

BioBit 暑期培训

Handwritten signature or initials, possibly "gib", written in white ink on a dark background.

Handwritten text, possibly "BIB", written in white ink on a dark background.

230822



0822 钱斌 DNA数字信息存储 北大CAB

冻干粉保存

00 → C

3.4nm

01 → G

10.5 ↑ base

10 → A

21 bit

11 → T

"空间体积小于是基"

dNTP ~ 330g/mol

$$\frac{1g}{330g/mol} \times \frac{1}{2} \times NA = 0.91 \times 10^{21} \text{ base}$$

← ×2

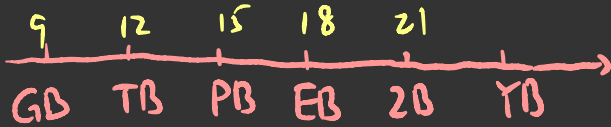
1.8 zeta bits

÷8

225 ExaBytes

1g 是很大的量 DNA通常是ng级

3 K MByte



2023 126.33 ZB whole world

30.02 ZB in China

25 188.03

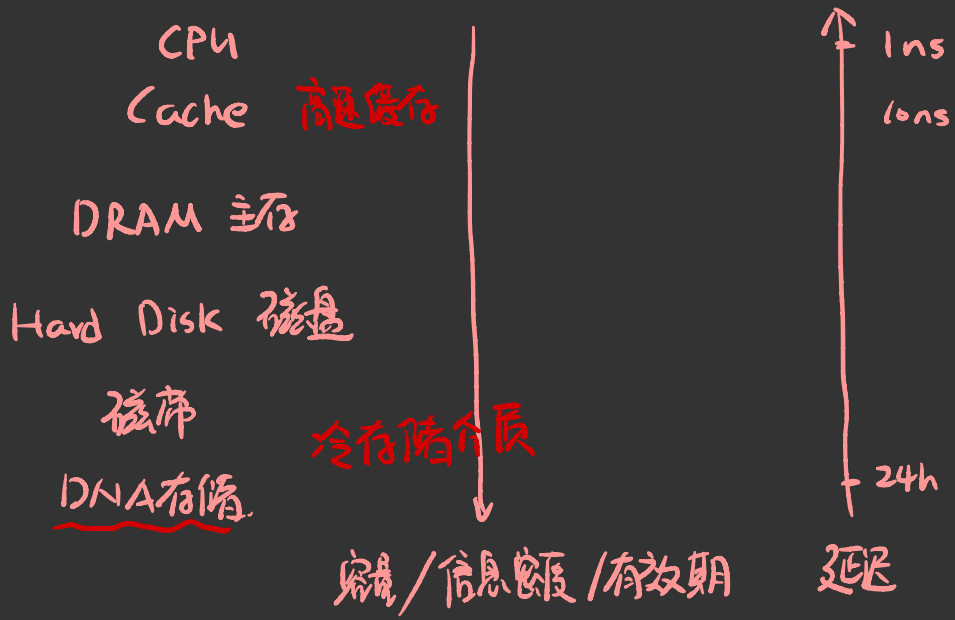
指数增长

27 284.30

硅储量不足

耗电量 3%-13% (2023)

✓ 脱氧核苷酸平均分子量是葡萄糖的3倍



✓ 阿里云的机房
在杭州附近的一个
湖中，从而确保冷却

每5年倒换一次，防止介质损坏

E3级存储中心：~3.8万个机柜 1.68亿度电/年。 5年高 磁干扰 (安全问题)

1988：首个DNA存储 // 2012：Chrunch Oiglo DNA存储
cell population 存一个5帧动画。 霍夫曼编码

2013：Huffman 存储

近几年: 存储介质、实际应用

2021年6月: «DNA存储白皮书»
微软等企业发布

Q&A: DNA存储能否落地?

2019年3月 1w\$ ATCG合成模块

8.4 hr

✓ 哈佛 MIT 出了一套完整
DNA存储装置.

存5个字母 hello

存储

共10个小时

测序模块 (纳米孔测序) 36 min

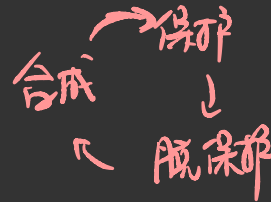
目前 DNA storage: 200 MB 存储量, 10B/s - 3kB/s 写 28\$/kB

1GB/s 读 \$0.4-4/GB

合成、测序技术能否推动! → 液相法技术能否落地

Q&A: 技术难题是什么?

DNA合成技术: 化学循环合成



Microsoft: 柱式合成 酶合成技术

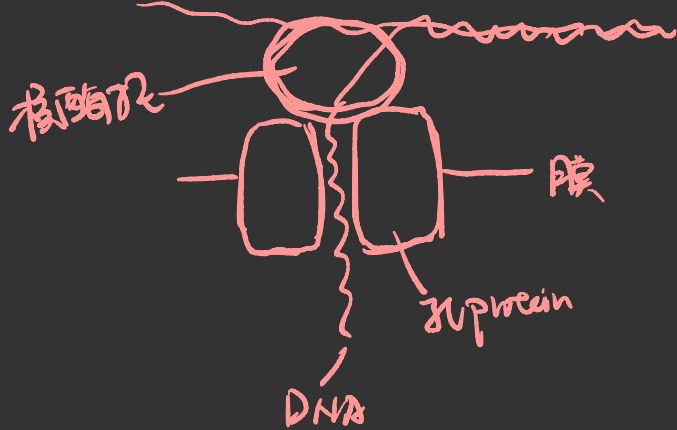
阵列合成技术

阵列上反应发生的密度决定
这项技术的成本

- 光刻法: 用光激活合成活性 光学分辨率 亚微米级
- 热敏法
- 电化学法: 微电极控制 \rightsquigarrow 最密反应腔 $25\text{M}/\text{cm}^2$
- 喷墨打印: 喷出极少量碱基合成 \rightsquigarrow 主流技术 \checkmark

测序

- 一代 Sanger
 - 二代 illumina 非 de-novo 合成 边识别边合 短片段
 - 三代 单分子 边识别边合
 - PacBio 查 长片段测序 贵很多
 - 纳米孔 Roche
- 实时! 边识别边合 \rightsquigarrow 不需组装. 序列相关性保留



错误率较二代高10-100倍

内外侧电压不同

查 液相复制 + 荧光标记



DNA 拼接

活字印刷

Catalog 公司 16G13

wikipedia 数据存档

武汉病毒所

编码: 信道编码器 → 调制写入单元 ^{合成、组装}

纠错码 + 冗余码

事物不均 - mismatch / indel

噪声 → 信道
衰减 降解 稳定性

DNA 断裂

* 物理冗余

↳ 在截取时各有三份

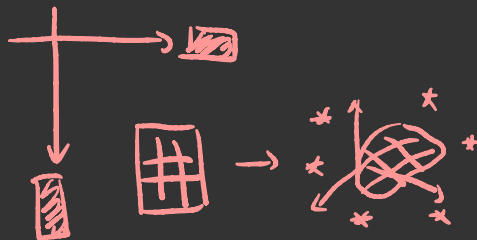


* 逻辑冗余 Reed Solomon Code (2015, Grass)
冗余位.

↓
测序解码

查

分组、每组加3位 RS 码



喷泉码 DNA fountain (2017) Dropout 错误!

三个数据包 \rightarrow XOR \rightarrow 丢失后仍可恢复
(只需要一定量数据包)
冗余度最低

成功率 0.0001% 效率当年最高. Best one!

图像压缩 (2022) 等.

关于 Indel : 纳米孔测序易发 Indel



进孔速率不稳定 \checkmark 电流变化产生 Indel

2023年钱老师文章.

长片段纠错编码体系

GC含量. 避免 poly A/C ...

避免生物活性. 避免酶结合. 避免 seq blacklist

Random access

数据随机访问

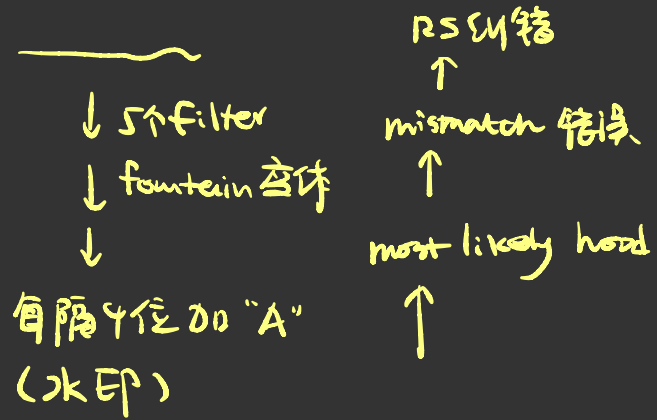
引物 → 索引序列 → PCR

- ✓ 磁盘上有物理位置 (移动磁头)
- ✓ 溶液中需要随机分子碰撞

100μL 中至少 10 个 copy 的分子

- ✓ 2021 年 将 DNA 包在硅质小球中 (SiO₂ 膜?)
↳ 从而物理上精准提取

- ✓ Gene 纳米颗粒 → 文件系统 可的基因组文件逻辑结构
(钱老师工作)



将DNA束缚在表面时，人为提高 ΔE ，测序精度 \uparrow

文件锁定、册序、重命名 也陆续发布。

lgDNA 很难达到 EB 级存储

(可能体积信息密度更重要！)

物理冗余 拷贝数、分子数
逻辑冗余

✓ 隐私便携性 (封装)

DNA SiO_2 球 \rightarrow 3D 打印热塑材料 \rightarrow 兔子

(2020, Nat Biotech) 塑料兔子。

眼镜片 \rightarrow

溶液中不能长期存。会断！

✓ 细菌基因组存储

嗜盐菌中的整合酶.

(钱老师工作)

↳ 插入基因组.

天大 酵母人工染色体.

保留时间不久



2000代 三个月

10^{-10} 突变率 保真度高

(2013, Adv Sci)

枯草芽孢杆菌 \rightarrow 芽孢状态封装 DNA

Olgo Pool

体内

max 2-3 PB/mL

max 1TB/mL

50 kb info DNA

✓ 分子加密 or 销毁

↓
DNA 信息存储

↳ 依赖于 DNA 索引

(2019 奖 Nat comm)

折纸不做信息存储

转化 \rightarrow 多层信息存储 95°C 加热破坏索引.

✓ 痕迹追踪 (DNA 稳定不易坏)

DNA barcode 标签! (2020 Science) → 芽孢喷在表面.

(2020 Nat Comm) → 将DNA喷在表面. 使用电信号差异大的短序列

✓ 动态写入 & 擦除

2021, 电信号 ⇒ 数字信号

(Nat Chem Biol) → 诱导 CRISPR 系统

2021. 双颗粒信号可视化

做扁码. 虽然短序列但可以用 Nanopore 测序.

0822 陈昊荣

北大微电子

普渡 DNA Origami

MIT 博士后 合成生物学



"Cello"

MIT 导师的工作

硬件描述语言

↓ Parsing

↓ Logic

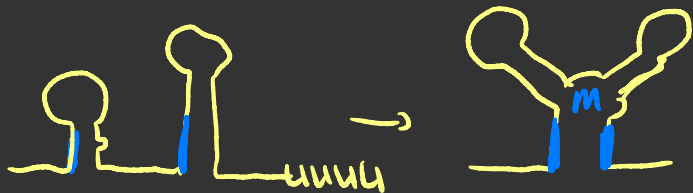
↓ Gate assignment

↓ DNA

DH10β 中 表达

How to make general?

RNA world!



转录

Riboswitch

天然存在



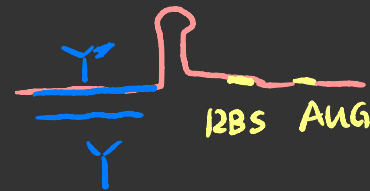
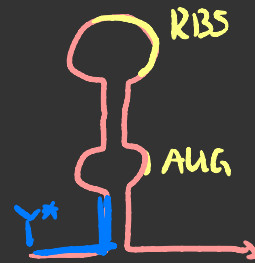
用配体的结合力
翻转构象

RNA 可以感知小分子!

~ 40nt 的筛选库
Library pool

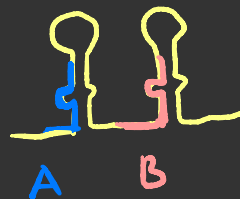
SELEX

筛 → ✓



哈佛

成体系的核酸
元件



OR Gate



AND Gate



一种化学本质是DNA的
 RNA酶，可以定向地切
 割RNA分子

DNAzyme
 Mg^{2+} 即可切

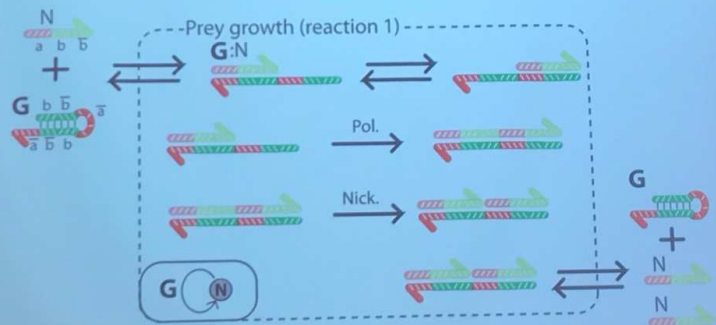
不切割其它分子

NOT Gate
A AND (NOT B) ✓

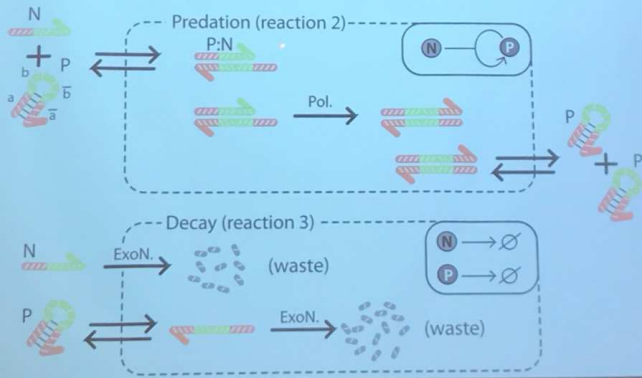
可振荡!

In vitro computing

in vitro replication



Construct predator-prey dynamic system



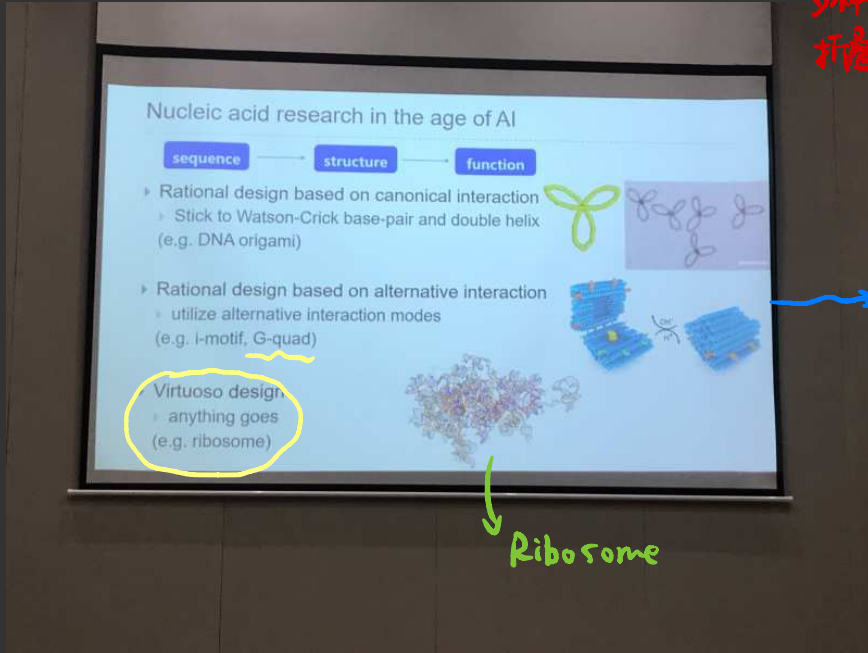
Sequence → structure → function

多核酸
打叠技术

DNA brick

DNA origami

ss-origami 单链DNA



有更复杂

M-fold ✓

0.3 off-ratio

Ribosome

0822 马大程 川大 → 清华 → MIT → 莱斯 → 之江
博士后

CRISPR: 脱靶、靶向范围窄 // AAVS: 容量小

↓
PAM 限制

Cas9 同源 protein - 32个

PAM区不保守

⇓
嵌合 Cas9 ← 互作短胚移植
(SaCas9)

与DNA类似的?

