

DNA-based Cryptography

Ashish Gehani, Thomas H. LaBean
and John H. Reif

A. Gehani, Thomas H. LaBean, and John H. Reif, DNA-based Cryptography, Proc. DNA Based Computers V: Cambridge, MA, June 14-16, 1999.

A. Gehani, Thomas H. LaBean, and John H. Reif, DNA-based Cryptography, chapter in "Aspects of Molecular Computing", Springer Verlag series in Natural Computing (edited by N. Jonoska, G. Paun and G. Rozenberg) LNCS 2950 Festschrift, Springer, pp. 167-188, (2004).

Biotechnological Methods (e.g., recombinant DNA) have been developed for a wide class of operations on DNA and RNA strands

DNA Computing: makes use of such biotechnological methods for doing computation

- Uses DNA as a medium for ultra-scale computation
- Comprehensive survey of Reif [R98]
- Use of DNA Computing for solution of combinatorial search problems:
 - Hamiltonian path problem [Adleman94]
 - Data Encryption Standard (DES)
[Boneh, et al 95] [Adleman, et al 96]

But ultimately limited by volume requirements, which may grow exponentially with input size.

DNA Storage of Data

- A medium for *ultra-compact information storage*: large amounts of data that can be stored in compact volume.
- Vastly exceeds storage capacities of conventional electronic, magnetic, optical media.
- A *gram* of DNA contains 10^{21} DNA bases
= *10^8 tera-bytes*.
- A few grams of DNA may hold *all data stored in world*.
- Most DNA computing techniques are applied at concentrations of 5 grams of DNA per liter of water.

DNA Data Bases:

- A “wet” data base of *biological data*
 - *natural DNA* obtained from biological sources may be recoded using nonstandard bases [Landweber,Lipton97], to allow for subsequent DNA Computing.
- DNA containing data obtained from more conventional *binary storage media*.
 - *input and output of the DNA data* can be moved to conventional binary storage media by *DNA chip arrays*
 - binary data may be *encoded* in DNA strands by use of an alphabet of short oligonucleotide sequences.
- Associative Searches within DNA databases:
 - methods for fast associative searches within DNA databases using hybridization [Baum95]
 - [Reif95] data base join operations and various massively parallel operations on the DNA data

Cryptography

Data security and cryptography are *critical* to computing data base applications.

Plaintext: non-encrypted form of message

Encryption: process of scrambling plaintext message; transformation to encrypted message (*cipher text*).

Example:

- fixed codebook provides an initial mapping from characters in the finite plaintext alphabet to a finite alphabet of codewords,
- then a sophisticated algorithm depending on a key may be applied to further encrypt the message.

Decryption: the reverse process of transforming the encrypted message back to the original plaintext message.

Cryptosystem: a method for both encryption and decryption of data.

Unbreakable cryptosystem: one for which successful cryptanalysis is not possible.

Our *Unbreakable* DNA *Cryptography Method:*

DNA-based, molecular cryptography system

- Plaintext message data encoded in DNA strands by use of a (publicly known) alphabet of short oligonucleotide sequences.
- Based on *one-time-pads* that are in principle *unbreakable* and may be practical for DNA:
 - Practical applications of cryptographic systems based on one-time-pads are *limited in conventional electronic media*, by the size of the one-time-pad.
 - DNA provides a much more *compact storage media*, and an extremely small amount of DNA suffices even for huge one-time-pads.

Our DNA *one-time-pad encryption scheme:*

- a *substitution method* using libraries of distinct pads, each of which defines a specific, randomly generated, pair-wise mapping
- an *XOR scheme* utilizing molecular computation and indexed, random key strings

Applications of DNA-based Cryptography Systems

- the encryption of (recoded) *natural DNA*
- the encryption of DNA encoding *binary data*.

Methods for *2D data input and output*:

- use of *chip-based DNA micro-array* technology
- transform between conventional binary storage media via (photo-sensitive and/or photo-emitting) DNA chip arrays

DNA Steganography Systems:

Clelland CT, Risca V, Bancroft C (1999) Hiding messages in DNA microdots. Nature 399:533–534

- ***Secretly tag*** the input DNA
- Then ***disguise*** it (without further modifications) within collections of other DNA.
- Original plaintext is ***not actually encrypted***
- Very appealing due to ***simplicity***.
- But we shall show it can be decrypted.

Example:

- DNA plaintext messages are appended with one or more secret keys
- Resulting appended DNA strands are hidden by mixing them within many other irrelevant DNA strands (e.g., randomly constructed DNA strands).
- Can be stored as amplifiable microdots

Our Decryption Results for *DNA Steganography Systems*:

- **Potential Limitations of these DNA Steganography methods:**
 - **Show proposed DNA steganography systems can be *broken*, with some assumptions on information theoretic entropy of plaintext messages.**
- **We also discuss various modified DNA steganography systems which appear to have *improved security*.**

Organization of Talk

- ❖ *Introduction* of DNA Computing and cryptography terminology, and results.
- ❖ *Unbreakable DNA cryptosystems* using randomly assembled *one-time pads*.
- ❖ Example of a *DNA cryptosystem for two dimensional images*, using a DNA chip for I/O and also using a randomly assembled one-time pad.
- ❖ *DNA Steganography Techniques:*
 - show that they can be *broken* with some modest assumptions on the entropy of the plaintext, even if they employ perfectly random one-time pads.
 - Provide possible improvements
- ❖ Conclusions

Cryptosystems Using Random One-Time Pads

Use *secret codebook* to convert short segments of plaintext messages to encrypted text:

- Must be *random* codebook
- Codebook can be used only *once*

In secret, assemble a large one-time-pad in the form of a DNA strand:

- randomly assembled from short oligonucleotide sequences,
- isolated, and cloned.

