

Genomic encryption of digital data stored in synthetic DNA

Robert N. Grass*^[a], Reinhard Heckel^[b], Christophe Dessimoz^[c], Wendelin J. Stark^[a]

[a] Prof. Dr. R.N. Grass, Prof. Dr. W. J. Stark
Department of Chemistry and Applied Biosciences
ETH Zurich
Vladimir-Prelog-Weg 1, 8093 Zurich, Switzerland
E-mail: rograss@ethz.ch

[b] Prof. Dr. R. Heckel
Department of Electrical and Computer Engineering
Rice University
6100 Main Street, Houston, TX 77005 Houston, United States

[c] Prof. Dr. C. Dessimoz
Department of Computational Biology and Center for Integrative Genomics
University of Lausanne, Lausanne, Switzerland
Swiss Institute of Bioinformatics, Lausanne, Switzerland
Department of Genetics, Evolution & Environment and Department of Computer Science
University College London, UK

Supporting information for this article is given via a link at the end of the document.

Abstract: Today, we can read human genomes and store digital data robustly in synthetic DNA. Here we report a strategy to intertwine these two technologies to enable the secure storage of valuable information in synthetic DNA, protected with personalized keys. We show that genetic short tandem repeats (STRs) contain sufficient entropy to generate strong encryption keys, and that only one technology, DNA sequencing, is required to simultaneously read key and data. Using this approach, we experimentally generated 80 bit strong keys from human DNA, and used such a key to encrypt 17 kB of digital information stored in synthetic DNA. Finally, the decrypted information was recovered perfectly from a single massively parallel sequencing run. Paste your abstract here. Please note it may not exceed 250 words. It may include up to three cited (non-numerical) references.

Due to its high theoretical data density of 455 exabyte per gram^[1] and its high stability,^[2] DNA has recently been proposed as a capable digital data storage medium. Poems, books, music, images and whole operating systems have already been stored in and successfully retrieved from synthetic DNA.^[3] Another advantage of DNA as a technical data storage substrate is that by having the same properties as natural DNA, it can be read using the same high-throughput “next-generation sequencing” (NGS) platforms. As a result, it is now possible to combine natural and synthetic DNA for storage and reading. To illustrate such an avenue, here, we demonstrate a data storage scheme with biometric authentication entirely based in DNA. By utilizing a user’s genomic short tandem repeat (STR) profile to generate a personalized cryptographic key, with an entropy of at least 80 bits, combined with AES-256 symmetrical encryption and Reed-Solomon error correction coding,^[2] this scheme achieves a biometrically encrypted and long-term storage of delicate information in synthetic DNA. To illustrate the performance of the scheme, we encrypt and store a paper of Alan M. Turing,^[4] which was originally kept classified for over 60 years, and demonstrate successful recovery without information loss.

Prior to the days of modern encryption technologies, secret and personal messages had to be physically hidden to avoid unauthorized access. With the development of the mathematical tools of one-way functions, which are cheap to evaluate but

prohibitively expensive to invert,^[5] a secret message can be encrypted using of a key, so that the encrypted message becomes useless to anyone who does not have access to the (secret) key. As a result, only the key has to be kept secret or private. Such encryption technologies are now at the core of our digital life as we utilize private keys (passwords) to access valuable information. For the encryption to work, the key has to have sufficient variability (entropy) so that it cannot easily be guessed by experimentation. Currently, keys with an entropy of at least 128 bit are regarded as safe, and it is envisioned that keys of 256 bit cannot be forged by current or future technologies, if the data is encrypted according to the Advanced Encryption Standard (AES).^[6]

However, it is difficult to memorize such complex keys, which are equivalent to 32 random alphanumeric values. As a result, utilized keys are often shorter, and easy to guess, rendering the encrypted information vulnerable. A possible solution to this problem is offered by biometrics, where measurable and differentiating features of individuals are utilized to generate a numeric encryption key. Examples thereof, termed biocryptography,^[7] are fingerprint-, eye- and face-scanners, which have most recently been integrated into consumer electronics, such as cell-phones and laptops. While currently possible with low-cost devices, the measurement of these personal features is imprecise and the amount of distinguishing features (entropy) is limited, resulting in relatively weak keys and making such biometric keys unfeasible for the encryption of highly valuable data. As examples, the fingerprint readers utilized in current smartphones have a false acceptance rate (FAR) of 1/50'000, which is equivalent to a key entropy of ca. 16 bit and Apple’s face recognition is advertised^[8] with a FAR of 1/1'000'000, ca 20 bit.

In this paper we explore a potential alternative solution by using personal, genetic information instead of the resulting phenotypes (fingerprint, iris, face-features etc.) to generate biometric keys.^[9] While the reading of genetic information is currently certainly still more complex than the reading of the features of fingerprints, the field of DNA sequencing is rapidly advancing, and Zaaier *et al.*^[10] have recently shown that the

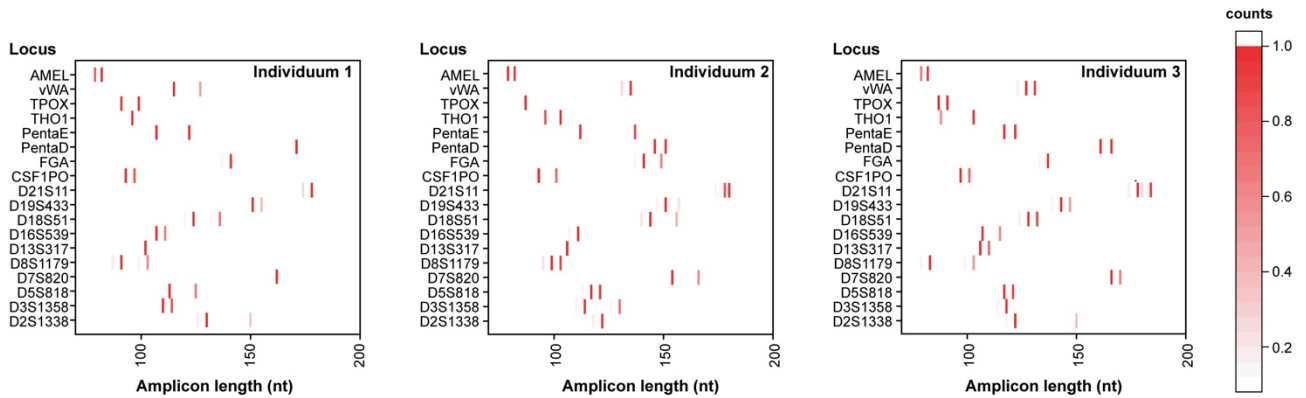


Figure 1. Sequencing of STR profiles. Short tandem repeat (STR) profiles obtained by next generation sequencing from three individuals using a primer mix previously reported by Kim *et al.*^[11] The sequencing reads are sorted by the corresponding primer sequences for each STR marker, and the lengths of the sequenced amplicons are represented here as relative counts per locus to give a representation that reflects capillary electropherograms. The corresponding plots show one or two different STR alleles per marker, and in some cases (e.g. vWA for individuals 2 and 3) STR stutter.^[12]

sequencing time required to identify humans from buccal swab samples by nanopore sequencing SNPs is within several minutes. While much current work on genotyping focuses on single nucleotide polymorphisms (SNPs), as these are more directly available from shotgun sequencing experiments, we chose to derive the genomic biometric information from short tandem repeats (STRs). This choice is motivated by the long tradition of STRs in forensics (since the early 90s).^[13] The distribution of various genotypes within the population have been characterized to some certainty^[14] and are readily available (for example at the website strbase.nist.gov run by the US National Institute of Standards and Technology (NIST)). While STR profiling is traditionally performed by PCR followed by capillary electrophoresis,^[15] several NGS-based protocols have recently been developed for the task,^[9c, 16] and have been made available to specialized laboratories. As we had difficulties in obtaining a commercial kit (Illumina's Forenseq prep kit has a starting price of 16'000 USD, and the advertised Promega Powseq Auto/Mito/Y is still in the Prototype stage^[17]) we decided on a NGS STR analysis utilizing a recently published multiplex PCR primer mix,^[11] which generates relatively short STR amplicons (77-210bp). A preliminary experiment with the buccal-swab DNA of three individuals revealed that the amplicons could be sequenced

successfully and all included 17 forensic autosomal STRs and amelogenin markers could be read.

As visible in Figure 1, the STR profiles of the three unrelated individuals differ markedly, and it is therefore conceivable that these STR profiles could be utilized as an access key. From a more statistical standpoint, the probabilities of the individual STR profiles (two genomic alleles per STR marker) can be compared with statistical datasets. Here, we compared with the widely utilized NIST reference 1036 US population dataset (revised version from July 2017 and the therein reported probability of identity (PI_{STR})^[14b, 18]), which ranges from 1.45 % to 50% per STR marker. Taking the 17 chosen STR markers and amelogenin, the worst-case probability that two non-related individuals have the same STR profile can be calculated under the assumptions of marker independence and an unstructured population^[19]. For the selected markers ($PI_{17+AM} = \prod(PI_{STR})$) this equals to 6.9×10^{-22} , which is about one in a trillion human populations. As STR profiles can be measured exact and essentially error-free,^[20] this low number can be considered as the false acceptance rate (FAR) of the biometric measurement, and is significantly lower than the false acceptance rate of currently utilized biometric technologies (face recognition and fingerprint analysis), making it exceptionally

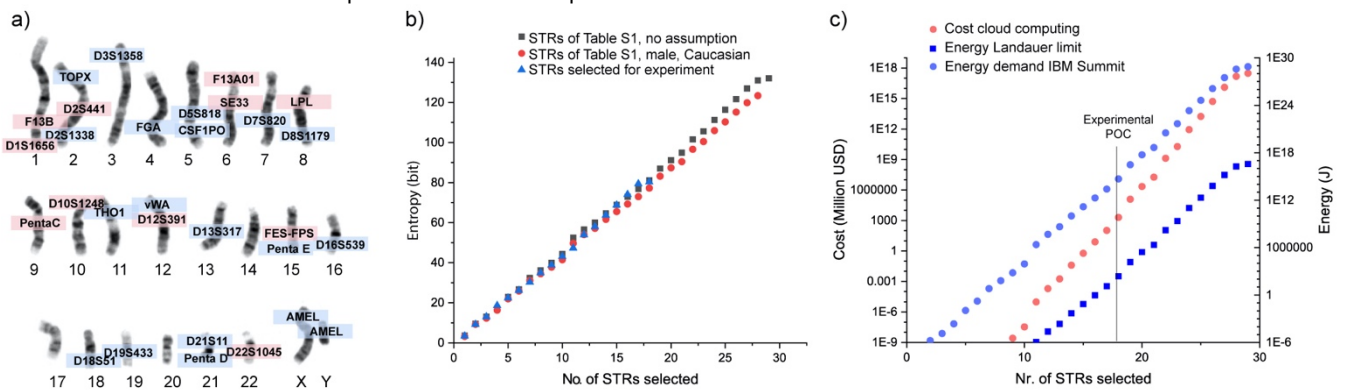


Figure 2. Entropy of an STR profile. a) Position of the 29 STR loci of Table S1 within a human genome. The STRs marked in blue are accessible to next generation sequencing via a literature primer mix and were utilized in the experimental proof of concept (POC).^[11] The karyotype was generated from an open access image from the National Institutes of Health.^[21] b) The entropy (stochastic information content) of the individual STR markers is computed from the probability distribution and the entropy of several markers is additive as the markers can be considered fully independent of each other (see supporting information). c) The Landauer limit gives the thermodynamically minimal amount of energy to delete a bit,^[22] and is seen as a theoretical lower energy limit for computational efforts (blue squares). Current supercomputing infrastructure is several orders of magnitude less energy effective (blue dots) and the current cost of large-scale cloud computing (red dots) is inhibiting, even for relatively low key strengths. The proof of concept experiments utilize 17 STR markers and amelogenin.

