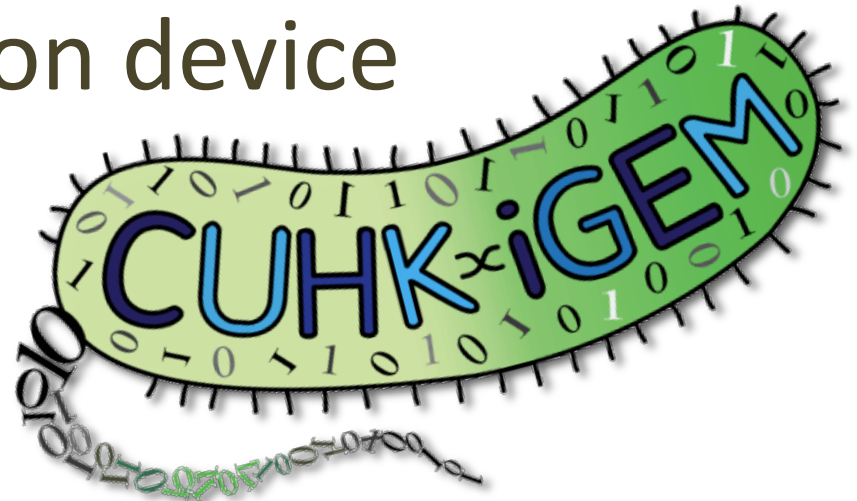
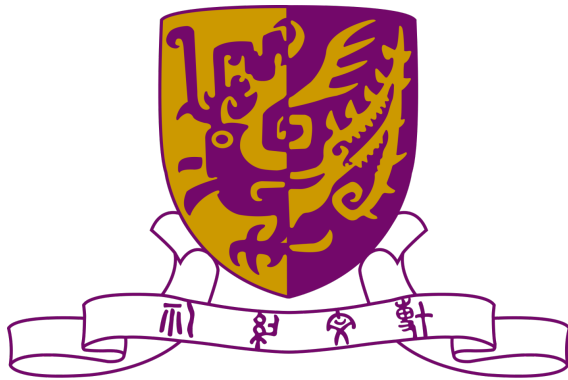


The Chinese University of Hong Kong – iGEM 2010

Bacterial based storage and
encryption device



Bacterial based information storage device

- Bancroft's group (2001)

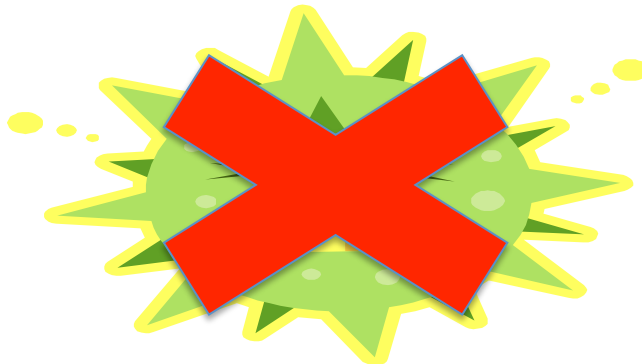
Mount Sinai School of Medicine

- Yachie's group (2007)

Keio University



However..



This year, The CUHK...

- True, massively parallel bacterial storage system

It is not the only thing we did..!!



In Addition...

- Encryption module with DNA shuffling system
 - Rci system
- The data proof-read
 - Checksum
- Strategy deal with synthesis/sequencing difficulties
 - Homopolymer, repetitive sequence

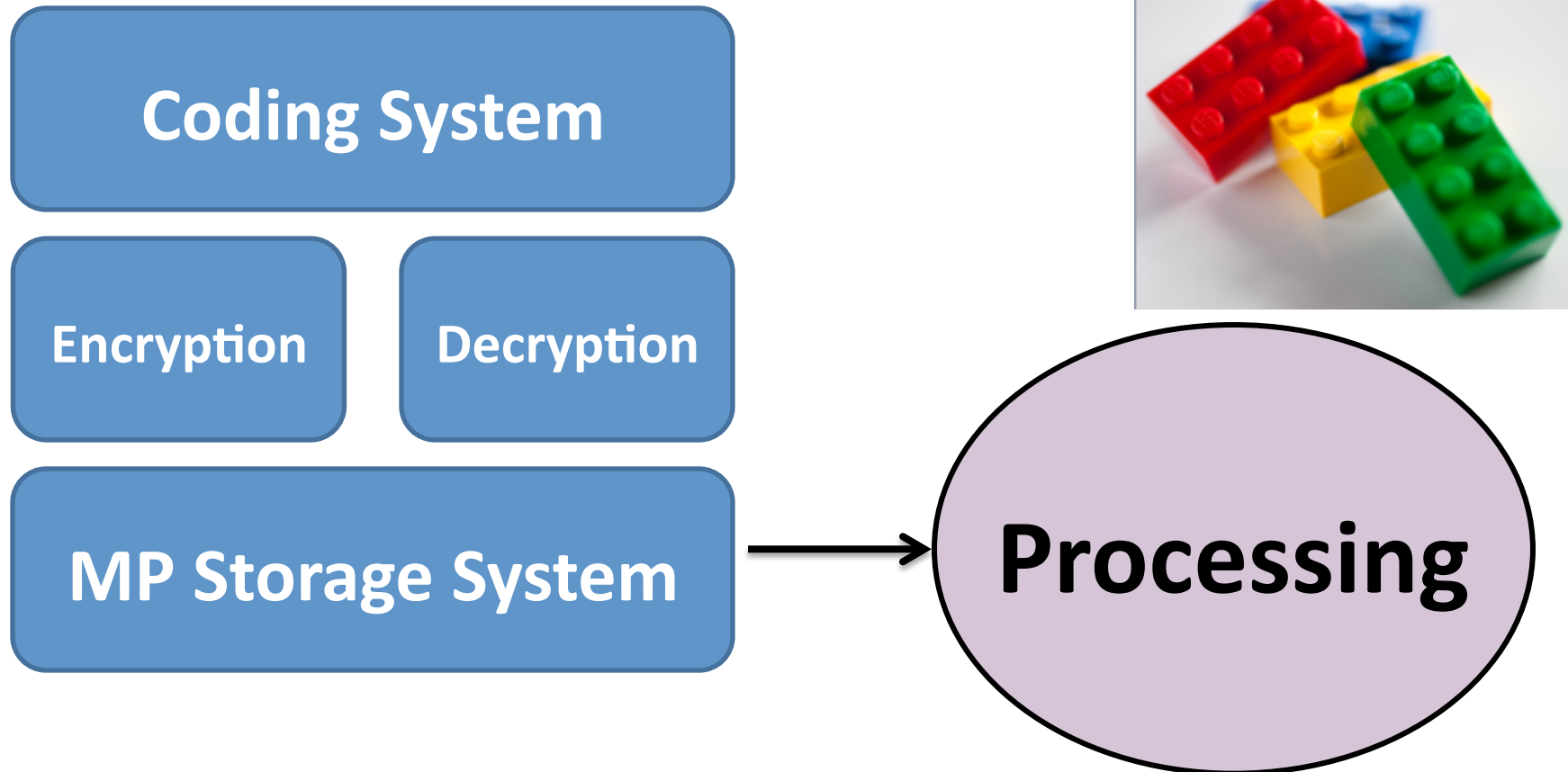




A LIVING DATA STORAGE SYSTEM

Basic infrastructure of the system

Module-based



Coding System

Encryption

Decryption

MP Storage System

Coding System

Encoding table

Quaternary Number System

DNA sequence

Compression

Use Numbers to Represent the Letters

From ASCII Table	Change to Quaternary Numbers
• i = 105	• 105 → 1221
• G = 71	• 71 → 0113
• E = 69	• 69 → 0111
• M = 77	• 77 → 0131

iGEM → 1221011301110131

Use “A, T, C & G” to Represent the Numbers

- 0 = A
- 1 = T
- 2 = C
- 3 = G

iGEM → 12211011301110131
→ ATCTATTGATTTATGT

Enter your text here:

24 chars

iGEM is very interesting

Original message input

Quaternary Encoding:

96 chars

122110131011103102001221130302001312121113021321020012211232131012111302121113031310122112321213

Converted to Quaternary number

DNA Encoding:

96 bp

TCCTTATGTATTTAGTACAATCCTTGAGACAATGTCTCTTTGACTGCTACAATCCTTCGCTGTATCTTTGACTCTTTGAGTGTATCCTTCGCTCTG

Converted to DNA sequence

Coding System

Encryption

Decryption

MP Storage System

Coding - Compression

- **DEFLATE — a compression algorithm**

1. Can reduce the **homopolymer** and **repetitive regions**

2. Can store more information

Homopolymer

Vector Insert Size

(~200 bp will be used by shufflon system)



1000 bp

Minimum Number of bacteria required for storing compressed DNA encoded message:

1

Enter your text here:

650 chars

[illegible]

Quaternary Encoding:

2600 chars

[illegible]

DNA Encoding:

2600 bp

Handwriting practice lines with a blue top line and a red bottom line. The lines are empty for writing practice.

