGACGAGTGACTCTGTACAGTACGTACG}$
 ATACGTACGTACGTC G (length $1 + 115 + 1 = 117$)

$F''_1 = \text{T ATAGTCGTACGTACGTACGTACGTACGTACGTACTGTACAGAGTCAC TCG}$
 $\text{TCATCGATACTCACAGCATGCTGCGTAGCAGCGTATCTCGCTGCGAGATG}$
 ATACGTACGTACGAG C (length 117)

3 Decoding

Decoding is simply the reverse of encoding, starting with sequenced DNA fragments F''_i of length 117 nt. Reverse complementation during the DNA sequencing procedure (e.g. during PCR reactions) can be identified for subsequent reversal by observing whether fragments start with A|T and end with C|G, or start G|C and end T|A. With these two ‘orientation’ nt removed, the remaining 115 nt of each segment can be split into the first 100 ‘message’ nt and the remaining 15 ‘indexing’ nt. The indexing nt can be decoded to determine the file identifier ID and the position index $i\beta$ and hence i , and errors may be detected by testing the parity trit P . Position indexing information permits the reconstruction of the DNA-encoded file, which can then be converted to base-3 using the reverse of the encoding table in step (1.4) above and then to the original bytes using the given Huffman code. Precise details of the decoding procedure are left as an exercise for the reader. Errors introduced during DNA synthesis, storage or sequencing could lead to various artefacts, particularly nt insertion, deletion or substitution. Recovery of information from fragments with such errors may be possible (details left as exercise)—we have not found this to be necessary due to the large numbers of perfectly-sequenced fragments available via high-throughput sequencing.