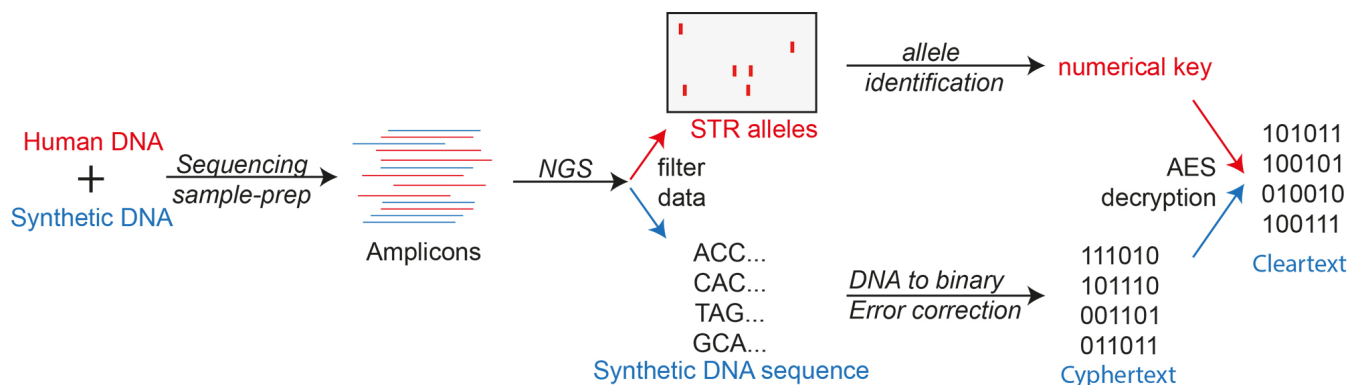
interesting in protecting highly valuable information. This is, of course, also the reason why DNA analysis utilizing STR profiles is regarded as the gold standard in forensic analysis. For our task of securely encrypting data, it is however not only relevant how rare a given STR profile is within a population, it is also highly relevant how secure the key is against a brute force attack (key guessing). For keys that are chosen uniformly at random from a set of keys, security against a brute force attack is quantified by the length of the key (in bit), as this determines the average number of attempts required to brute force guess the key. Since we are using STR markers as keys, we have to account for the markers not being uniformly distributed, so a brute force approach to guess the key would start with the most probable key. In general, the average amount of guesses required is determined by the entropy of the key (see supporting information).

Even if only considering the 17 forensic autosomal markers accessible via the PCR primers in our experiments (Table S1), the resulting STR key entropy is 80 bits strong. To put this into context, the encryption keys of the RC5 cypher were recovered by the distributed computing project distributed.net in 1997 for a 57 bit key, in 2002 for a 64 bit key, and the current challenge to break a 72 bit key is still ongoing. The entropy of the biometric key could be increased straightforwardly by using more STR markers and/or the addition of SNPs. The usage of all 29 independent STR markers included in the 1036 NIST table would enable a key strength of ca. 132 bit, and also pre-knowledge about an individual (e.g. Caucasian, male), would only slightly reduce the key strength (Figure 2).

To display the feasibility of the proposed approach as well as a potential use-case, we chose to protect and store sensitive scientific information. A manuscript written by Alan M. Turing in ca. 1941 was chosen for this purpose. The manuscript "Paper on Statistics of Repetitions" can be considered as a mathematical basis for the breaking of the enigma code, which is considered a key factor in bringing the Second World War to an end. Following declassification of the document after more than 60 years of classified storage, the paper was recently typeset and is available on arXiv,^[4b] and we chose to directly store the document in its LaTeX format, exactly as available from arXiv (digital cleartext; 17kB).

The STR markers of one of the individuals of Figure 1 were translated to a numerical value (see supporting information) and hashed using PBKDF2^[23] to generate a fixed-length key (256 bit). To encrypt the cleartext data, the digital bitstream and key were fed to an AES implementation (in Matlab, see supporting information) and then translated to 1426 DNA sequences, each 159 nucleotides long. This transformation was performed according to our previously published scheme and includes concatenated Reed Solomon error correction capabilities and constant amplification primers.^[2] The DNA sequences were then synthesized using an array technology by Customarray to yield 3.4 µg of synthetic DNA.

A key feature of using data stored in DNA in conjunction with STR encryption keys is that in the reading/decryption procedure only one technology is required to simultaneously read the cyphertext and the encryption key. During sequencing preparation (Scheme 1), the DNA of a buccal swab of the key individual is mixed with the synthetic DNA pool, and the information of both data sources is read concurrently within the same sequencing run. Using appropriate primer signatures (see supporting info) and knowledge of the synthetic DNA sequence length, the sequences corresponding to the cyphertext can be identified in the sequencing data. The embedded Reed Solomon error correction code further enables the correction of DNA synthesis, storage and sequencing errors. This cyphertext data is only useful, and can only be deciphered, if the STR marker data found in the same sequencing data computes the correct decryption key. In our experiment using the Turing paper and the DNA buccal swab of Individuum 1, a sequencing run of 2.5 million 150nt pair end reads, read the synthetic DNA in 70 fold coverage, and each of the 17 STR markers and amelogenin could be identified in the sequencing data at least 100 times. The individual errors within the synthetic DNA could be fully resolved by the decoding step resulting in 0 bit errors, the STR profile of the individual agreed perfectly with the STR profile of the same individual recorded for encryption (see supporting information). As a result, the cyphertext could be deciphered, yielding a perfect reconstitution of the original file.



Scheme 1. Data readout and decryption. Human DNA and synthetic DNA are mixed during sequencing preparation reactions using appropriate primer pairs to yield diverse amplicons. Following Next Generation Sequencing the data is filtered according to the individual primer sequences present in the data. Sequences containing STR primers are utilized to generate the numerical decryption key, and sequences of the expected synthetic DNA length are used to compute the digital cyphertext, thereby using a previously established DNA error correction and DNA to digital conversion scheme.^[2] Only if the correct numerical key is fed into the AES decryption process, can the resulting cleartext file be interpreted. If a wrong numerical key is fed to the AES decryption process, the resulting file yields essentially no information about the original file and resembles random data.

The above analysis and experiments show that STR profiles can be used to encrypt digital information encoded in DNA, and that there is an intrinsic advantage of having the encryption key and encrypted data present in the same medium. As a potential use case we foresee a specifically designed DNA sequencer, which is loaded with the mixed human and synthetic DNA samples, performs the decryption process and yields the deciphered file, if the correct human DNA is supplied. The device could also use the raw data to judge the age/freshness of the DNA sample (e.g. by measuring DNA degradation markers^[24] in the STR reads) to impede data recovery using non-authentic material (e.g. shed DNA).

In addition, the usage of STRs might be especially interesting in storing encrypted information for long time frames, as in contrast to the phenotypes currently used for encryption, the genotype is inherited. As such, it would be easier for a close family member to guess the encryption key and read the information than for a stranger. The entropic advantage of various close family members displayed in Supplementary Figure S5 clearly shows that parents and siblings have a large enough advantage so that the key could be guessed with standard IT infrastructure, whereas the anticipated effort would be too large for more distant relatives (e.g. cousins). Using the STR profile of the experiment as an example (17 STR markers), a cousin would have to solve a 76 bit problem, whereas a sibling only 52 bits. In terms of computational time and cost on a current p3.16xlarge system using the effort assumptions of Fig. 2 this equates to 8 hours and 75 USD for the sibling and about 10'000 years and 1 billion USD for the cousin.

For ubiquitous digital data storage, DNA read and write is currently too expensive.^[3b] While there are several efforts ongoing to change this in the future,^[3c] DNA storage may however already be useful today for highly valuable, or very private information. For this, the format offers unprecedented compactness (easy to hide), high data stability,^[2] does not have to be copied/resaved,^[3a] and as shown above, intrinsic possibilities for biometric protection. Our analysis shows that STRs (which are already forensically applied) carry enough entropy to be used as cryptographic keys, and brute-force attacks on such keys are beyond the current computing capabilities, and therefore are extremely unlikely to result in information exposure - even accounting for the foreseeable continuing increase in computation speed.

Acknowledgements

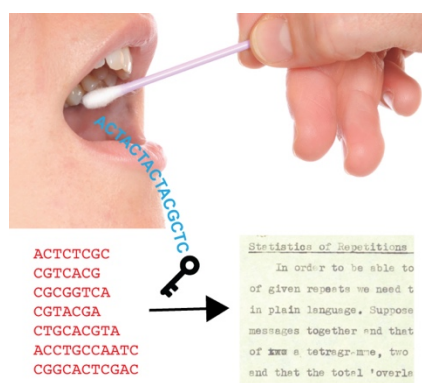
We thank the team from Microsynth for their support and ETH Zurich for funding.

Keywords: Synthetic DNA • biometrics • STR profiles • Next Generation Sequencing

- [1] G. M. Church, Y. Gao, S. Kosuri, *Science* **2012**, 337, 1628-1628.
 [2] R. N. Grass, R. Heckel, M. Puddu, D. Paunescu, W. J. Stark, *Angew. Chem. Int. Edit.* **2015**, 54, 2552-2555; *Angew. Chem.* **2015**, 127, 2582-2586.
 [3] a) N. Goldman, P. Bertone, S. Y. Chen, C. Dessimoz, E. M. LeProust, B. Sipos, E. Birney, *Nature* **2013**, 494, 77-80; b) Y. Erlich, D. Zielinski, *Science* **2017**, 355, 950-953; c) L. Organick, S. D. Ang, Y. J. Chen, R. Lopez, S. Yekhanin, K. Makarychev, M. Z. Racz, G. Kamath, P. Gopalan, B. Nguyen, C. N. Takahashi, S. Newman, H. Y. Parker, C. Rashtchian,