⇒
NLP 预训练范式:

预训练阶段:

微调阶段

雷帕霉素: 免疫抑制剂
(不可临床)

PROTAC 改造!

microRNA 与人工启动子

内分泌 多输入“与门”

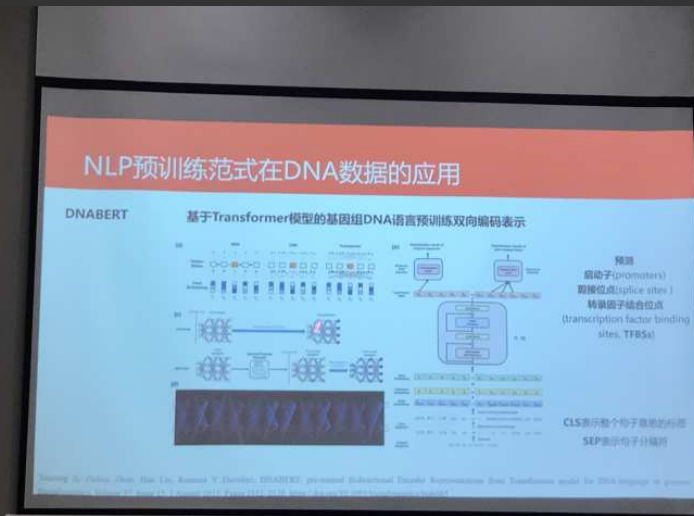
(2016 Not comm) microRNA toggle switch

看看合作的点在哪?

↑

当时漏掉信息

现在看似乎就是空
后序的工作。



二维分子结构 \rightarrow Graph? 空间关系如何取用?

CGT 细胞基因治疗

230822 傅帅

BT+IT 融合时代

Biological Computer ✓
 \rightarrow 碳基替代硅基

Computational Biology?
 \rightarrow AlphaFold2

Bioinformatics?
 \rightarrow 数据获取、加工、存储、分析

Biological computation 信息处理能力

Qian: 如何从模型中提取知识?
钱鹏昊提问

ation x

↓
没有明确回答

《星际穿越》

《火星救援》

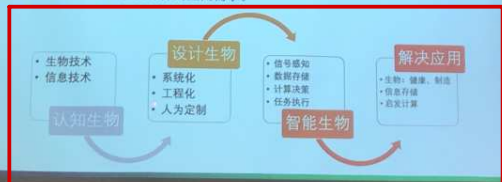
DBTL 循环

Design → Build → Test → Learn

► 什么是生物计算? / What is Biological Computation

生物计算:

是指通过人工智能和生物技术破译核酸序列中的遗传语言规律, 通过系统化、工程化方法, 重新设计和构造生物分子、细胞和生物体, 使其具备人为定制、超自然状态的逻辑运算能力, 实现智能化信号感知、数据存储、计算决策与任务执行能力, 用于解决生命、信息存储与启发计算的应用需求。



明确的研究范式

泰州国家园区

22年来传化

生物技术产业新势力



大力支持 + 强监督

“挤压销售泡沫”

反腐

合成生物学

工程化流程

技术分类

60%-70%

产业驱动力

新材料 & 先进疗法 ATMP

AI崛起

AI+医疗 商业模式

医疗器械类

AI制药企业 “共85家”

非医疗器械类

(北京 深圳 居多)

(截至23年7月)

医药产业关键技术

传统化学

产品改性

纺织 纸 塑料

.....

竹林式发展

纺织助剂 打破垄断

降本20%

无氟防水 冲锋衣

智能工厂

柔性供应链

05年入股新安化工

S: P CI

降毛利 双循环

物流

生产基地

公路港

"货运网"

信息平台

交易平台

服务平台

奥的斯电梯

产业集群

600+亿营收

农业

pǔ yáng

石斛 兰花

玉米 大豆

种曲

科技城

一城一产 全球科学家

生活 工作

创新药 基因治疗 合成生物学

400w方 - 200w方 产业
200w方

浅湾智谷

初只是一小片
园区名

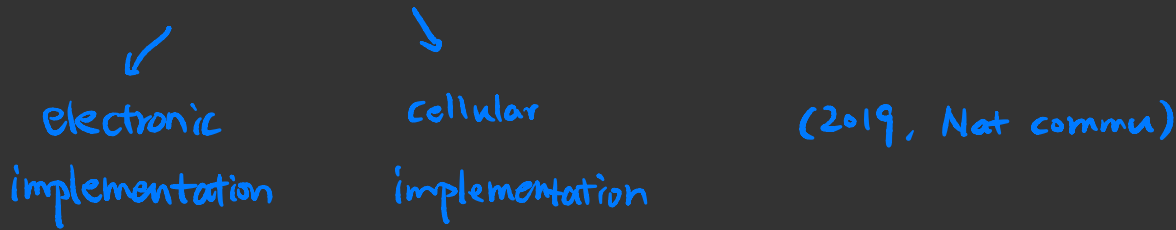
500强

1500多人

0823 Angel Cellular computing with contextual dependencies
towards environmental applications

@ Angel G Moreno angel.goni@upm.es

Turing machine $\left\{ \begin{array}{l} \text{符号位 (操作符)} \\ \text{输入内容} \end{array} \right\} \Rightarrow \text{输出.}$



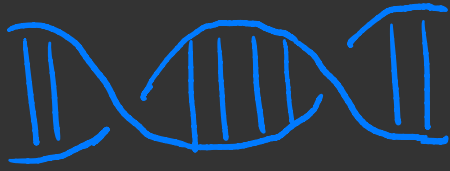
"Combinatorial logic"

"Boolean logic gate"



(Wang, 2011, Nat Commu)

10.1038/ncomms1516 & doi



(Nielsen A 2016, Science)

3 input \Rightarrow 1 output

"Contextual dependencies"

My python environment is so messy.

Predictions are stopped by non-linearities

(Huseyin, 2021, Nat comm)

Different Plasmid, Different Host



"different result"

Parts 的描述是因环境而异的.

Predicting compatibility

Pseudomonas putida 假单胞菌

Understanding



Using

(Nikel, P, 2015, mBio)

(Garcia, 2017, eLife) → How to use noise?

2023. pBLAM 基因组插入工具

(Gorfi-Moreno, 2017, ACS synthetic biology) 6(7). 1359-1369

0824 王宝俊 大老师!

石油基原料 → 合成化学工业 → 产品 (20世纪)

生物基原料 → 合成生物学 → 产品 (21世纪)

DNA → 生物元件 → ^{模块}生物器件 → ^{系统}基因线路 → 细胞、组织

合成生物学 ⇒ 工程生物学

模块化、标准化、层次化

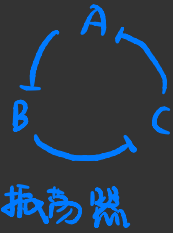
DIBTL循环

第三次生物技术革命

BioTech Giants = Microsoft Cambridge / Redmond
做合成 做DNA存储

Orthogonal split inteins 蛋白胶水
(后面有图)

基础研究 or 应用研究



digital-like circuits \Rightarrow genetic logic gate

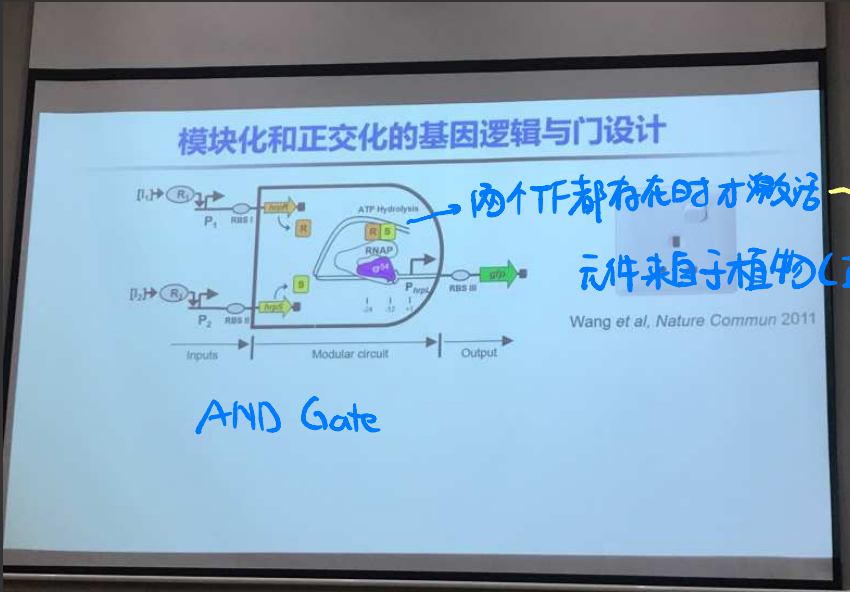
- 标准化的逻辑门库
- 电子设计自动化辅助工具.

帝国理工博士

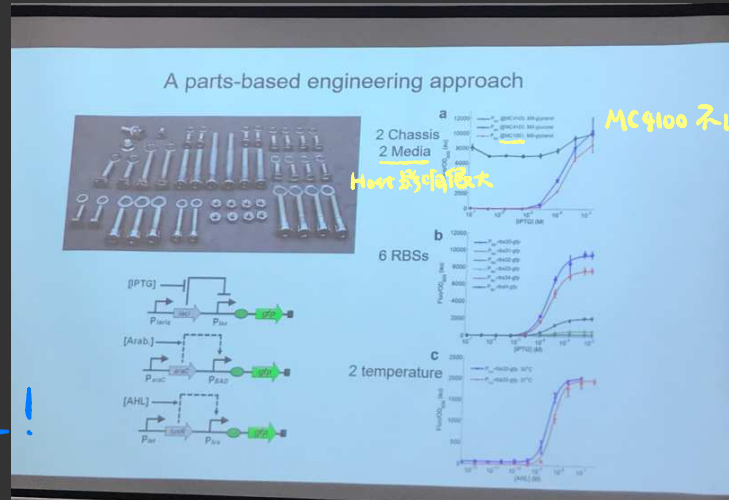
/ˈædəm/ (2007-2011)
Adam, Krishna /ˈkrɪʃnə/

Q: 对比一下不同层次的Gate? 如何选择?

DNA RNA protein
(陈洪荣组) (内分泌组) 有没有什么 principle?



Forward engineer the AND Gate!

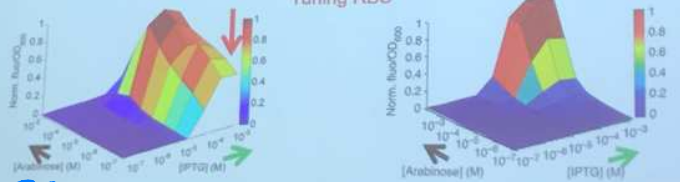


Forward engineer the genetic AND gate



strong → weak

Tuning RBS



(a) AND gate constructed in trial and error

(b) Forward engineered AND gate

In *E. coli* MC1061, M9-glycerol, 30°C

没底更大和IPTG不能拟合图例

模块拼合



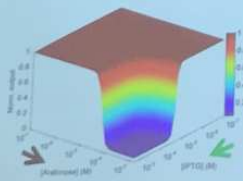
Combinatorial NAND gate – modelling & characterisation



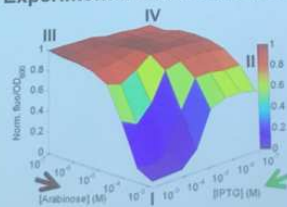
$$f([I]) = k_1 \left(a + \frac{([I]/K)^n}{1 + ([I]/K)^n} \right) \cdot \frac{([R]/K_R)^{n_1} \cdot ([S]/K_S)^{n_2}}{1 + ([R]/K_R)^{n_1} + ([S]/K_S)^{n_2}} [G]_{\max}$$

$$k_2 \left(a_2 + \frac{K_2^{n_2}}{K_2^{n_2} + [R_2]^{n_2}} \right)$$

Model prediction



Experimental characterisation



Q: 拼合时需要关注什么问题

如何标准化?

接口

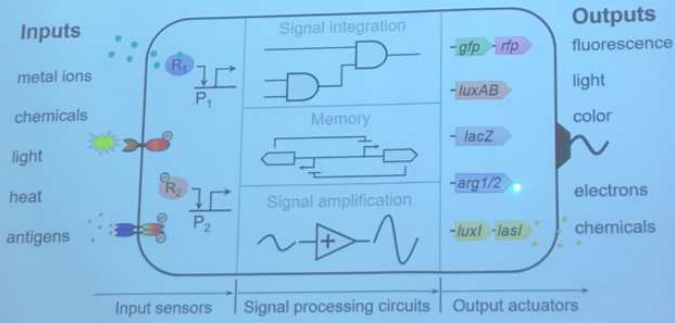
主要用实验测量表征

也许存在

隔离序列等方法

来标准化

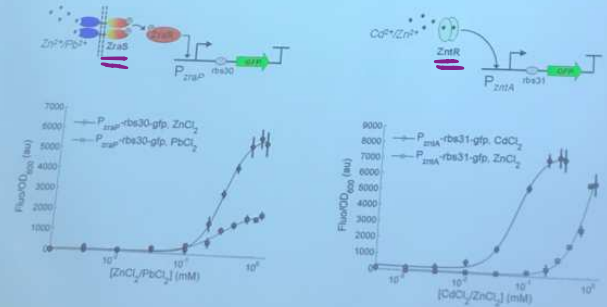
Application example: synthetic cellular biosensors



Modular, Programmable, Multiplex

Wang et al, Biosensors & Bioelectronics, 2013

Non-selective zinc, cadmium cellular sensors



Wang et al, Biosensors & Bioelectronics, 2013

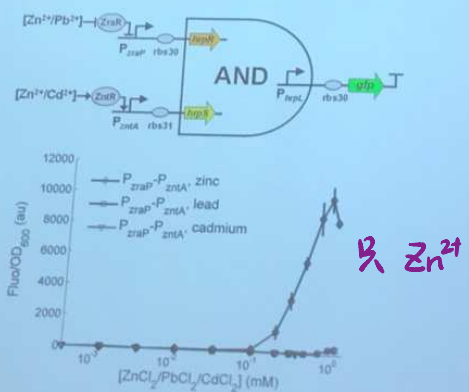
用自然界存在的感受模块做感应。

↳ 不够特异。

↓
用AND Gate 提高特异性

Q: 非线性如何检测?

Solution: selective sensor using an AND gated filter



只 Zn²⁺

Wang et al, Biosensors & Bioelectronics, 2013

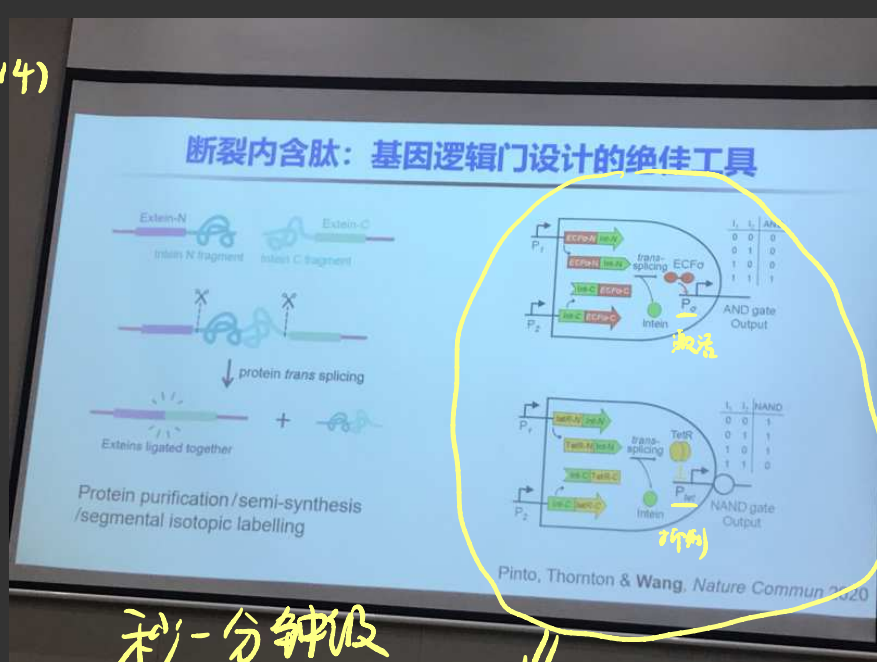
Liu Y et al, Nat Commun 5, 3393 (2014)

国内癌细胞检测并定向消杀

断裂内含肽

不依赖 cofactor、无残留、自切割

纯化、合成领域中广泛使用



秒-分钟级



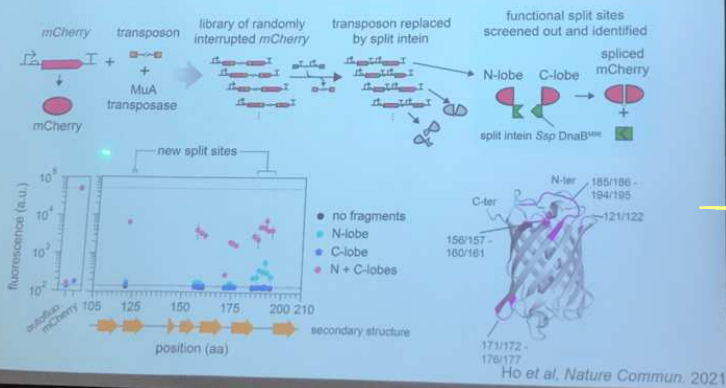
可用内含肽搭建
与门、与非门

15对正内含肽被发现

↓ NEXT

如何在有限断裂?