Compressed DNA Encoding:

60 bp

TGCAGTCCAACGAAGTATGTAATTCCAGTCAAAGAAAAAAAAAATGGAATCGGTTGGTAG

The length and repetitive sequence is greatly reduced

Convert

Repetitive Regions

Vector Insert Size

(~200 bp will be used by shufflon system)

A horizontal line representing a 1000 bp DNA fragment. A single vertical tick mark on the line indicates a restriction site. A label '1000 bp' is at the right end of the line.

Minimum Number of bacteria

required for storing compressed
DNA encoded message:

1

Enter your text here:

672 chars

[illegible]

Quaternary Encoding:

2688 chars

[illegible]

DNA Encoding:

2688 bp

[illegible]

Compressed DNA Encoding:

69 bp

TGCAGTCCTGAGAAGCGGTTGGAATGTCATGCTTCAGTGCCTAATCACAAAAGGCTCCGGGATAAGCT

Convert

Coding System

Encryption

Decryption

MP Storage System

Encryption

- Provide DNA variation
- DNA Shuffling system
- Examples:

Homologous recombination
RACHITT

Characters!!!



Fragmentation of message

- Larger than the maximum vector insertion size
 - Limitation of current DNA synthesis technology
- Split the message into different parts

How do you deal with the problem of positioning?