- K. Stewart, G. Gupta, R. Carlson, J. Mulligan, D. Carnean, G. Seelig, L. Ceze, K. Strauss, *Nat. Biotechnol.* **2018**, 36, 242-248; d) M. Blawat, K. Gaedke, I. Hutter, X.-M. Chen, B. Turczyk, S. Invreso, B. W. Pruitt, G. M. Church, *Procedia Comput. Sci.* **2016**, 80, 1011-1022.
 [4] a) A. M. Turing, *UK National Archives, HW 25/38 c. 1941*; b) A. M. Turing, *arXiv.org* **2015**, 1505.04715; c) S. Zabell, *Cryptologia* **2012**, 36, 191-214.
 [5] L. Lamport, *Commun. ACM* **1981**, 24, 770-772.
 [6] O. Dunkelmann, N. Keller, A. Shamir, *J. Cryptol.* **2015**, 28, 397-422.
 [7] K. Xi, J. Hu, *Bio-Cryptography*, Springer, Berlin, Heidelberg, **2010**.
 [8] "About Face ID advanced technology". Available at <https://support.apple.com/en-us/HT208108>, accessed 25.01.2019.
 [9] a) M. Hashiyada, *Tohoku J. Exp. Med.* **2004**, 204, 109-117; b) J. M. Butler, *Philos. Trans. Royal Soc. B* **2015**, 370, 20140252; c) K. B. Gettings, K. M. Kiesler, S. A. Faith, E. Montano, C. H. Baker, B. A. Young, R. A. Guerrieri, P. M. Vallone, *Forensic Sci. Int. Genet.* **2016**, 21, 15-21; d) S. T. Amin, M. Saeb, S. El-Gindi, *Proceedings of the Second IASTED International Conference on Computational Intelligence* **2006**, 120-125.
 [10] S. Zaaier, A. Gordon, D. Speyer, R. Piccone, S. C. Groen, Y. Erlich, *Elife* **2017**, 6, e27798.
 [11] E. H. Kim, H. Y. Lee, I. S. Yang, S. E. Jung, W. I. Yang, K. J. Shin, *Forensic Sci. Int. Genet.* **2016**, 22, 1-7.
 [12] C. Brookes, J. A. Bright, S. Harbison, J. Buckleton, *Forensic Sci. Int. Genet.* **2012**, 6, 58-63.
 [13] a) D. Tautz, *Nucleic Acids Res.* **1989**, 17, 6463-6471; b) A. J. Jeffreys, V. Wilson, S. L. Thein, *Nature* **1985**, 314, 67-73; c) A. Edwards, A. Civitello, H. A. Hammond, C. T. Caskey, *Am. J. Hum. Genet.* **1991**, 49, 746-756.
 [14] a) B. Budowle, B. Shea, S. Niezgodna, R. Chakraborty, *J. Forensic Sci.* **2001**, 46, 453-489; b) C. R. Hill, D. L. Duewer, M. C. Kline, M. D. Coble, J. M. Butler, *Forensic Sci. Int. Genet.* **2013**, 7, E82-E83.
 [15] a) C. P. Kimpton, P. Gill, A. Walton, A. Urquhart, E. S. Millican, M. Adams, *PCR Methods Appl.* **1993**, 3, 13-22; b) K. Lazaruk, P. S. Walsh, F. Oaks, D. Gilbert, B. B. Rosenblum, S. Menchen, D. Scheibler, H. M. Wenz, C. Holt, J. Wallin, *Electrophoresis* **1998**, 19, 86-93.
 [16] a) F. Guo, Y. S. Zhou, F. Liu, J. Yu, H. Song, H. Y. Shen, B. Zhao, F. Jia, G. W. Hou, X. H. Jiang, *Forensic Sci. Int. Genet.* **2016**, 23, 111-120; b) C. Borsting, N. Morling, *Forensic Sci. Int. Genet.* **2015**, 18, 78-89.
 [17] E. A. Montano, J. M. Bush, A. M. Garver, M. M. Larjani, S. M. Wiechman, C. H. Baker, M. R. Wilson, R. A. Guerrieri, E. A. Benzinger, D. N. Gehres, M. L. Dickens, *Forensic Sci. Int. Genet.* **2018**, 32, 26-32.
 [18] C. R. Steffen, M. D. Coble, K. B. Gettings, P. M. Vallone, *Forensic Sci. Int. Genet.* **2017**, 31, E36-E40.
 [19] E. Guichoux, L. Lagache, S. Wagner, P. Chaumeil, P. Leger, O. Lepais, C. Lepoittevin, T. Malausa, E. Revardel, F. Salin, R. J. Petit, *Mol. Ecol. Resour.* **2011**, 11, 591-611.
 [20] S. Salceda, A. Barican, J. Buscaino, B. Goldman, J. Klevenberg, M. Kuhn, D. Lehto, F. Lin, P. Nguyen, C. Park, F. Pearson, R. Pittaro, S. Salodkar, R. Schueren, C. Smith, C. Troup, D. Tsou, M. Vangbo, J. Wunderle, D. King, *Forensic Sci. Int. Genet.* **2017**, 28, 21-34.
 [21] "Talking Glossary of Genetic Terms, National Institutes of Health. National Human Genome Research Institute.". Available at <https://www.genome.gov/glossary/>, accessed 28.01.2019.
 [22] R. Landauer, *IBM J. Res. Dev.* **1961**, 5, 183-191.
 [23] B. Kaliski, *IETF* **2000**, 10.17487/RFC12898.
 [24] a) F. F. Wang, L. L. Wang, C. Briggs, E. Sicinska, S. M. Gaston, H. Mamon, M. H. Kulke, R. Zamponi, M. Loda, E. Maher, S. Ogino, C. S. Fuchs, J. Li, C. Hader, G. M. Makrigiorgos, *J Mol Diagn* **2007**, 9, 441-451; b) N. Ludyga, B. Grunwald, O. Azimzadeh, S. Englert, H. Hoffer, S. Tapio, M. Aubele, *Virchows Arch* **2012**, 460, 131-140; c) T. J. Anchordoquy, M. C. Molina, *Cell Preserv Technol* **2007**, 5, 180-188.

Entry for the Table of Contents



Inherited features from an individual's genome are used as an encryption key to protect digital data stored in synthetic DNA. Although encryption key and cyphertext originate from very different sources, the chemical nature of the information is identical and can consequently be analyzed with the same reading technology.

Materials and Methods

STR profiles of three laboratory members were measured following the method recently published by Kim *et al.* ^[1]. In detail, buccal swabs were collected (Isohelix, UK) and the DNA was extracted/purified with a commercial kit (Nucleospin tissue, Machery-Nagel, DE) following instructions to yield 2.4 ± 1.4 ng/ μ l DNA. PCR primers were ordered from Microsynth (CH) on a genomics scale and desalted in 100 nmol/ml solutions. A primer mix was prepared by taking 5-10 μ l of each primer pair (according to Table 1 in Kim *et al.*^[1]) and adding 40 μ l of water. Individual DNA samples were amplified via qPCR (10 μ l Kapa SYBR Fast qPCR Master mix, 6 μ l primer mix, 4 μ l sample; 10 minutes at 94°C activation followed by 22 cycles of 59°C 90 sec, 72°C 60 sec, 94°C 20 sec), purified by gel electrophoresis and extracted from the gel via a commercial kit (HighPure PCR Product Purification Kit, Roche, DE). Sample preparation (Illumina TruSeq amplicon library) and sequencing was performed by the company Microsynth (Illumina MiSeq 500 K amplicon reads (2*250 v2)) to yield 904944 past filter reads with an average length of 153 bp.

For every sample, the reads were analyzed with a Matlab script, searching for sequences containing corresponding forward and reverse primers of the individual STR markers ^[1], and collecting normalized sequence length distributions for every marker (See Figure 1, Matlab script and Excel tables available from the authors upon reasonable request). In disagreement with the original paper ^[1], where coverage was reported as very homogeneous (Figure 2 in Kim *et al.*), the relative coverage of the various markers in our experiments was quite inhomogeneous (Figure S1). However, the minimal coverage of 149 reads (for D19S433) was more than sufficient to analyze the STR profile for all markers.

Encryption key generation

For the generation of the cryptographic key, the data of Figure 1 had to be translated to a numerical key: The amplicon length for every STR marker corresponds to a specific genomic allele (2.2 to 43.3), equal to the number of pattern repeats and variants, but the possible alleles are not the same for all STR markers (e.g., CSF1PO can take the values 8.0, 9.0, 10.0, 11.0, 12.0, 13.0, 14.0, etc., while THO1 can take the values 5.0, 6.0, 7.0, 8.0, 9.0, 9.3, 10.0, 11.0 etc.). We therefore assigned every reported allele to an integer number (CSF1PO 8.0 \rightarrow 1, 9.0 \rightarrow 2, 10.0 \rightarrow 3 etc. THO1: 5.0 \rightarrow 1, 6.0 \rightarrow 2, 7.0 \rightarrow 3 etc.) to generate a numerical key consisting of integers. These integers (1...27) were put together in a fixed order (specified in Table S2) to give a numerical key. For individuum 1 this key is
61134366636665649810111145111113136944684712.

Since most cryptographic functions require keys of a fixed key length, the alphanumeric key was fed to a key stretching function, PBKDF2 (Implemented by Parvez Anandam at <http://anandam.name/pbkdf2/>). While the resulting key is 256 bits long, the entropy of the key is only 80 bits, as discussed in the main body of the paper (see Figure 2 of main manuscript). If however, the attacker does not have any prior information on how the key is generated, it would inherit the security of a 256bit long key. For individuum 1 the stretched key (in hex encoding with a salt of 0 and 10000 iterations, 32 bit key size) calculates to:

ed3ddc3e957c7e7df5e9bea414c459a596b30be457cbf9d4838097e9c171ef76.

Digital file selection, encryption and storage in DNA

The LaTeX file downloadable from arXiv (<https://arxiv.org/format/1505.04715>) contains a tex file of 17 kB. The file was padded according to the Public Key Cryptography Standards (PKCS)#7 so that length of the file could be divided by 16. This resulting byte-vector was fed into a validated AES routine ^[2] in Matlab (implemented by Stepan Matejka, Revision 1.1.0, 2011/10/12, ecb mode) in conjunction with the key derived above for individuum 1 resulting in a cyphertext. This cyphertext was translated to DNA sequences and redundancy was added for correcting errors in the DNA according to the scheme described in Grass *et al.* 2015 ^[3] resulting in 1426 sequences,

each 117 nt long. In order to guarantee direct compatibility with the PCR primers utilized for the STR amplification, we decided to utilize one of the STR PCR primers as a handle for the synthetic DNA, and the TPOX primers were selected for this. As a consequence, each of the 1426 synthetic DNA sequences had the following format:

5' CAGAACAGGCACTTAGGGAAC--Data--GCAAATAAACGCTGACAAGGA 3'

The resulting 159 nt long DNA sequences were ordered from Customarray (USA) on a 12K format chip, and delivered as an 80 µl solution (Tris-EDTA buffer) containing 43.1 ng/µl DNA.

File and key reading, decoding and decryption

Two sequencing pre-preparation methods were performed and resulted in very similar results: individual PCR amplification of synthetic and human DNA, and amplification of the two information sources together.

Separate amplification:

The as-delivered synthetic DNA encoding for the encrypted manuscript was diluted with water 1:10 and amplified by qPCR for 22 cycles with the primer mix and cycling specifics given above. Separately thereto a fresh buccal sample of individual 1 was collected with a buccal swab (Isohelix, UK), and the DNA therein was extracted and purified (Nucleospin tissue, Machery-Nagel, DE) and eluted into 100 µl of supplied buffer. The resulting DNA solution was diluted 1:6 with water and amplified by qPCR for 22 cycles with the primer mix and cycling specifics given further above. The PCR products of both the synthetic DNA and human DNA sample were purified by gel electrophoresis, extracted from the gel (HighPure PCR Product Purification Kit, Roche, DE) and mixed in a ratio 1:18 prior to library assembly.

Co-amplification:

The as-delivered synthetic DNA was diluted with water 1:10, the purified human DNA (see above) was diluted with water 1:6 and the two DNA solutions were mixed in a ratio of 7:1. This mixture was amplified by qPCR using the same primer mix and cycling specifics as given above. The resulting amplicons were purified by gel electrophoresis and extracted from the gel utilizing the same method prior to library assembly. *Library assembly and sequencing.*

Illumina TruSeq amplicon library preparation was performed for the separate and co-amplified sample individually and the two indexed experiments were sequenced together by the company Microsynth (Illumina MiSeq, v2 micro, 2x150bp) to yield 9369162 past filter reads with an average length of 147 bp and a 95.2% Q30 quality score.