Intein-assisted bisection mapping (IBM) to find split sites for split intein insertion



查一下位点怎么系

优先

系统的方法

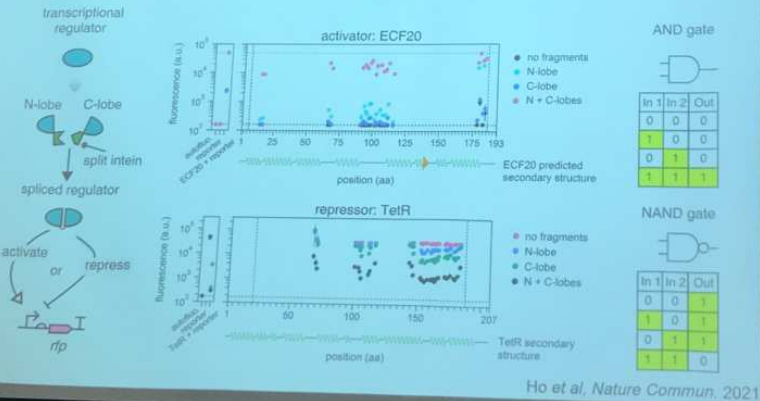
Q: 真核中 work?

全细胞反应?

荧光强度比较好? 查

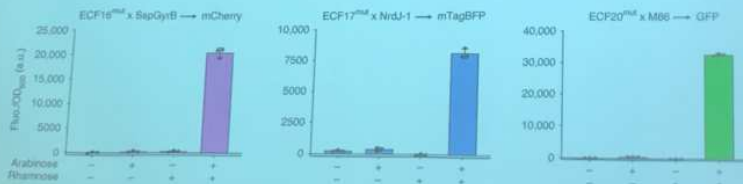
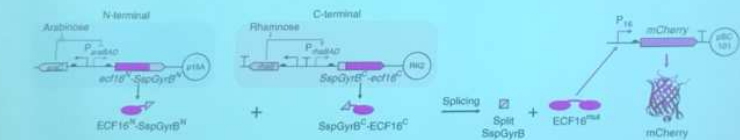
没能直接问出来, 但下来查证一下
看看是否有细胞应用潜力

Any protein with screen-able function can be turned into logic gates for bio-computation



提供了系统化的
断裂位点筛选方法

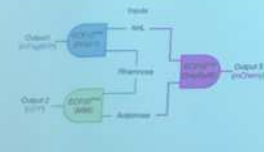
基于断裂内含肽的正交逻辑与门设计



Split inteins can be coupled to transcription factors to engineer logic gates

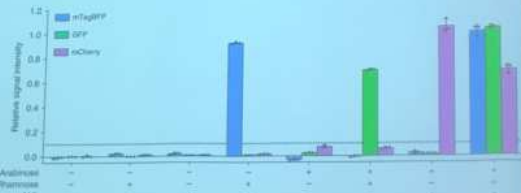
基于断裂内含肽的复杂逻辑线路设计—多输入组合区分

3-input 3-output majority voting circuit



Input 1 (Arabidose)	Input 2 (Rhamnose)	Input 3 (AHL)	Output 1 (mCherry)	Output 2 (mTagBFP)	Output 3 (GFP)
0	0	0	0	0	0
0	1	0	0	0	0
0	0	1	0	0	0
0	1	1	0	0	0
1	0	0	0	1	0
1	1	0	0	1	0
1	0	1	0	1	0
1	1	1	1	1	1

Experimental data match the predicted behavior



Split inteins enable gene circuits scalability and complexity

正交性不错。

Analogue circuits — "transcriptional amplifiers"

模拟电路

模拟线路 — 生物放大器的设计

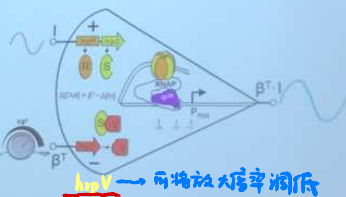
Imperial College London

First ever biological amplifier created by Imperial scientists

24 July 2014



Modular transcriptional amplifier for ultrasensitive chemical sensors

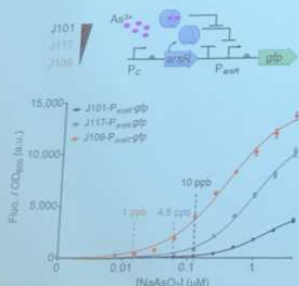


Wang et al, Nuclear Acids Res, 2014

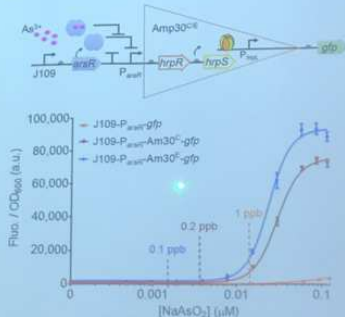
要搞清楚这几个基因的生物学原理

多层放大器级联的超敏感砷污染传感器

1 Tuning receptor density

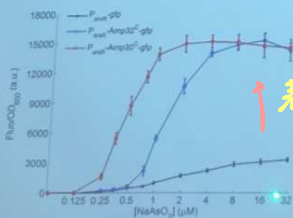
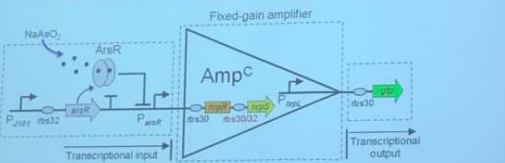


2 Amplifying transcriptional input



Wang et al, Nucleic Acids Res, 2015, 43:1955

具有转录信号放大功能的砷污染传感器

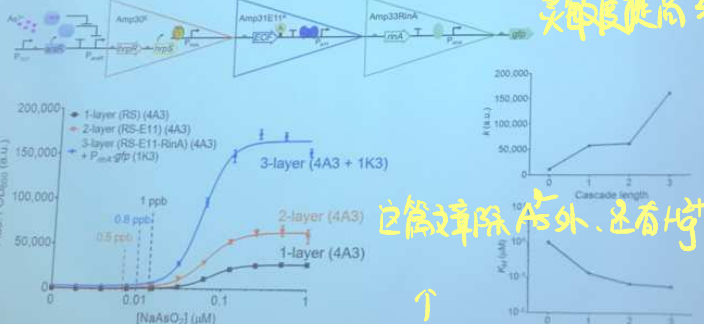


表达量放大

放大关系成倍率 (0-20倍)

Wang et al, Nucleic Acids Res, 2014

多层放大器级联的超敏感砷污染传感器



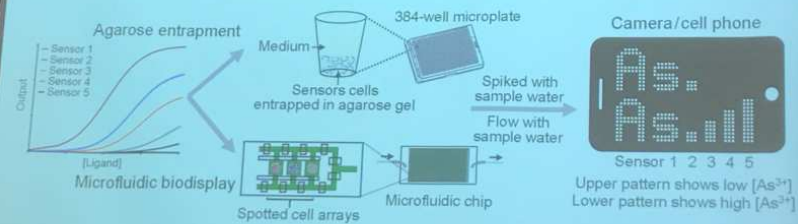
灵敏度提高3个数量级

这篇论文外，还有诗

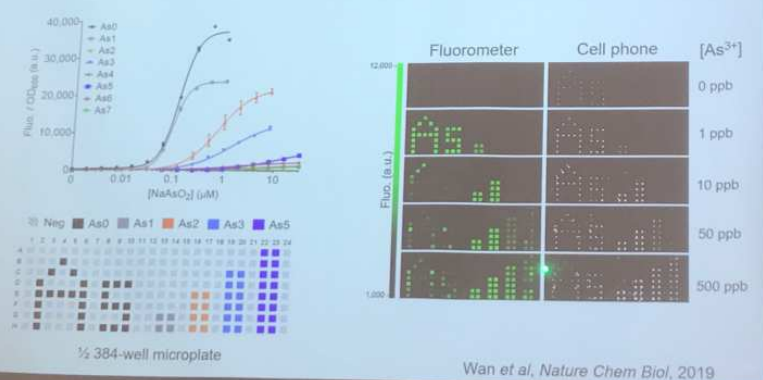
Wan et al, Nature Chem Biol 2019

搬到生产上!

Field application-oriented sensor arrays



Agarose gel-encapsulated sensor array for arsenic



成果应用：改造大肠杆菌灵敏感知饮用水砷污染



Wan, ... & Wang*, Nature Chem Biol 2019
Wan, ... & Wang*, Nature Commun. 2020

- 首次建立一套完整的全细胞传感器优化方案
- 设计出目前最敏感的神汞重金属污染微生物传感器 (检测极限提高2个数量级, 达0.01 ppb)

孟加拉国超一半以上人均饮用水中

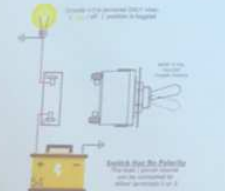
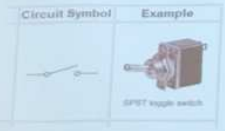
As³⁺超标

"Nat chem bio"

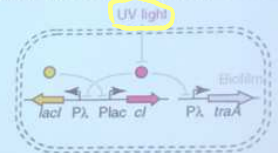
Engineer memory circuits — data storage in living cells

拨动开关 (toggle switch)

控制电路的闭合与断开



作为双稳态基因开关使工程细菌形成生物膜



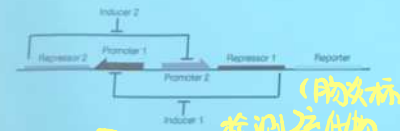
生物
⇄
工程



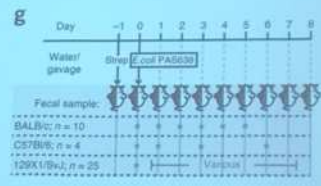
无DNA损伤 DNA损伤

最早的 toggle switch
Collins et al. PNAS, 2004

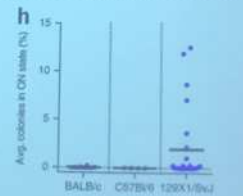
A genetic toggle switch in *E. coli* for recording gut inflammation signal



(防止病原体)
检测硫化物



可逆性
之间以上



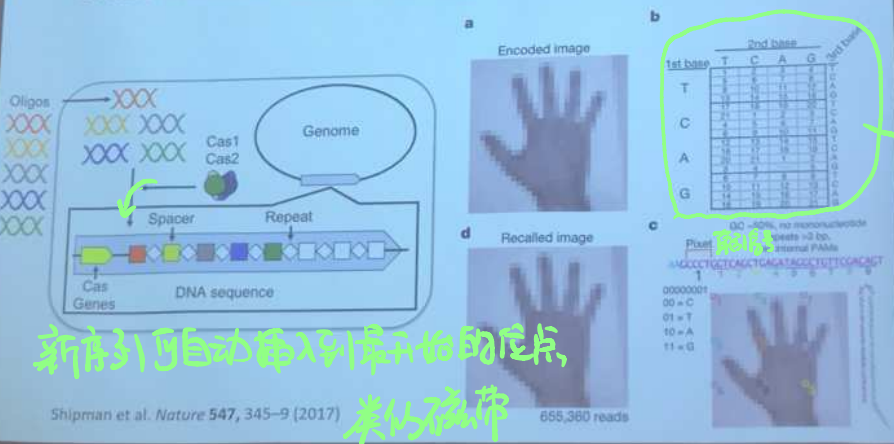
Riglar et al. Nature Biotech 35, 653-8 (2017)

可逆的 toggle switch?

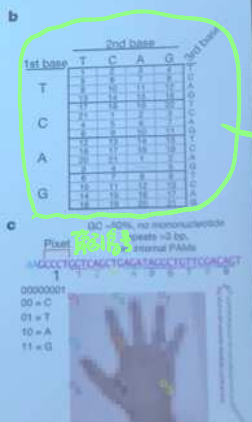
Angel 的报告有提到 (后面PT)

甚至存在用相同信号调控状态切换的 toggle switch.

CRISPR-based memory storage in the genome of cells



新序列会自动插入到最开始的位点，
类似磁带



三个碱基 有一个灰度

Cello (Nielsen 2016 Science)
自动化设计软件

SynBio trends in the next decade

Present

- Trial-and-error
- Building - limiting step
- Labour intensive
- Small scale systems
- Single cell organisms
- Static, module focused
- Lab proofs-of-concept

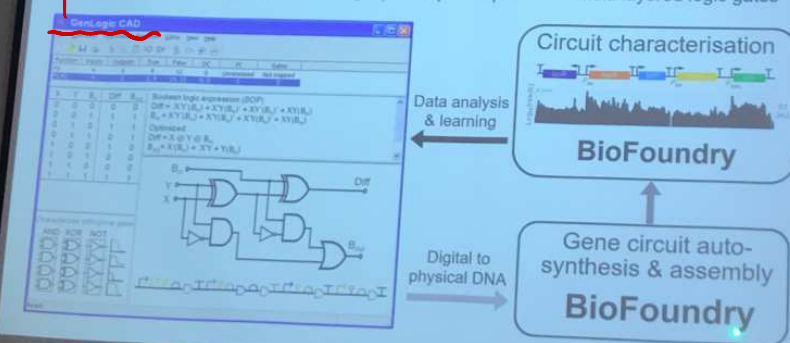
Future 2021-2030

- Predictive & automated
- Data interpretation -limiting step
- Machine intensive
- Large scale complex systems
- Multicellular organisms
- Dynamic, system focused
- Real-world applications

假想的模式图，不存在真实工具

基因线路设计自动化

Design program to automate matching input-output responses of multi-layered logic gates



Analogy of Bio-IC design to ARM processor design

ARM develop and license CPU cores

+

ARM development tools (EDA software)

→

ARM core designs licensees



IDEs [DS-5](#) | [Keil MDK](#)
Compilers [Arm Compiler 6](#)
Boards [Juno Arm Platform](#)

Apple (A6 – A12 chips)
Huawei (HiSilicon)
Samsung, Intel ...

Bio-IC design IP cores
Libraries of genetic logic blocks

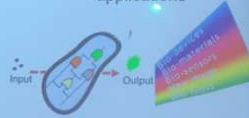
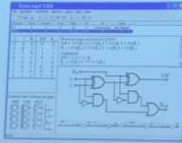
+

Bio-CAD software

→

Bio-circuits/ASICs
for bio-computing/control applications

Characterised orthogonal gate libraries



0824 Angel

Introduction – the Biocomputation Lab



I am a Computer Scientist by training, and I am still doing computer science, but using living biological systems as hardware / software. I got to know about DNA Computing in my undergrad and started researching on Synthetic Biology for my PhD. After postdoc-ing with Martyn Amos (MMU, UK) and Victor de Lorenzo (CNB, Spain), I established my first group at Newcastle University (UK). Back to Madrid in 2020 to join the CBBG and launch the Biocomputation Lab.



[From top left] Lorea Alejandre (biochemist), Matthew Crowther (computer scientist), Ana Valero (biotechnologist), Elena Rodríguez (lab technician), Coral García (lab manager), Angeles Hueso (microbiologist / molecular biologists), Lewis Grozinger (computer scientist), Jesús Miró-Bueno (theoretical physicist), Paula Múgica (biotechnologist), Juan Méndez (microbiologist), metabolic engineer), Juan Rico (genetic engineer)

Introduction – the Biocomputation Lab

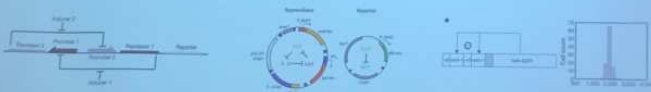
- Engineered biological circuits** / Designing biocircuits to perform predefined functions
Genetic and metabolic networks are like cascades of events triggered by certain inputs, such as environmental signals, to deliver outputs. We design and implement synthetic networks that process inputs into outputs according to predefined rules encoded in the genetics of microbes. This set of rules is equivalent to an algorithm. In other words, we program biocomputations.
- Biological complexity: from context to evolution** / Understanding how context affect performance
Genetic circuits are much more than DNA sequences and they do not guarantee the same behaviour in different organisms. Their performance is influenced by how the host context interacts with the circuits, something we call contextual dependencies. Our aim is to characterise these. Evolution is the best dependency!
- Mathematical modelling beyond experiments** / Towards understanding of biological systems
How does ribosome jamming alter translation? What are the effects of transcription-translation coupling of circuits? Can multicellular systems be effectively modulated? We use mathematical modelling and computational simulations to better understand how biological systems work.
- Standards, automation, and data handling** / Committed to synthetic biology standardisation
Our laboratory uses and has contributed to the Standard European Vector Architecture (SEVA) and the Synthetic Biology Open Language (SBOL). We are developing methods to visualise and edit design information in a user-friendly manner, and we use automation machinery.

Introduction – Synthetic Biology

[Definition #1] "Synthetic biology is a) the design and construction of new biological parts, devices and systems and b) the re-design of existing natural biological systems for useful purposes." SyntheticBiology.org

[Definition #2] "Synthetic biology [is] ... the design and construction of novel artificial biological pathways, organisms or devices, or the redesign of existing natural biological systems." UK, Royal Society

[Definition #3] "A technology that makes the design and construction of biological systems easier". Drew Endy



Gardner, T. S., Cantor, C. R., & Collins, J. J. (2000). Construction of a genetic toggle switch in *Escherichia coli*. *Nature*, 403(6767), 339-342.