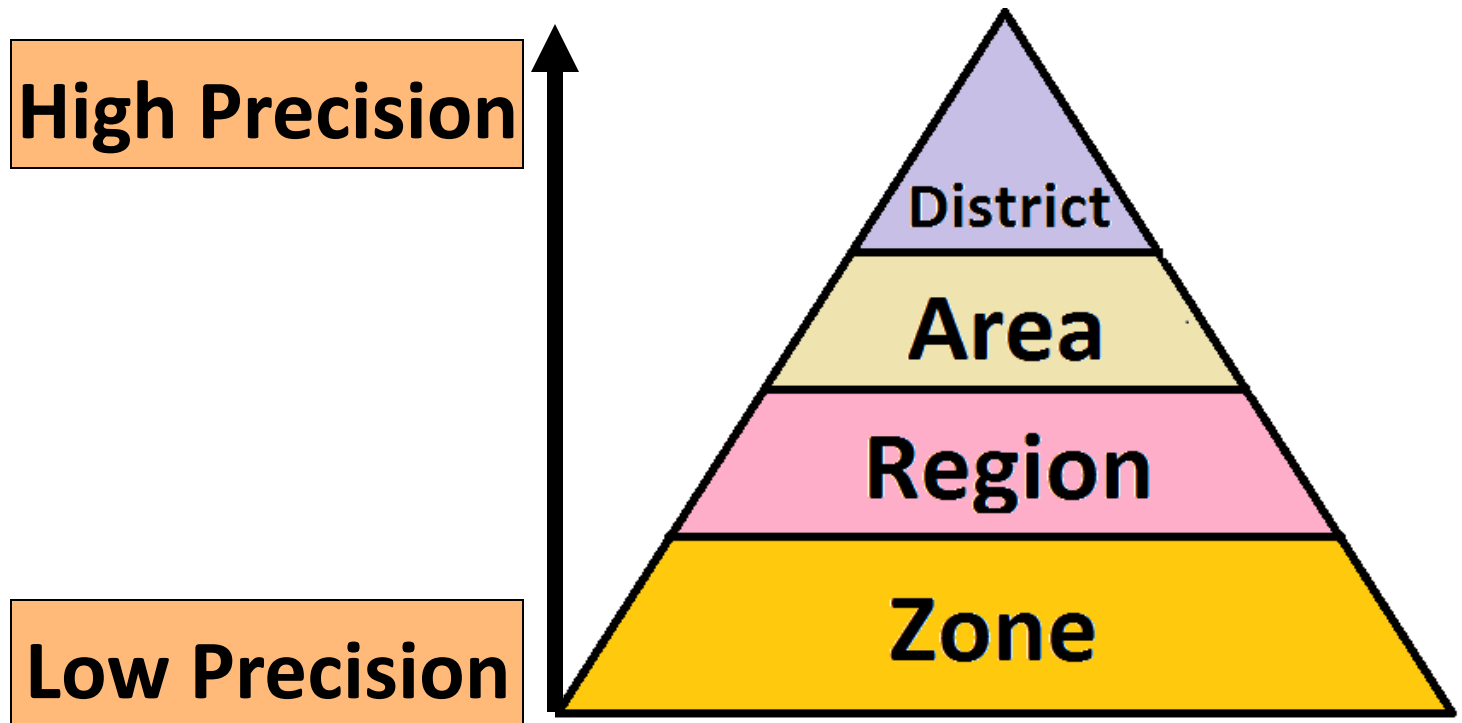


Postal Code

Storage – Massively parallel

Header – Locate particular data fragment of the message

Analogy to the hard disk : 4 address units



Header of
2nd
fragment



AGAT	AGAC	AGTA	AGAG
------	------	------	------

Header of
1st
fragment



AGAT	AGAC	AGAG	AGCT
------	------	------	------

Header of
3rd
fragment

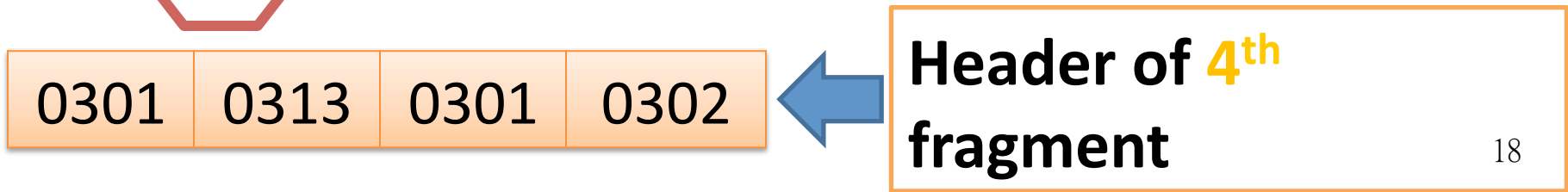
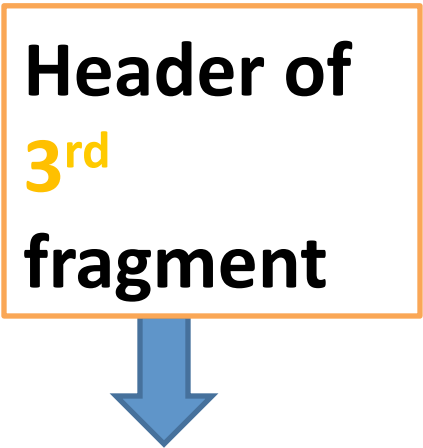
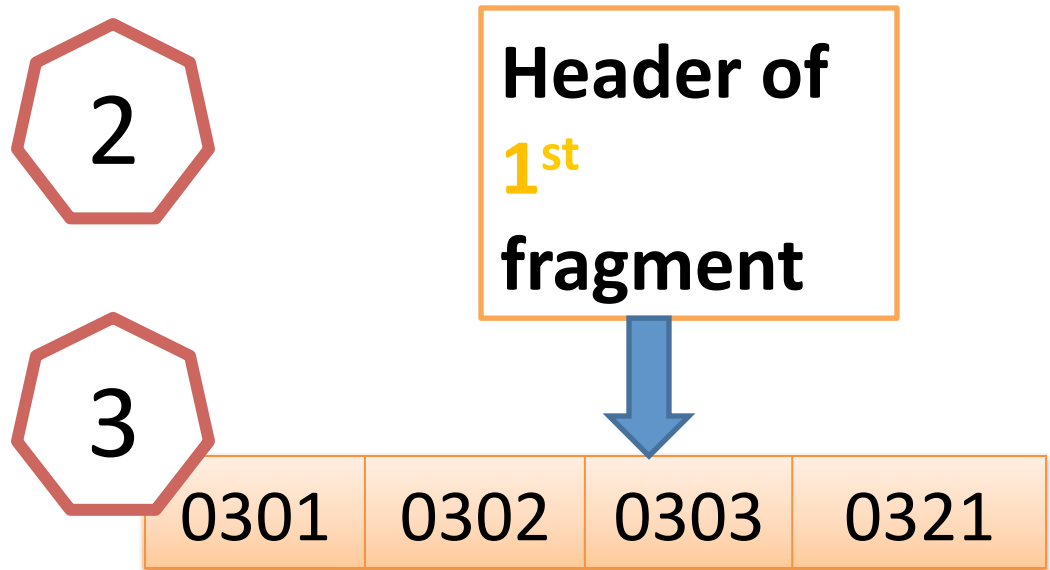
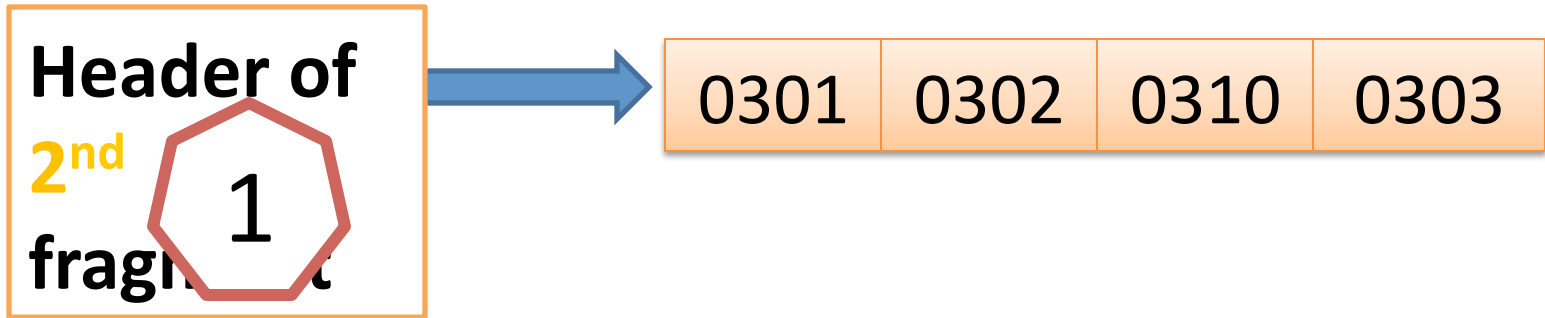


AGAT	AGAC	AGTA	AGTG
------	------	------	------

AGAT	AGTG	AGAT	AGAC
------	------	------	------



Header of
4th
fragment

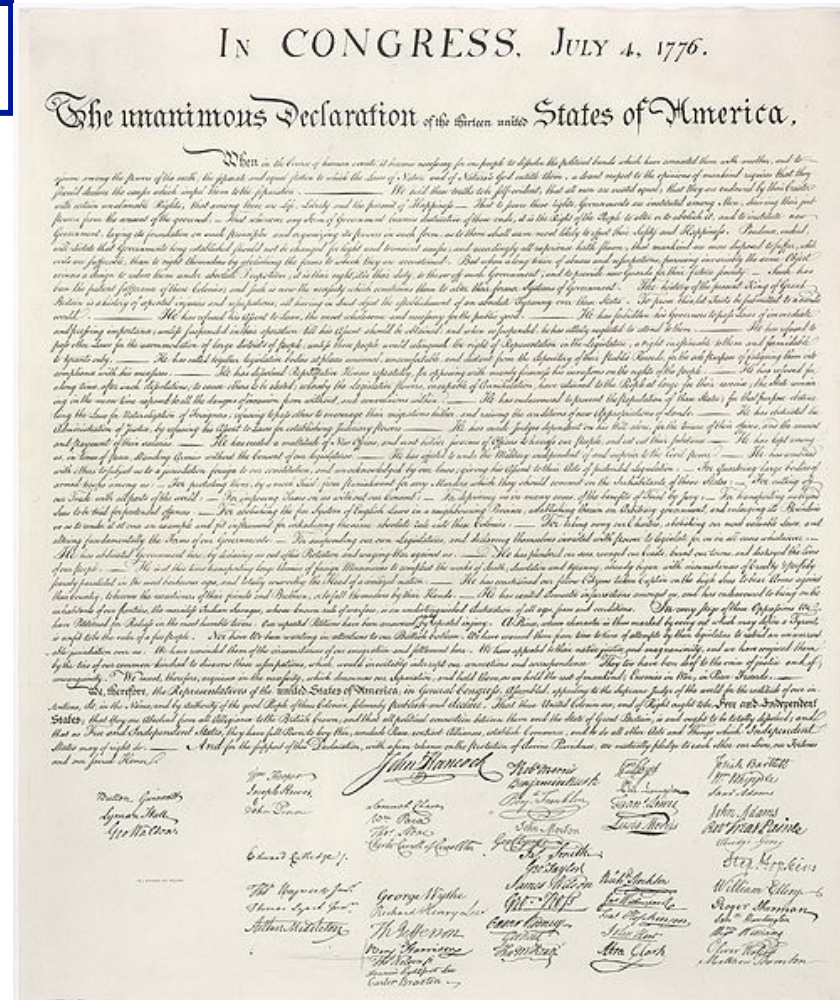


Capacity

If insertion size per cell is 1kb.....

The United States Declaration of Independence

Only
18 cells!!!



Capacity

1 gram of cells consists of~ **10 Million** cells.

The United States Declaration of
Independence requires **18cells.....**

Each fragment will have at least
500,000 copies!!!