One-time-pad *shared in advance* by both the sender and receiver of the secret message:

- requires initial communication of one-time-pad between sender and receiver
- facilitated by compact nature of DNA

Our DNA Cryptosystem Using Substitution

A. Gehani, Thomas H. LaBean, and John H. Reif, DNA-based Cryptography, Proc. DNA Based Computers V: Cambridge, MA, June 14-16, 1999.

A. Gehani, Thomas H. LaBean, and John H. Reif, DNA-based Cryptography, chapter in "Aspects of Molecular Computing", Springer Verlag series in Natural Computing (edited by N. Jonoska, G. Paun and G. Rozenberg) LNCS 2950 Festschrift, Springer, pp. 167-188, (2004).

Our DNA Cryptosystem Using Substitution

Substitution one-time-pad encryption:

- a substitution method using libraries of distinct pads, each of which defines a specific, randomly generated, pair-wise mapping.
- The decryption is done by similar methods.

Input:

plaintext binary message of length n ,
partitioned into plaintext words of fixed length,

Substitution One-time-pad:

a table randomly mapping all possible strings of plaintext words into cipher words of fixed length, such that there is a unique reverse mapping.

Encryption:

by substituting each i th block of the plaintext with the cipher word given by the table, and is decrypted by reversing these substitutions.

DNA Implementation of Substitution One-time-pad Encryption:

- ***plaintext*** messages:
 - one test tube of short DNA strands
- ***encrypted*** messages:
 - another test tube of different short DNA strands

Encryption by substitution:

- maps these in a random yet reversible way
- plaintext is converted to cipher strands and plaintext strands are removed

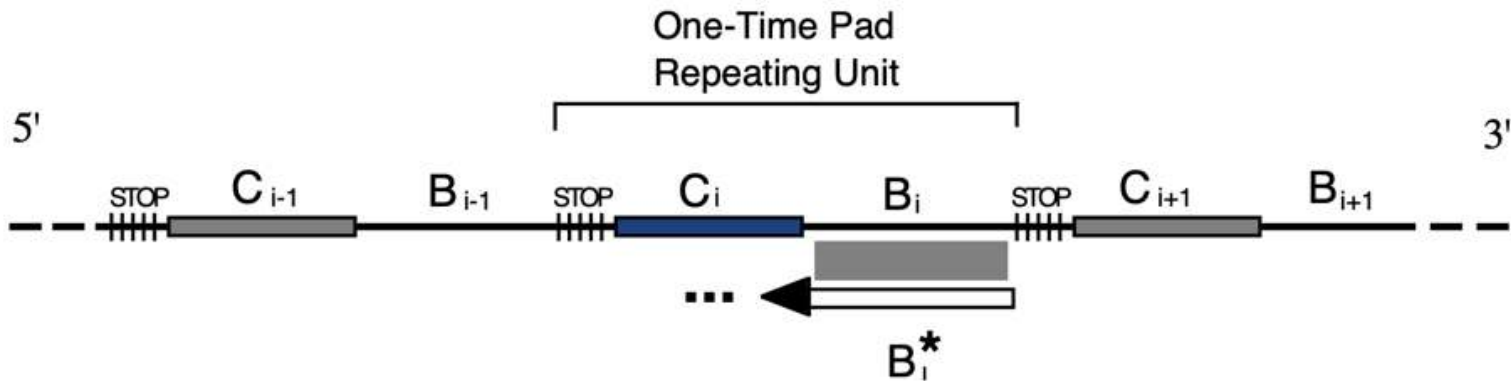
DNA Substitution one-time pads:

- use long DNA pads containing many segments:
 - each segment contains a cipher word followed by a plaintext word.
- ***cipher word:*** acts as a hybridization site for binding of a primer
- cipher word is appended with a ***plaintext word*** to produce ***word-pairs***.

These ***word-pair DNA strands*** used as a lookup table in conversion of plaintext into cipher text.

One-time-pad DNA Sequence:

- Length n with $d = n/(L_1 + L_2 + L_3)$ repeating units:



For each $i = 1, \dots, d$ the *ith Repeating Unit* is made up of:

- B_i = a cipher word of length $L_1 = c_1 \log n$
- C_i = a plaintext word length $L_2 = c_2 \log n$
 - Each sequence pair uniquely associates a plaintext word with a cipher word.
- Polymerase "stopper" sequence of length $L_3 = c_3$,

Generation of oligonucleotides corresponding to the plaintext/cipher word-pair strands:

- B_i^* used as polymerase primer, is *extended* with polymerase by specific attachment of plaintext word C_i .

Stopper sequence prohibits extension of growing DNA strand beyond boundary of paired plaintext word.

**Word-pair strands are essentially:
*a lookup table for a random
codebook.***

***Feasibility* depends upon:**

- size of the lexicon;
- number of possible pads available;
- size, complexity, and frequency of message transmissions.

<u><i>Parameter</i></u>	<u><i>Range</i></u>
Lexicon size	10,000 – 250,000 words
Word size	8 – 24 bases
Message size	5 – 30% of lexicon size
Pad diversity	10^6 - 10^8

***Pad diversity:* total number of random pads generated during a single pad construction experiment.**

Codebook Libraries:

- previous gene library construction projects [LK93, LB97]
- used in DNA word encoding methods used in DNA computation [DMGFS96, DMGFS98, DMRGF+97, FTCSC97, GDNMF97, GFBCL+96, HGL98, M96].

