

Computational Design of Biological Systems by Automatic Methods



Alfonso JARAMILLO

Synth-Bio group Programme Epigenomique CNRS-Genopole-UEVE & Ecole Polytechnique

http://synth-bio.org



- Molecular Genetics in the post-genomic era
- Design of molecular parts (I): Macromolecules
- Design of molecular parts (II): Networks
- Design of molecular parts (III): Cells

Molecular Genetics in the post-genomic era

- Can we understand complex genetic systems as a combination of molecular parts?
 - Proteins, RNAs, Genetic circuits, Metabolic circuits, Genomes...
- Approach 1: build a list of parts, construct computational models and test their predictions against experimental data.
 Systems Biology
- Approach 2: design, construct & validate synthetic systems from molecular parts
 - Synthetic Biology



Introd

Enabling breakthroughs in a postgenomic era

- Advances in computing power
- Internet
- Genomic sequencing
- Crystal structures of proteins
- High through-put technologies



Productivity improvements in DNA sequencing and synthesis, compared with Moore's law Oct 2002, Log scale











Introd

Understanding complex genetic systems as a combination of molecular parts

- Approach: design, construct & validate synthetic systems from molecular parts
 - Problem: Genetic Engineering has been around for more than 30 years and its technology does not scale with current molecular part lists.
 - Solution: Make biology more engineerable
- How we could facilitate the engineering of genetic systems?
 - By embracing engineering paradigms
 - Abstraction, Modularization and standardization
 - By developing computational design methods that apply our knowledge





Path 1: the construction of engineered DNA, which allows manipulation at every level of the natural hierarchy.
Path 2: the use of engineered DNA to produce novel nanostructures.
Path 3: the development of nonstandard amino acids and base pairs, which can then be assembled into foldamers and DNA analogs.

Path 4: the creation of alternative genetic systems.

Path 5: producing minimal genomes (synthetic chromosomes) and transplanting them into prokaryotic hosts.

Path 6: adding new functions to living organisms by manipulating cell machinery.Path 7: the fusion of proteins to produce assemblies with novel functions.

Path 8: the use of peptide synthesis to create programmable building blocks that can assemble further into functional protein



Design principles of SB

- Decoupling Design & Fabrication
 - Rules insulating design process from details of fabrication
 - Enable parts, device, and system designers to work together
 - VLSI electronics (1970s)
- Abstraction
 - Insulate relevant characteristics from overwhelming detail
 - Simple components that can be used in combination
 - From Physics to Electrical Engineering (1900s)
- Standardization of Components
 - Predictable performance
 - Off-the-shelf
 - Mechanical Engineering (1800s) & the manufacturing revolution (e.g. Henry Ford)

Abstraction levels





Off-the-shelf biological parts and devices







Application: hydrogen production

- Hydrogen is considered the energy carrier of the future
- Use cyanobacteria for photoproduction of hydrogen:
 - Solar energy is inexpensive
 - Production is clean and sustainable

The problem:

- Photosynthesis can produce hydrogen (hydrogenase)
- Photosynthesis produces oxygen
- Oxygen inhibits hydrogen production by hydrogenase!

Options:

- 1) Use resistent hydrogenase (NiFe), less efficient
- 2) Use efficient hydrogense (Fe), fight inhibition







Oxygen consumption





Fighting inhibition



Tune photosynthesis so production and consumption match





BioModularH2 project



- Abstraction
 - -Parts -Devices
- Oxygen consumption and sensing Regulation

H₂ production

Specification

-Systems

- Modularity
- Simulation
- Optimization



Multi-scale computational design

103

235

Macromolecules



 $+\sum_{nonbonded}\frac{q_i q_j}{4\pi Dr_{ij}} + \varepsilon_{ij} \left[\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{1/2} - 2\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{6/2} \right]$

Biological networks

$$\begin{aligned} \frac{d[Y]}{dt} &= \alpha \frac{1}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma \\ \frac{d[Y]}{dt} &= \alpha \frac{\left(\frac{[U]}{K}\right)^n}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma \end{aligned}$$

Genomic background

 $\max \mu = cv_{\frac{dA}{dt} = -v_1 - v_2 + v_3 + b_1}$ st: $Sv = b_{\frac{dB}{dt} = v_1 + v_4 - b_2}$ $v_{\min} \le v \le v_{\max}$





- De novo design of proteins:
- DESIGNER
- PROTDES

De novo design of:

- Transcriptional networks
 - GENETDES
 - ASMPARTS
- RNA networks
 - RNADES
- *De novo* design of metabolic pathways by retro-biosynthesis
 - DESHARKY
- Network inference from microarray & proteomics resp.
 - INFERGENE

Multi-scale computational design

103

235

Macromolecules



 $+\sum_{nonbonded}\frac{q_i q_j}{4\pi Dr_{ij}} + \varepsilon_{ij} \left[\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{1/2} - 2\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{6/2} \right]$

Biological networks

$$\frac{d[Y]}{dt} = \alpha \frac{1}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$
$$\frac{d[Y]}{dt} = \alpha \frac{\left(\frac{[U]}{K}\right)^n}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$

Genomic background

- $\max \mu = cv$
- st: Sv = b
- $v_{\min} \le v \le v_{\max}$





- De novo design of proteins:
- DESIGNER
- PROTDES

De novo design of:

- Transcriptional networks
 - GENETDES
 - ASMPARTS
- RNA networks
 - RNADES
- *De novo* design of metabolic pathways by retro-biosynthesis
 - DESHARKY
- Network inference from microarray & proteomics resp.
 - INFERGENE



How can we design protein structure and function?



Physical model at atomic scale

 $E = \sum_{bonds} K_b (b - b_o)^2 + \sum_{angles} K_{\theta} (\theta - \theta_o)^2 + \sum_{torsions} K \phi (1 + \cos(n\phi - \delta))$ $+ \sum_{impropers} K_{\phi} (\phi - \phi_o)^2 + \sum_{Urey-Bradley} K_{UB} (r_{1,3} - r_{1,3,o}) + \sum_{nonbonded} \frac{q_i q_j}{4\pi D r_{ij}} + \varepsilon_{ij} \left[\left(\frac{R_{\min,ij}}{r_{ij}} \right)^{1/2} - 2 \left(\frac{R_{\min,ij}}{r_{ij}} \right)^{6/2} \right]$









Macromo

Main challenges in protein design

The main challenges in protein design require methodological advances.

- Model unfolded state
 Syst & Synth Biol 2009
- Model folded state

Syst & Synth Biol 2009. PROTDES software

Implicit solvation

Biophys J. 2005

Side-chain and backbone flexibility

Proteins 2009

 Combinatorial explosion J Comput Chem 2008







Macromol

Design & Construction of Parts

Synthetic Protein Scaffolds



Design of a New Fold





Baker's group, (Science 2003)



Designed protein



Blue computationally designed, red x-ray structure

RMSD 1.17A



Redesign of natural protein domains





Macromol

New Molecular Recognition



Design of new sensor proteins

Redesign 5 periplasmic binding proteins (PBP) to bind trinitrotoluene (TNT), L-lactate or serotonin in place of the wild-type sugar or amino-acid ligands Hellinga's group,

Macromol (Nature 2003) Serotonin TNT ∟-Lactate Ribose NHa Tryptamine D-Lactate TNB NHa⁺ -Tryptophan Pyruvate 2.4-DNT 2.6-DNT 28 closed open



Design of new sensor proteins









[Macromo]

RDX biosensor

 RDX: a very common explosive used in landmines.

- Traditionall landmine detection mechanism (metal detection) useless in plastic made ones. Has more than 80% of false positives.
- Use *Pseudomonas putida* as a specific RDX detector.
- Use ribose binding protein and redesign binding site for specificity towards RDX.

Collaboration with Prof. V. de Lorenzo (CNB-CSIC).







Macromo

Design of MHC-I inhibitors

- Find sequences of 9 residues long binding to MHC- I
- Minimize the binding energy between the MHC-I and the peptide



- All with binding
- 131% of reference binding
- Less than 55% identity with known peptides
- 3 peptides recognized by the TCR





Computational Redesign of Endonuclease





Crystal Structure of the DES Enzyme-DNA Complex



Electro-density map of the redesigned region:

gray: computational designed model

Superposition: salmon: design model cyan: crystal structure



Macromol

De Novo Design of Novel Enzymes



Macromol

Catalytic site design






Catalytic site design

• Find sequence that both folds AND has activity: two-objective problem





Aacromo]

Experimental validation of min. energy designs

- Thermostable choristmate mutase with restricted AI/KE alphabet
 Coll. Profs. Hilvert (ETH-Zurich), Wodak (Toronto) & Karplus (Harvard &ISIS)
- Design of 10 peptide sequences of 9 residues long binding to MHC- I

J Biol Chem 2003





Redesign thioredoxin by grafting esterase activity on p-nitrophenyl acetate while preserving original function. Promiscuous enzyme design.