Coding System

Encryption

Decryption

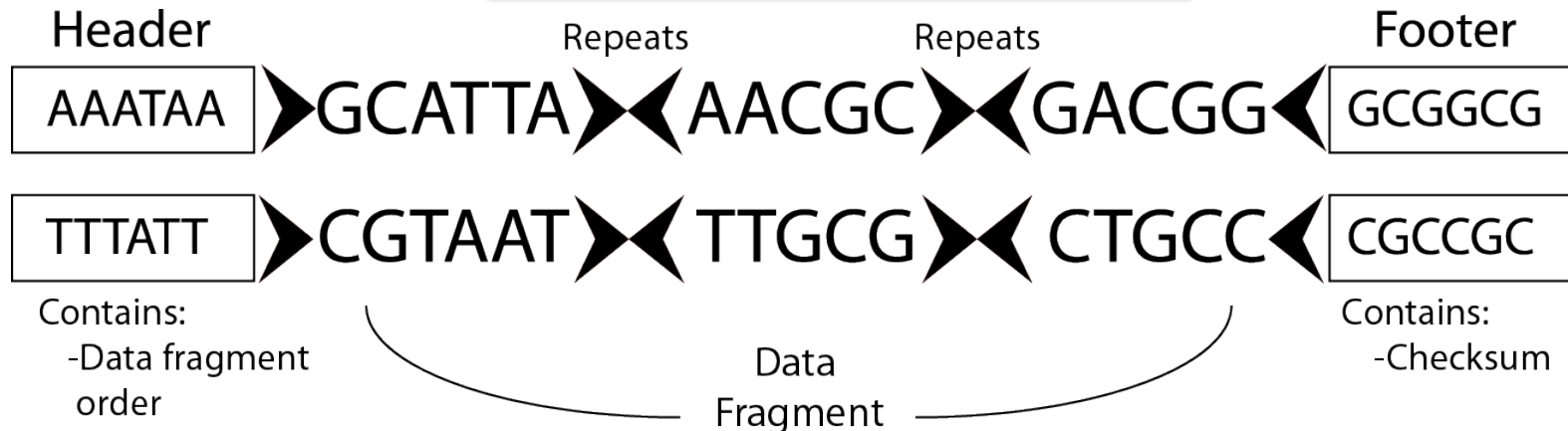
MP Storage System

Decryption

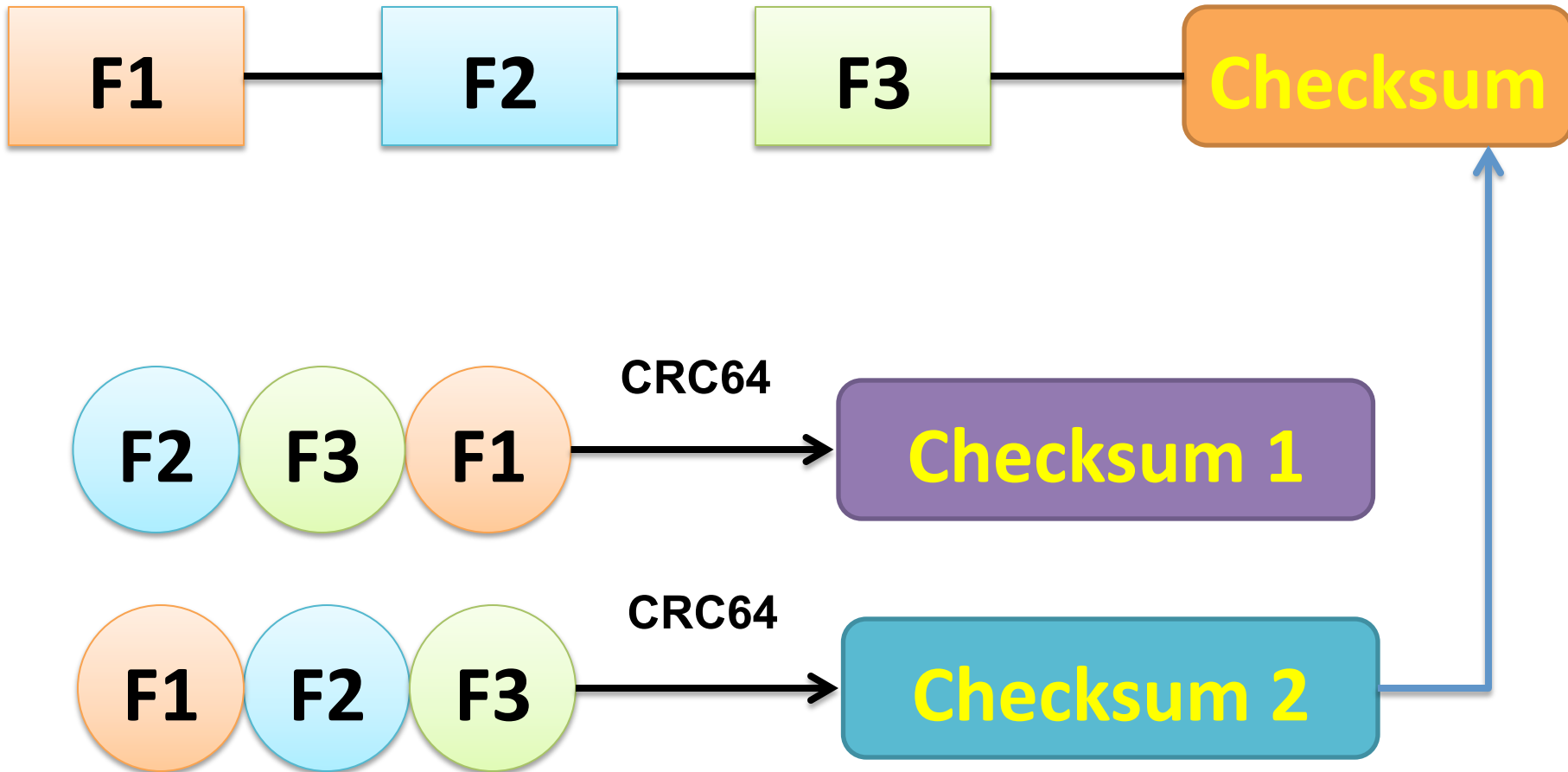
sequencing

Identification of repeat, message, checksum

Checksum system



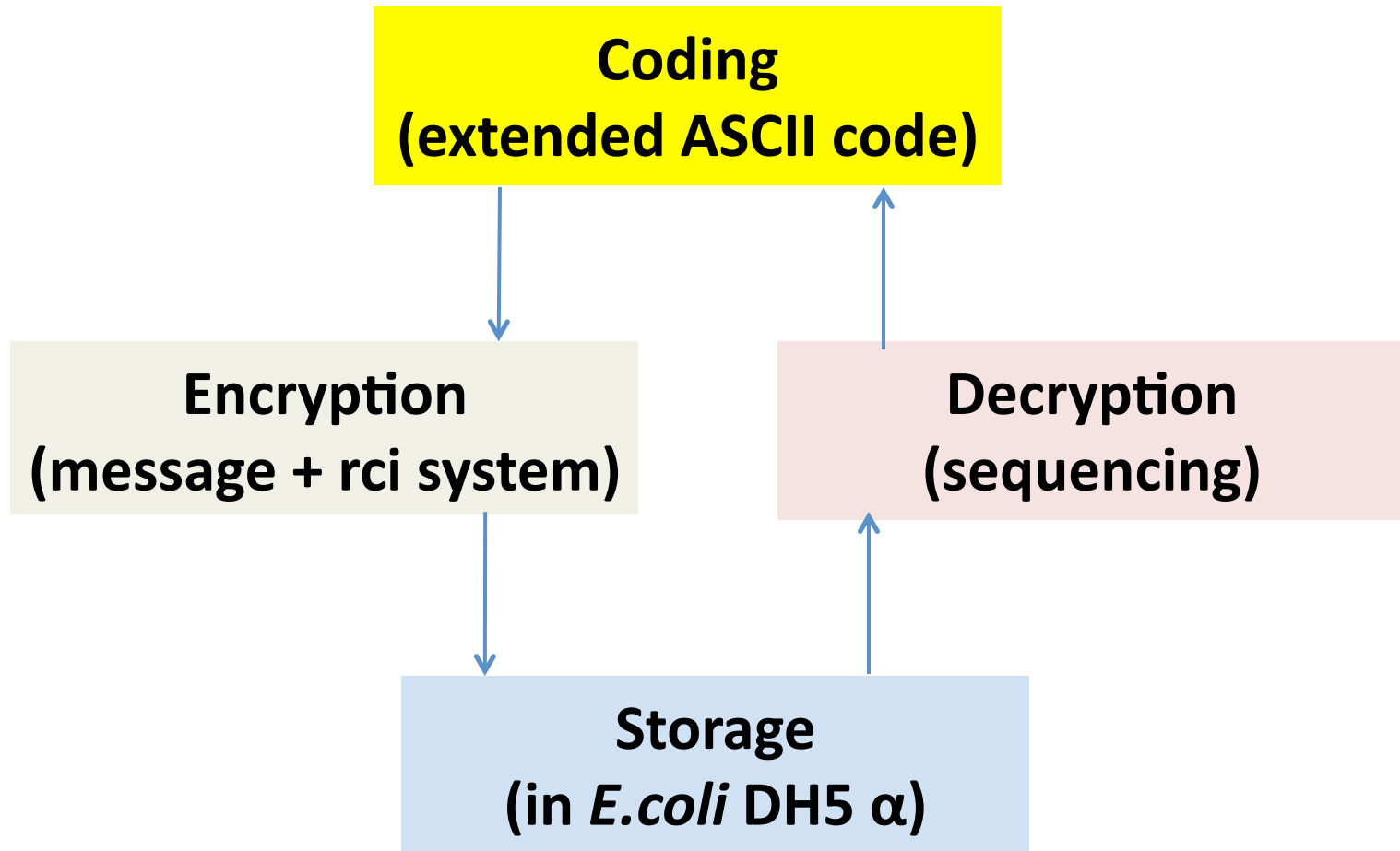
Checksum Mechanism



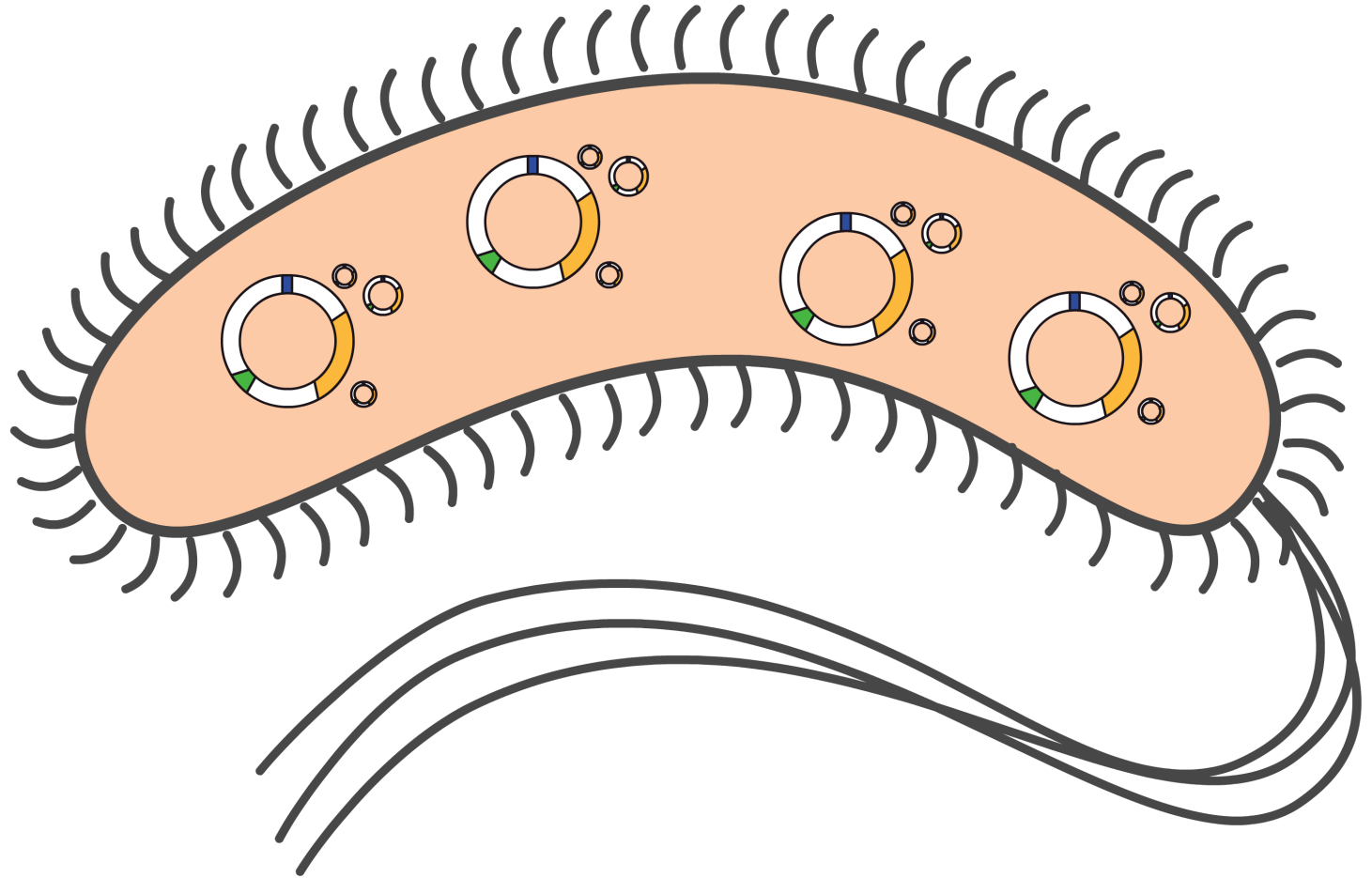


WET LAB TEAM

To prove our concept...



4. The host cell duplicates the data storage vectors, which helps to ensure data integrity by redundancy.



Message

- we must learn to live together as brothers or perish together as **tools**

<<code form Dr. Martin Luther King, Jr., a prominent leader in the African American civil rights movement >>

eg. “tools”

DNA encoding:

TGTATCGGTC
GGTCGATGAG
(20bp)

Enter your text here:

70 chars