Use *two distinct lexicons* of sequence words:

- for cipher words
- for plaintext words.

Can *generate lexicons* by normal DNA synthesis methods:

- utilize sequence *randomization at specific positions* in sequence words.

Example:

For $N = A+C+G+T$, $R = A+G$, and $Y = C+T$:

RNNYRNRRYN produces

$2 \times 4 \times 4 \times 2 \times 2 \times 4 \times 2 \times 2 \times 2 \times 4 = 16,384$ possible sequences.

Methods for Construction of DNA one-time pads.

- (1) *Random assembly* of one-time pads in solution (e.g. on a synthesis column).
 - *Difficult to achieve both full coverage* and yet still avoiding possible *conflicts by repetition* of plaintext and/or cipher words.
 - can set c_1 and c_2 large so probability of repeated words on pad of length n is small, but coverage is be reduced.
- (2) Use of *DNA chip technology* for random assembly of one-time pads

Advantages:

- currently commercially available (Affymetrix) chemical methods for construction of custom variants are well developed.
- direct control of coverage and repetitions

***DNA chip Method* for Construction of DNA one-time pads.**

- an array of immobilized DNA strands,
- multiple copies of a single sequence are grouped together in a microscopic pixel.
- optically addressable
- known technology for synthesis of distinct DNA sequences at each (optically addressable) site of the array.
- combinatorial synthesis conducted in parallel at thousands of locations:
 - For preparation of oligonucleotides of length L , the 4^L sequences are synthesized in $4n$ chemical reactions.

Examples:

- 65,000 sequences of length 8 use 32 synthesis cycles
- 1.67×10^7 sequences of length 10 use 48 cycles

DNA Chip Method for

Construction of DNA One-time pads

- *plaintext and cipher pairs* constructed:
- *nearly complete coverage* of the lexicon on each pad,
- *nearly unique word mapping* between plaintext and cipher pairs.
- resulting cipher word, plaintext word pairs can be assembled together in random order (with possible repetitions) on a long DNA strand by a number of known methods:
 - blunt end ligation
 - hybridization assembly with complemented pairs [Adleman97]
- Cloning or PCR used to *amplify* the resulting one-time pad.

XOR One-time-pad **(Vernam Cipher) Cryptosystem**

One-time-pad S:

a sequence of independently distributed random bits

M: a plaintext binary message of n bits

• Encrypted bits:

$$C_i = M_i \text{ XOR } S_i \text{ for } i = 1, \dots, n.$$

XOR: given two Boolean inputs, yields 0 if the inputs are the same, and otherwise is 1.

• Decrypted bits:

Use commutative property of XOR:

$$\begin{aligned} C_i \text{ XOR } S_i &= (M_i \text{ XOR } S_i) \text{ XOR } S_i \\ &= M_i \text{ XOR } (S_i \text{ XOR } S_i) \\ &= M_i. \end{aligned}$$

DNA Implementation of XOR One-time-pad Cryptosystem

- *plaintext messages:*
one test tube of short DNA strands
- *encrypted messages:*
another test tube of different short DNA strands

Encryption by XOR One-time-pad:

- maps these in a random yet reversible way
- plaintext is converted to cipher strands and plaintext strands are removed

For *efficient* DNA encoding:

- use *modular base 4* (DNA has four nucleotides)
- Encryption:
addition of one-time-pad elements modulo 4
- Decryption:
subtract one-time-pad elements modulo 4

Details of DNA Implementation of XOR One-time-pad Cryptosystem

- Each plaintext message has appended unique *prefix index tag* of length L_0 indexing it.
- Each of one-time-pad DNA sequence has appended unique *prefix index tag* of same length L_0 , forming *complements* of plaintext message tags.
- Use Recombinant DNA techniques (annealing and ligation) to *concatenate into a single DNA strand* each corresponding pair of a plaintext message and a one-time-pad sequence
- These are *encyphered by bit-wise XOR computation*:
 - fragments of the plaintext are converted to cipher strands using the one-time-pad DNA sequences, and
 - plaintext strands are removed.

Reverse decryption is similar: use commutative property of bit-wise XOR operation.

DNA Computing Methods to effect bit-wise XOR on Vectors.

Can adapt DNA Computing methods for *binary addition*:

- similar to bit-wise XOR computation
- can *disable carry-sums* logic to do XOR

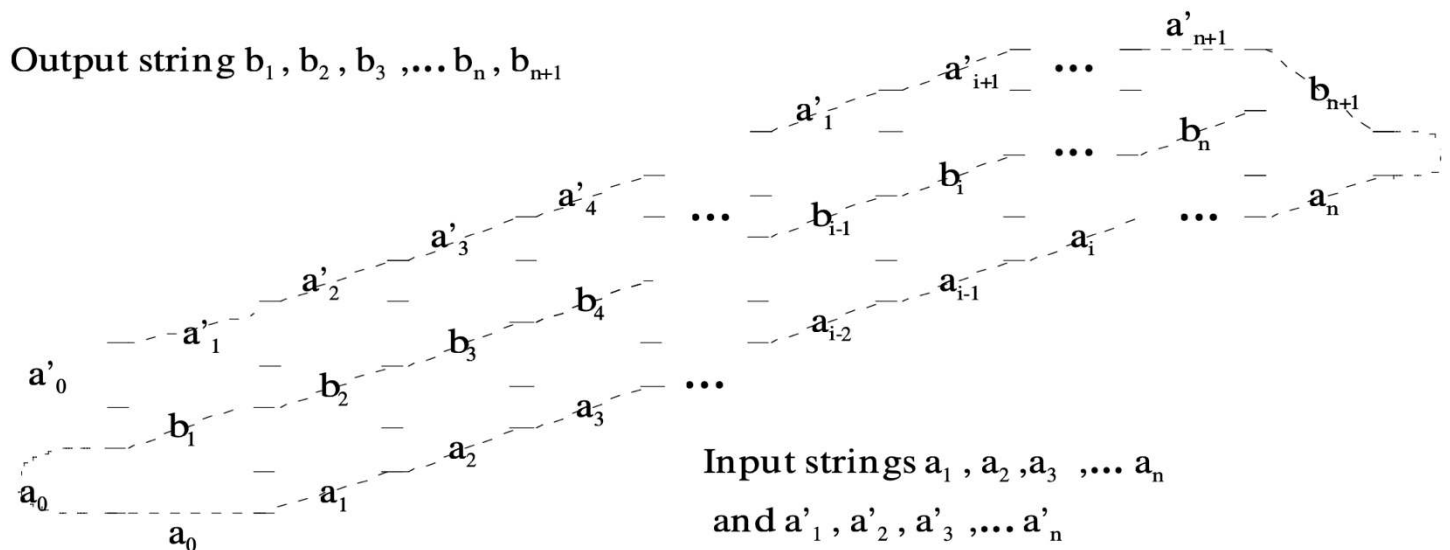
DNA Computing techniques for Integer Addition:

- (1) [Guarnieri, Fliss, and Bancroft 96] first DNA Computing addition operations (on single bits).
- (2) [Rubin et al 98, OGB97,LKSR97,GPZ97] permit chaining on n bits.
- (3) Addition by *Self Assembly* of DNA tiles
[Reif,97][LaBean, et al,2000]

XOR by Self Assembly of DNA tiles

Thomas H. LaBean, Erik Winfree, and John H. Reif, Experimental Progress in Computation by Self-Assembly of DNA Tilings, Proceeding of DNA Based Computers V: Cambridge, MA, June 14-16, 1999. Published in DIMACS Series in Discrete Mathematics and Theoretical Computer Science, Volume 54, edited by Erik Winfree and D.K. Gifford, American Mathematical Society, Providence, RI, pp. 123-140, (2000).

Chengde Mao, Thomas H. LaBean, John H. Reif, Nadrian C. Seeman, Logical Computation Using Algorithmic Self-Assembly of DNA Triple-Crossover Molecules, Nature, vol. 407, pp. 493-495 (Sept. 28 2000); C. Erratum: Nature 408, 750-750 (2000).



XOR by Self Assembly of DNA tiles

[LaBean, et al,2000]

[1] For each bit M_i of the message, construct sequence a_i that represents the i th bit.

[2] Scaffold strands for binary inputs to the XOR:

- Using linkers, assemble message M 's n bits into scaffold strand sequence $a_1 a_2 \dots a_n$,**
- One-time-pad is a further portion of the scaffold strand $a'_1 a'_2 \dots a'_n$ which are created from random inputs**

[3] Add output tiles; annealing give self assembly of the tiling.

XOR by Self Assembly of DNA tiles cont.

[4] adding ligase yeilds reporter strand

$$\mathbf{R = a_1 a_2 \dots a_n . a'_1 a'_2 \dots a'_n . b_1 b_2 \dots b_n}$$

where $b_i = a_i \text{ XOR } a'_i$, for $i = 1, \dots, n$.

[5] reporter strand is extracted by melting away the tiles' smaller sequences, and purifying.

➤ **contains concatenation of:**

input message, encryption key, ciphertext

[6] Using a marker sequence:

ciphertext can be excised and separated based on its length being half that of remaining sequence.

[7] Ciphertext can be stored in a compact form

DNA Cryptosystem for 2D Images

Our DNA 1-Time Pad Cryptosystem consists of:

- Data set to be encrypted: *2-dimensional image*
- *DNA Chip* bearing immobilized DNA strands:
 - addressable array of nucleotide sequences immobilized s.t. multiple copies of single sequence grouped together in a microscopic pixel.
- *Library of one-time pads* encoded on long DNA strand

DNA Cryptosystem for 2D Images

using:

- **DNA Chip and**
- **Randomly Assembled One-Time Pad**

Encryption and Decryption of 2D images recorded on microscopic arrays of a DNA chip:



Message

Encrypted

Decrypted

Simulated fluorescence microscopy of DNA I/O chip.

