

Embedding DNA-Based Natural Language in Microbes for the Benefit of Future Researchers

Heqian Zhang¹, Jiaquan Huang¹, Xiaoyu Wang¹, Zhizeng Gao², Song Meng³, Hang Li⁴, Shanshan Zhou⁵, Shang Wang⁶, Shan Wang⁷, Xunyou Yan⁸, Xinwei Yang⁹, Xiaoluo Huang^{10*}, Zhiwei Qin^{1*}

¹Center for Biological Science and Technology, Advanced Institute of Natural Sciences, Beijing Normal University, Zhuhai, Guangdong, 519087, China. ²School of Marine Sciences, Sun Yat-sen University, Zhuhai, Guangdong, 519000, China. Southern Laboratory of Ocean Science and Engineering, Zhuhai, Guangdong, 519000, China. ³State Key Laboratory of Drug Research & Natural Products Research Center, Shanghai Institute of Materia Medica, Chinese Academy of Sciences, Shanghai, 201203, China. University of Chinese Academy of Sciences, Beijing, 100049, China. Zhongshan Institute for Drug Discovery, Shanghai Institute of Materia Medica, Chinese Academy of Sciences, Zhongshan, Guangdong, 528400, China. ⁴School of Pharmaceutical Sciences, Sun Yat-sen University, Guangzhou, Guangdong, 510006, China. ⁵State Key Laboratory of Microbial Resources, Institute of Microbiology, Chinese Academy of Sciences, Beijing, 100101, China. ⁶College of Plant Science, Jilin University, Changchun, Jilin, 130062, China. ⁷State Key Laboratory of Microbial Technology, Shandong University, Qingdao, Shandong, 266237, China. ⁸College of Life Science, Langfang Normal University, Langfang, Hebei, 065000, China. ⁹School of Pharmaceutical Sciences, Zhengzhou University, Zhengzhou, Henan, 450001, China. ¹⁰Shenzhen Key Laboratory of Synthetic Genomics, Guangdong Provincial Key Laboratory of Synthetic Genomics, Key Laboratory of Quantitative Synthetic Biology, Shenzhen Institute of Synthetic Biology, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, Guangdong, 518055, China.

Correspondence: Zhiwei Qin, Xiaoluo Huang
E-mail: z.qin@bnu.edu.cn; huangxl@siat.ac.cn

Material and methods

Encoding the text data in a single DNA strand

The algorithm is accessible on the Storage-D platform (<http://storage.dailab.xyz:16666/>). During the encoding process, two 900 base pair (bp) fragments were generated. To ensure synthesis feasibility and sequence stability, homopolymer stretches were limited to fewer than six nucleotides, and the GC content was targeted to range between 40% and 60%. Each 900 bp fragment was augmented with a 2-byte Reed-Solomon (RS) error correction code, following the platform's protocols. These fragments were then ligated to form a contiguous 1.8 kilobase (kb) sequence. To facilitate specific amplification, a set of 20-base pair random access primers was appended to the concatenated sequence, designed using the primer design tool integrated into the Storage-D platform.

General experiment for bacterial transformation

The DNA sequence was synthesized by GENEWIZ (Suzhou, China) and subsequently cloned into the *Streptomyces* integrative plasmid pIB139 via homologous recombination at the NdeI and EcoRI restriction sites. This plasmid was then transformed into *E. coli* ET12567/pUZ8002 and cultured on Luria-Bertani (LB) agar plates containing apramycin (25 µg/mL), kanamycin (50 µg/mL), and chloramphenicol (25 µg/mL) at 37 °C overnight. Colonies were isolated and verified through sequencing.

The PCR reaction was set up in a 50 µL system containing 25 µL 2× Phanta Max Buffer, 1 µL dNTP, 1 µL DNA Polymerase, 1 µL each of forward and reverse primers, 1 µL template DNA, and 20 µL ddH₂O. The PCR conditions were: initial denaturation at 95 °C for 3 minutes; 35 cycles of 98 °C for 10 seconds, 60 °C for 30 seconds, and 72 °C for 2 minutes; followed by a final extension at 72 °C for 7 minutes.

The confirmed colony was grown in fresh LB liquid medium with antibiotics under constant shaking (200 rpm) until the culture reached the exponential phase (OD₆₀₀ ~0.8). The culture broth was then harvested by centrifugation at 5,000 rpm for 10 minutes and washed three times with LB medium without antibiotics.

To obtain *Streptomyces* spores, strain G7 was cultured on MS solid medium (20 g mannitol, 20 g soya flour, 20 g agar, and 1 L tap water) at 30 °C for 15 days. Spores were collected by gently washing the surface of the strain with 20% glycerol using a sterile cotton bud. The spores were washed three times with tryptic soy broth (TSB)

medium at 12,000 rpm for 5 minutes to remove glycerol, then resuspended in TSB and incubated at 50 °C for 10 minutes before cooling to room temperature.