we must learn to live together as brothers or perish together as tools

Our message (70 characters)

Quaternary Encoding:

280 chars

13131211020012311311130313100200123012112011302123202001310123302001230122113121211020013101233
121312111310122012111302020012011303020012021302123313101220121113021303020012331302020013001211
1302122113031220020013101233121312111310122012111302020012011303020013101233123312301303

DNA Encoding:

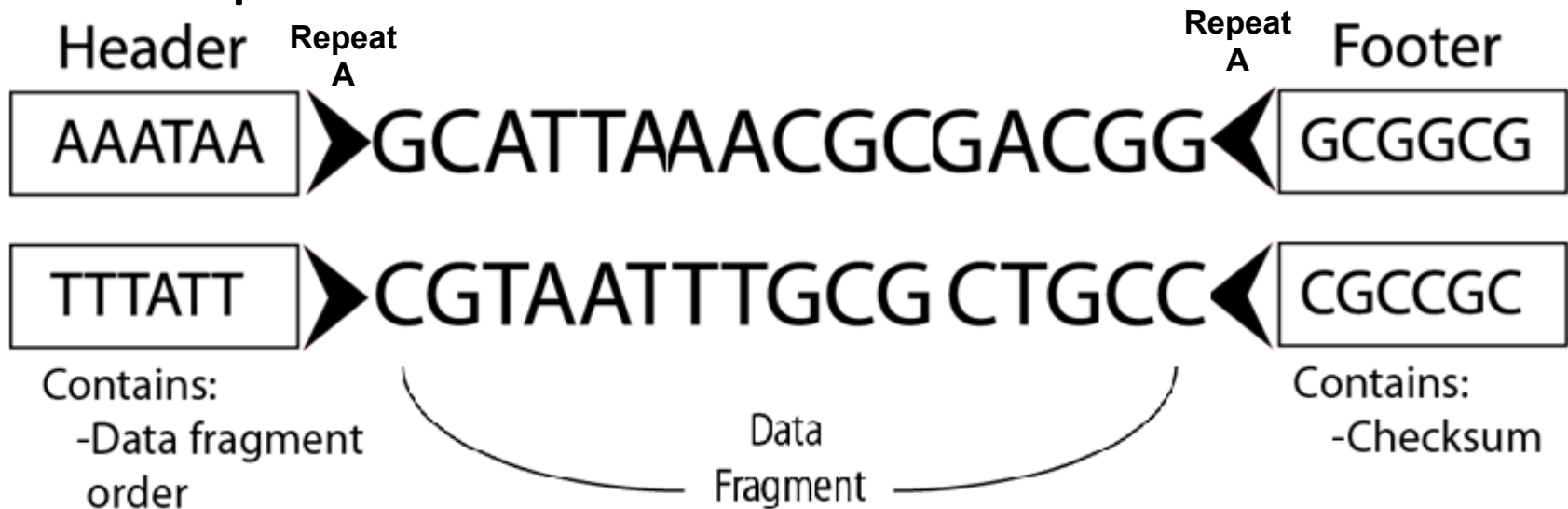
280 bp

TGTGTCTTACAATCGTTGTTGAGTGTAACAATCGATCTTTCATTGACTCGCACAATGTATCGGACAATCGATCCTTGTCTC
TTACAATGTATCGGTCTGTCTTTGTATCCATCTTTGACACAATCATTGAGACAATCACTGACTCGGTGTATCCATCTTTGAC
TGAGACAATCGGTGACACAATGAATCTTTGACTCCTTGAGTCCAACAATGTATCGGTCTGTCTTTGTATCCATCTTTGACA
CAATCATTGAGACAATGTATCGGTGCGTCGATGAG