Coll. Prof. Sánchez-Ruiz (Granada, Spain)





Aacromo]

Experimental validation of min. energy designs

- Thermostable choristmate mutase with restricted AI/KE alphabet
 Coll. Profs. Hilvert (ETH-Zurich), Wodak (Toronto) & Karplus (Harvard &ISIS)
- Design of 10 peptide sequences of 9 residues long binding to MHC- I

J Biol Chem 2003





Redesign thioredoxin by grafting esterase activity on p-nitrophenyl acetate while preserving original function. Promiscuous enzyme design.

Coll. Prof. Sánchez-Ruiz (Granada, Spain)





Macromo

Enzyme with minimal aminoacid alphabet

Chorismate mutase with restricted AI/KE alphabet?



 Redesign of helix 1 using hydrophobic /hydrophilic patterning.





Macromo

Enzyme with minimal aminoacid alphabet

Chorismate mutase with restricted AI/KE alphabet?



- Redesign of helix 1 using hydrophobic /hydrophilic patterning.
- We did two parallel experiments:
 - Computational protein design—
 - In vivo directed evolution Coll. Profs. Hilvert (ETH-Zurich)



Coll. Wodak (Toronto) & Karplus

(Harvard & ISIS)

Status and

Methodology Computational Protein design

State space: $\{x_i\}$ = Aminoacid side-chain at residue i

Mathematical description: Calculation of folding free energy $exp(-G/kT) = \int d(solute)d(solvent)exp(-E/kT) = \int d(solute)exp(-(E + E_{solv}^{eff})/kT)^{n}$ $M_{ij}(x_{ij}) = \sum_{x_{ij}} exp(-E_{ij}^{pair}(x_{ij},x_{j})/RT) exp(-E_{ij}^{sing}(x_{ij})/RT)\prod_{k\neq j}M_{ki}(x_{ij})$ Belief propagation



Optimisation: Inverse folding problem Heuristic: Monte Carlo Simulated Annealing Exact: Branch & Bound $\{x_i\}$



Challenges:

NP-Hard problem of large size: 10^{200} , new approaches are needed We use a physical model with almost no fitted parameters

Multi-scale computational design

90

Macromolecules



 $+\sum_{nonbounded}\frac{q_i q_j}{4\pi Dr_{ij}} + \varepsilon_{ij} \left[\left(\frac{R_{\min,ij}}{r_{ij}}\frac{1}{J}^2 - 2\left(\frac{R_{\min,ij}}{r_{ij}}\frac{1}{J}^6\right)\right] \right]$

Biological networks

$$\frac{d[Y]}{dt} = \alpha \frac{1}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$
$$\frac{d[Y]}{dt} = \alpha \frac{\left(\frac{[U]}{K}\right)^n}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$

Genomic background

$$\max \mu = cv$$

st :
$$Sv = b$$

$$v_{\min} \le v \le v_{\max}$$





- *De novo* design of proteins:
- DESIGNER
- PROTDES

De novo design of:

- Transcriptional networks
 - GENETDES
 - ASMPARTS
- RNA networks
 - RNADES
- *De novo* design of metabolic pathways by retro-biosynthesis
 - DESHARKY
- Network inference from microarray & proteomics resp.
 - INFERGENE



Networks

Can we design protein networks with targeted behaviour?

























Elowitz & Leibler. 2000. Nature 403:335-8







Networks

Atkinson oscillator

Development of Genetic Circuitry Exhibiting Toggle Switch or Oscillatory Behavior

in Escherichia coli

Mariette R. Atkinson et al. Cell, Vol. 113, 597-607, May 30, 2003,





Hasty oscillator



J Stricker *et al. Nature*, (2008) doi:10.1038/ nature07389





Networks

Can we design such protein networks in an automated way?