Equal volumes of *E. coli* and G7 spores were mixed, centrifuged at 12,000 rpm for 5 minutes to remove the supernatant, resuspended in 200 µL of TSB, and plated onto MS agar supplemented with 10 mM Mg²⁺. After incubation at 30 °C for 20 hours, the plates were overlaid with apramycin (25 µg/mL) and nalidixic acid (25 µg/mL) and cultured for 5 days at 30 °C. Colonies were isolated and purified on MS agar with antibiotics for three cycles over 12 days at 30 °C. Finally, the colonies were cultured on MS agar without antibiotics for 5 days at 30 °C and identified by Sanger sequencing.

The information of natural language and its encoded DNA sequence

This is *Streptomyces ginsengnesis* G7 that was discovered in 2020 in Jilin Province China and demonstrated the ability to produce natural products such as lydicamycins which inhibit weed growth, and other structurally unknown compounds which inhibit the growth of fungal pathogens in crops.

```
GTGAACGCATGGTCACCTACCGTATCAATGCCAAATTTCCAAGCGGCGAAGACCCGGATTGGCCAGA
CCGCGAGGGGTAAGACACGCCCTCGCCCGAAGACTCAAAAACCCAGCACACGAGAATTCCAAGAC
TCCAAACTACGAGACTGCGAACTGCCGAACCGTAAAACACACGAAGAGACGGACAGGCCAGCCTAC
GAGACGGCGAGATAACCAAAAAGTCAATACAAAAAATTACCAGCATAAGAGATAACCAATACTCGA
ACGCCCCGAGACGGCGAAAAACCGACCCCGCCCTCAAGAAGAAATACGACCAGCCCAGCGTGCCAAA
TGCCCAATTCGAAAAAACCCGACAGGCCAGAAAGAAAACACGCGAAGACCCGACTTGCCAAACCC
ACCAAGCGGCGAAAAACCGAACAGTCGAAAGATCGAGGGATCCAACACTCGAAGGGTCCAATATCC
GAGCACTCAGACACGCCAGACGCCGAAGCTACCGATTGGCCGAAAACCCGAGGCGCAGAAACCCCG
AGGGGCCGATATTCCAAGCCGGAAGCATCCGACATTCCTCGCGTCCAGACCGAAAAATCACCAGAC
GCCGAAACATCCGACCGGCCTCCGCGCAAAAACCCCAAACAGCGAGAAAACCAAAGGCCGAGACAG
CCAGCCCCCGAATTCAAAGATACGCATCCACCCGAAATGGCCAAAGTTCGAAACTACAAACACGCGG
ATAGCCCAGCCGCGAAGCTACCAACCTGCCACCATTCAAACACACCAGCACTCGAGACCACGAGAT
AACGAGACATCCAAAACCGAGCGCACCAGCGGTCCAGCATTGCAAAATGCCAATACACGAATACA
CCACTTGCAAAAACACCAGACGGCCAATATCCGGATTGCCGAGAATACAAACAGTTAGTAGGGGA
TAAGGTAGTCAGGATGGATGGTAGCCGTTAAACCAGCGCCTTTGGGTATTGAATTTTAGGACTCCT
GGTAGGATAGAGGGAAGGGTACATGCGGGCAAGGAACGAACCCCATATGTTTAAATGGTCTTATGG
TATTTGACCCCTCTGGTATTAGTTACACCAGGGATAATGGGTACCTGGGATCCAGGTAATTTGCCGTC
AGAGTATAAAGTTTTAGTCGATCCAGGATCCCTGGTCTCCTGCGTTGAGTCCAACATGCAAGTTAG
GACCCGGAATGAAATGTCTAAAAGTTTCTAAGGTATTAACCCAACGGAGTGATCAGTATAACCACGA
ACGGACCTACGTGGAGCGATGTGGACCAGCCAGTCAGCCCTGATGGTGGCAGTCCCACGTGGAAAC
ATGCGCATCAGGATTCCTGGTCACGAGGTTGTTTGCGTAGATATATTTAGGGTATATGGAGTAAAC
TAAACCAGATAATCAGATTACGCCTAAACCAGTACCGAAGAAGTCAAGGTTGTTTGGACTCAAGAAC
```

CTCACGGTGCGACCGCACGAACCATCAATCGCATACTCCCACAATTAATCCATCGCTCGCACCAACCA
TCCCACGCTCGAACCATCCAATACTCGCACGATCGATCCCTCCAACAAATACTCGCACGCACCCTCCA
TCCATCCAACACTCCCATACACCCACGCTTACTTAATCGATCGCTCCCTTAATCCATTAATCACACACA
CCCACAATCGCACCCCTCGCACACTCGATCCCTCGCACAATCCAACACTCGCTCAATCGCTCGATCGCA
CACTCCAATAATCCCTCACTCCCTCAATCAAACCCACCAACAATCGCACCATCCCTCGCTTAATAATC
GCTCGATGAATCTCATTCCGTCTCCCCGTCTAAC

NOTE. Sequences highlighted in green are designated for PCR amplification.