DNA Encoding(280bp)₂₆

Structure of message

- Repeat A sequence in natural shufflon system has the highest inversion frequency
- 19bp



Parts designed



DNA
Message

Message gene template (438bp)

Synthesized DNA

K361000



Rci site-specific recombinase (1155bp)

Synthesized DNA (rci gene sequence of *E. coli* (strain: K-12))

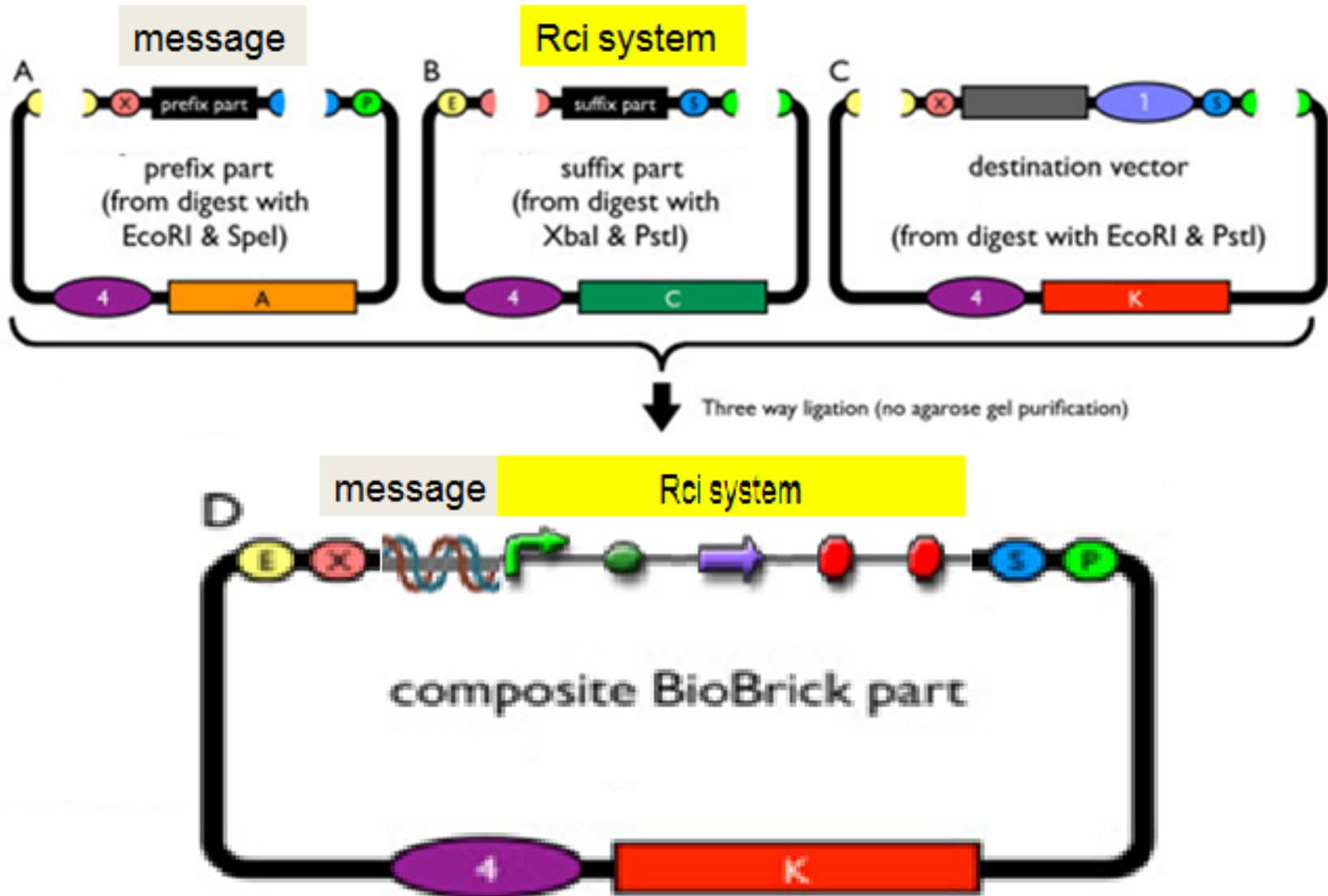
LacI
R0010 B0034 K361000 B0012 B0011



Rci system (1484bp)

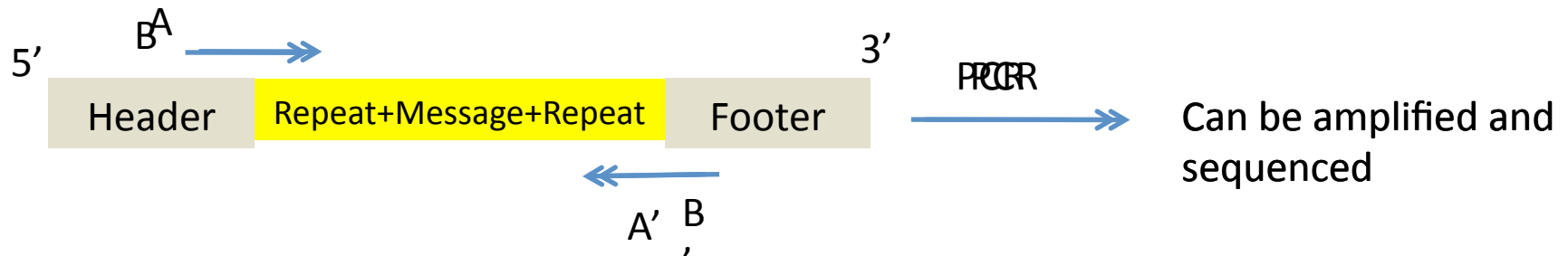
- lac promoter
- ribosome binding site
- rci gene
- double terminator