Summer of the

Networks

Automatic design gene networks

- We are going to use a coarse-grained description at the protein level
- Focus on transcription regulation
- Combinatorial optimisation



45

5:1111-010 5:1111-010

Networks

Methodology to design gene networks

- We are going to use a coarse-grained description at the protein level
- Focus on transcription regulation
- Combinatorial optimisation





Networks

GENETDES software



GENETDES software

Rodrigo et al. Bioinformatics 2007

Evolve and optimise network using a targeted time-course to construct a **score**





Networks

GENETDES software



GENETDES software

Rodrigo et al. **Bioinformatics 2007**

Evolve and optimise network using a targeted time-course to construct a **score**





Extend to combinatorial assembly of SBML models

48



Networks

Input u1	Input u2	<u>Output y</u>
0	0	0
0	1	0
0	1	0
Ι	0	0
1	1	1

500













Networks

New developments: constructing the circuits by assembling



Workflow



BBa_F2620

30C₆HSL → PoPS Receiver

Mechanism & Function

Input Compatibility*

CHSL 500 ---- C. HSL - 30C HSL 400 C,HSL CHSL 30C,HSL

___C_HSL 200 ---- 0, HSL

Part Compatibility (qualitative)

Nucleotides: 0 / 6.6xNtucleotides cell-1 s-1

Polymerases: 0 / 1.5E-1xNt RNAP cell-1

pSB3K3 and pSB1A2

E0240, E0430 and E0434

(Nt = downstream transcript length)

___ 800

Ě 300-

100

Chassis:

Plasmids:

Devices:

5 0-18

A transcription factor (LuxR) that is active in the presence of a cell-cell signaling molecule (3OC₆HSL) is controlled by a regulated operator (PLtetO-1). Device input is 3OC6HSL. Device output is PoPS from a LuxR-regulated operator. If used in a cell containing TetR then a second input such as aTc can be used to produce a Boolean AND function.





PLMO-1 RBS luxR Term. Plus.B





Dynamic Performance*



Chassis: MG1655 *Equipment: PE Victor3 multi-well fluorimeter **Equipment: BD FACScan cytometer



From biological discovery to an engineered device



The device is re-engineered using standardised biological parts



Asmparts: *in silico* assembly of parts



$$\frac{d}{dt}[mRNA] = POPS - \delta[mRNA]$$

$$\frac{d}{dt}[Protein] = RIPS - \beta[Protein]$$

$$POPS = POPS_0 + \frac{\alpha_0 + \alpha(\frac{[TF]}{K})^n}{1 + (\frac{[TF]}{K})^n}$$

 $RIPS = \lambda[mRNA]$

 $POPS_0 = \eta POPS$



Asmparts: in silico assembly of parts



$$\frac{d}{dt}[mRNA] = POPS - \delta[mRNA]$$

mRNA degradation constant

$$\frac{d}{dt}[Protein] = RIPS - \beta[Protein]$$

transcription rate in presence of TF

protein degradation constant

basal transcription rate

$$POPS = POPS_0 + \frac{\alpha_0 + \alpha(\frac{[TF]}{K})}{1 + (\frac{[TF]}{K})}$$

Hill coefficient

regulatory coefficient

$$RIPS = \lambda mRNA$$
 Ribosome binding constant

$$POPS_0 = \eta POPS$$
 1-transcription termination
efficiency



assembly of standard model parts



Methodology computational gene network design

State space: $\{x_i\}$ = Concentration/number of molecule i



Heuristic: Monte Carlo Simulated Annealing



Challenges:

Avoid solving the dynamics by using analytical approximations for perturbations Adapt it to stochastic processes and discrete events (*e.g.* signalling)

GENETDES software Rodrigo et al. Bioinformatics 2007



Naturally RNA-based gene regulation systems





RNA-based Synthetic Biology




FJ Isaacs et al., Nature Biotechnology, 2004

Computational Design of Riboswitches



We can use combinatorial optimisation to stabilise an unbound active/inactive ribozyme and to destabilise a bound inactive/active conformation. Several logic gates can be created.

Breaker's group 2005



Networks

RNA-based digital devices

A Functional composition of an RNA device



B Signal integration (SI) schemes





Smolke's group (Science 2008)



Networks

RNA-based digital devices



Automatic design with nucleic acids

Example of biological integrated circuits by using RNA



Multi-scale problem:

- Scale 1: Extend Genetdes to use generalised reactions in a modular way (Genetdes++).
- Macroscopic scale, governed by chemical reactions among several species.
- Scale 2: Obtain the nucleotide sequence that will produce a given reaction (RNAdes): Inverse folding problem
- Microscopic, controlled by statistical physics.



Methodology computational RNA network design

State space: Scale 1: $\{x_i\}$ = Concentration of molecule i

Scale 2: $\{x_i\}$ = Nucleotide at residue position i



Challenges:

Kinetic modelling considering secondary structure Generalise inverse folding to "inverse kinetics" problem

State of

Multi-scale computational design

90.