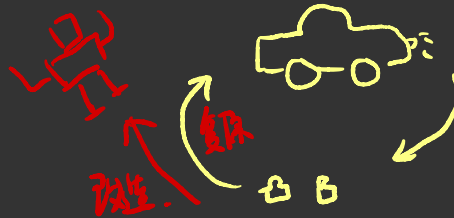
Elowitz, M. B., & Leibler, S. (2000). A synthetic oscillatory network of transcriptional regulators. *Nature*, 403(6767), 335-338.

Beekel, A., & Serrano, L. (2000). Engineering stability in gene networks by autoregulation. *Nature*, 405(6786), 590-593.

Synthetic vs System

Bottom-Up

Top-Down

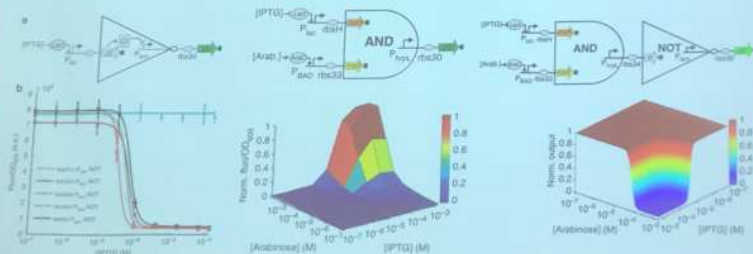


Introduction – Biocomputing / programming cells



Synthetic biology: Engineering Escherichia coli to see light. Levskaya et al. Nature 438, 441-442 (2005)

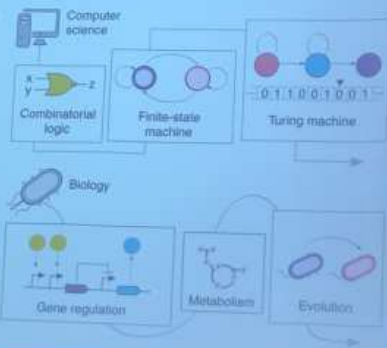
NOT / AND / NAND logic gates



Wang, B. Engineering modular and orthogonal genetic logic gates for robust digital-like synthetic biology. Nature Communications (2011) DOI: 10.1038/ncomms1516

Introduction – Biocomputing

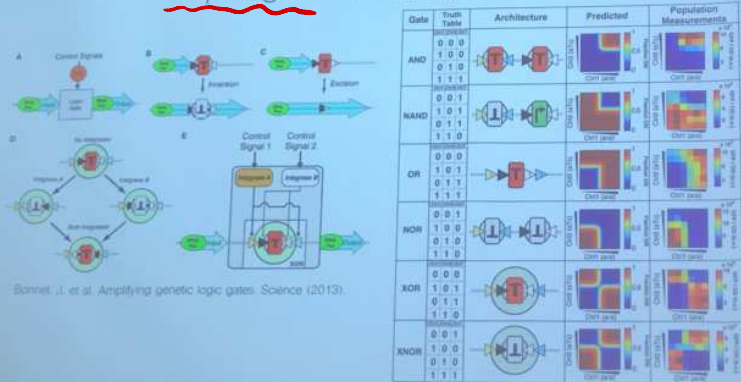
- Not all models of computation are equally powerful, that is, they are not equivalent in terms of the range of computations they can describe.
- Combinatorial logic is extremely weak.
- Two extremely powerful models are the Lambda calculus and the Turing Machine.



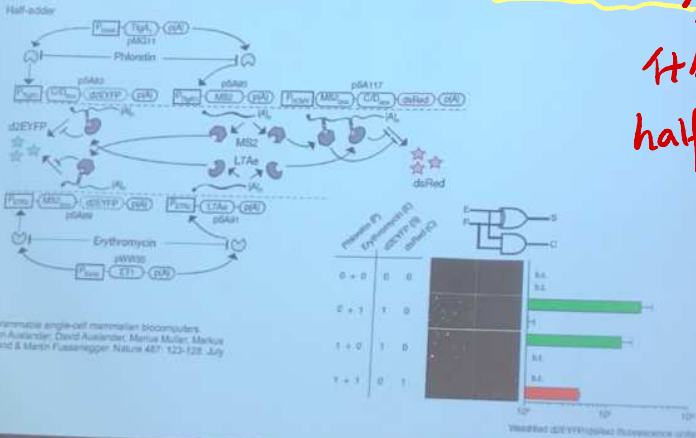
Diogenes, L., Anos, M., Gorskiowski, T. E., Carbonell, P., Opatow, D. A., Sood, P., & Sood-Alexander, A. (2019). Pathways to cellular supremacy in biocomputing. Nature Communications, 10(1), 1-11

DNA/RNA Logic Gate!

Flip logic – memory circuits

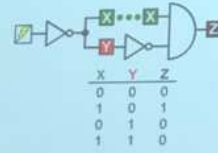
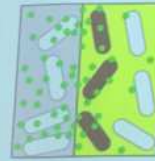


Mammalian logic circuits – half adder

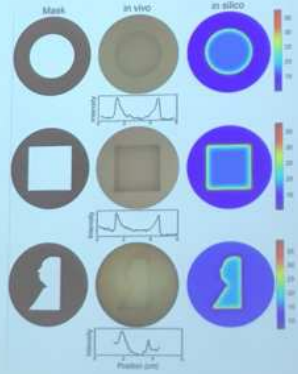


查一下
什么叫做
half adder

The edge detector (single cell + populations)

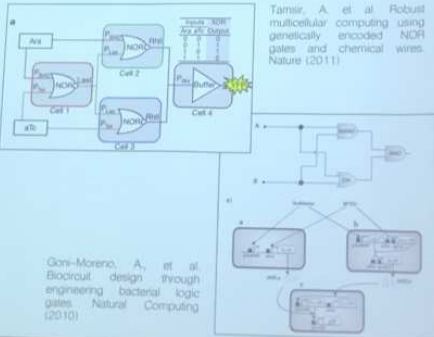


A synthetic genetic edge detection program. Jeffrey J. Tabor, Howard M. Salis, Zachary Booth Simpson, Aaron A. Chevalier, Anamir Levskaya, Edward M. Marcotte, Christopher A. Voigt and Andrew D. Ellington. Cell, 137: 1272-1281 (2009)

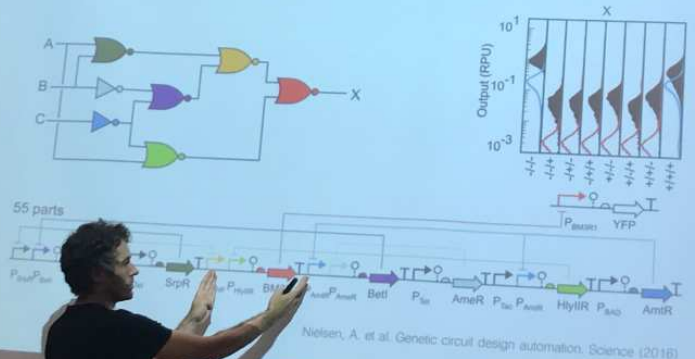


Multicellular distributed computations

However, for certain computations, algorithms exist in which the ordering of computational steps may be relaxed, or abandoned entirely. These algorithms display a level of concurrency.



Large-scale logic circuits



Predictive Biology

纷乱的复杂

体系混乱的复杂

Complicated is different than complex

An engine is complicated. A complicated system has a direct and clear cause and effect relationship. Its elements interact in a predictable way. For each action, there is a proportional reaction. Its problems are often difficult to solve but can be solved with rules and processes after a rational and specialized assessment.

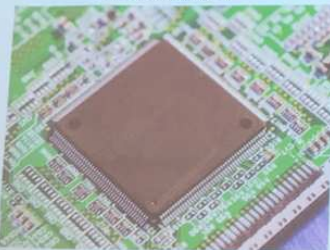


A complex system, in turn, is composed of elements that interact with each other exhibiting a dynamic and adaptive behavior. The relationships are more important than the elements themselves. Interactions are what matter. These interactions can give rise to unpredictable behaviors, with no single identifiable cause. Small actions can lead to huge reactions, while great interventions may prove ineffective.

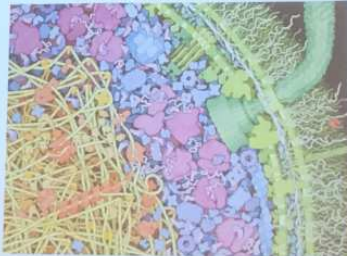
Complicated is different than complex

We can't control [complex] systems [permanently] or figure them out. But we can dance with them! —Donella Meadows

Chips are complicated



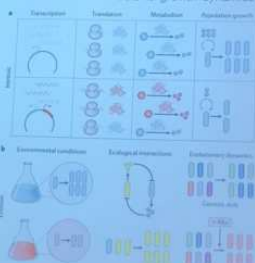
Living systems are complex



Predictive biology as a challenge

Predictive biology is the next great chapter in synthetic and systems biology, particularly for microorganisms. Tasks that once seemed infeasible are increasingly being realized such as designing and implementing intricate synthetic gene circuits that perform complex sensing and actuation functions, and assembling multi-species bacterial communities with specific, predefined compositions. These achievements have been made possible by the integration of diverse expertise across biology, physics and engineering, resulting in an emerging, quantitative understanding of biological design. As ever-expanding multi-omic data sets become available, their potential utility in transforming theory into practice remains firmly rooted in the underlying quantitative principles that govern biological systems.

Factors that contribute to growth dynamics



Loosekin, A. J. & Collins, J. J. (2020) Predictive biology: modeling, understanding, and harnessing microbial complexity. *Nature Reviews Microbiology*, 18(1), 527-532

Predictive biology as a challenge

Complex dynamics in engineered and natural populations

- Systematic characterization of defined ecological modules tested in diverse environments
- Implement high-throughput experimental techniques to increase the parameter space of every isolate-by-environment interaction

Increasingly large data sets

- Minimize 'fishing expeditions', where appropriate, by validating insights from large quantitative data sets with specific testable hypotheses through mathematical models and controlled experiments
- Computational tools that facilitate integrating diverse levels of information into predictive models
- Centralized and accessible parameter reporting

Accurate parameter estimates

- Establish an appropriate level of abstractness with which to characterize a given system
- Match parameter definitions with experimental estimation
- Use conditions that more closely mimic the process of interest to improve model accuracy
- Improved technical tools to standardize parameter estimates
- Efforts to use consistent terminology to facilitate cross-literature data compilation and review

Accounting for stochastic and deterministic evolution in complex populations

- Integration with machine learning and systematic analyses into evolutionary constraints
- Expand sequencing breadth and depth to explore evolutionary responses at lower detection limits

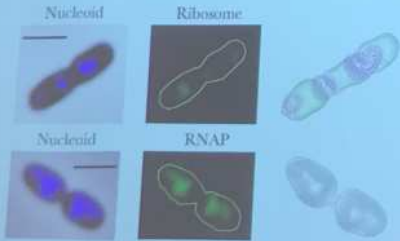
Translating in vitro predictions to in vivo outcomes

- Implement experimental conditions that are biologically relevant
- Involvement of animal models and other experimental platforms that can better simulate natural environments of interest
- Incorporate physiological conditions into modeling analysis
- Advances in culturing techniques
- Establish in situ quantification methods

Loosekin, A. J. & Collins, J. J. (2020) Predictive biology: modeling, understanding and harnessing microbial complexity. *Nature Reviews Microbiology*, 18(1), 527-532

Measure with image analysis

- Using an image analysis software (CellShape) we measured the correlation between nucleoid, ribosome and RNAP areas.
- Prediction: localization important & RNA/TF diffusion should play a role



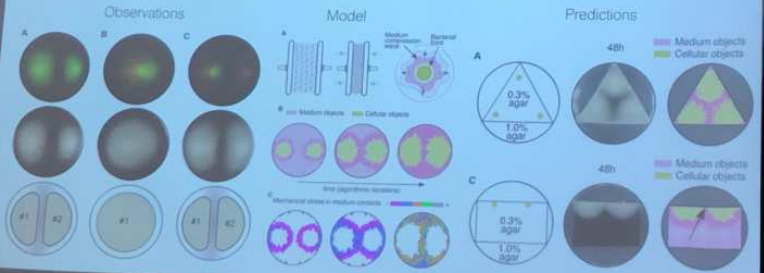
Kim, J., Goh-Moreno, A., Calles, B. and de Lorenzo, V. (2018). Spatial organization of the gene expression hardware in *Pseudomonas putida*. *Environ Microbiol*, 21, 1645–1658.



Predicting populations with agent-based models (I)

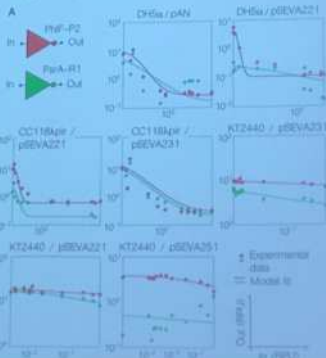
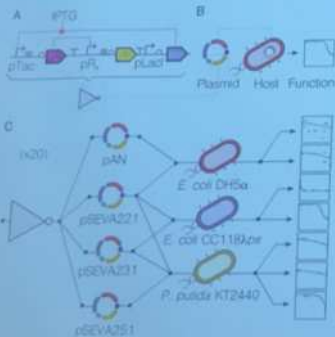
Predictions suggested that the medium is compressed in the direction of the bacterial front motion. This phenomenon generates what was termed a *compression wave* that goes through the medium preceding the swimming population and that determines the visible high-level pattern.

Espejo, D. R., Martínez-García, E., De Lorenzo, V., & Goh-Moreno, A. (2019). Physical forces shape group identity of swimming *Pseudomonas putida* cells. *Frontiers in microbiology*, 7, 1437.



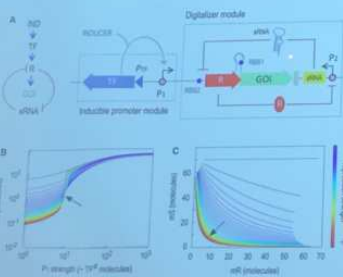
Predictions are stopped by nonlinearities

Two gates change performance in different ways. Very difficult to build predictions!



Nelson, T. et al. Contextual dependencies impact the reusability of genetic inventory. *Nature Communications* (2021).

Predicting dynamics with ODEs (the digitalizer)



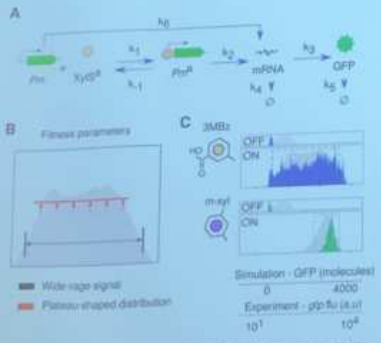
Bottom left graph. Each line is a single simulation that measures the level of mR (y axis) while the concentration of TF molecules increases (x axis). The colour of the lines goes from dark red (both repressions very strong*) to dark blue (both repressions very weak).

Bottom right graph. Same as before, but plotting mS vs mR. We see the relationship between the two mRNA species in the system. The stronger the repression forces (both) are, the more exclusive this relationship is.

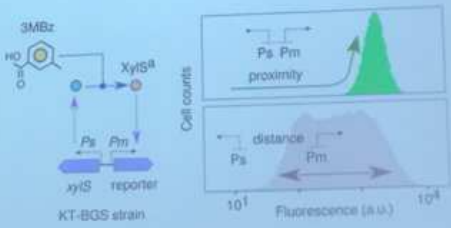


Calles, B., Goh-Moreno, A., & de Lorenzo, V. (2019). Digitizing heterologous gene expression in Gram-negative bacteria with a portable ON/OFF module. *Molecular systems biology*, 15(12), e877.

Predicting dynamics with SSA

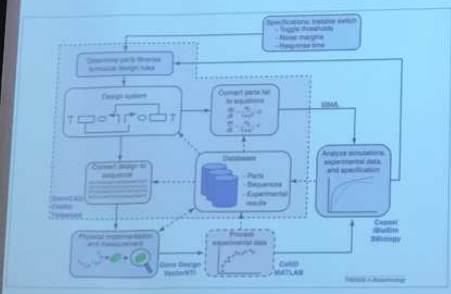


Transcriptional noise was predicted to depend on the intracellular physical distance between regulator source (where XylS is produced) and the target promoter. Experiments with engineered bacteria in which the distance is minimized or enlarged confirmed the predicted effects of source/ target proximity on noise patterns.



Guth-Mohr, A., Benedetti, I., Kim, J., & de Lorenzo, V. (2017). Decoupling of gene expression noise and static dynamics of transcription factor-promoter interact. ACS synthetic biology, 6(7), 1369-1380.

Computer Aided Design (CAD) / Aided Modelling (CAM)



"The aim of synthetic biology is to make genetic systems more amenable to engineering, which has naturally led to the development of computer-aided design (CAD) tools. Experimentalists still primarily rely on project-specific ad hoc workflows instead of domain-specific tools, which suggests that CAD tools are lagging behind the front line of the field. Here, we discuss the scientific hurdles that have limited the productivity gains anticipated from existing tools. We argue that the real value of efforts to develop CAD tools is the formalization of genetic design rules that determine the complex relationships between genotype and phenotype"

Lee, M. W., Shinar, B. W., Bai, D. A., & Pedraza, J. (2012). Genetic design automation: engineering biology in scientific research. Trends in biotechnology, 30(2), 120-126.