Integration of message to rci system

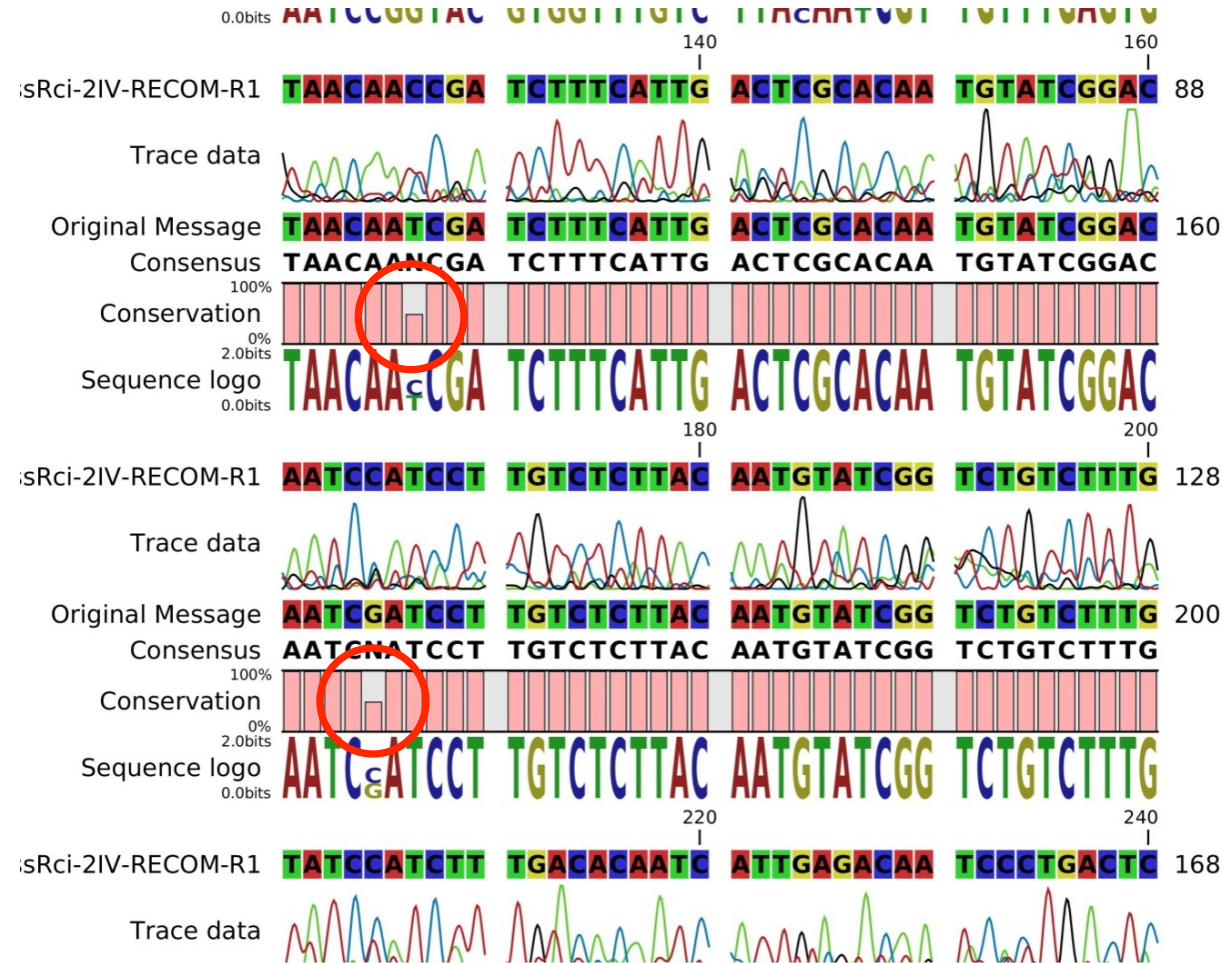


Expectation

- Repeat sequence + Message + Repeat sequence
- There should be two scenarios:
 1. Inversion of message
 2. No change of original message
- **Two** sets of primers are used



Results



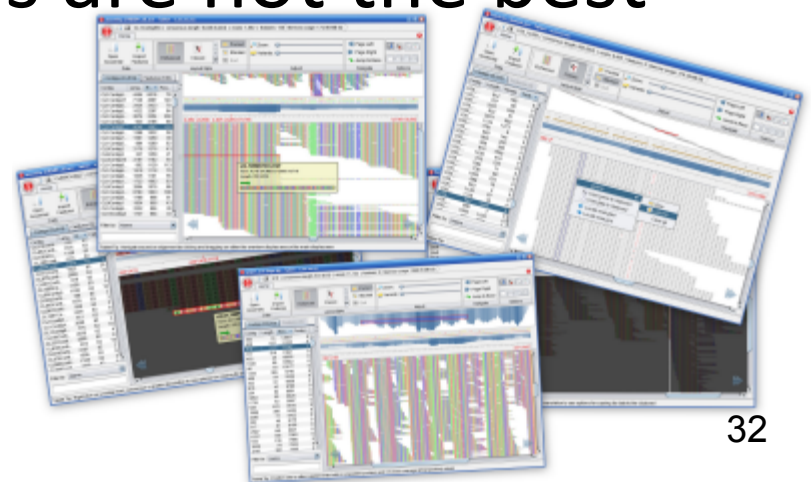
- Inverted and original message were found

- No loss of DNA

Checksum and high throughput sequencing!!!

High throughput sequencing

- **Massively parallel** sequencing process
- **Multiple copies** of sequencing products (reads) that can cover a particular message stored within the DNA
- Enable us to perform a **majority voting** on bases for which qualities are not the best

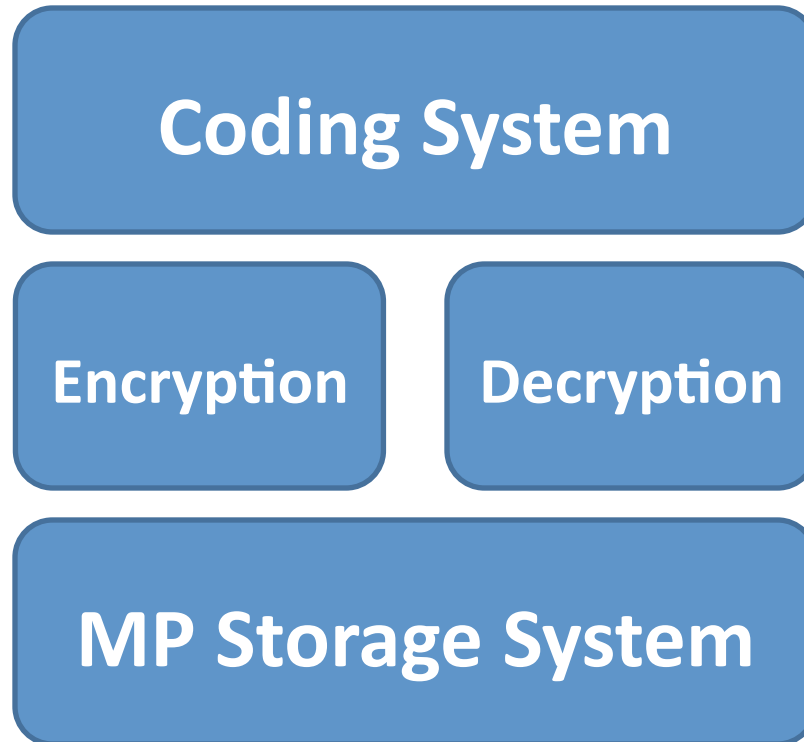




FUTURE PERSPECTIVES

To summarize...

- Infrastructure of our system



- Experimental proof

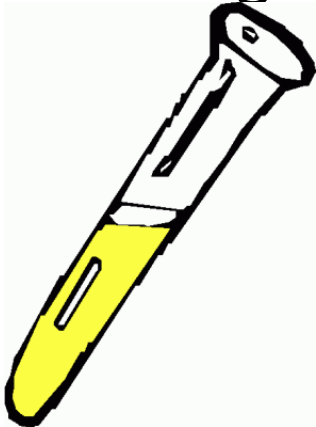
Bio-hard disk

	Storage
Hard disk	2000GB
1 gram E.coli	900,000GB



Therefore....

1 gram(wet weight) of E.coli



= 450

2 TB hard disk



Rapid & Specific access

- Parallel storage system

Insert *Header & Footer* in every message fragment



Design *specific probe* corresponding to Header



Pick up particular message from pool of data

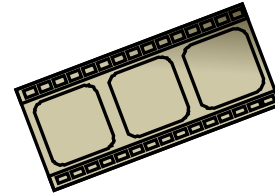
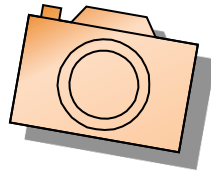
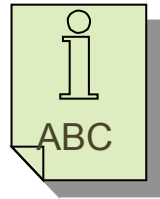


Targeted sequencing



Future Application

- Can store text, images, music, movies.....



- Insertion of barcodes *into synthetic organisms* as a part of current safety protocols to distinguish between synthetic & natural organisms
- Store additional information:

Copyrights

Safety protocols

Designers of the organisms

Acknowledgement



The Chinese University
of Hong Kong



Further Information

- Gyohda, A. & Komano, T. (2000). Purification and Characterization of the R64 Shufflon-Specific Recombinase. J. Bacteriol., 182 (10), 2787-92.
- Gyohda, A., Zhu, S., Furuya, N. & Komano, T. (2005). Asymmetry of Shufflon-specific Recombination Sites in Plasmid R65 Inhibits Recombination between Direct sfx Sequences. J. Biol. Chem., 281 (30), 20772-9.