SUPPLEMENTARY INFORMATION

doi:10.1038/nature11875

FastQC Report

Tue 4 Sep 2012
ain.bam

Summary

Basic Statistics

Per base sequence quality

Per sequence quality scores

Per base sequence content

Per base GC content

Per sequence GC content

Per base N content

Sequence Length Distribution

Sequence Duplication Levels

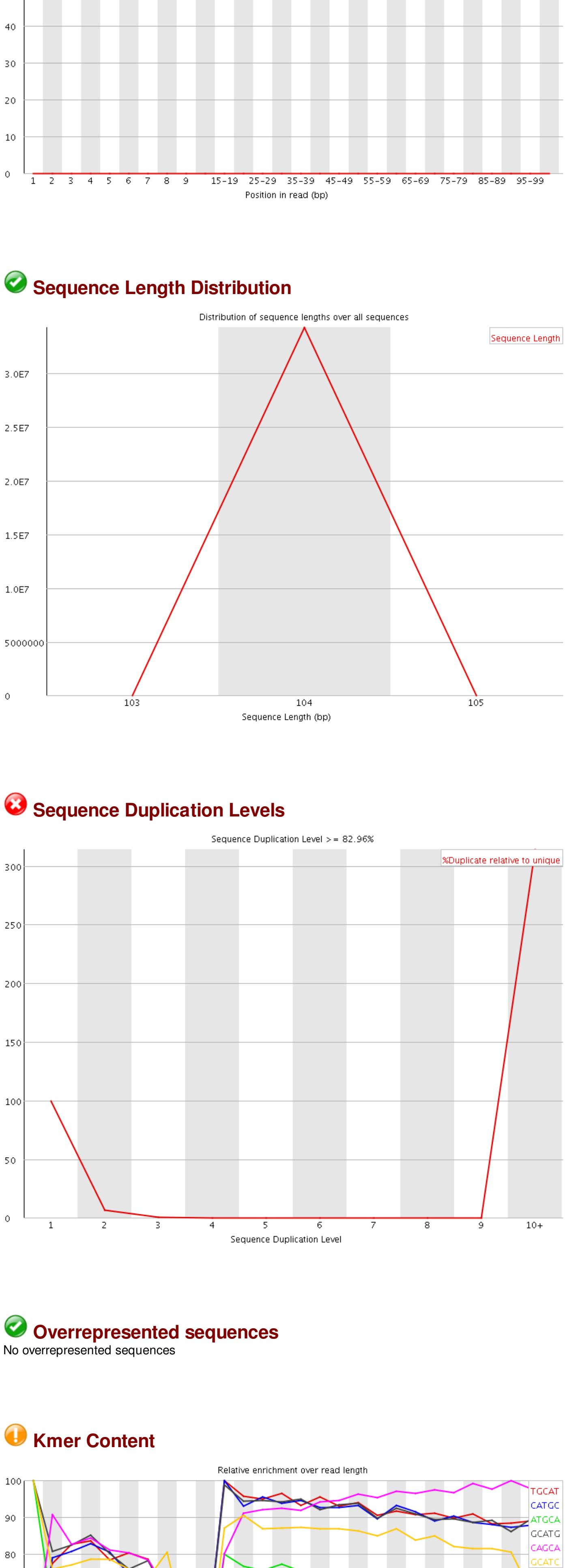
Overrepresented sequences

Kmer Content

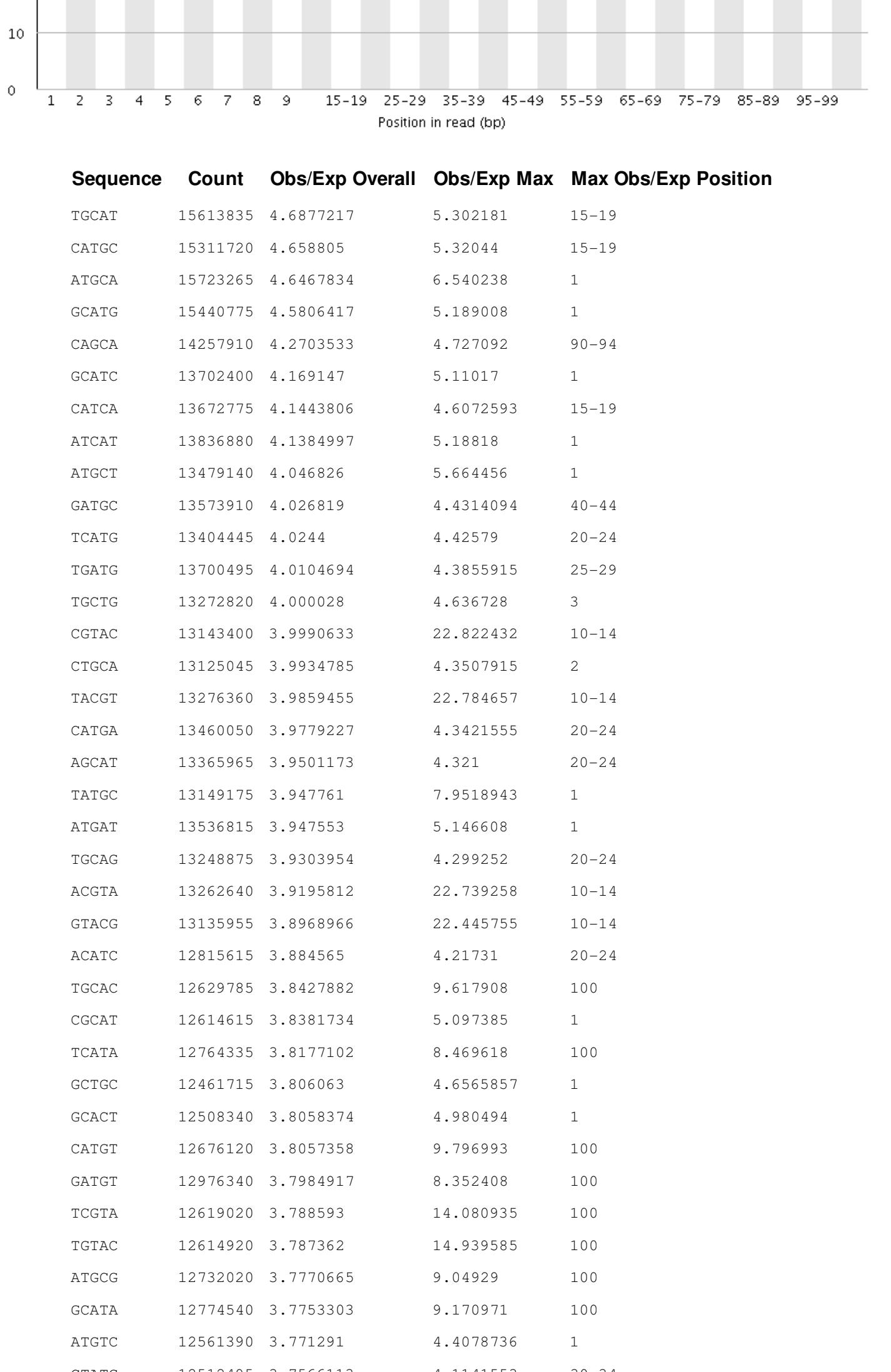
Basic Statistics

Measure	Value
Filename	ain.bam
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	34252010
Filtered Sequences	0
Sequence length	104
%GC	49