Initialization and Message Input

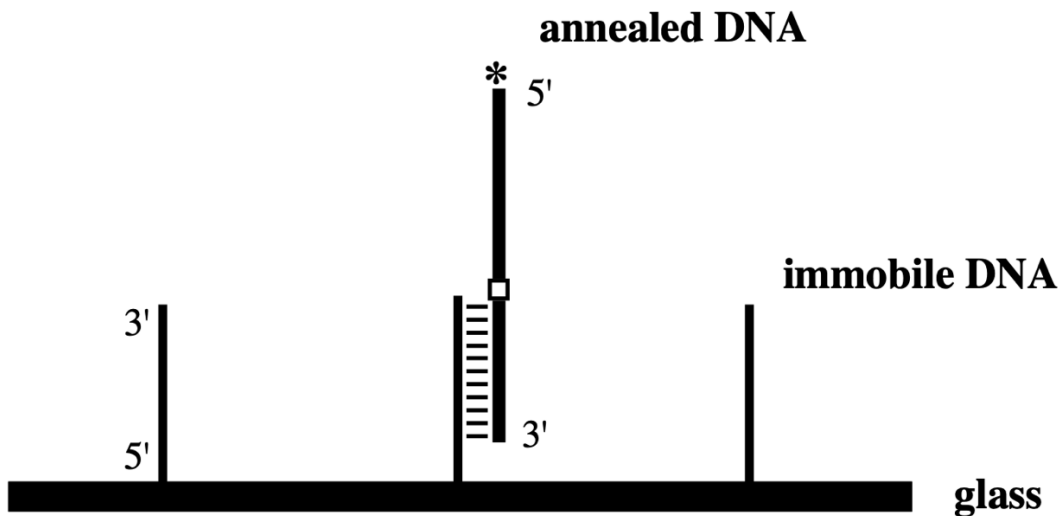
- **Fluorescent-labeled, word-pair DNA strands are prepared from a substitution pad codebook**
- **These are annealed specifically to their sequence complements at unique sites (pixels) on the DNA chip.**
- **The message information is transferred to a photo mask with transparent (white) and opaque (black) regions:**



Message Input to DNA Chip

Initialization and Message Input

- *Immobile DNA strands* are located on the glass substrate of the chip in a sequence addressable grid.
- *Word-pair strands* prepared from random substitution pad:
 - The 5' (*unannealed*) end carries a cipher word
 - The 3' (*annealed*) end carries a plaintext word.
 - Contains a *photo-cleavable* base analog between two sequence words.

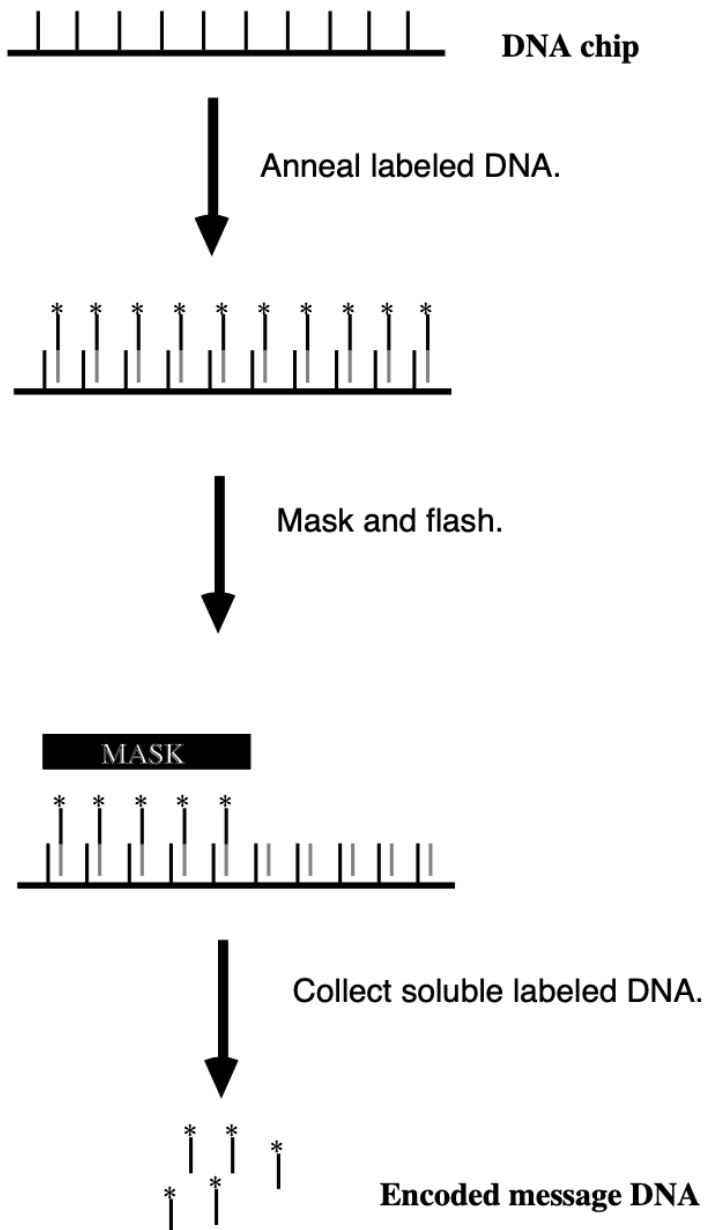


Annealed DNA:

- a fluorescent label on its 5' end (asterisk);
- codebook-matching sequence word
- a photo-labile base (white square) capable of cleaving the DNA backbone; and
- a chip-matching word (base-paired to immobile strands)

Encryption:

Encryption Scheme



[1] start with DNA chip displaying sequences *complementary* to plaintext lexicon.

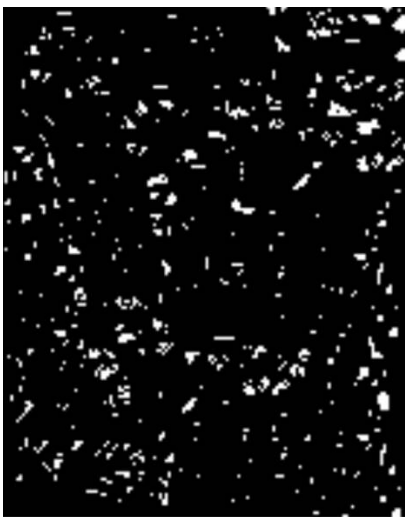
[2] fluorescent-labeled word-pair strands from one-time-pad *annealed* to chip at pixel with complement to plaintext 3' end.