Decoding and decryption

For file decoding, only sequences of the expected length (159 bp) were considered, and decoded with the Reed-Solomon error correction code ^[3]. In both cases (separate amplification and co-amplification) the cyphertext could be decoded without a single bit-error.

Decryption key generation

During file decryption we chose to use shorter Illumina reads (2x150 bp instead of 2x250bp used during encryption key reading) and some STR amplicons had expected length larger than 150 bp. As a consequence, sequencing data was first filtered for appropriate STR primers, and then stitched utilizing a Matlab script. As both the amplicons encoding for the synthetic DNA, as well as amplicons derived from the genetic marker TPOX had the same primer sequences (see file

selection and translation above), only sequences containing the TPOX primers, and at least five copies of the TPOX repetitive motive (AATG) were considered for marker analysis (See Fig S2). From here on, the STR marker alleles and key derivatization precisely followed the description of the key derivatization procedure above. For the separately amplified sample, the STR profile of individuum 1 did not differ from the STR profile recorded during the encryption key generation step in any allele (Fig. S3). For the STR profile extracted from the co-amplified sample, the STR profile of individuum 1 was complete, with the exception of the amelogenin, where only 1 copy read was insufficient for the correct allele identification. As shown in Figure S4 the reason for this was not that it was not possible to amplify and sequence the synthetic DNA together with the human DNA sample to generate the STR and file amplicons in a single amplification reaction, but that due to the mixing ratio chosen, the STR markers were rather underrepresented in the sample (average of 551 reads per marker compared to an average of 10132 reads per marker for separate amplification). Also, the relative coverage between the individual markers was slightly higher if the STR amplicons were generated in the presence of the synthetic DNA sample (see Fig S4). Optimization of the primer mix volumes (ratio of individual marker primer, and mixing ratio optimization between human and synthetic DNA will in future allow the reading of file and key with significantly less total Illumina amplicon reads. As the sample preparation via PCR only results in a marginal cost and effort (especially, if compared to the following library preparation and sequencing), the more conservative route of separate amplification is considered optimal.

Supplementary Information Text

Entropy of an STR marker

Typically, the security of a key is quantified by the length of the key. The reason for that is that the length of the key is a measure for the difficulty to guess the key, provided the key is drawn uniformly at random from all possible keys. Under this assumption, the average number of trials to guess a key is at least $2^{(L-2)}$, where L is the number of bits of the key. For example, a key with 128 bits on average requires at least 2^{126} trials to guess the key, which is an infeasible task for today's computers due to the large computational complexity, and the reason such a key is considered secure.

In our setup, the key is the signature of the STR markers, but the markers are not uniformly distributed. In order to quantify the security of such a key, we are again interested in the average number of trials required to guess the signature or key. This number is at least $2^{(E-2)}$, where E is the entropy of the signature.

The entropy of a discrete random variable which takes on m different values is defined as

$$E = - \sum_{i=1}^m p_i \log_2 p_i$$

where p_i are the probabilities of the random variable taking on its i-th value.

Note that the entropy is maximized if the random variable is uniform distributed, i.e., if the p_i are equal for all i. If the p_i are not equal the entropy is lower. As an example, consider a fair dice, which has equal probability for each number. Its entropy can be computed as 2.58. In contrast, a dice with the non-uniform probabilities specified in the table below only has an entropy of 2.37.

Nr (i)	Probability (p_i)	$-p_i \log_2 p_i$
1	0.05	0.21
2	0.1	0.33
3	0.1	0.33
4	0.2	0.46
5	0.25	0.5
6	0.3	0.52
		Entropy = <u>2.37 bit</u>

For every STR marker the NIST data gives the required probabilities for every allele. Due to the diploid nature of our genome, every STR marker gives two reads, and the two reads are independent of each other but there is no notion of order of the two reads. This is equivalent to

throwing two dice, and viewing the outcome as a set. For example, the outcomes that the first dice is 2 and the second is 5 and the outcome that the first dice is 5 and the second 2 are equivalent, if we view the outcome as a set. Let p_{ij} be the probability that we observe the set (i,j) with $i < j$. Then $p_{ii} = p_i * p_j$ and $p_{ij} = 2 * p_i * p_j$ where p_i is the probability of a read taking on the value i . With those probabilities, the entropy for the above dice experiment can be computed as:

Result i,j	Probability (p_{ij})	$-p_{ij} \log_2 p_{ij}$	Result i,j	Probability ($p_{i.}$)	$-p_{ij} \log_2 p_{ij}$
1,1	0.025	0.022	3,3	0.01	0.066
1,2	0.01	0.066	3,4	0.04	0.186
1,3	0.01	0.066	3,5	0.05	0.216
1,4	0.022	0.113	3,6	0.06	0.244
1,5	0.025	0.133	4,4	0.04	0.186
1,6	0.03	0.152	4,5	0.1	0.332
2,2	0.01	0.066	4,6	0.12	0.37
2,3	0.02	0.113	5,5	0.0625	0.25
2,4	0.04	0.186	5,6	0.15	0.41
2,5	0.05	0.216	6,6	0.09	0.313
2,6	0.06	0.244			
				Entropy=	3.95 bit

As can be seen from the calculation above, the entropy if the outcome of the experiment is a set is lower than if we obtain two independent draws of the dices and can distinguish the events say 2,5 and 5,2.

With this approach, the entropy of every diploid STR marker can be calculated from the NIST population data, and ranges from ca. 1 bit for AMEL (male/female) to 8.1 bit for SE33, with an average of 4.7 bit per diploid STR marker.

Entropy of an STR profile

The entropy of individual variables is additive, if the variables are independent of each other. As STR markers are inherited, it may be assumed that they follow the Mendelian law of independent assortment,^[4] especially if the markers are spaced well apart of each other on the chromosome (lower chance of linkage), or on different chromosomes, which completely inhibits genetic linkage and renders the markers fully independent.

For the 29 STR makers in the NIST 1036 tables, only the marker D6S1043 was not included due to the reported linkage to SE33. All other used STR markers (Table S3) are either the only STR marker on a given chromosome (no linkage possible) or the independence of these markers has been proven in the forensics literature.^[5]

For the 18 markers (17 STR marker and amelogenin) used in the experiment (Table S4), three pairs (TPOX, D2S1338; D5S818, CSF1PO and PentaD, D21S11) lie on the same chromosome, and as for the discussion above, the genetic independence of these markers has been discussed in detail in the forensics literature with the result that they may be considered as independent (non-linked with recombination probabilities of $> 10\%$).^[5] As a result of this analysis, it can be safely assumed that the 18 markers are full independent (non-linked), and that the entropies of the STR profile can be calculated by the sum of the entropies of the individual markers (Table S3). If more STR markers are used, the risk of linkage increases (as the new STR markers have to be placed nearer to existing markers), and the amount of entropy that can be added it not without limits. Still, it may be expected that the introduction of additional markers on Chromosomes 1, 6, 9, 10, 14, 17 & 20 would bring an additional several independent markers, and at an assumed average entropy per STR of 4.7 bit, this would result in an additional 33 bit.

Calculation of computational efforts for brute force attack

As a theoretical lower limit of energy required to run through all of the passwords, the Landauer limit may be considered. The Landauer limit represents the minimal thermodynamic energy required for a computation. Under the minimal assumption that only one computation is required to check a key, this gives the minimal energy required for such a check. ($L = k T \ln(2)$).

To calculate the energy required for a modern supercomputer, the energy demand of the computer, the computations per seconds (FLOPS) and computations per key guess are required. For the currently fastest supercomputer (IBM Summit), computational data and power are given as 122.3 petaFLOPS^[6] and 13 MW.^[7] The number of FLOPS required per key is not known, and the supercomputer does not have optimal architecture for this purpose (integer instead of floating point operations), but a relatively conservative assumption is an equivalent of 1000 flops per key.^[8]

To calculate the cost of trying a key on current large scale cloud computing infrastructure, the following assumptions were made:

- Cost per hour p3.16xlarge computing time (Amazon, 3-year Reserved, as of Nov 2018): 9 USD.
- Computational speed of one p3.16xlarge unit: 60'000 million SHA-256 hashes per second.^[9]
- Assumption: Test of an AES key is as computationally as demanding as evaluating one SHA-256 hash. This can be seen as a minimal value, as in the current scheme every key derived from a STR profile has to be hashed (key-stretched) using 10'000 rounds of PBKDF2, causing a significant computational effort for testing a single STR profile combination.

Entropy of close relatives

Due to the Mendelian inheritance of STR marker, close relatives have an advantage in guessing their relatives STR profiles, as for every marker every descendent will inherit one allele from each parent. As a result, a direct descendent only has to guess one allele of his parents, knowing that the other allele is equivalent to one of his. This can be extended to other close relatives, resulting in the probabilities of sharing alleles given in Table S5.

If the relative shares both alleles, he does not have to perform any guess, if he shares one, he has to guess which of the two he shares plus he has to guess the second allele, if he shares none the entropy of guessing the correct allele is calculated as derived above.

As an example, for a specific marker, full siblings have the following entropy of guessing their siblings alleles, knowing their own profile:

$$E_s = 0.25 \sum(-p_{ij} \log_2(p_{ij})) + 0.5 (1 + \sum(p_i \log_2(p_i))) + 0.25 \times 0.$$

SI References

- [1] E. H. Kim, H. Y. Lee, I. S. Yang, S. E. Jung, W. I. Yang, K. J. Shin, *Forensic Sci. Int. Genet.* **2016**, 22, 1-7.
- [2] J. Kepner, V. Gadepally, B. Hancock, P. Michaleas, E. Michel, M. Varia, *High Performance Extreme Computing Conference (HPEC), 2017 IEEE* **2015**, 10.1109/HPEC.2015.7322470.
- [3] R. N. Grass, R. Heckel, M. Puddu, D. Paunescu, W. J. Stark, *Angew. Chem. Int. Edit.* **2015**, 54, 2552-2555.
- [4] H. H. Li, U. B. Gyllensten, X. F. Cui, R. K. Saiki, H. A. Erlich, N. Arnheim, *Nature* **1988**, 335, 414-417.
- [5] a) T. Tamura, M. Osawa, E. Ochiai, T. Suzuki, T. Nakamura, *Legal Med-Tokyo* **2015**, 17, 320-325; b) C. Phillips, L. Fernandez-Formoso, M. Garcia-Magarinos, L. Porras, T. Tvedebrink, J. Amigo, M. Fondevila, A. Gomez-Tato, J. Alvarez-Dios, A. Freire-Aradas,

- A. Gomez-Carballa, A. Mosquera-Miguel, A. Carracedo, M. V. Lareu, *Forensic Sci. Int. Genet.* **2011**, 5, 155-169; c) K. L. O'Connor, A. O. Tillmar, *Forensic Sci. Int. Genet.* **2012**, 6, 840-844.
- [6] "Summit supercomputer ranked fastest computer in the world". Available at <https://www.energy.gov/articles/summit-supercomputer-ranked-fastest-computer-world>, accessed 25.01.2019.
- [7] L. Zhiye, "US dethrones China with IBM Summit supercomputer". Available at <https://www.tomshardware.com/news/us-supercomputer-china-top500-summit,37367.html>, accessed 25.01.2019.
- [8] M. Arora, "How secure is AES against brute force attacks?". Available at https://www.eetimes.com/document.asp?doc_id=1279619, accessed 25.01.2019.
- [9] D. Stamat, "AWS EC2: P2 vs P3 instances". Available at <https://blog.iron.io/aws-p2-vs-p3-instances/>, accessed 25.01.2019.
- [10] a) C. R. Hill, D. L. Duewer, M. C. Kline, M. D. Coble, J. M. Butler, *Forensic Sci. Int. Genet.* **2013**, 7, E82-E83; b) C. R. Steffen, M. D. Coble, K. B. Gettings, P. M. Vallone, *Forensic Sci. Int. Genet.* **2017**, 31, E36-E40.
- [11] C. Phillips, D. Ballard, P. Gill, D. S. Court, A. Carracedo, M. V. Lareu, *Forensic Sci. Int. Genet.* **2012**, 6, 354-365.
- [12] P. Hatzler-Grubwieser, B. Berger, D. Niederwieser, M. Steinlechner, *Forensic Sci. Int. Genet.* **2012**, 6, E50-E51.