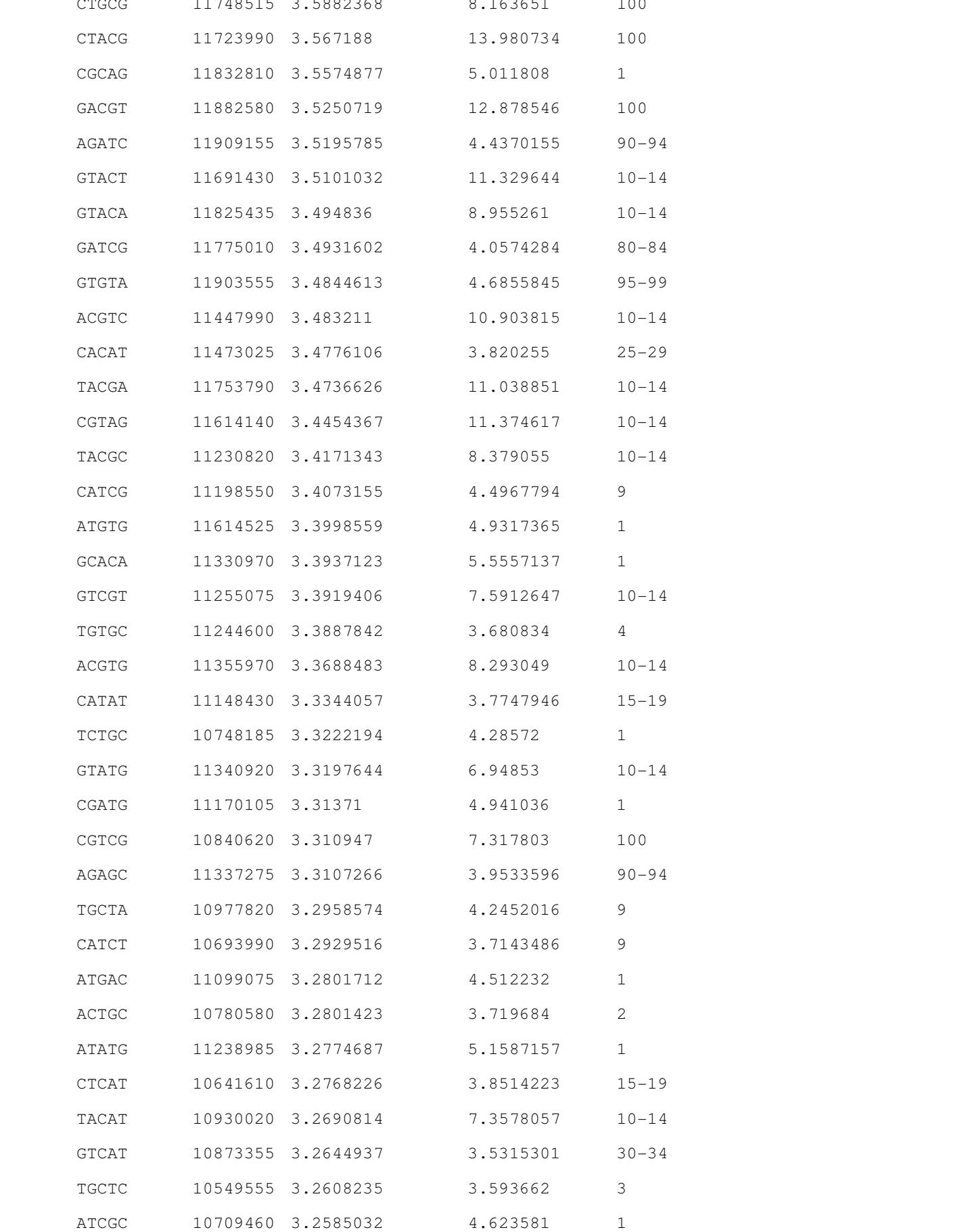
Per base sequence quality



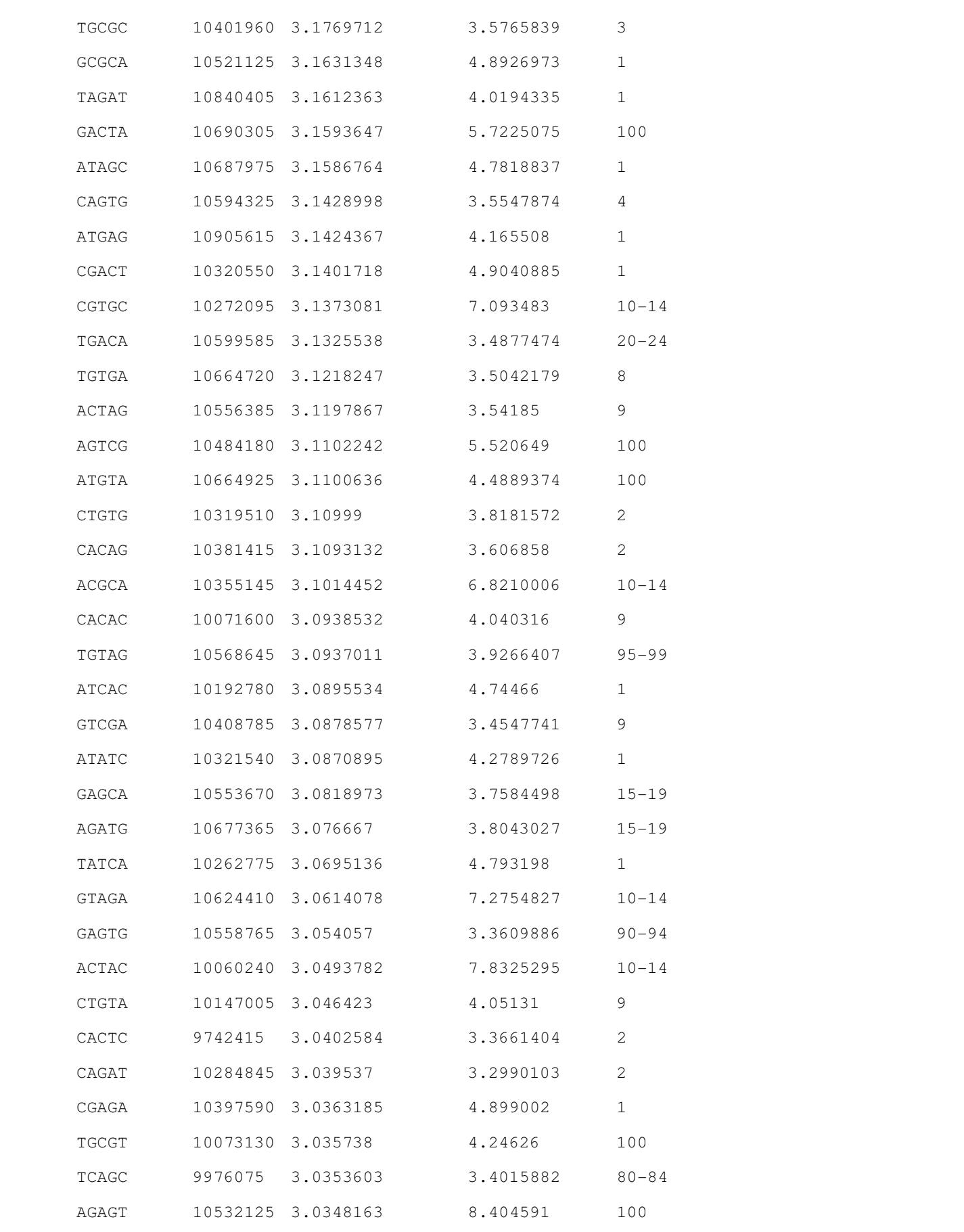
Per sequence quality scores



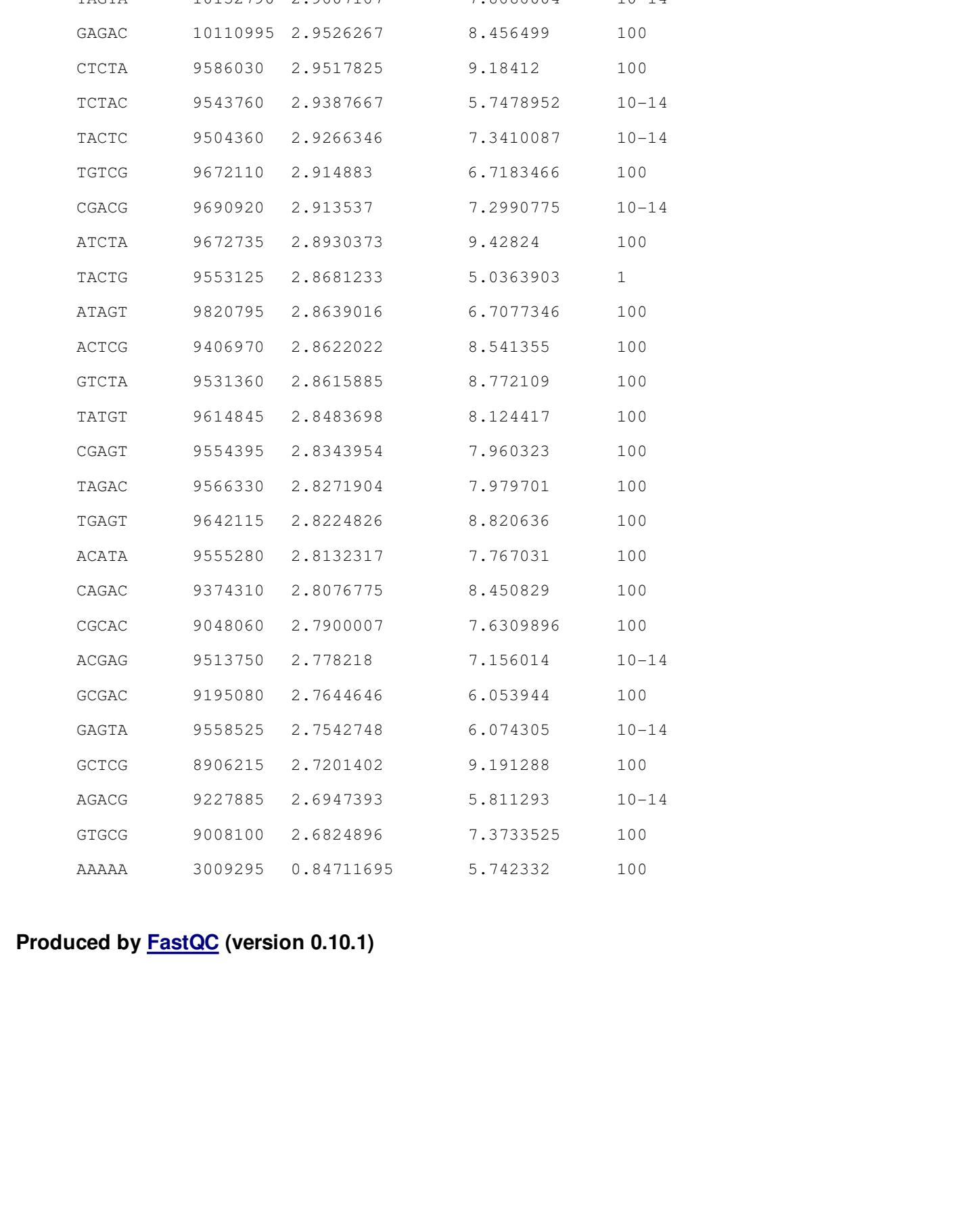
Per base sequence content



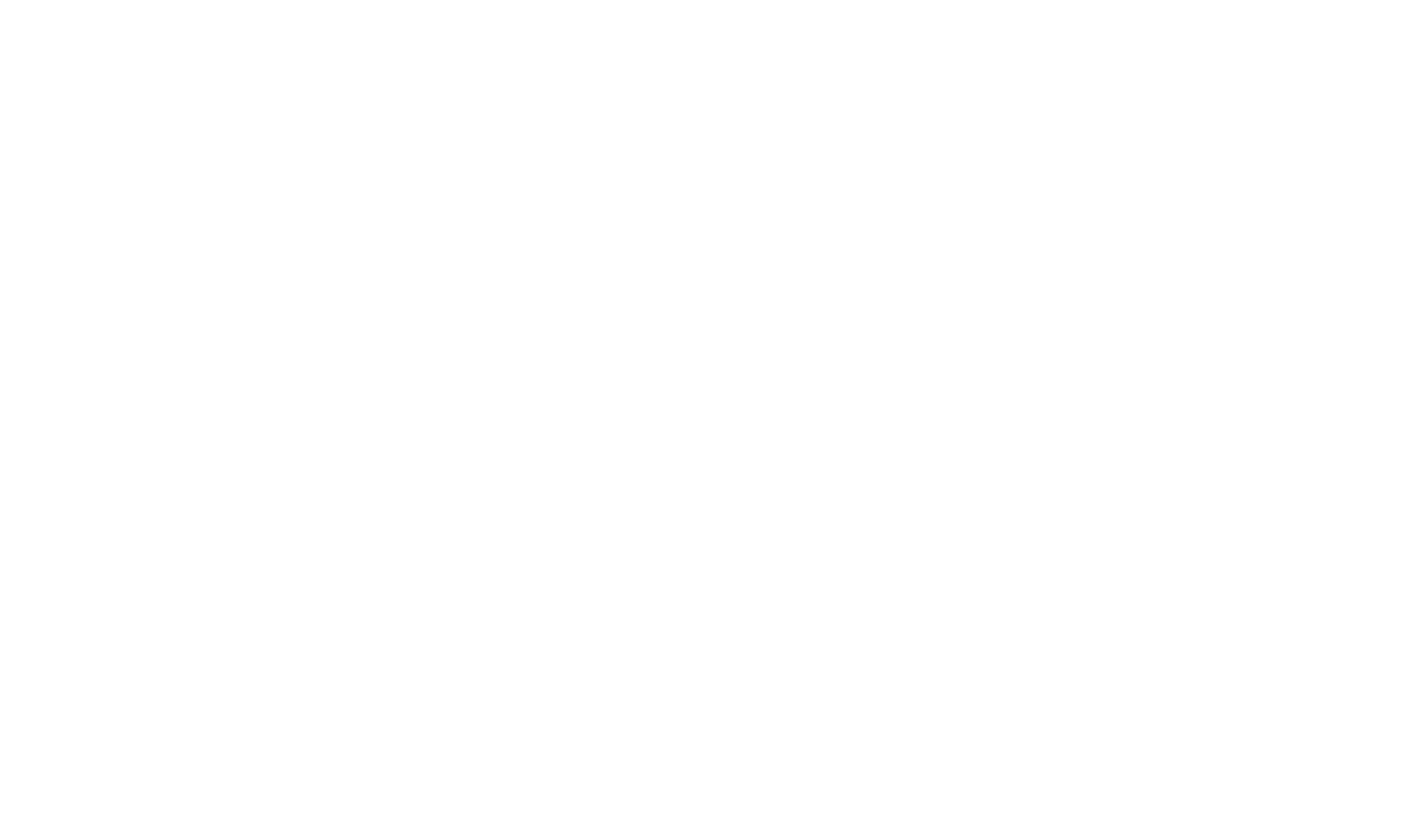
Per base GC content



Per sequence GC content



Per base N content



Sequence Length Distribution

Sequence Duplication Levels

Overrepresented sequences

No overrepresented sequences

Kmer Content

Produced by FastQC (version 0.10.1)

FastQC version 0.10.1

Copyright (c) 2012 EMBL-EBI

http://www.bioinformatics.babraham.ac.uk/projects/fastqc

FastQC is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

FastQC is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program. If not, see <http://www.gnu.org/licenses/>.

FastQC was developed at the European Bioinformatics Institute (EMBL-EBI) in the UK.

FastQC version 0.10.1 was released on 2012-09-04.

FastQC is part of the QIIME project.

FastQC is part of the Galaxy project.

FastQC is part of the Bioconductor project.

FastQC is part of the BioPython project.