[3] mask protects some pixels from a light-flash. At unprotected regions, DNA *cleaved* between plaintext & cipher words.

[4] cipher word strands, still labeled with fluorophore at 5' ends, are collected and transmitted as *encrypted message*.

Encryption of the Plaintext Message:

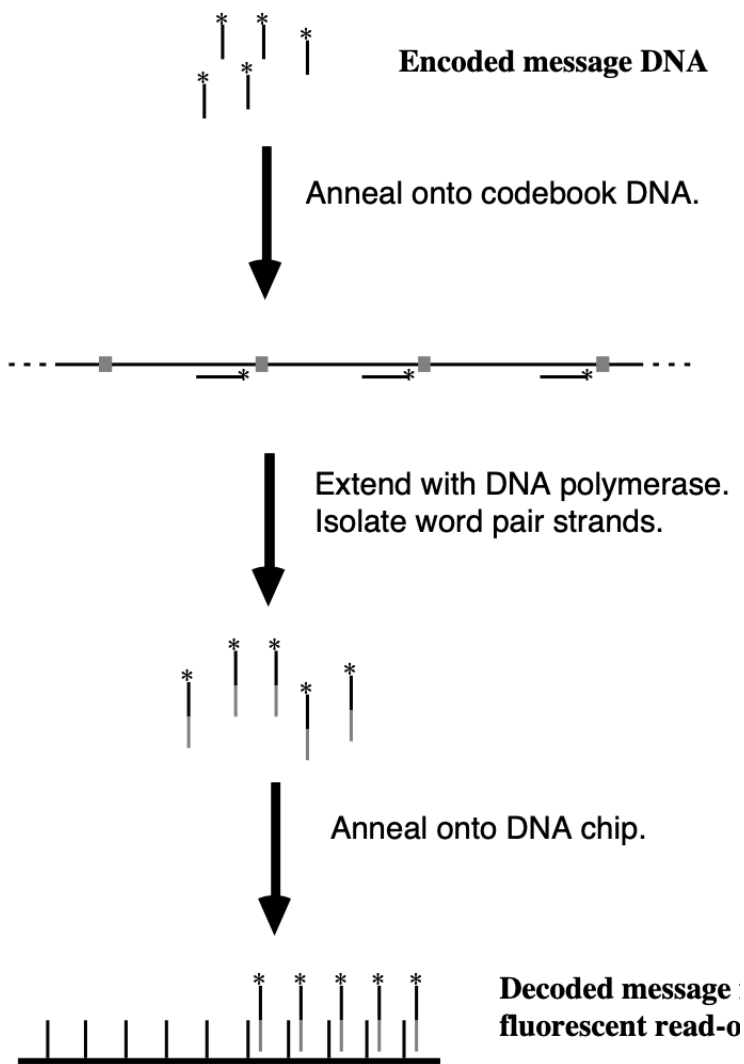
- Following a light-flash of mask-protected chip, *annealed* oligonucleotides beneath transparent mask pixels are *cleaved* at a photo-labile position:
 - their 5' sections are dissociated from annealed 3' section and collected in solution.
- This test tube of strands is *encrypted message*.
- Annealed oligos *beneath opaque mask* are unaffected by light-flash and can be *washed off chip*.
- If *encrypted message oligos are reannealed* onto a (washed) DNA chip, message information would be *unreadable*:



Simulated DNA Chip Read-Out of Encrypted Message

Decryption

Decryption Scheme



[1] word-pair strands constructed, *appending* cipher word with proper plaintext word, by polymerase extension or lop-sided PCR using cipher words as primer and one-time-pad as template.

[2] cipher strands *bind* to their specific locations on the pad and are appended with their plaintext partner.

[3] binding reformed *word-pair strands* to DNA chip and reading message by fluorescent microscopy.

Decryption of the Message

- use the fluorescent labeled oligos as primers in one-way (lopsided) PCR with the same one-time codebook which was used to prepare the initial word-pair oligos.
- When word-pair PCR product is bound to the same DNA chip, the decrypted message is revealed:



Decrypted Message

Simulated DNA Chip Read-Out of Decrypted Message

Steganography:

a class of techniques that hide secret messages within other messages:

- plaintext is not actually encrypted but is instead disguised or hidden within other data.

Historical examples:

- use of grills that mask out all of an image except the secret message,
- micro-photographs placed within larger images
- invisible inks, etc.

Advantages:

- it is very appealing due to its *simplicity*.

Disadvantages:

- Cryptography literature generally consider conventional steganography methods to have *low security*:
 - All known steganography methods have been often broken in practice [Kahn67] and [Schneier96]

DNA Steganography Techniques:

Clelland CT, Risca V, Bancroft C (1999) Hiding messages in DNA microdots. Nature 399:533–534

- Take one or more input DNA strands (considered to be the plaintext message)
- Append to them one or more randomly constructed “*secret key*” strands.
- Resulting “*tagged plaintext*” DNA strands are *hidden* by mixing them within many other additional “distracter” DNA strands which might also be constructed by random assembly.

Decryption:

- Given *knowledge* of the “*secret key*” strands,
- Resolution of DNA strands can be decrypted by a number of possible known recombinant DNA separation methods:
 - plaintext message strands *separated out* by hybridization with complements of “secret key” strands can be placed in solid support on magnetic beads or on a prepared surface.
 - These separation steps may combined with amplification steps and/or PCR

DNA Steganography Systems:

Clelland CT, Risca V, Bancroft C (1999) Hiding messages in DNA microdots. Nature 399:533–534

(1) Encoding rule: A novel encoding method is used instead of the traditional binary encoding. Nucleotides are used as quaternary code and each letter is denoted by three nucleotides. For example, the letter A is denoted by CGA, the letter B is denoted by CCA, etc.