Supplementary Figures and Tables

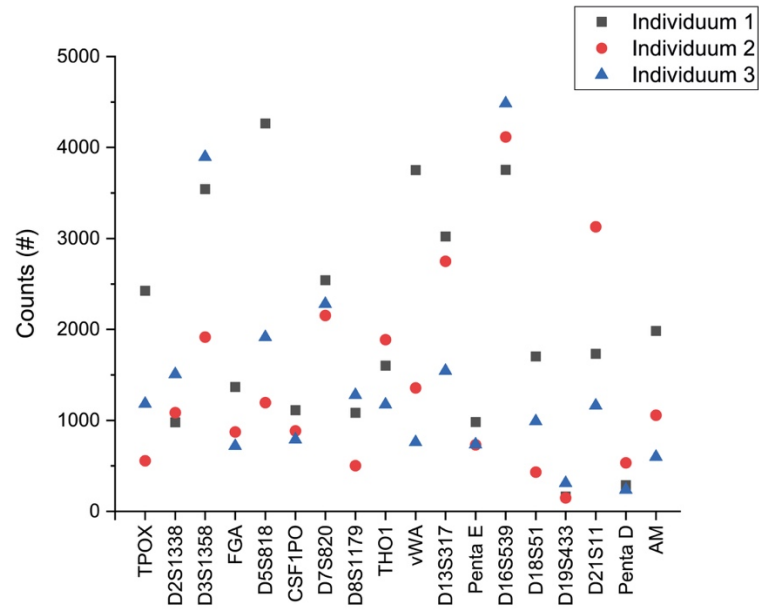


Figure S1. Sequencing coverage of the individual STR markers for three different individuals during initial key generation experiments (see Figure 1 in main text).

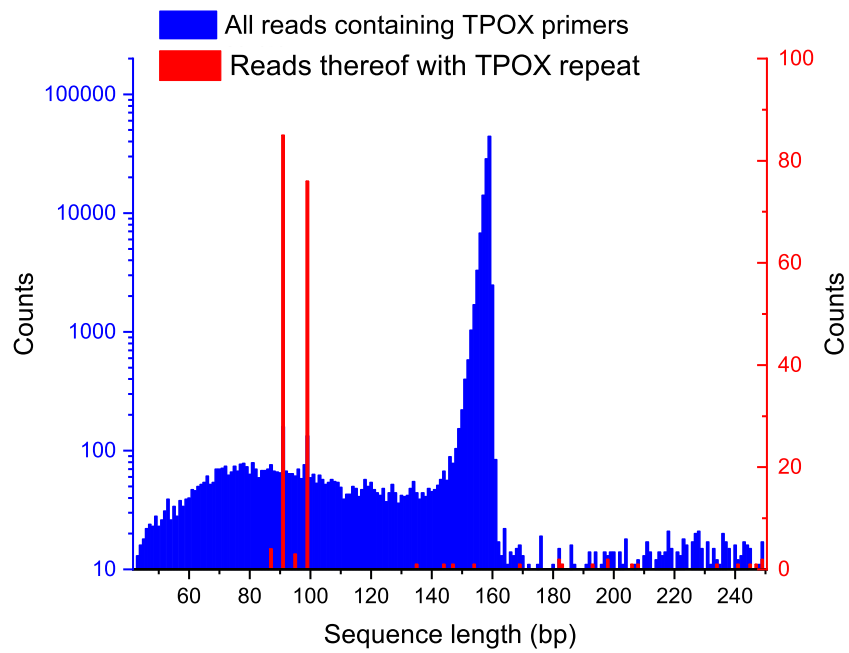


Figure S2. Analysis of the sequences starting and ending with TPOX primer sequences. The sequences in red additionally contain at least five copies of the TPOX repeat (AATG). The blue data shows the presence of the synthetic DNA amplicons (expected length of 159 bp), and the red data represents the alleles 6 and 8 for the TPOX STR marker of individual 1.

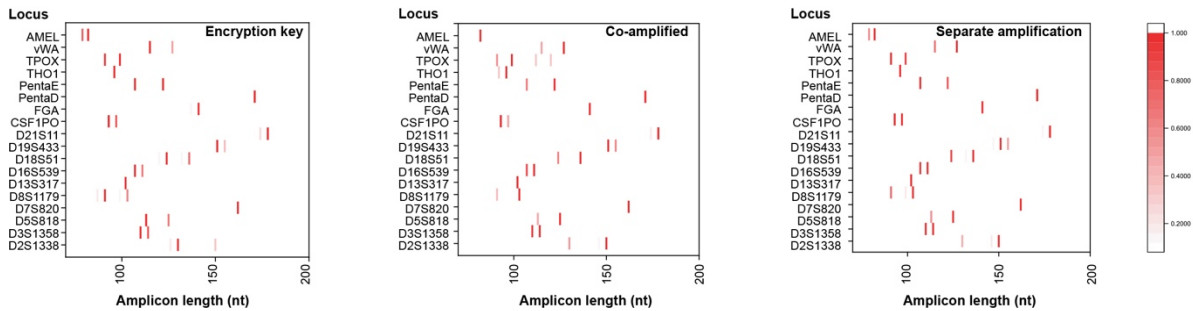


Figure S3. STR marker amplicon length profiles for individual 1 measured with three different procedures, either in the absence of synthetic DNA (left), co-amplified and co-sequenced together with the synthetic DNA (middle) and separately amplified, but co-sequenced with the synthetic DNA (right).

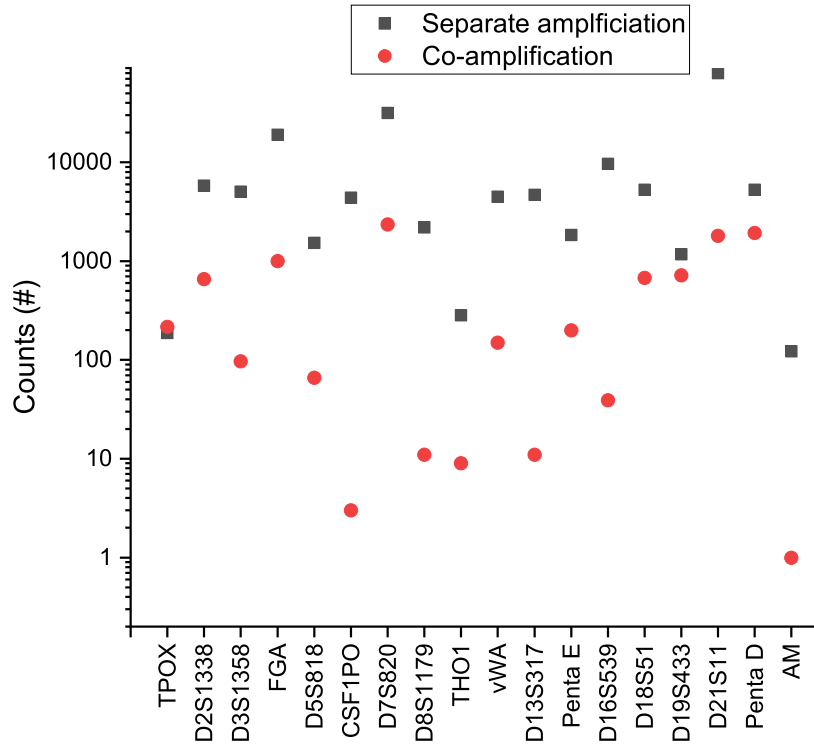


Figure S4. Coverage (counts) of the individual STR markers read during the decryption stage in the presence of the synthetic DNA.

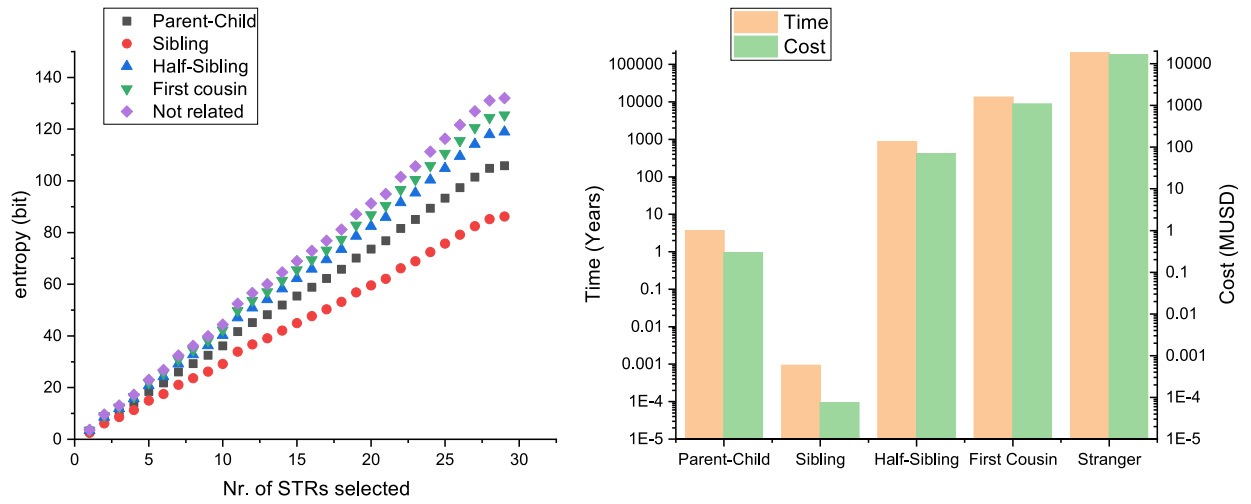


Figure S5. Entropy of the STR markers of Table S1 with pre-knowledge of the STR profile of a close relative, and the estimated time and costs to guess the relatives STR profile.

Table S1. Independent STR markers selected from the NIST 1036 US population dataset ^[10]. Markers highlighted in blue are used in the experimental work.