**If you would like to know more about our project,
you are welcome to visit our Wiki page:**

http://2010.igem.org/Team:Hong_Kong-CUHK



Q&A

Inversion frequency

1. types of 19-bp repeat sequences
(repeat-a > repeat-d > repeat-b or repeat-c)
2. distance between repeat sequences
(distance increases, frequency increases)
3. DNA sequences surrounding the repeat sequences (symmetric repeat sequence increase frequency)

Inversion frequency

4. presence of HU protein

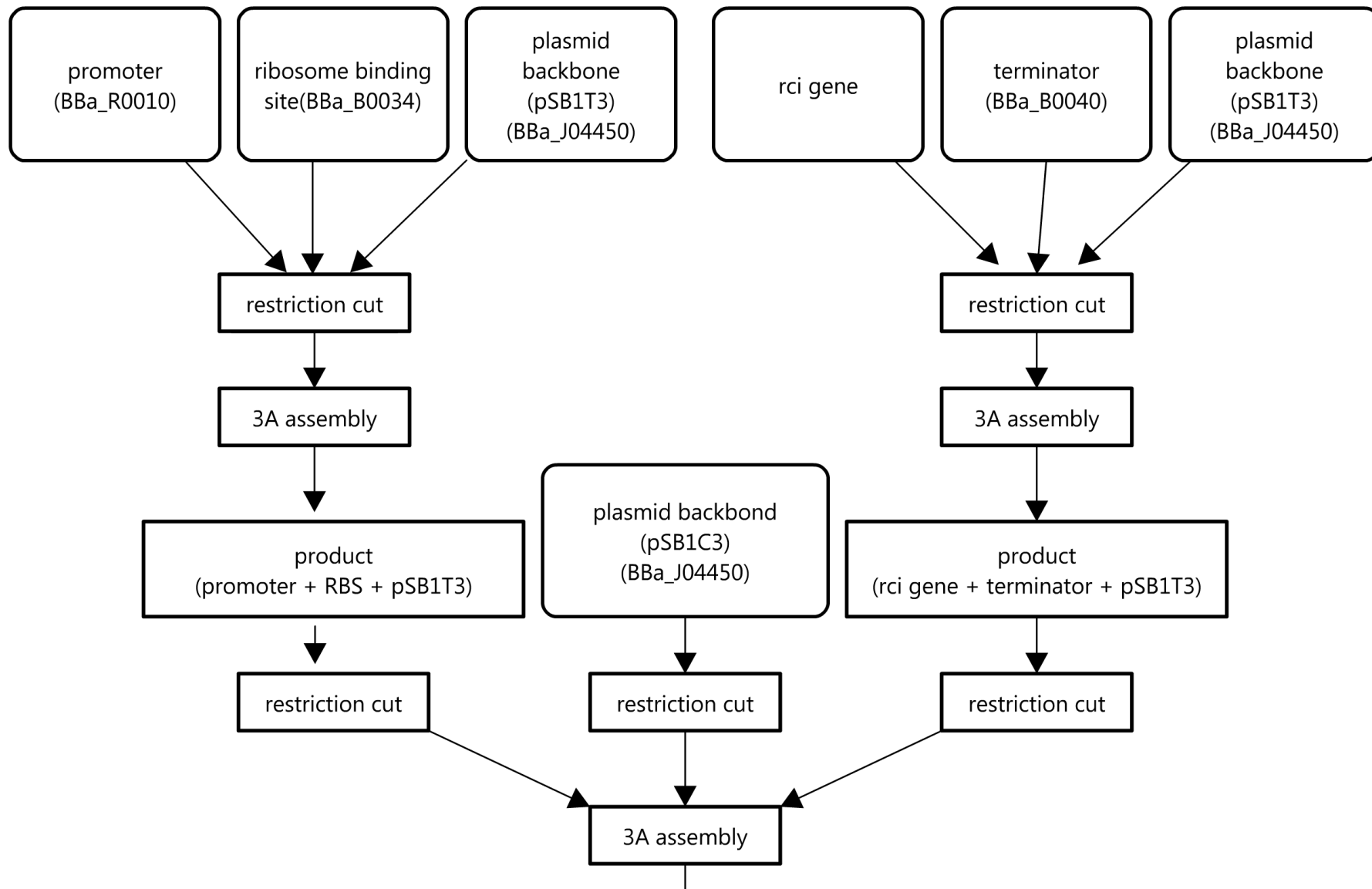
(binding of HU protein to DNA might facilitate assembly and/or stabilization of the Rci-DNA complex at the recombination sites, increases frequency)

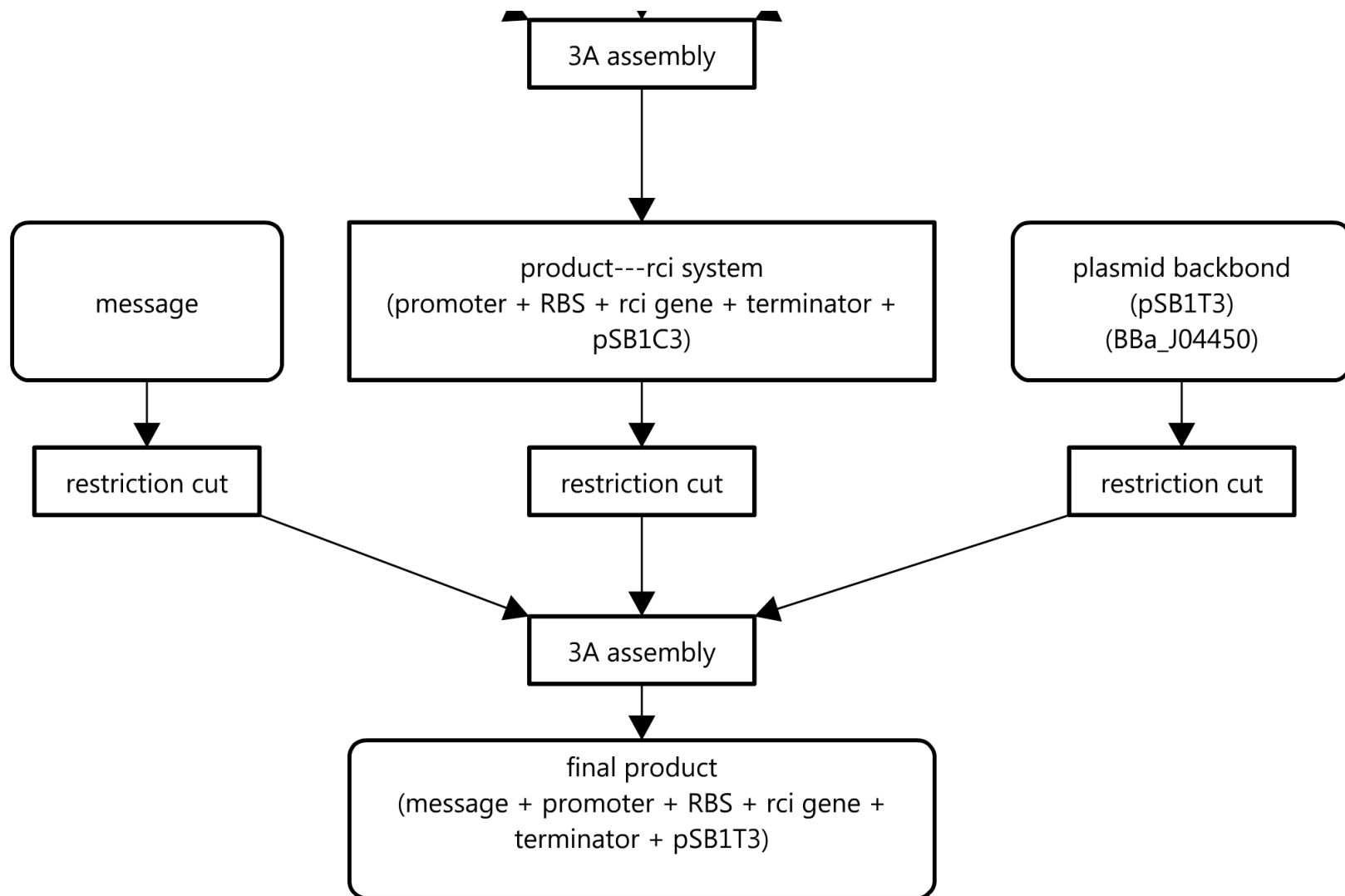
5. extent of DNA supercoiling

(Inhibition of DNA supercoiling → decrease Rci activity → decrease inversion frequency)

To avoid mutation

- Reduce reproductive cycle
- Provide favorable condition
- Move on to eukaryotes, make use of eukaryotes' proofreading system(more sophisticated DNA repair system)







PACIFIC
BIOSCIENCES™

Pacific Bioscience

- Real time
- Read Length : 1000 - 10000bp
- Single Molecule Sequencing
- 30 minutes sequencing process