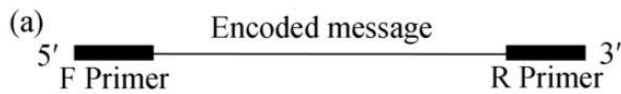
(2) Synthesizing secret-message DNA: The secret message is encoded into DNA sequence according to the above code. For instance, AB is encoded as CCGCCA. After coding, they synthesized a secret- message DNA oligodeoxynucleotide containing an encoded message 69 nucleotides long flanked by forward and reverse PCR primers, each 20 nucleotides long. Thus, the secrete-message DNA is prepared.

(3) Hiding message: They prepared concealing DNA that is physically similar to the secret-message DNA by sonicating human DNA to roughly 50 to 150 nucleotide pairs (average size) and denaturing it. The secret-message DNA and concealing DNA were mixed and attached on a piece of paper using common adhesives to form colorless microdots. Then the paper containing microdots can be posted by general mail service.

(4) Read Message: The sharing secrets for the sender and the receiver are encoding rule and primers. After the receiver gets the paper, he can easily find the microdots. Since the intended receiver had gotten the primers and encoding rule through a secure way, he could amplify the secret-message DNA by perform PCR on DNA microdots, sequence it and retrieve the message (plain- text) according to the encoding rule.

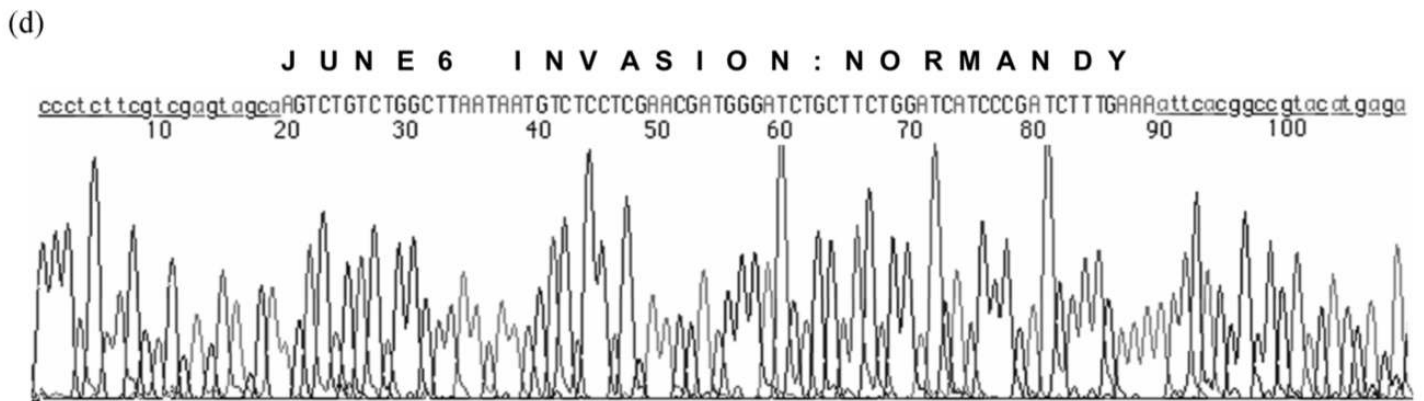
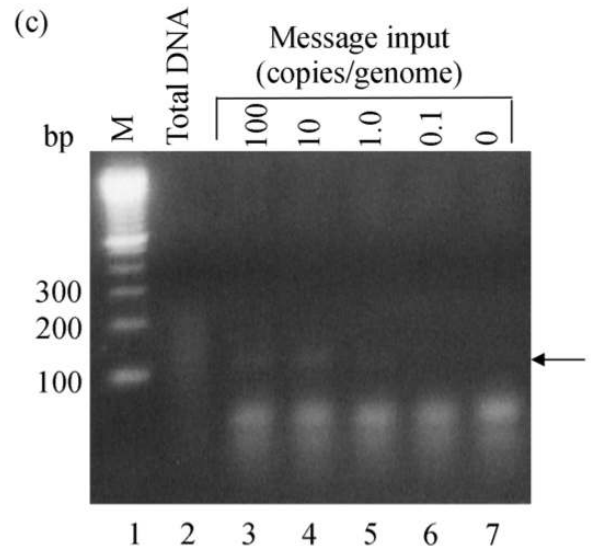
DNA Steganography Systems:

Clelland CT, Risca V, Bancroft C (1999) Hiding messages in DNA microdots. Nature 399:533–534



(b) Encryption key

A=CGA	K=AAG	U=CTG	0=ACT
B=CCA	L=TGC	V=CCT	1=ACC
C=GTT	M=TCC	W=CCG	2=TAG
D=TTG	N=TCT	X=CTA	3=GCA
E=GGC	O=GGA	Y=AAA	4=GAG
F=GGT	P=GTG	Z=CTT	5=AGA
G=TTT	Q=AAC	=ATA	6=TTA
H=CGC	R=TCA	=TCG	7=ACA
I=ATG	S=ACG	=GAT	8=AGG
J=AGT	T=TTC	=GCT	9=GCG



Example of Clelland et al steganography method:

- Secret-message synthesizing process is shown
- Text encoding rule.
- PCR result.
- Secret-message DNA and corresponding message (plaintext).