Macromolecules

- $$\begin{split} E &= \sum_{bonds} K_b (b b_o)^2 + \sum_{angles} K_\theta (\theta \theta_o)^2 + \sum_{torsions} K\phi (1 + \cos(n\phi \delta)) \\ &+ \sum_{impropers} K_{\phi} (\varphi \varphi_o)^2 + \sum_{Urey-Bradley} K_{UB} (r_{1,3} r_{1,3,o})^2 + \end{split}$$
- $+\sum_{nonbonded}\frac{q_iq_j}{4\pi Dr_{ij}} + \varepsilon_{ij}\left[\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{12} 2\left(\frac{R_{\min,ij}}{r_{ij}}\right)^{6}\right]$

Biological networks

$$\frac{d[Y]}{dt} = \alpha \frac{1}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$
$$\frac{d[Y]}{dt} = \alpha \frac{\left(\frac{[U]}{K}\right)^n}{1 + \left(\frac{[U]}{K}\right)^n} - \beta[Y] + \gamma$$

Cells: Genomic background max $\mu = cv$ st : Sv = b

 $v_{\min} \leq v \leq v_{\max}$



De novo design of proteins:

- DESIGNER
- PROTDES

De novo design of:

- Transcriptional networks
 - GENETDES
 - ASMPARTS
- RNA networks
 - RNADES
- *De novo* design of metabolic pathways by retro-biosynthesis
 - DESHARKY
- Network inference from microarray & proteomics resp.
 - INFERGENE

Computational designs in a genomic background

We obtained a ODE model for the global transcription network of *E. coli*:



We developed a methodology for the automatic design of metabolic pathways:





In silico genome evolution and design

- Evolution moves:
 - Add/remove TF or enzyme
 - Replace promoter
 - Add/remove operon
 - Modify kinetic parameters
- Biological part models (Asmparts)
- Desharky to move in metabolic space
- Fitness/scoring function:
 - Use chassis model to estimate cell growth
 - Cost/benefit model:

Cells

- Expressing genes is decremental to growth
- Expressing "useful" pathways contributes to growth
- FBA for fast metabolic reactions, ODEs for slow transcriptional ones.





Methodology genome-scale modelling

State space: $\{x_i\}$ = Concentration of metabolite i

 $\{y_i\}$ = Concentration of transcription factor i

Mathematical description:

Where v_i are the cell metabolic fluxes, c their contributions to the growth rate, S the stoichiometry matrix, and b the uptake fluxes

$$\frac{\mathrm{d}}{\mathrm{dt}}y_i = a_i + \sum_{j \in TF} b_{ij}y_j + \sum_{j \in TF} \sum_{k \in TF} b_{ijk}y_jy_k - \delta_i y_i,$$
$$\frac{\mathrm{d}}{\mathrm{dt}}x_i = \sum S_{ij}v_j - b_i$$

Objective function:	$\max \mu = cv$	
Steady state assumption	Subject to $S_{\mathcal{V}} = b$	
	$v_{\min} \le v \le v_{\max}$	

Optimisation:

Exact: Linear Programming $\{v_i\}$ Heuristic: Monte Carlo Simulated Annealing to evolve the ODEs



Challenges:

Couple transcription regulation to metabolic reactions Integrate discrete events (*e.g.* signalling).





Synth-Bio group

- Guillermo Rodrigo PhD Student (co-supervised IBMCP, Spain) (2006)
- Javier Carrera PhD Student (co-supervised IBMCP, Spain) (2007)
- Filipe Pinto PhD Student (co-supervised IBMC, Portugal) (2009)
- Daniel Camsund PhD Student (co-supervised Uppsala, Sweeden) (2009)
- Boris Kirov PhD Student (2008)
- Thomas Landrain PhD Student (2008)
- Bogdana Barlacu Technician (2008)
- Vijai Singh Postdoc (2009)
- Mariel Montesinos
 Administrative assist. & project management (2007)

Recent members:

- Pablo Tortosa EMBO postdoc (2004-2007)
- Maria Suarez Postdoc (2006-2009)

http://synth-bio.org

Postdoc openings!!!

Funding

ATIGE	Genopole/UEVE	2008-20
SynthBioClock	CNRS IPCB	2008
TARPOL	FP7	2008-20
BioModularH2	FP6 NEST	2007-20
Solar ethanol	IFCPAR/CEFIPRA	2008-20
Aide projets EU	Ile-de-France	2008-20
Emergence	FP6 NEST	2006-20
Laccase design	Alliance (Columbia)	2006-20
Glucaric acid p.	MIT-France 76	