The Synthetic Biology Open Language (SBOL)

The SBOL data standard is a data exchange representation for synthetic biology designs. Its goal is to improve the efficiency of data exchange and reproducibility of synthetic biology research. SBOL introduces a standardized format for the electronic exchange of information on the structural and functional aspects of biological designs. The standard has been designed to support the explicit and unambiguous description of biological designs by means of a well defined data model. The standard further describes the rules and best practices on how to use this data model and populate it with relevant design details. SBOL uses existing Semantic Web practices and resources, such as Uniform Resource Identifiers (URIs) and ontologies, to unambiguously identify and define genetic design elements. The definition of the data model and associated format, the rules on the addition of data within the format and the representation of this in electronic data files are intended to make the SBOL standard a useful means of promoting global data exchange between laboratories and between software programs.



<https://sbolstandard.org>

A community effort. Current specification version 3.0.1. You can find it here:

<https://sbolstandard.org/docs/SBOL3.0.1.pdf>

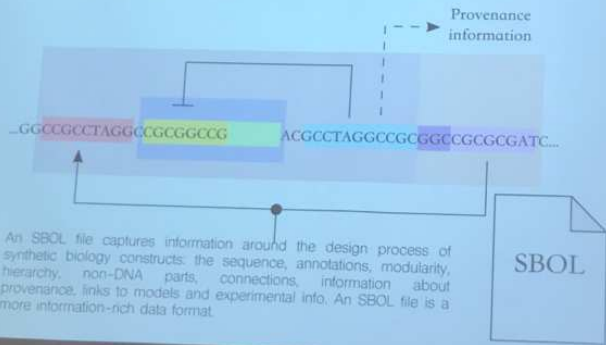
The Synthetic Biology Open Language (SBOL)

...GCCCGCTAGGCGCGGCCGCGCGAACGECTAGGCCGCGCGCGCGATC...

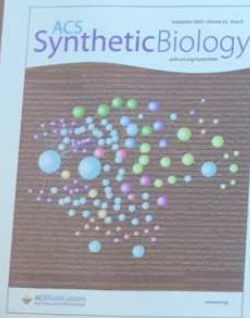
A GenBank file encodes a sequence with annotations. That is, it captures information on what are specific regions on that sequence and other information the user wants to input concerning sequence annotations. This type of format is much more useful if you need to identify different regions, parts or sub-sequences.



The Synthetic Biology Open Language (SBOL)



Networks to visualise/analyse SBOL information

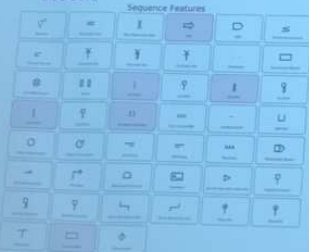


"As genetic circuits become more sophisticated, the size and complexity of data about their designs increase. The data captured goes beyond genetic sequences alone; information about circuit modularity and functional details improves comprehension, performance analysis, and design automation techniques. However, new data types expose new challenges around the accessibility, visualization, and usability of design data (and metadata). Here, we present a method to transform circuit designs into networks and showcase its potential to enhance the utility of design data. Since networks are dynamic structures, initial graphs can be interactively shaped into subnetworks of relevant information based on requirements such as the hierarchy of biological parts or interactions between entities. Additionally, several visual changes can be applied, such as coloring or clustering nodes based on types (e.g., genes or promoters), resulting in easier comprehension from a user perspective."

Crowther, M., Wast, A., & Gori-Moretti, A. (2012). A network approach to genetic circuit design. *ACS Synthetic Biology*, 1(8), 3058-3066

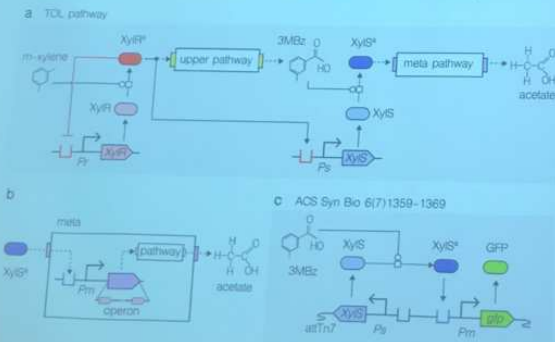
SBOL Visual

SBOL
VISUAL



People who are engineering biological organisms often find it useful to communicate in diagrams, both about the structure of the nucleic acid sequences that they are engineering and about the functional relationships between sequence features and other molecular species. Some typical practices and conventions have begun to emerge for such diagrams. SBOL Visual aims to organize and systematize such conventions in order to produce a coherent language for expressing the structure and function of genetic designs. At the same time, we aim to make this language simple and easy to use, allowing a high degree of flexibility and freedom in how such diagrams are organized, presented, and styled—in particular, it should be readily possible to create diagrams both by hand and with a wide variety of software programs.

SBOL Visual – example 1



Beck, J., Nguyen, T., Gorochowski, T. E., Gori-Moretti, A., Scott-Brown, J., McLaughlin, J. A., & Wast, A. (2018). Communicating structure and function in synthetic biology diagrams. *ACS synthetic biology*, 8(6), 1818-1825

The Systems Biology Markup Language (SBML)

<https://sbml.org>

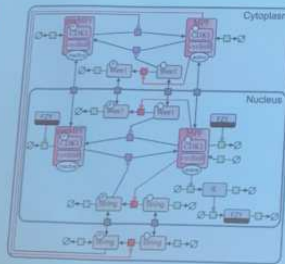


The starting point is an appreciation that computational modeling of biological systems is no longer a fringe activity—it's a requirement for us to make sense of our vast and ever-expanding quantities of data. This reality is acknowledged and reinforced by a vast increase over the past two decades in the number of journals, books and articles having computational and systems-biology emphases.

At the level of software, a different format is needed for quantifying a model to the point where it can be simulated and analyzed. That's where the Systems Biology Markup Language (SBML) comes in.

Simply put, SBML is a machine-readable format for representing models. It's oriented towards describing systems where biological entities are involved in, and modified by, processes that occur over time. An example of this is a network of biochemical reactions, SBML's framework is suitable for representing models commonly found in research on a number of topics, including cell signaling pathways, metabolic pathways, biochemical reactions, gene regulation, and many others.

The Systems Biology Graphical Notation (SBGN)



"Welcome to the global portal for documentation, news, and other information about the Systems Biology Graphical Notation (SBGN) project, an effort to standardise the graphical notation used in maps of biological processes"

<https://sbgn.github.io>

Trink, V., Le Novère, N., Wollersdorf, O., & Wolkowicz, D. (2018). Quick start for creating effective and impactful biological pathways using the Systems Biology Graphical Notation: Plus computational biology, 14(2), e1005740.

The Standard European Vector Architecture (SEVA)

<https://seva-plasmids.com>



The Standard European Vector Architecture (SEVA) platform is a web-based resource and a material clone repository to assist the choice of optimal plasmid vectors for de-constructing and re-constructing complex prokaryotic phenotypes.

The SEVA database (SEVA-DB) is a resource for implementation of a standard for physical assembly of vector plasmids and for their non-ambiguous nomenclature as well as the index for a repository of functional sequences and actual constructs available to the community. The database was designed to simplify the choice of a given vector for the sake of specific applications, in such a way the user can easily decide the best configuration of replication origins, antibiotic resistance and business segments.



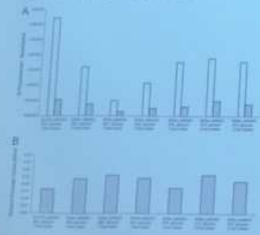
The Standard European Vector Architecture (SEVA)



Rules for naming each of the 4 positions available for composing a complete SEVA code. The first position identifies the antibiotic resistance. The second position is the origin of replication (a sole numeric code 1 to 9 and then a capital letter is added). In case of either variants or addition(s) of either the antibiotic marker or a second replication origin, a lower case letter is then inserted next. The third position is the cargo, which can be mono-function (named 1 to n) or variants/multi-function thereof (number followed by a capital letter). Finally, the fourth position is for the gadgets, which are designated by lowercase Greek letters (α to ω).

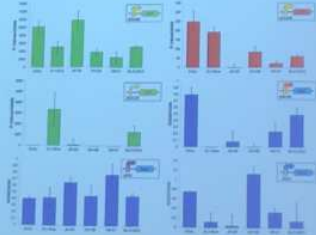
Metrology (the biggest challenge!) – RPU example

RPU: Reference/Relative Promoter Units. A single reporter protein fluctuates over time, and changes from condition to condition—from host to host. Using a second one (always the same) and correlate units gives more stable measurements.



Kelly, J. R., Rubin, A. J., Davis, J. H., An-Popstein, C. M., Cumber, J., Choi, M. J., & Endy, D. (2009). Measuring the activity of *BioBrick* promoters using an *in vivo* reference standard. *Journal of biological engineering* 3(1), 1-13.

However, not all promoter-reporter systems fluctuate in the same way from host to host, so the reference may drastically change measurements in specific cases.



Vilanova, C., Tanner, K., Dorado-Morales, P., Villoslada, P., Chapman, D., Fiala, A., & Forster, M. (2015). Standards not that standard. *Journal of biological engineering* 9(1), 1-4.

Standards – to sum up...

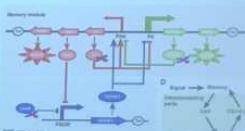
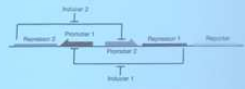
- SBOL: We use the Synthetic Biology Open Language to represent information and capture details beyond DNA sequences and annotations. We use SBOL for designing and sharing.
- SBML: We use the Systems Biology Markup Language to share mathematical models. To run them in different software tools. We use SBML for the analysis and prediction.
- SBGN: We use the Systems Biology Graphical Notation to illustrate models. This representation corresponds to SBML models.
- SEVA: We use the Standard European Vector Architecture to build plasmid vectors with known parts, sequences and functions. We share functions that are well characterised.
- RPUs: We use reference promoter units to move away from "arbitrary units" and calibrate our measurements in absolute units. We share experimental data using RPUs in order to compare and reproduce.
- Chassis: This is still a challenge... What is a chassis?

de Lencastre, J., Andrade, A., & Salgado, M. (2011). For the sake of the Bioscience, please visit a Synthetic Biology Chassis at: www.Bioscience.20.41-21

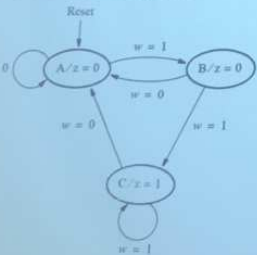
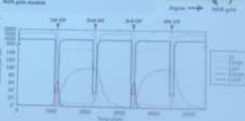
Sequential Logic – States

The output of stateful computations depends not only on the current input, but also on the current state. State Machines are more powerful than Combinatorial Logic.

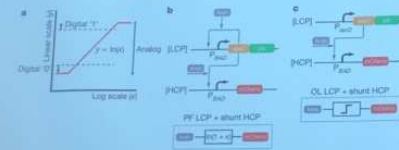
Gardner, T., Caribé, C., Collins, J. Construction of a genetic logic switch in *Saccharomyces cerevisiae*. *Nature* (2003).



Lin, C., Liu, X., Ni, M., Huang, Y., Huang, Q., Huang, L., & Doyagi, Q. (2018). Synthesizing a novel genetic sequential logic circuit: a push-on switch. *Molecular systems biology*, 14(1), 260.



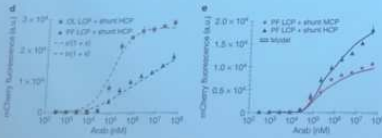
Analogue Biocomputing (in/out range pop. level)



This work shows variations in dynamic ranges in relation to network motif (PF vs OL).

PF returns wider dynamic range (i.e., more values at the output). Here, the OL returns a more digital-like behaviour (at the population level).

Computing here is done by turning inputs into outputs, in the analogue domain!



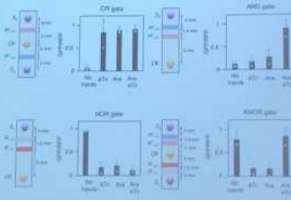
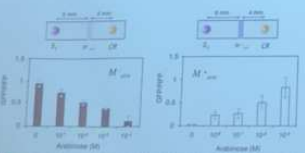
Daniel, R., Rubens, J. R., Sengupta, R., & Liu, T. K. (2013). Synthetic analog computation in living cells. *Nature*, 497(7451), 819–822.

Analogue Biocomputing (in/out range pop. level)

By introducing (extracellular) space as another parameter, analogue signalling is used to process information. S: sender strain; CR: reporter strain; M: modulation cells; CS: carrying signal.

Modulating space in between S, CR and M cells; the reporter performed different functions. The genetic programs are the same, but processing is analogue.

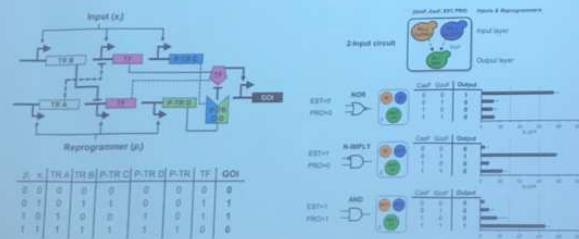
M: The more Ara, the less CS goes through (block).
M': The more Ara, the more CS goes through.



Choi, S., Srinivasan, P., & Mazza, J. (2021). 3D printed multicellular devices performing digital and analogue computation. *Nature* (comm), 3(1), 1–10.

Reprogrammable biocomputers

Individual cells (left) can let inputs (x) go through (ID) if the reprogramming signal (p) is 0. They invert (NOT) the input if the reprogramming signal is 1. By mixing different cells in a chamber, the program is reconfigured dynamically (right): same cells, same genetic circuits, different reprogramming signals to build a NOR an N-IMPLY or an AND function. Key concept: there are inputs, outputs and reprogramming signals (three different types).

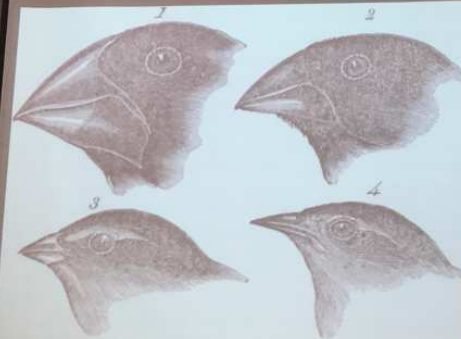


Canadell, O., Dru-Vaquenas, N., Mogen-Diaz, S., de Nadal, E., Maza, J., & Paves, F. (2022). Implementing reconfigurable biological computation with distributed multicellular consortia. *Nucleic Acids Research*.

Biocomputing beyond Turing?

"A priori it is not obvious that every function which we could intuitively regard as computable by an algorithm can be computed using a Turing Machine. Church, Turing and many other people have spent a great deal of time gathering evidence for the Church-Turing thesis, and in sixty years no evidence to the contrary has been found. Nevertheless it is possible that in the future we will discover in Nature a process that computes a function not computable on a Turing Machine. It would be wonderful if that ever happened, because we could then harness that process to help us perform new computations which could not be performed before. Of course we would also need to overhaul the definition of computability, and with it, computer science."

Nelson M. A. & Chung I. (2002)
Quantum computation and quantum information

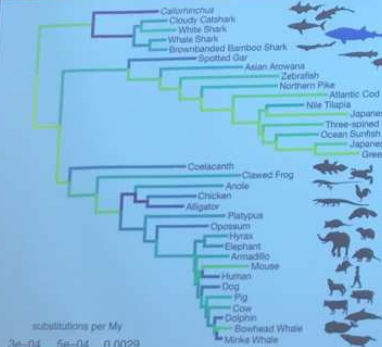


1. *Geospiza magnirostris*.
2. *Geospiza parvula*.

3. *Geospiza fortis*.
4. *Certhidea olivacea*.

Is this not the solution to a problem? (no matter how the solution is found). Is this not a computation? There are inputs, outputs and rules. What about those rules?

Darwin's finches or Galapagos finches Darwin, 1845. Journal of zoology and the natural history and geology of the countries visited during the voyage of H.M.S. Beagle round the world, under the Command of Capt. Fitz. Roy, Vol. 28, plates 1.

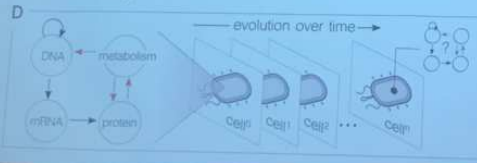


Dolphins and cows have more in common than dolphins and sharks. Their physics are solutions to a problem. Is evolution then a set of algorithmic rules to turn inputs into outputs? A computation then?

Tan, M., Redmond, A. K., Dooly, H., Noss, R., Sato, K., Kuraku, S., ... & Read, T. (2015). The whale shark genome reveals patterns of vertebrate gene family evolution. *Genome*, 58, 465-484.