Cryptanalysis of DNA

Steganography Systems:

DNA steganography system's *security is entirely dependent* on degree that message DNA strands are *indistinguishable* from “distracter” DNA strands.

Cryptanalysis Assumptions:

- no knowledge of the “secret key” strands
- secret tags are indistinguishable from “distracter” DNA strands.
- plaintext is not initially compressed, and comes from a source (e.g., English or natural DNA) with Shannon information theoretic entropy $E_s > 1$
- the “distracter” DNA strands are constructed by random assembly

Then: DNA Steganography System can be *broken:*

- Original plaintext portion of “tagged plaintext” DNA strands are *distinguishable* from “distracter” DNA strands, and

Shannon (information theoretic)

Entropy E_s

- Provides a measure of the factor that a source can be *compressed* without loss of information.

Examples:

- Most images have entropy nearly 4
- English and other natural language text has entropy about 3
- Computer programs have entropy about 5
- Most natural DNA have entropy range 1.2 to 2

Lossless Data Compression [Lempel-Ziv 77]

Ziv, Jacob; Lempel, Abraham (May 1977). "A Universal Algorithm for Sequential Data Compression". IEEE Transactions on Information Theory. 23 (3): 337–343. CiteSeerX 10.1.1.118.8921 . doi:10.1109/TIT.1977.1055714

Input: text string of length n with entropy E_s

[1] Form a *dictionary* D of the $d = n/L$ most frequently occurring subsequences of length at least $L = E_s \log_2 n$ in the known source distribution.

[2] In place of subsequences of the input text matching with elements of the dictionary D , *substitute their indices* in the dictionary D .

Cryptanalysis of a DNA Steganography System:

Input: test tube T containing:

a mixture of “tagged plaintext” DNA strands mixed with a high concentration of “distracter” DNA strands, of length n .

- form a *dictionary* D of the $d = n/L$ most frequently occurring subsequences of length at least $L = E_s \log_2 n$ in the known plaintext source distribution.
- Give procedure for *separating* out plaintext message strands by repeated rounds of hybridization with complements of elements of D.

$r(T) =$ *ratio of concentration* of “distracter” DNA strands to “tagged plaintext” DNA strands.

On each *round of separation*:

- form a new test tube $F(T)$ with expected $r(F(T))$ considerably *reduced* from the previous ratio $r(T)$.

Separation Procedure for Extracting Secret Message:

[1] Pour a fraction $s = 1/2$ of volume of current test tube T into a test tube T_1 and pour remaining fraction $1-s$ of T into test tube T_2 .

[2] Choose a *random text phrase* x in D (not previously considered in a prior trial), and using Watson-Crick complement of x , do a *separation* on test tube T_2 , yielding a new test tube T_3 whose contents are only DNA strands containing phrase x .

[3] Pour contents of test tubes T_1 and T_3 into a new test tube $F(T)$.

- Ratio $r(F(T))$ of “distracter” DNA strands to plaintext DNA will expect to *decrease* from original ratio $r(T)$ by a constant factor $c < 1$
- After $O(\log(r/r'))$ repeated rounds of this process, ratio of concentration in test tube T will expect to *decrease* from initially $r = r(T)$ to any given smaller ratio r' .

Another cryptanalysis technique for breaking steganographic systems:

Use “*hints*” that *disambiguate plaintext*.

Example:

- wish to make secret the DNA of an individual (e.g., the President).
- use an improved steganography system: “distracter” DNA strands (mixed with DNA of an individual) are from a similar genetic pool.
- Steganography system may often be *broken* by use of distinguishing “*hints*” concerning DNA of the individual
 - e.g., the individual might have a particular set of observable expressed gene sequences (e.g., for baldness, etc.).
- *Hints* may allow for subsequent identification of the full secret DNA:
 - use of a series of separation steps with complement of portions of known gene sequences.

Improved DNA Steganography Systems with Enhanced Security:

Idea: make it more difficult to *distinguish* probability distribution of plaintext source from that of “distracter” DNA strands.

(1) Mimicking Distribution of “Distracter” DNA:

- use improved construction of the set of “distracter” DNA strands, so distribution better mimics the plaintext source distribution
- construct the “distracter” DNA strands by random assembly from elements of Lempel-Ziv dictionary.
- *Drawback:* Cryptanalysis using “hints” that disambiguate plaintext.

Improved DNA Steganography Systems with Enhanced Security:

Idea: make it more difficult to *distinguish* probability distribution of plaintext source from that of “distracter” DNA strands.

(2) Compression of Plaintext.

- Recode the plaintext using a universal lossless compression algorithm (e.g., Lempel-Ziv 77].
- Resulting distribution of the recoded plaintext approximates a universal distribution, so uniformly random assembled distracter sequences may suffice to provide improved security.
- *Drawbacks:*
 - Unlike conventional steganography methods, plaintext messages need to be preprocessed.
 - May still be decryptable.

Conclusion

Presented an initial investigation of DNA-based methods for Cryptosystems.

- **Main Results** for DNA one-time-pads cryptosystems:
 - Gave DNA substitution and XOR methods based on one-time-pads that are in principle *unbreakable*.
 - Gave an implementation of our DNA cryptography methods including *2D input/output*.
- Results for *DNA Steganography*:
 - A class of proposed DNA steganography methods offer only limited security; can be *broken* with some reasonable assumptions on entropy of plaintext messages.
 - Modified DNA steganography systems may have some *improved security*.

Open Problem

Show whether DNA steganography systems with natural DNA plaintext input can or cannot be made to be *unbreakable*.