	STR	Chromosome ^[a]	CODIS	Entropy ^[b]	Entropy (Cauc.)
1	F13B	1	-	3.63	3.3
2	D1S1656	1	2017	5.96	6.05
3	TPOX	2	Core	3.49	2.99
4	D2S441	2	2017	4.15	4.06
5	D2S1338	2	2017	5.77	5.6
6	D3S1358	3	Core	3.79	3.96
7	FGA	4	Core	5.62	5.25
8	D5S818	5	Core	3.75	3.32
9	CSF1PO	5	Core	3.76	3.32
10	F13A01	6	-	4.44	3.73
11	SE33	6	-	8.08	8.13
12	D7S820	7	Core	4.15	4.36
13	LPL	8	-	3.42	3.17
14	D8S1179	8	Core	4.61	4.5
15	Penta C	9	-	4.3	3.87
16	D10S1248	10	2017	4.04	3.73
17	THO1	11	Core	3.87	3.74
18	vWA	12	Core	4.35	4.26
19	D12S391	12	2017	5.88	5.95
20	D13S317	13	Core	4.18	4.15
21	FES-FPS	15	-	3.68	3.08
22	Penta E	15	-	6.63	6.15
23	D16S539	16	Core	4.08	3.84
24	D18S51	18	Core	5.7	5.51
25	D19S433	19	2017	5.01	4.32
26	D21S11	21	Core	5.31	4.93
27	Penta D	21	-	5.29	4.64
28	D22S1045	22	2017	4.09	3.55
29	AMEL	X Y	-	1	1

[a] Markers on different chromosomes are fully independent. For markers on the same chromosome, the non-linkage has been discussed and proven in the forensics literature, with the exception of the linkage between SE33 and D6S1043 (see Phillips *et al.* 2012^[11]) [b] Calculated per diploid genome, in bit, details see above.

Table S2. Number of possible alleles reported in literature^[10] for each marker, number of allele possibilities per marker for a diploid genome, and the integer values derived from the STR profile for individual 1. These integer values are utilized to generate the numerical key.

Marker	Possible alleles per haploid genome	Possibilities per diploid genome	Integer key Individuum 1
D2S1338	13	169	6/11
D3S1358	11	121	3/4
D5S818	9	81	3/6
D7S820	11	121	6/6
D8S1179	11	121	3/6
D13S317	8	64	6/6
D16S539	9	81	5/6
D18S51	22	484	4/9
D19S433	16	256	8/10
D21S11	27	729	11/11
CSF1PO	9	81	4/5
FGA	27	729	11/11
PentaD	16	256	13/13
PentaE	23	529	6/9
THO1	8	64	4/4
TPOX	10	100	6/8
vWA	11	121	4/7
AM	-	2	1/2
Total possibilities		1.4064E+38	

Table S3. STRs of the NIST 1036 tables, their position in the chromosome as well as reasoning for independence of the individual STR marker.

STR	Chrom o-some arm	Approx Pos. (Mb) ¹	Entropy All	Entropy Cauc	Independence (reasoning)
F13B	1q	197.04	3.63	3.26	include (<i>not linked**</i> , $Rc^*(F13B:D1S1656) = 0.34$)
D1S1656	1q	230.91	5.96	6.05	include (<i>not linked**</i> , $Rc^*(F13B:D1S1656) = 0.34$)
TPOX	2p	1.49	3.49	2.99	include (<i>not linked**</i> , $Rc^*(TPOX:D2S441) > 0.40$)
D2S441	2p	68.24	4.15	4.06	include (<i>not linked**</i> , $Rc^*(TPOX:D2S441) > 0.40$)
D2S1338	2q	218.88	5.77	5.60	include (<i>separate chromosome arm</i>)
D3S1358	3	45.58	3.79	3.95	include (<i>separate chromosome</i>)
FGA	4	155.51	5.62	5.25	include (<i>separate chromosome</i>)
D5S818	5q	123.11	3.75	3.32	include (<i>not linked**</i> , $Rc^*(CSF1PO:D5S818) = 0.25$)
CSF1PO	5q	149.46	3.76	3.32	include (<i>not linked**</i> , $Rc^*(CSF1PO:D5S818) = 0.25$)
F13A01	6p		4.44	3.73	include (<i>separate arm to SE33</i>)
SE33	6q	88.99	8.08	8.13	include (<i>separate chromosome</i>)
D6S1043	6q	92.45	5.66	4.96	exclude (<i>linked** to SE33</i> , $Rc(SE33:D6S1043)=0.044$)
D7S820	7q	83.79	4.15	4.36	include (<i>separate chromosome</i>)
LPL	8p		3.42	3.17	include (<i>separate arm</i>)
D8S1179	8q	125.91	4.61	4.50	include (<i>separate chromosome</i>)
Penta C	9p		4.30	3.87	include (<i>separate chromosome</i>)
D10S1248	10	2.24	4.04	3.73	include (<i>separate chromosome</i>)
THO1	11	2.19	3.87	3.74	include (<i>separate chromosome</i>)
vWA	12p	6.09	4.35	4.26	include (<i>not linked**</i> , $Rc^*(vWA:D12S391) = 0.117$)
D12S391	12p	12.45	5.88	5.94	include (<i>not linked**</i> , $Rc^*(vWA:D12S391) = 0.117$)
D13S317	13	82.72	4.18	4.15	include (<i>separate chromosome</i>)
FESFPS	15q		3.68	3.08	include (<i>not linked**</i> , $Rc^*(Penta E:FES-FPS) = 0.181$)
Penta E	15q	97.37	6.63	6.15	include (<i>not linked**</i> , $Rc^*(Penta E:FES-FPS) = 0.181$)
D16S539	16	84.94	4.08	3.84	include (<i>separate chromosome</i>)
D18S51	18	60.95	5.70	5.51	include (<i>separate chromosome</i>)
D19S433	19	30.42	5.01	4.32	include (<i>separate chromosome</i>)
D21S11	21q	20.55	5.31	4.93	include (<i>not linked**</i> , $Rc^*(Penta D:D21S11) > 0.3$)
Penta D	21q	45.06	5.29	4.64	include (<i>not linked**</i> , $Rc^*(Penta D:D21S11) > 0.3$)
D22S1045	22	37.54	4.09	3.55	include (<i>separate chromosome</i>)
AMEL	X & Y		1	1	include (<i>separate chromosome</i>)
Total			132.0 bit	124.5 bit	

* Rc = Recombination rate from Kosambi mapping function.^[11]

** Non-linkage has been shown for **Rc recombination fractions for ~0.12,^[11] as found for

vWA:D12S391, and various studies have shown marker independence for this relatively close STR pair for non-close relatives.^[5]

** Only included profiles (**Bold values**) used for sum

Table S4. STR Markers used in experimental work and their history in the forensic analysis.

	Chromosome	Codis Core Locus	New FBI core locus	European Locus*	Entropy per diploid genome (bit)
TPOX	2	YES	YES		3.49
D2S1338	2		YES		5.77
D3S1358	3	YES	YES	YES	3.78
FGA	4	YES	YES	YES	5.62
D5S818	5	YES	YES		3.75
CSF1PO	5	YES	YES		3.76
D7S820	7	YES	YES		4.15
D8S1179	8	YES	YES	YES	4.61
THO1	11	YES	YES	YES	3.87
vWA	12	YES	YES	YES	4.35
D13S317	13	YES	YES		4.18
PentaE	15				6.63
D16S539	16	YES	YES		4.08
D18S51	18	YES	YES	YES	5.70
D19S433	19		YES		5.01
D21S11	21	YES	YES	YES	5.31
PentaD	21				5.29
AMEL	Y&X		YES	YES	1
Total Entropy					80.4 bit

* European Standard Set of Loci and new ESS loci.^[12]

Table S5. Probability of having a given number of STR alleles shared between close relatives

Relationship	0 alleles	1 allele	2 alleles
Parent-child	0	1	0
Full siblings	1/4	1/2	1/4
Half siblings	1/2	1/2	0
Grandparent-grandchild	1/2	1/2	0
Uncle-Nephew	1/2	1/2	0
First cousins	3/4	1/4	0
Entropy per marker	$\sum(-p_{ij}\log_2(p_{ij}))$	$1 + \sum(p_i\log_2(p_i))$	0

Supplementary Materials for

Genomic encryption of digital data stored in synthetic DNA

Robert N. Grass^{1*}, Reinhard Heckel², Christophe Dessimoz^{3,4,5,6,7}, Wendelin J. Stark¹

Correspondence to: robert.grass@chem.ethz.ch

Supplementary Text

Entropy of an STR marker

Typically, the security of a key is quantified by the length of the key. The reason for that is that the length of the key is a measure for the difficulty to guess the key, provided the key is drawn uniformly at random from all possible keys. Under this assumption, the average number of trials to guess a key is at least $2^{(L-2)}$, where L is the number of bits of the key. For example, a key with 128 bits on average requires at least 2^{126} trials to guess the key, which is an infeasible task for today's computers due to the large computational complexity, and the reason such a key is considered secure.

In our setup, the key is the signature of the STR markers, but the markers are not uniformly distributed. In order to quantify the security of such a key, we are again interested in the average number of trials required to guess the signature or key. This number is at least $2^{(E-2)}$, where E is the entropy of the signature.

The entropy of a discrete random variable which takes on m different values is defined as

$$E = - \sum_{i=1}^m p_i \log_2 p_i$$

where p_i are the probabilities of the random variable taking on its i-th value.

Note that the entropy is maximized if the random variable is uniform distributed, i.e., if the p_i are equal for all i. If the p_i are not equal the entropy is lower. As an example, consider a fair dice, which has equal probability for each number. Its entropy can be computed as 2.58. In contrast, a dice with the non-uniform probabilities specified in the table below only has an entropy of 2.37.

Nr (i)	Probability (p_i)	$-p_i \log_2 p_i$
1	0.05	0.21
2	0.1	0.33
3	0.1	0.33
4	0.2	0.46
5	0.25	0.5
6	0.3	0.52
Entropy =		2.37 bit

For every STR marker the NIST data gives the required probabilities for every allele. Due to the diploid nature of our genome, every STR marker gives two reads, and the two reads are independent of each other but there is no notion of order of the two reads. This is equivalent to throwing two dice, and viewing the outcome as a set. For example, the outcomes that the first dice is 2 and the second is 5 and the outcome that the first dice is 5 and the second 2 are equivalent, if we view the outcome as a set. Let p_{ii} be the probability that we observe the set (i,j) with $i < j$. Then $p_{ii} = p_i * p_j$ and $p_{ij} = 2 p_i * p_j$ where p_i is the probability of a read taking on the value i.

With those probabilities, the entropy for the above dice experiment can be computed as:

Result i,j	Probability (p _{ij})	-p _{ij} log ₂ p _{ij}	Result i,j	Probability (p _i)	-p _{ij} log ₂ p _{ij}
1,1	0.025	0.022	3,3	0.01	0.066
1,2	0.01	0.066	3,4	0.04	0.186
1,3	0.01	0.066	3,5	0.05	0.216
1,4	0.022	0.113	3,6	0.06	0.244
1,5	0.025	0.133	4,4	0.04	0.186
1,6	0.03	0.152	4,5	0.1	0.332
2,2	0.01	0.066	4,6	0.12	0.37
2,3	0.02	0.113	5,5	0.0625	0.25
2,4	0.04	0.186	5,6	0.15	0.41
2,5	0.05	0.216	6,6	0.09	0.313
2,6	0.06	0.244			
				Entropy= 3.95 bit	

As can be seen from the calculation above, the entropy if the outcome of the experiment is a set is lower than if we obtain two independent draws of the dices and can distinguish the events say 2,5 and 5,2.

With this approach, the entropy of every diploid STR marker can be calculated from the NIST population data, and ranges from ca. 1 bit for AMEL (male/female) to 8.1 bit for SE33, with an average of 4.7 bit per diploid STR marker.

Entropy of an STR profile

The entropy of individual variables is additive, if the variables are independent of each other. As STR markers are inherited, it may be assumed that they follow the Mendelian law of independent assortment (*I*), especially if the markers are spaced well apart of each other on the chromosome (lower chance of linkage), or on different chromosomes, which completely inhibits genetic linkage and renders the markers fully independent.