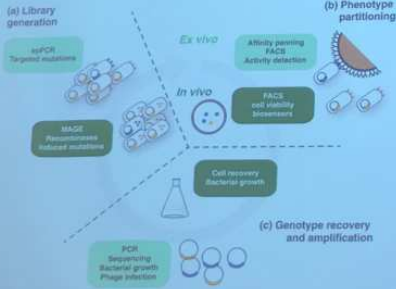
Motivation: Evolutionary Computing vs. Evolution

- Evolutionary computing is a field within computer science that develops and applies algorithms inspired by Darwinian evolution.
- Living systems evolve.
- Both computer scientists and living systems use evolutionary algorithms to generate algorithmic solutions.
- However, the rational engineering of autonomous evolutionary computing in living cells is still an overarching challenge.
- A potential initial step would be to expand on the standard representation of the central dogma of molecular biology.



Groeninge, L., Amici, M., Garofanelli, T., E. Carbonell, P., Ovaron, D. A., Saut, R., ... & Carlini-Moreno, A. (2019). Pathways to cellular superscience in biocomputing. *Communications*, 10(11), 1-11. Nature

Side [2]: using natural evolution



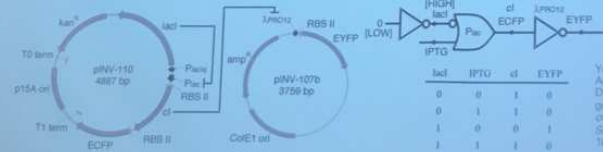
- Directed evolution for optimizing and engineering novel functions in both nucleic acids and proteins.
- Similarities with Darwinian evolution: genetic and phenotypic diversity, selection...
- Those simple steps hide a huge number of methodologies and procedures: library quality, evolutionary landscapes...
- Rational design vs. directed evolution (full knowledge of the system vs. gaps)
- More than 20 years of successful directed evolution (> than synbio)

Pedro Tizei et al. Selection platform for directed evolution in synthetic biology. *Biochem. Soc. Trans.* (2016) 44, 1156-1175

Side [2]: using natural evolution (Logic example)

Evolved the connection between gates (the ci repressor in this example). Warning: directed evolution needs manual steps (i.e., selection).

"We propose a combined rational and evolutionary design strategy for constructing genetic regulatory circuits, an approach that allows the engineer to fine-tune the biochemical parameters of the networks experimentally in vivo. By applying directed evolution to genes comprising a simple genetic circuit, we demonstrate that a nonfunctional circuit containing improperly matched components can evolve rapidly into a functional one"



Yokobayashi, Y., Weiss, R. & Arnold, F. H. (2002). Directed evolution of a genetic circuit. *Proceedings of the National Academy of Sciences*, 99(26), 16587-16591

CV: Paper Production

Virus Machines

Using Virus Machines to Compute Pairing Functions, *IJNS*, 2023

Generating, Computing and Recognizing with Virus Machines, *TCS*, 2023

No Virus Machines

A Discrete Representation of the Second Fundamental Form, *Mathematics*, 2022

A Protocol for Solutions to DP-Complete Problems Through Tissue Membrane Systems, *Mathematics*, 2023

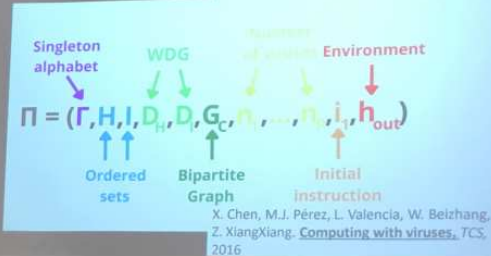
CV: Presentations

Conferences

- IWINAC 2022
Spain
- CMC2022
Italy
- (A)CMC2023
China

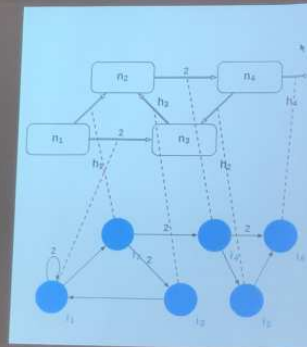
Best Paper Award

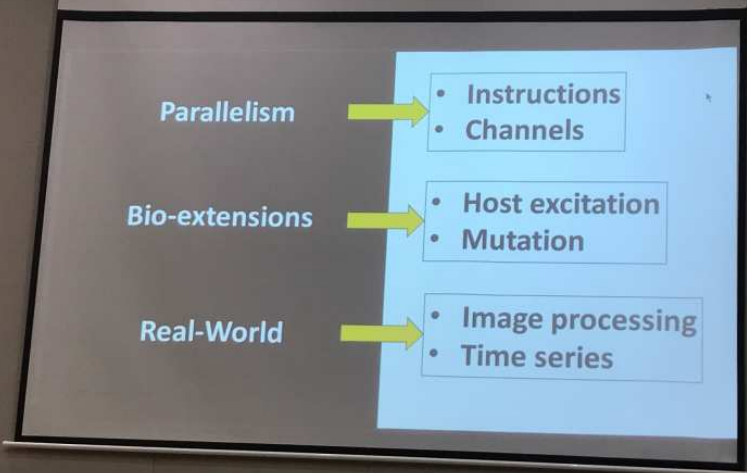
Virus Machine of degree (p,q)



Example

A Virus Machine of degree (4,6)





2 My Research Contents: 研究背景

脉冲神经P系统 (Spiking neural P systems, SNP systems)

脉冲神经P系统是神经型膜计算模型的主要形式，是一类受生物神经元以脉冲的方式传递信息的机制启发设计的神经计算模型，抽象于生物神经元发放脉冲的现象。

2 My Research Contents: 研究动机

Peng et al. Nonlinear Spiking Neural P Systems, International Journal of Neural Systems, 2020.

非线性脉冲神经P系统 (Nonlinear spiking neural P systems, NSNP systems)
(为脉冲神经P系统引入非线性机制构建了非线性脉冲神经P系统)

Motivator: 利用SNP systems的机制开发一个循环神经网络(RNN)模型

从NSNP systems出发构建这个RNN模型

```

    graph LR
      NSNP([NSNP]) --> RNN([RNN])
      RNN --> Emotion([情感分类])
  
```

RNN – Recurrent neural networks (循环神经网络)



一个度为 $m \geq 1$ 的非线性脉冲神经网络系统 (简称NSNP系统) 的定义为:

$$\Pi = (O, \sigma_1, \sigma_2, \dots, \sigma_m, \text{syn}, \text{in}, \text{out})$$

- (1) $O = \{a\}$ 表示一个单字母表 (a 称为脉冲);
- (2) $\sigma_1, \sigma_2, \dots, \sigma_m$ 是 m 个神经元, 形式为:

$$\sigma_i = (u_i, R_i), 1 \leq i \leq m$$

其中:

- (a) $u_i \in R^+$ 是神经元 σ_i 所含脉冲的初始值, 表示神经元 σ_i 的初始状态;
- (b) R_i 是脉冲规则的有限集合, 形式为 $T \cap a^{\theta(u_i)} \rightarrow a^{f(u_i)}$; d , 其中 $T \in R^+$ 是点火阈值, $g(u_i)$ 是线性或非线性的函数, $f(u_i)$ 是非线性函数, $g(u_i) \geq f(u_i) \geq 0$, 并且 $d \geq 0$;
- (3) $\text{syn} \subseteq \{1, 2, \dots, m\} \times \{1, 2, \dots, m\}$, 并且对于所有 $(i, j) \in \text{syn}$ 有 $i \neq j$, 其中 $1 \leq i, j \leq m$;
- (4) in, out 分别表示了系统的输入神经元和输出神经元。

彭宏教授

罗晓晖教授

西华大学

先荣豪

袁艳萍

图像

MAP

“非线性脉冲P系统”

演化计算方法

(DE PSO等)

0825 卷春波 cb.lou@siat.ac.cn

中科院深圳先进院合成所

物理 (北大本博 理论生物物理中心 (QB))

饶议、施一公 Bio2000

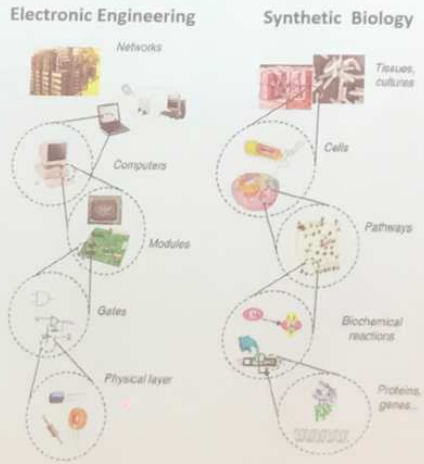
很好糊复合管 } 清华 暑期学校

设计

↕ 在电子领域可分开

制造

合成生物学理想的研究方式

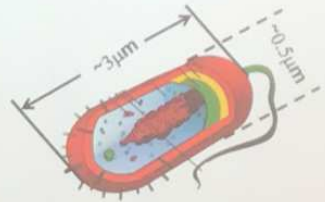
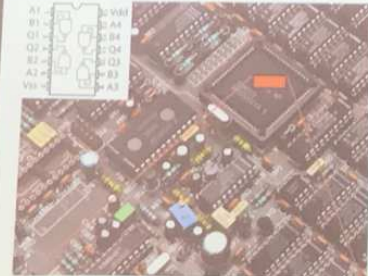


参数化的电子元件
 电阻: (R)
 电感: (L)
 电容: (C)
 晶体管:
 ($\beta, P_{CM}, F_D, V_{O}, I_c$)

Adrianantoandro et al. Mol Sys Biol. 2006

经典的图

生物系统的复杂性



The key cytoplasmic factors in E.coli

Name	Diameter	Number
Chromosome	~300nm	1-2
RNA	~20nm	1,000-10,000
Polemerase		
Ribosome	~40nm	20,000-70,000
Protein	~2nm	2,400,000
mRNA	~10nm	6,000
tRNA	~5nm	60,000-600,000
ATP	~5nm	10,000,000,000



10亿-百亿的体量

两个基因元件设计的核心科学问题：

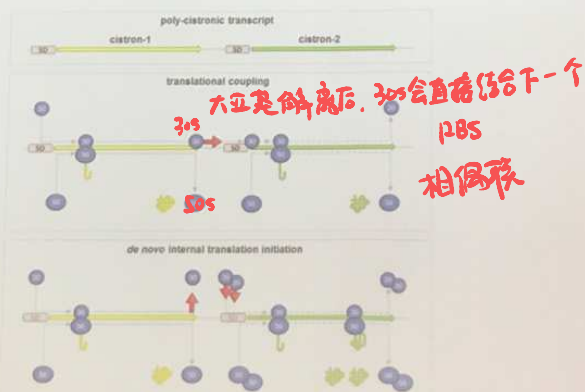
- 1, 模块化：解决顺式干涉作用
- 2, 正交化：解决反式干涉作用

复杂的体系： $\xrightarrow{\hspace{2cm}}$ 模块化
解偶合、绝缘化

SD序列 (RBS)

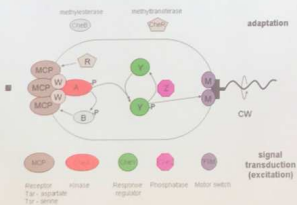
生物复杂性实例一：翻译偶联

► Translational Coupling in prokaryotic operon

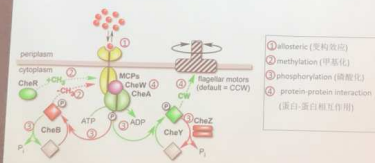


实例二：细菌趋化效应

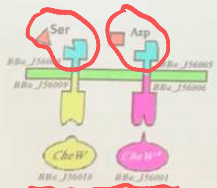
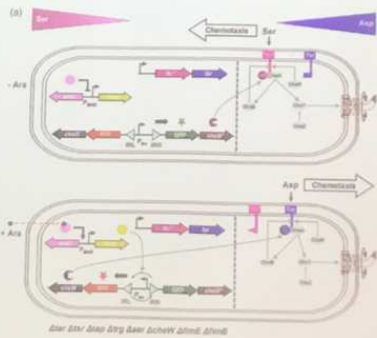
一个大肠杆菌的记忆、思考和追求——Tu Yuhai (IBM)
“Nano-brain”——Thorsten Mascher (TU, Germany)



实例二：细菌趋化效应



2006年UCSF的iGEM项目



Mol. Biol. (2011), 406, 215-227

工程化改造 cheW 的感应接收模块
转化中的 Sensor

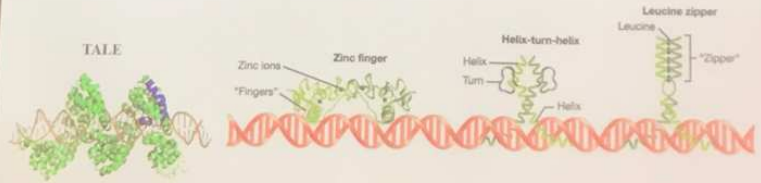
20年前的工作

基因元件的设计原则

设计原则

➢ 简化原则（模块化）

{以转录因子(TF)与operators相互识别为例}



基因元件的设计原则 (LacI)



蛋白质-DNA识别“密码”



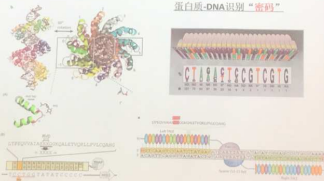
2007年，中科院iGEM项目

Zhan, J., Ding, S., Ma, R., Ma, X., Gu, X., Zhou, Y., & Liu, H. (2007). Molecular systems biology. 4(1), 204.

一个aa对应多个
碱基

Zinc Finger
3个半碱基

基因元件的设计原则 (TALE)



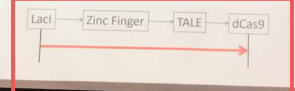
蛋白质-DNA识别“密码”

一个aa \Leftrightarrow 一个碱基

基因元件的设计原则

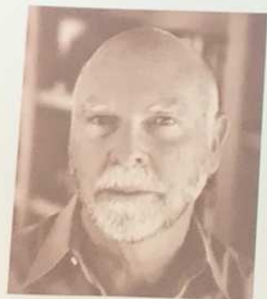


设计原则：简化原则（模块化）

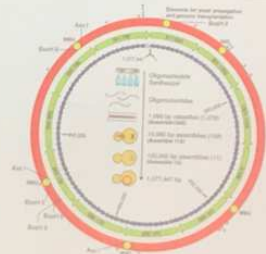
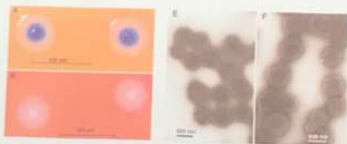


逐步模块化

合成生物学发展的关键问题



J. Craig Venter



四种基本密码子(A, T, C, G)

科学疯子

2008年 合成第一个基因组 2010年文章

14年采访: 合成生物学最大的挑战?

不仅是DNA合成价格, 更是在不清楚功能时做设计.

DOI: 10.1089

本课程的主要内容:

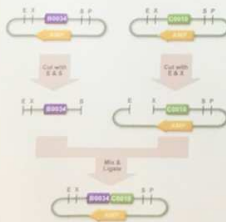
- 1, iGEM等工程化标准.. "Tom Nedd" 计算机
- 2, DNA组装技术 "Jun Andy" 数学
- 3, 基因调控线路设计
- 4, 次级代谢基因簇的设计
- 5, DNA折纸的设计
- 6, 信号转导网络的设计

BioBricks-I



Synthetic Biology
based on standard parts

Biobrick assembly standard

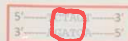


BioBrick standard assembly prefix and suffix sequences

Prefix: 5' GAATTC GCGGCCGC T TCTAGAG 3'
EcoRI NotI XbaI

Suffix: 5' T ACTAGT A GCGGCCGC CTGCAG 3'
SpeI NotI PstI

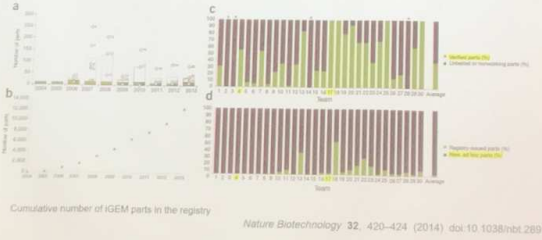
XbaI/SpeI scar



切割后末端相同

Registry of Standard Biological Parts

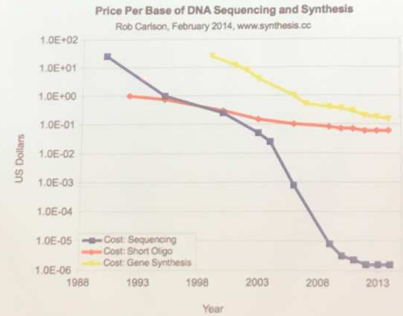
<http://parts.igem.org/Catalog>



DNA组装技术与基因组合成

Method	Mechanism	PCR-free	Scar-free	parallelity	Size (bp)
BioBricks	RE	Yes	NO	single	10-10k
BglBricks	RE	Yes	NO	single	10-10k
Golden Gate	Type IIs RE	Yes	NO	multiple	10-10k
Infusion	Overlap	NO	Yes	multiple	10-10k
SLIC	Overlap	NO	Yes	multiple	1k-10k
USER	Overlap	NO	Yes	multiple	10-10k
Gibson	Overlap	NO	Yes	multiple	1-500k
RedET	Recombination	NO	Yes	single	1k-250k
Yeast TAR	Recombination	NO	Yes	multiple	1k-1M
OLMA	Type IIs RE / others	YES	YES	multiple	10-250k

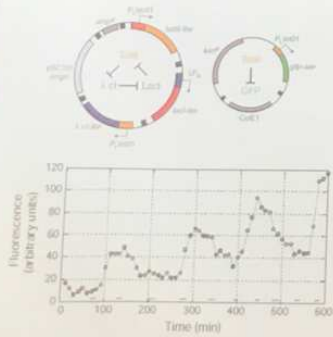
基因“读”和“写”的价格趋势



gBlock 是一个基因合成的公司

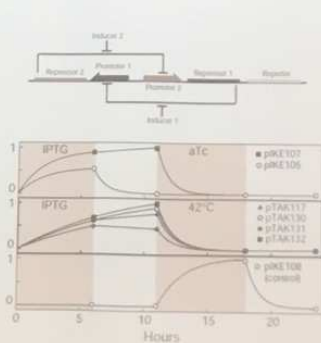
基因线路的标志性工作

Repressilator (振荡器)



Elowitz & Leibler Nature (2000) 403:335-338

Toggle switch (双稳开关)



Gardner & Collins Nature (2000) 403:339-342

实例1: 生物振荡器设计

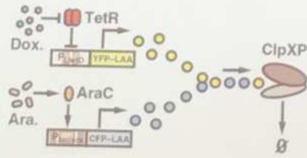


Elowitz and Leibler, Nature, 2000, 403:335-338

Elowitz & Leibler Nature (2000) 403:335-338

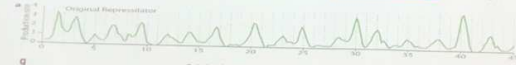
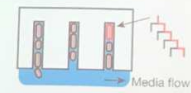
ssrA tag的解耦合与“未解耦合”

Correlated signaling through coupled degradation



Cookson et al., Mol. Syst. Biol., 2011, 7:1-11

2016年改进



Original repressor (INDL332)



Integrated repressor (LPT25)



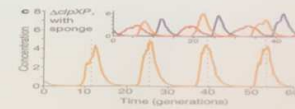
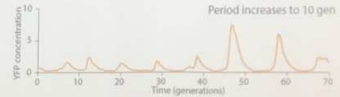
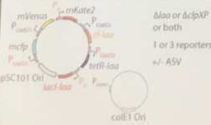
Integrated repressor in $\Delta clpXP$ with P_{ssrA} -mCherry-ssv (LPT5-4)

Time →
Potvin-Trottier, Laurent, et al. Nature 538 7626 (2016): 514-517.

Removing degradation



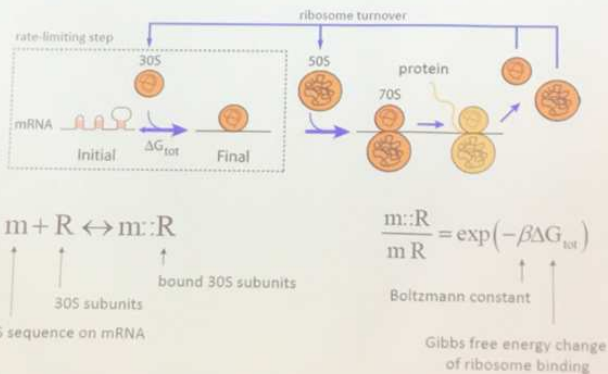
Without degradation but with titration sponge



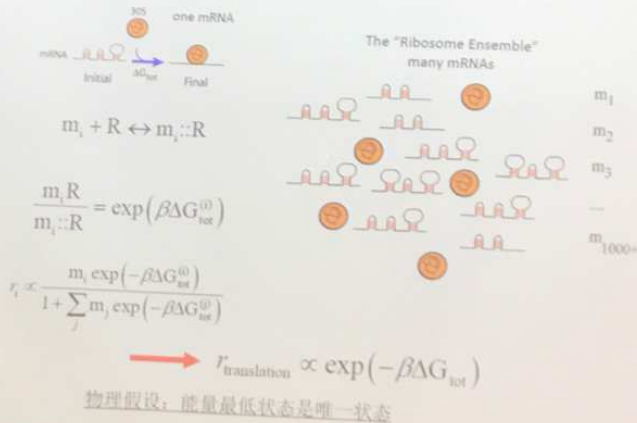
Potvin-Trottier, Laurent, et al. Nature 538.7626 (2016): 514-517.

RBS 计算器的生物物理模型

Translation is a multi-step process
Translation initiation is often the rate-limiting step



RBS 计算器的热力学模型 (续)



150+ Predictions Tested

RBS Calculator_{v1.1}

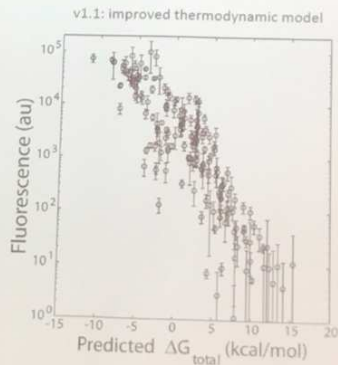
tunable control of the translation initiation rate

$$r_{translation} \propto \exp(-\beta \Delta G_{tot})$$

Rational control of translation initiation over a 100,000-fold scale

$R^2 = 0.80$

$$\beta^{-1} = 0.45 \pm 0.05 \text{ mol/kcal}$$



Salis, Methods in Enzymology, 2011
Salis, Mirsky, & Voigt, Nature Biotechnology, 2009

正向设计的RBS元件更准确

Systematic metabolic pathway optimization requires three ingredients:

1. the ability to quantitatively control enzyme expression

We've developed a way to control bacterial enzyme expression across a 100,000-fold scale

RBS Calculator_{v1.1}

rational control over the translation initiation rate

ribosome binding site protein coding sequence

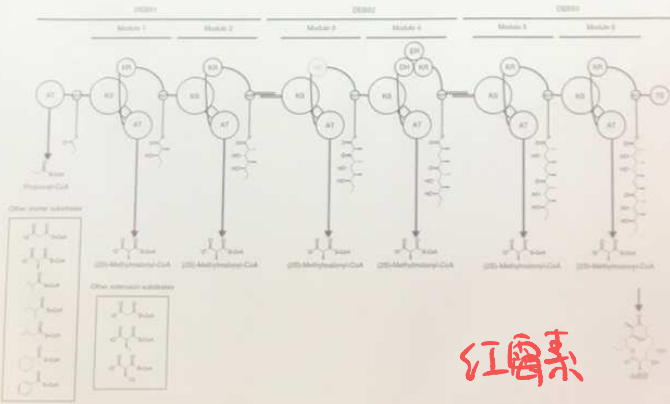


1064 au (proportional scale from 1 to 100,000)

Salis, Methods in Enzymology, 2011
Salis, Mirsky, & Voigt, Nature Biotechnology, 2009

为教不多的
靠谱的
预测软件

Polyketide synthase (PKS)

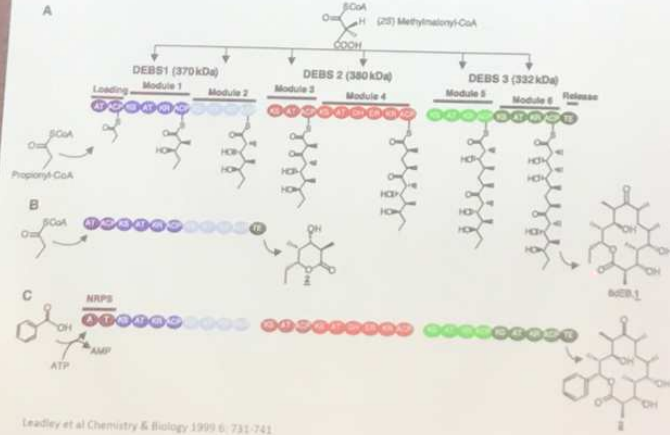


Leadley et al Chemistry & Biology 1999 6: 731-741

红霉素

大多数不成功

成功的例子



Leadley et al Chemistry & Biology 1999 6: 731-741

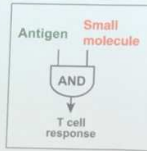
Building 2D structures with DNA bricks

Success rate of design:
 > 94%
 > 107 out of 114

Wen Yan and Feng He, Nature 2012, 485: 623-626



Logic gates for Cancer therapy



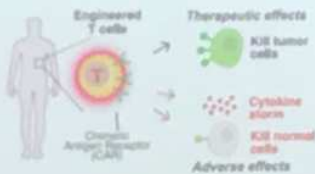
7.5 亿

Wu C, Sandoz A, Faganini T, M. Gaudin, A. S. Kim, W. A. Lim, Science 2008, 320: 100-103

魏泽西

CAR-T 免疫风暴

Logic gates for Cancer therapy



Wendell Lim (UCSF)

Conventional CAR design



Split CAR design



Wu, C. Y., Boydel, K. L., Paulsen, E. M., Ouyang, J., & Lim, W. A. (2015). Remote control of therapeutic T cells through a small molecule-gated chimeric receptor. *Science*, 350(6258), aab4077.

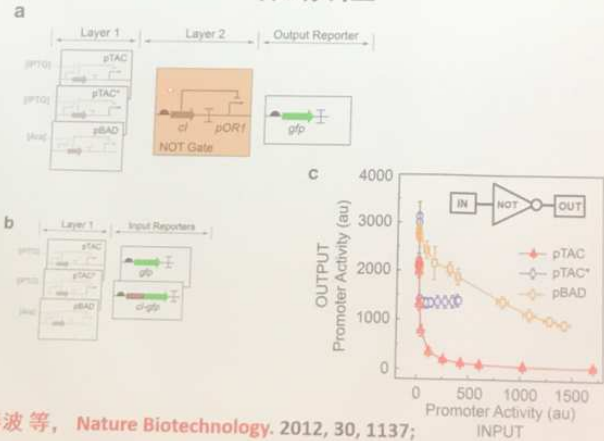
NEXT

老师自己的工作

五类人工调控元件的设计原则

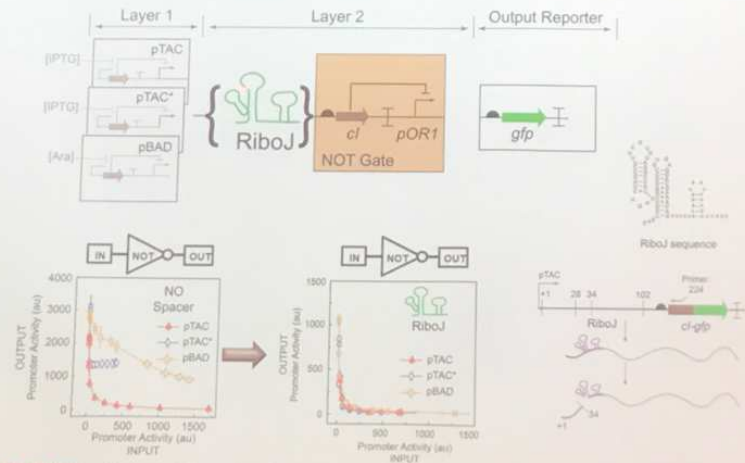
1. 启动子与RBS的模块化设计
2. 启动子与操作子的模块化设计
3. 协同性转录因子的模块化设计
4. 人工翻译激活元件的正交化设计
5. 人工细胞通讯系统的正交化设计

生物信号“转换函数”的周边序列依赖性



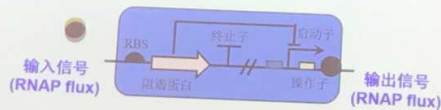
姜春波等, *Nature Biotechnology*. 2012, 30, 1137;

Predictable NOT gate transfer function



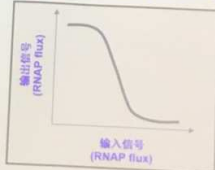
姜春波等, *Nature Biotechnology*. 2012, 30, 1137;

唯一信号载体: RNA聚合酶“流” (RNAP flux)



(生物元件的模型化、参数化)

- ✓ 可定量
- ✓ 可预测
- ✓ 可微调



细胞重编程的前提条件



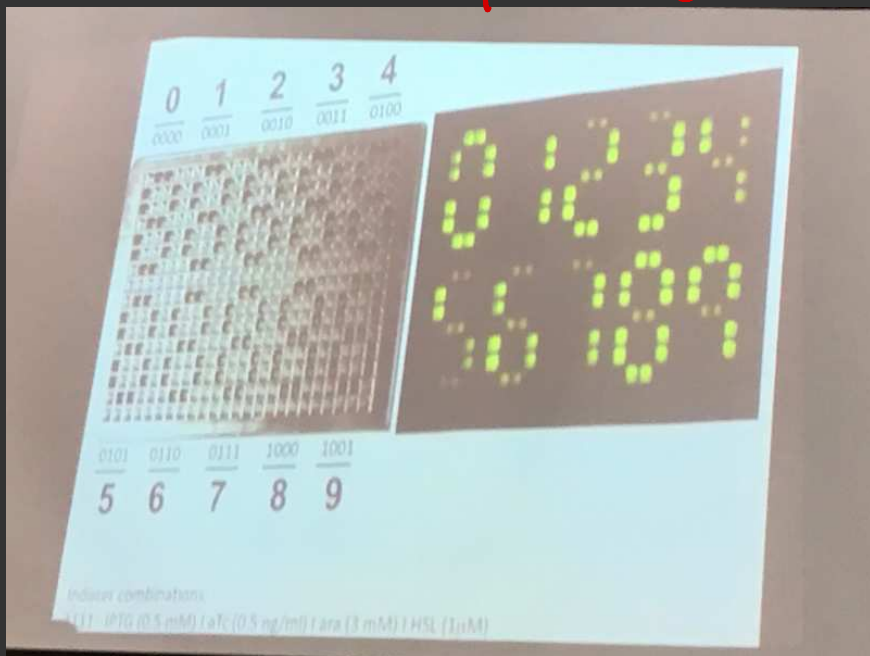
Nielsen, A. A., Der, B. S., Shi, J., & Voigt, C. A. (2016). Genetic circuit design automation. *Science*, 352(6281).

Voigt等, *Science*. 2016, 352, 6281;

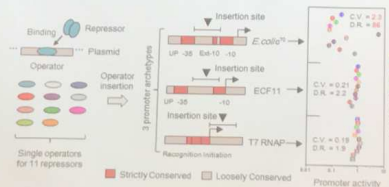
cello 公司
 ~\$2.6亿 } 融资
 ~¥10+亿 }

7段式数字显示

Q: 单cell or 多cell?

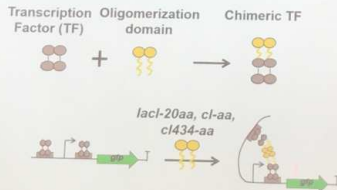


启动子“核心区”模块化



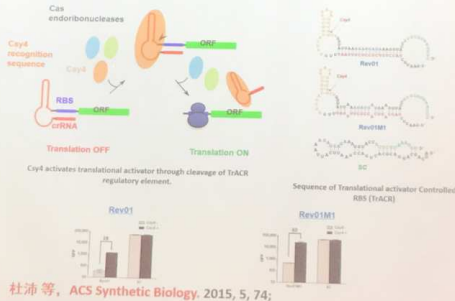
Zong, ... Lou[†], Nature Communications, 2017, 8, 8101

协同性转录因子的模块化设计



侯君然等, ACS Synthetic Biology, 2017, 7, 1188;

基于Cas6 RNA内切酶的人工翻译激活元件

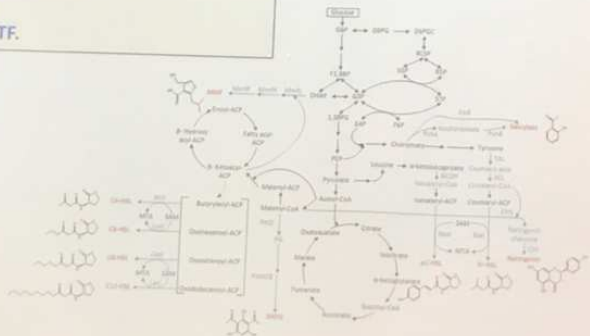


杜沛等, ACS Synthetic Biology, 2015, 5, 74;

更多正交型细胞间通讯系统

设计原则:

- diffusible molecules across membrane
- Synthesized from CCM
- <5 enzymes
- sensed by a TF.



元器件设计原则：正交化

- 已完成的元器件：
细胞通讯系统 (10对)；Csy4翻译激活元件 (3对)；T7-family RNAP (6个)；LacI阻遏元件 (5对)；TALE (26个)；TetR家族阻遏元件 (20对)；ECF sigma 激活元件 (20对)；sRNA (14对)；toehold-RNA (26对)；CRISPRi (任意多)，NF-kB (6对) 等等

- 尚待完成的元器件：
原核调控元件；双组分调控元件；磷酸化修饰元件；翻译机器元件；蛋白降解元件；DNA复制调控元件；感应蛋白元件；能量供应元件；RNA结合蛋白；RNA降解元件；...

真核调控元件 (几乎是空白!!)：启动子；激活因子；增强因子；泛素化修饰因子；磷酸化修饰元件；蛋白降解元件；RNA剪切元件；RNA结合元件；...

Ribo profiling

翻译速率 \rightarrow 表达量

复合物中的蛋白比例固定

理想的务实主义者.