For the 29 STR makers in the NIST 1036 tables, only the marker D6S1043 was not included due to the reported linkage to SE33. All other used STR markers (Table S2) are either the only STR marker on a given chromosome (no linkage possible) or the independence of these markers has been proven in the forensics literature (2-4).

For the 18 markers (17 STR marker and amelogenin) used in the experiment (Table S3), three pairs (TPOX, D2S1338; D5S818, CSF1PO and PentaD, D21S11) lie on the same chromosome, and as for the discussion above, the genetic independence of these markers has been discussed in detail in the forensics literature with the result that they may be considered as independent (non-linked with recombination probabilities of > 10%) (2-4).

As a result of this analysis, it can be safely assumed that the 18 markers are full independent (non-linked), and that the entropies of the STR profile can be calculated by the sum of the entropies of the individual markers (Table S2). If more STR markers are used, the risk of linkage increases (as the new STR markers have to be placed nearer to existing markers), and the amount of entropy that can be added it not without limits. Still, it may be expected that the introduction of additional markers on Chromosomes 1, 6, 9, 10, 14, 17 & 20 would bring an additional several

independent markers, and at an assumed average entropy per STR of 4.7 bit, this would result in an additional 33 bit.

Calculation of computational efforts for brute force attack

As a theoretical lower limit of energy required to run through all of the passwords, the Landauer limit may be considered. The Landauer limit represents the minimal thermodynamic energy required for a computation. Under the minimal assumption that only one computation is required to check a key, this gives the minimal energy required for such a check. ($L = k T \ln(2)$).

To calculate the energy required for a modern supercomputer, the energy demand of the computer, the computations per seconds (FLOPS) and computations per key guess are required. For the currently fastest supercomputer (IBM Summit), computational data and power are given as 122.3 petaFLOPS (5) and 13 MW (6). The number of FLOPS required per key is not known, and the supercomputer does not have optimal architecture for this purpose (integer instead of floating point operations), but a relatively conservative assumption is an equivalent of 1000 flops per key (7).

To calculate the cost of trying a key on current large scale cloud computing infrastructure, the following assumptions were made:

- Cost per hour p3.16xlarge computing time (Amazon, 3-year Reserved, as of Nov 2018): 9 USD
- Computational speed of one p3.16xlarge unit: 60'000 million SHA-256 hashes per second (8)
- Assumption: Test of an AES key is as computationally as demanding as evaluating one SHA-256 hash. This can be seen as a minimal value, as in the current scheme every key derived from a STR profile has to be hashed (key-stretched) using 10'000 rounds of PBKDF2, causing a significant computational effort for testing a single STR profile combination.

Entropy of close relatives

Due to the Mendelian inheritance of STR marker, close relatives have an advantage in guessing their relatives STR profiles, as for every marker every descendent will inherit one allele from each parent. As a result, a direct descendent only has to guess one allele of his parents, knowing that the other allele is equivalent to one of his. This can be extended to other close relatives, resulting in the probabilities of sharing alleles given in Table S4

If the relative shares both alleles, he does not have to perform any guess, if he shares one, he has to guess which of the two he shares plus he has to guess the second allele, if he shares none the entropy of guessing the correct allele is calculated as derived above.

As an example, for a specific marker, full siblings have the following entropy of guessing their siblings alleles, knowing their own profile:

$$E_s = 0.25 \sum(-p_{ij} \log_2(p_{ij})) + 0.5 (1 + \sum(p_i \log_2(p_i))) + 0.25 \times 0.$$

Supporting References

1. H. H. Li, U. B. Gyllensten, X. F. Cui, R. K. Saiki, H. A. Erlich, N. Arnheim, Amplification and Analysis of DNA-Sequences in Single Human-Sperm and Diploid-Cells. *Nature* **335**, 414-417 (1988).
2. T. Tamura, M. Osawa, E. Ochiai, T. Suzuki, T. Nakamura, Evaluation of advanced multiplex short tandem repeat systems in pairwise kinship analysis. *Legal Med-Tokyo* **17**, 320-325 (2015).
3. C. Phillips, L. Fernandez-Formoso, M. Garcia-Magarinos, L. Porras, T. Tvedebrink, J. Amigo, M. Fondevila, A. Gomez-Tato, J. Alvarez-Dios, A. Freire-Aradas, A. Gomez-Carballa, A. Mosquera-Miguel, A. Carracedo, M. V. Lareu, Analysis of global variability in 15 established and 5 new European Standard Set (ESS) STRs using the CEPH human genome diversity panel. *Forensic Sci. Int. Genet.* **5**, 155-169 (2011).
4. K. L. O'Connor, A. O. Tillmar, Effect of linkage between vWA and D12S391 in kinship analysis. *Forensic Sci. Int. Genet.* **6**, 840-844 (2012).
5. Summit supercomputer ranked fastest computer in the world. (2018) <https://www.energy.gov/articles/summit-supercomputer-ranked-fastest-computer-world>, accessed: 25.01.2019.
6. L. Zhiye, US dethrones China with IBM Summit supercomputer. (2018) <https://www.tomshardware.com/news/us-supercomputer-china-top500-summit,37367.html>, accessed: 25.01.2019.
7. M. Arora, How secure is AES against brute force attacks? (2012) https://www.eetimes.com/document.asp?doc_id=1279619, accessed: 25.01.2019.
8. D. Stamat, AWS EC2: P2 vs P3 instances. (2017) <https://blog.iron.io/aws-p2-vs-p3-instances/>, accessed: 25.01.2019.
9. C. R. Hill, D. L. Duewer, M. C. Kline, M. D. Coble, J. M. Butler, US population data for 29 autosomal STR loci. *Forensic Sci. Int. Genet.* **7**, E82-E83 (2013).
10. C. R. Steffen, M. D. Coble, K. B. Gettings, P. M. Vallone, Corrigendum to 'US Population Data for 29 Autosomal STR Loci' [*Forensic Sci. Int. Genet.* **7** (2013) e82-e83]. *Forensic Sci. Int. Genet.* **31**, E36-E40 (2017).
11. C. Phillips, D. Ballard, P. Gill, D. S. Court, A. Carracedo, M. V. Lareu, The recombination landscape around forensic STRs: Accurate measurement of genetic distances between syntenic STR pairs using HapMap high density SNP data. *Forensic Sci. Int. Genet.* **6**, 354-365 (2012).
12. P. Hatzler-Grubwieser, B. Berger, D. Niederwieser, M. Steinlechner, Allele frequencies and concordance study of 16 STR loci - including the new European Standard Set (ESS) loci - in an Austrian population sample. *Forensic Sci. Int. Genet.* **6**, E50-E51 (2012).

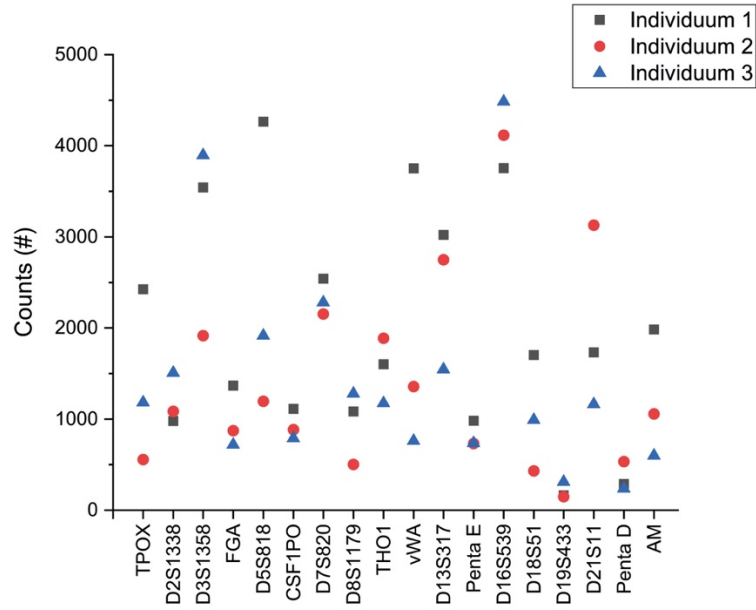


Figure S1. Sequencing coverage of the individual STR markers for three different individuals during initial key generation experiments (see Figure 1 in main text).

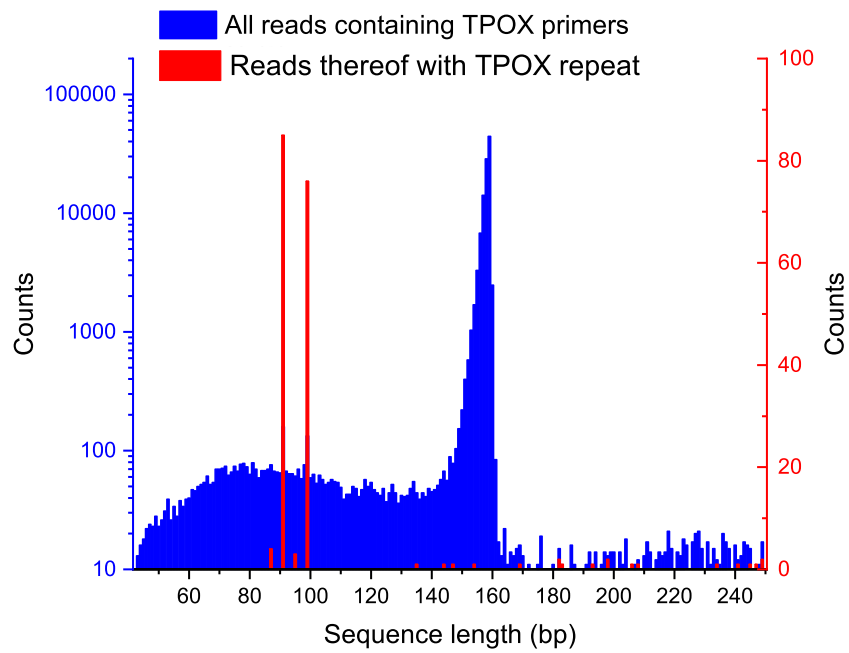


Figure S2. Analysis of the sequences starting and ending with TPOX primer sequences. The sequences in red additionally contain at least five copies of the TPOX repeat (AATG). The blue data shows the presence of the synthetic DNA amplicons (expected length of 159 bp), and the red data represents the alleles 6 and 8 for the TPOX STR marker of individual 1.

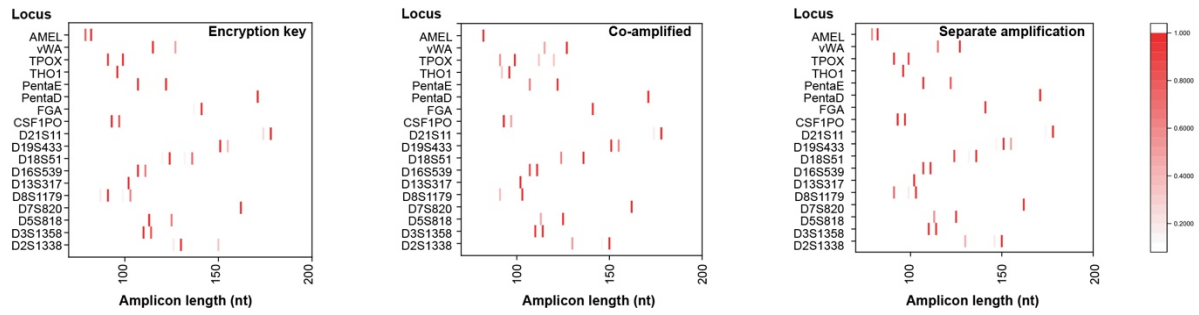


Figure S3. STR marker amplicon length profiles for individual 1 measured with three different procedures, either in the absence of synthetic DNA (left), co-amplified and co-sequenced together with the synthetic DNA (middle) and separately amplified, but co-sequenced with the synthetic DNA (right).

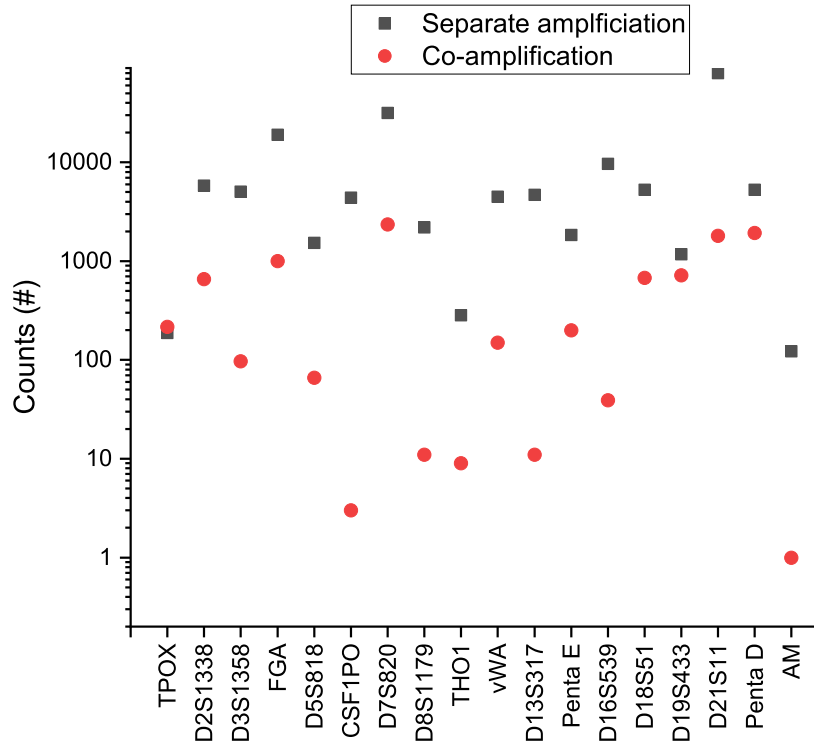


Figure S4. Coverage (counts) of the individual STR markers read during the decryption stage in the presence of the synthetic DNA.

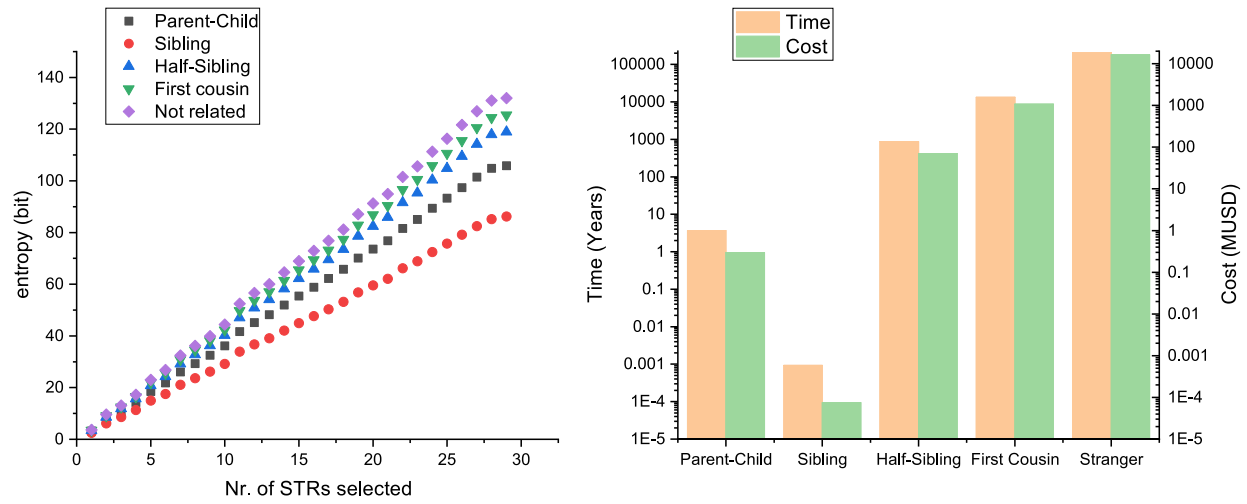


Figure S5. Entropy of the STR markers of Table 1 with pre-knowledge of the STR profile of a close relative, and the estimated time and costs to guess the relatives STR profile.

Table S1. Number of possible alleles reported in literature (9, 10) for each marker, number of allele possibilities per marker for a diploid genome, and the integer values derived from the STR profile for individuum 1. These integer values are utilized to generate the numerical key.

Marker	Possible alleles per haploid genome	Possibilities per diploid genome	Integer key Individuum 1
D2S1338	13	169	6/11
D3S1358	11	121	3/4
D5S818	9	81	3/6
D7S820	11	121	6/6
D8S1179	11	121	3/6
D13S317	8	64	6/6
D16S539	9	81	5/6
D18S51	22	484	4/9
D19S433	16	256	8/10
D21S11	27	729	11/11
CSF1PO	9	81	4/5
FGA	27	729	11/11
PentaD	16	256	13/13
PentaE	23	529	6/9
THO1	8	64	4/4
TPOX	10	100	6/8
vWA	11	121	4/7
AM	-	2	1/2
Total possibilities		1.4064E+38	

Table S2. STRs of the NIST 1036 tables, their position in the chromosome as well as reasoning for independence of the individual STR marker.

STR	Chromosome arm	Approx Pos. (Mb) ¹	Entropy All	Entropy Cauc	Independence (reasoning)
F13B	1q	197.04	3.63	3.26	include (<i>not linked</i> ^{**} , $Rc^*(F13B:D1S1656) = 0.34$)
D1S1656	1q	230.91	5.96	6.05	include (<i>not linked</i> ^{**} , $Rc^*(F13B:D1S1656) = 0.34$)
TPOX	2p	1.49	3.49	2.99	include (<i>not linked</i> ^{**} , $Rc^*(TPOX:D2S441) > 0.40$)
D2S441	2p	68.24	4.15	4.06	include (<i>not linked</i> ^{**} , $Rc^*(TPOX:D2S441) > 0.40$)
D2S1338	2q	218.88	5.77	5.60	include (<i>separate chromosome arm</i>)
D3S1358	3	45.58	3.79	3.95	include (<i>separate chromosome</i>)
FGA	4	155.51	5.62	5.25	include (<i>separate chromosome</i>)
D5S818	5q	123.11	3.75	3.32	include (<i>not linked</i> ^{**} , $Rc^*(CSF1PO:D5S818) = 0.25$)
CSF1PO	5q	149.46	3.76	3.32	include (<i>not linked</i> ^{**} , $Rc^*(CSF1PO:D5S818) = 0.25$)
F13A01	6p		4.44	3.73	include (<i>separate arm to SE33</i>)
SE33	6q	88.99	8.08	8.13	include (<i>separate chromosome</i>)
D6S1043	6q	92.45	<i>5.66</i>	<i>4.96</i>	exclude (<i>linked</i> ^{**} to SE33, $Rc(SE33:D6S1043)=0.044$)
D7S820	7q	83.79	4.15	4.36	include (<i>separate chromosome</i>)
LPL	8p		3.42	3.17	include (<i>separate arm</i>)
D8S1179	8q	125.91	4.61	4.50	include (<i>separate chromosome</i>)
Penta C	9p		4.30	3.87	include (<i>separate chromosome</i>)
D10S1248	10	2.24	4.04	3.73	include (<i>separate chromosome</i>)
THO1	11	2.19	3.87	3.74	include (<i>separate chromosome</i>)
vWA	12p	6.09	4.35	4.26	include (<i>not linked</i> ^{**} , $Rc^*(vWA:D12S391) = 0.117$)
D12S391	12p	12.45	5.88	5.94	include (<i>not linked</i> ^{**} , $Rc^*(vWA:D12S391) = 0.117$)
D13S317	13	82.72	4.18	4.15	include (<i>separate chromosome</i>)
FESFPS	15q		3.68	3.08	include (<i>not linked</i> ^{**} , $Rc^*(Penta E:FES-FPS) = 0.181$)
Penta E	15q	97.37	6.63	6.15	include (<i>not linked</i> ^{**} , $Rc^*(Penta E:FES-FPS) = 0.181$)
D16S539	16	84.94	4.08	3.84	include (<i>separate chromosome</i>)
D18S51	18	60.95	5.70	5.51	include (<i>separate chromosome</i>)
D19S433	19	30.42	5.01	4.32	include (<i>separate chromosome</i>)
D21S11	21q	20.55	5.31	4.93	include (<i>not linked</i> ^{**} , $Rc^*(Penta D:D21S11) > 0.3$)
Penta D	21q	45.06	5.29	4.64	include (<i>not linked</i> ^{**} , $Rc^*(Penta D:D21S11) > 0.3$)
D22S1045	22	37.54	4.09	3.55	include (<i>separate chromosome</i>)
AMEL	X & Y		1	1	include (<i>separate chromosome</i>)
Total					
Entropy			132.0 bit	124.5 bit	

* Rc = Recombination rate from Kosambi mapping function (11).

** Non-linkage has been shown for **Rc recombination fractions for ~0.12 (11), as found for vWA:D12S391, and various studies have shown marker independence for this relatively close STR pair for non-close relatives (2-4).

** Only included profiles (Bold values) used for sum

Table S3. STR Markers used in experimental work and their history in the forensic analysis.

	Chromosome	Codis Core Locus	New FBI core locus	European Locus*	Entropy per diploid genome (bit)
TPOX	2	YES	YES		3.49
D2S1338	2		YES		5.77
D3S1358	3	YES	YES	YES	3.78
FGA	4	YES	YES	YES	5.62
D5S818	5	YES	YES		3.75
CSF1PO	5	YES	YES		3.76
D7S820	7	YES	YES		4.15
D8S1179	8	YES	YES	YES	4.61
THO1	11	YES	YES	YES	3.87
vWA	12	YES	YES	YES	4.35
D13S317	13	YES	YES		4.18
PentaE	15				6.63
D16S539	16	YES	YES		4.08
D18S51	18	YES	YES	YES	5.70
D19S433	19		YES		5.01
D21S11	21	YES	YES	YES	5.31
PentaD	21				5.29
AMEL	Y&X		YES	YES	1
				Total	80.4 bit
				Entropy	

* European Standard Set of Loci and new ESS loci (12).

Table S4: Probability of having a given number of STR alleles shared between close relatives

Relationship	0 alleles	1 allele	2 alleles
Parent-child	0	1	0
Full siblings	1/4	1/2	1/4
Half siblings	1/2	1/2	0
Grandparent-grandchild	1/2	1/2	0
Uncle-Nephew	1/2	1/2	0
First cousins	3/4	1/4	0
Entropy per marker	$\sum(-p_{ij}\log_2(p_{ij}))$	$1 + \sum(p_i\log_2(p_i))$	0