Science partner 期刊 \leftarrow

之江实验室 ZHEJIANG LAB Intelligent Computing A SCIENCE PARTNER JOURNAL

Special Issue: Biological Computation

Call for Papers

Submission Deadline: December 15, 2023

Guest Editors:
Jonathan Cooper, University of Glasgow, UK
Angel Goffi-Moreno, Technical University of Madrid, Spain
Tom de Groot, Eindhoven University of Technology, Netherlands
Hanhuan Peng, SEU-ALLEN Jovel Center, China
Laura Na Uhi, University of Stuttgart, Germany

This special issue solicits original research, as well as review articles. Topics of interest include, but are not limited to:

- Cellular computing
- Cell-free computing systems
- DNA computing
- DNA data storage
- Physical computing
- Morphological computation
- Applications of biological computation
- Unconventional computing
- Synthetic biology tools to assist biocomputing

扫码了解详情或添加微信咨询

230825

宋理富

DNA 信息存储技术原理和应用

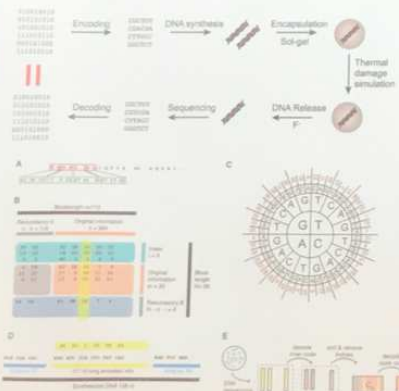
研究经历

年份	学历	单位	团队
2000-2007	本硕	山东大学生命科学院	许平教授
2007-2008	工作	华大基因BG1	
2008-2012	工作	天津工业生物技术研究所	
2012-2018	博士	德国汉堡工业大学	曾安平院士
2018-2019	工作	北京化工大学	曾安平院士
2019-2022	工作	天津大学	元英进院士
2022-至今	工作	天津工业生物技术研究所	系统生物技术中心

基因存储很厉害的单位

DNA 存储优势: 处理并行度高, 常温保存、功耗低、时间久
电流感性不敏感 → 数据安全度高
三维存储、信息密度好.

DNA信息存储研究进展介绍 3



Robust Chemical Preservation of Digital Information on DNA in Silica with Error-Correcting Codes*
 Arthur W. Coats¹, Richard Brinkley¹, Michelle Publico, David Phoenix, and Nicholas J. Leadley¹

Long-Term Memory

Abstract: Information and associated metadata can be preserved on DNA for centuries, but the ability to retrieve this information is limited by the degradation of the DNA over time. Here, we demonstrate that the use of error-correcting codes (ECCs) can significantly extend the lifetime of DNA-based data storage by protecting against errors caused by DNA damage. We show that ECCs can be used to protect digital information on DNA in silica, a robust storage medium, against errors caused by DNA damage. We demonstrate that ECCs can be used to protect digital information on DNA in silica, a robust storage medium, against errors caused by DNA damage. We demonstrate that ECCs can be used to protect digital information on DNA in silica, a robust storage medium, against errors caused by DNA damage.

DNA硅珠保护
 9.4度 2000年数据可靠

内码
 外码

Grass et al. 2015 ETH Zurich

DNA信息存储研究进展介绍 4

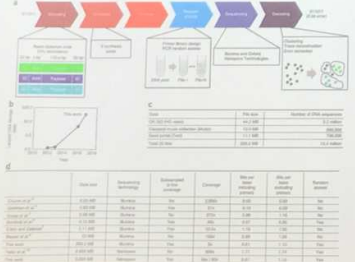
逻辑密度 1.57 bits/bp
 存储物理密度达到215PB/g

Erlich et al. 2017 Columbia University

喷泉码

DNA信息存储研究进展介绍 5

200 MB 数据规模随机访问



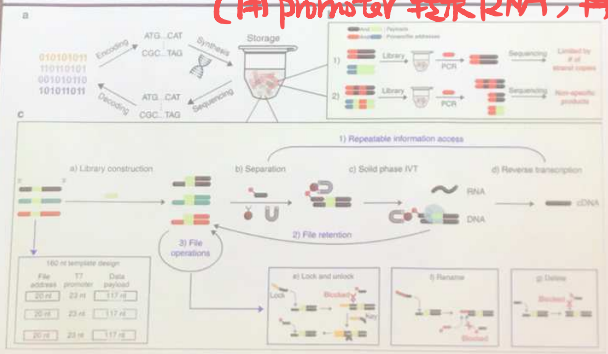
数据规模↑
 随机访问(引物)

Figure 1. Access to the DNA data storage medium and access time. The proposed system stores digital files (100 to 200 MB) on DNA. The resulting information can be retrieved by random access to individual data units. The resulting information can be retrieved by random access to individual data units. The resulting information can be retrieved by random access to individual data units.

Organick et al. 2018 Nature Biotechnology

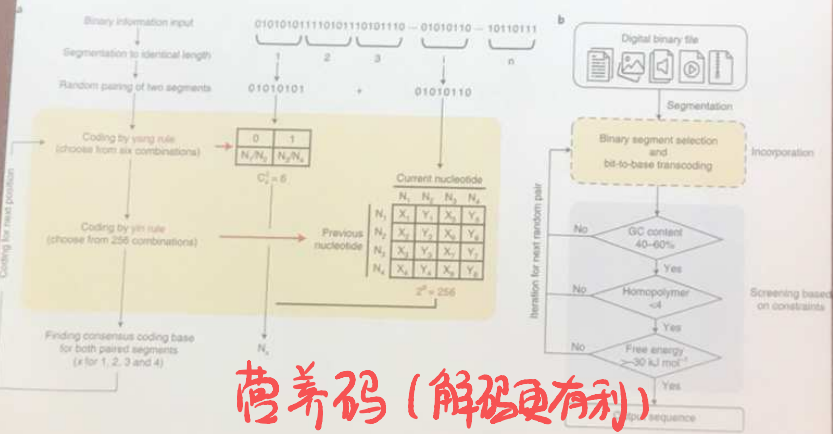
DNA信息存储研究进展介绍 6

测序时不消耗原DNA
 (A promoter 转录 RNA, 再cDNA后测序)



Lin, K.N., Vohler, K., Tuck, J.M. et al. Dynamic and scalable DNA-based information storage. *Nat Commun* 11, 2181 (2020). <https://doi.org/10.1038/s41467-020-16797-2>

DNA信息存储研究进展介绍 7



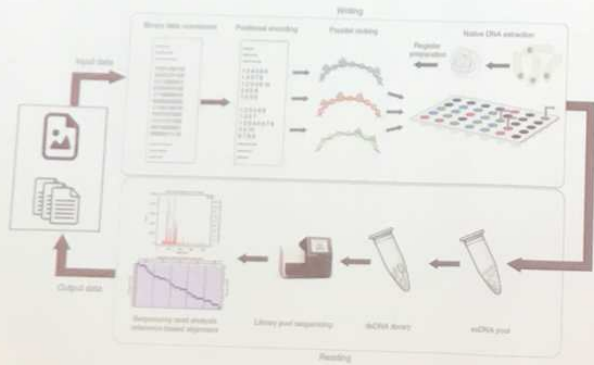
营养码 (解码更有利)

Ping, Z., Chen, S., Zhou, G. et al. Towards practical and robust DNA-based data archiving using the yin-yang coding system. *Nat Comput Sci* 2, 234-242 (2022). <https://doi.org/10.1038/s43588-021-00231-7>

喷泉码 解码依赖所有块

DNA信息存储研究进展介绍 8

DNA打孔卡技术

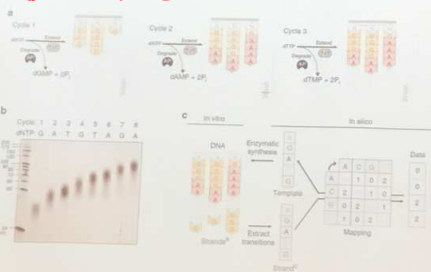


Tabatabaei, S. Kasra et al. (2020): DNA punch cards for storing data on native DNA sequences via enzymatic nicking. *Nature communications* 11 (1), p. 1742.

DNA信息存储研究进展介绍 9

非阻断型TdT酶法DNA合成的DNA信息存储

不好说



DNA信息存储研究进展介绍 10

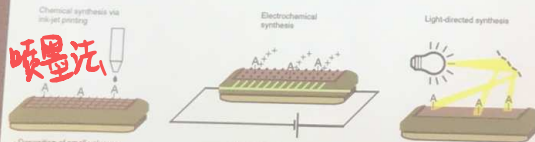
错误率高 通量低

高错误率的DNA合成技术用于DNA信息存储

喷泉码

化学

光控

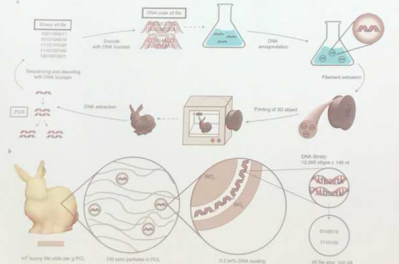


错误率低, 成本高

Grass, Robert N., et al. (2020): Low cost DNA data storage using photolithographic synthesis and advanced information reconstruction and error correction. *Nature communications* 11 (1), p. 5345.

DNA信息存储研究进展介绍 11

万物DNA存储



A DNA-of-things storage architecture to create materials with embedded memory.

DNA信息存储研究进展介绍 12

“人工染色体数据存储的关键概念验证”

-帝国理工 Tom Ellis



提出并验证了“人工染色体如何以一种稳定且廉价复制的方式用于数据存储

帝国理工 Tom Ellis "Self-replicating digital data storage with synthetic chromosomes" (2021, Rev 2021; 8: mwab086)

An artificial chromosome for data storage

Wenqiang Chen^{1,2}, Mengzhe Han^{1,2,3}, Jianping Zhou^{1,2,3}, Qi Ge¹, Panpan Wang¹,
Jinchun Zhang^{1,2}, Siyu Zhu^{1,2}, Lifu Song^{1,2} and **Tom Ellis**^{1,2,3*}

细胞膜阻碍
读写

进展介绍 13

清华大学

镜像 DNA 存储
天然中无降解酶 更稳定

DNA链式存储信道的错误类型重新定义



Original DNA sequence

G-T-A-A-G-A-C-G-T-C-T

Type 1. Substitution A to G

S1-G-T-A-A-G-G-C-G-T-C-T

Type 2. Indel insertion

S2-G-T-A-A-T-G-A-C-G-T-C-T

S3-G-T-A-A-A-C-G-T-C-T

G deletion

Type 3. DNA breaks

S4-G-T-A-A G-A-C-G-T-C-T

S5-G-T-A-A-G-A-C-G T-T-C-T

S6-G-T-A-A-G-A C-G-T-C-T

Type 4. DNA rearrangements

S7-G-T-A-A-T-T-C-T

S8-G-T-A-A-G-A-C-G-C-G-T-C-T

S9-G-T-A-A-G-A-G-A-C-G-T-C-T

DNA存储中可能的错误

断裂和重排是DNA链式存储最基本的错误类型

插入删除和替换错误也是特殊的“断裂重排”错误



DNA信息存储信道的多拷贝特征

写入—DNA合成 复制—PCR扩增



有大量尝试将拷贝数将为1的努力

读取—测序读取



多拷贝不可避免



多拷贝特性是关键!



传统基于聚类-多序列比对的内码解码过程



Low cost DNA data storage using photolithographic synthesis and advanced information reconstruction and error correction

无法处理DNA断裂和重排

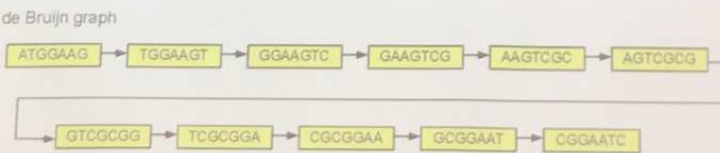
恢复耗时

De Bruijn Graph 德布莱英图的基本原理

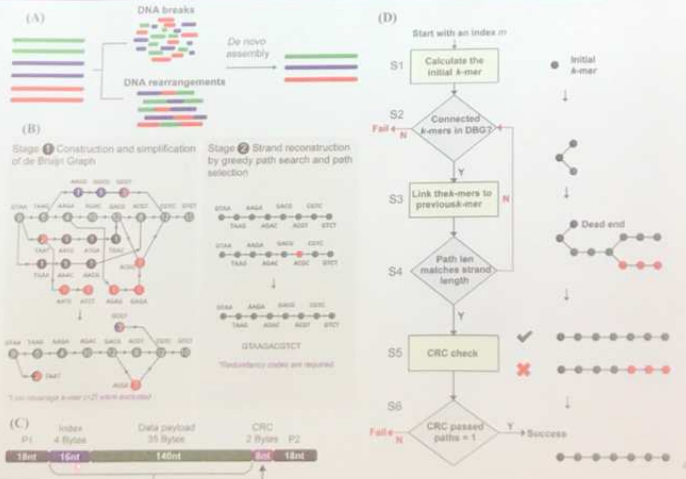
sequence **ATGGAAGTCGCGGAATC**

7mers

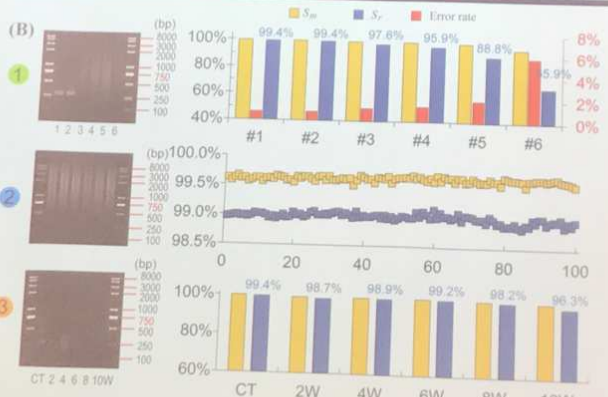
ATGGAAG
 TGGAAGT
 GGAAGTC
 GAAGTCG
 AAGTCGC
 AGTCGCG
 GTCGCGG
 TCGCGGA
 CGCGGAA
 GCGGAAT
 CGGAATC



基于图论的组装内码算法

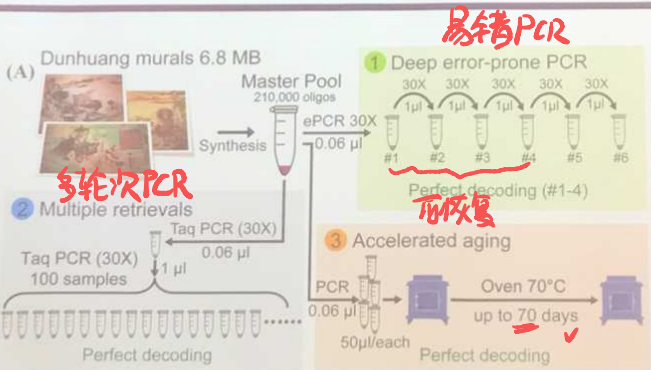


DNA信息存储的数据鲁棒技术

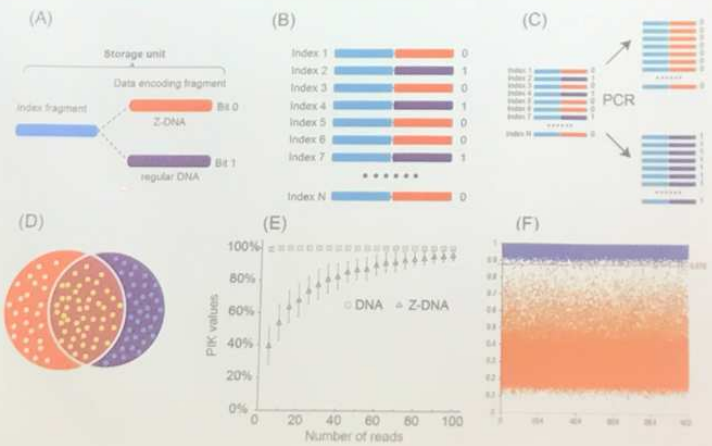


效果

DNA信息存储的数据鲁棒技术

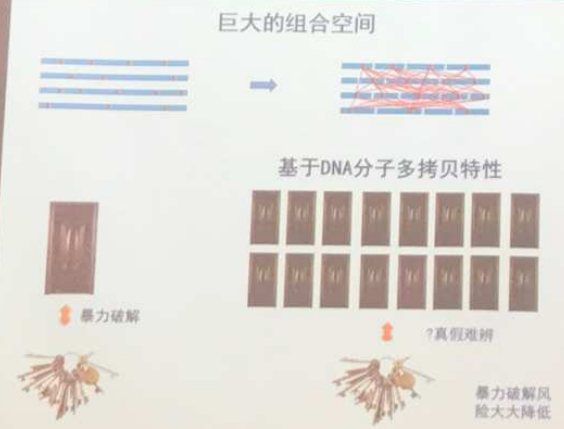


3) 基于Z-DNA的密钥存储技术



Z-DNA与普通DNA结构不同 PCR不可继承该信息

4) “影子保护”技术



5) 数据安全：“影子保护”技术

