

Data Integration and Mining for Synthetic Biology Design

September 2016



Goksel Misirli, and

Jennifer Hallinan, Matthew Pocock, Phillip Lord, James Alastair McLaughlin, Herbert Sauro, and Anil Wipat

**Interdisciplinary Computing and Complex BioSystems (ICOS) research group
School of Computing Science**



Engineering biological systems is challenging

$\begin{matrix} \text{C} & \text{G} \\ \text{A} & \text{T} \end{matrix} \rightarrow \text{NNNNNNNNNN} \rightarrow 4^{10} \sim 1 \text{ million solutions}$

Order and layout

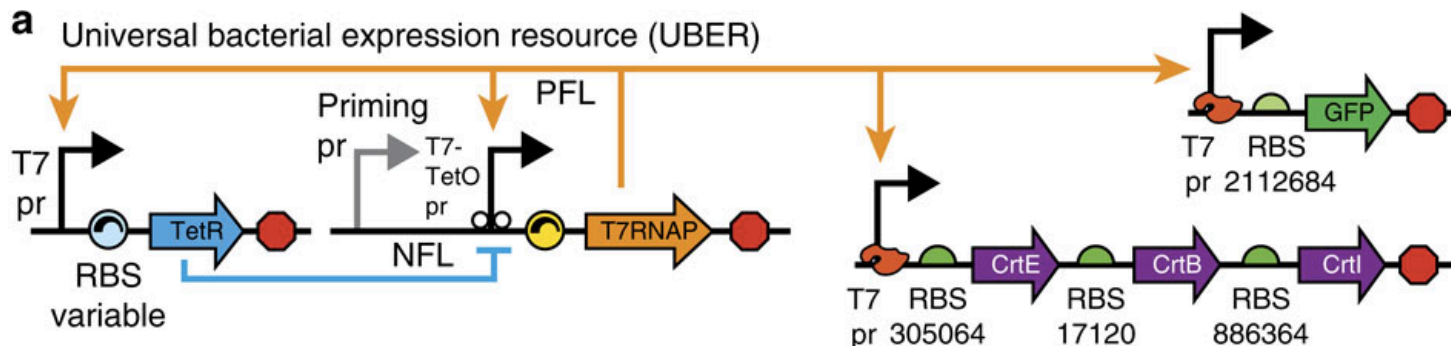
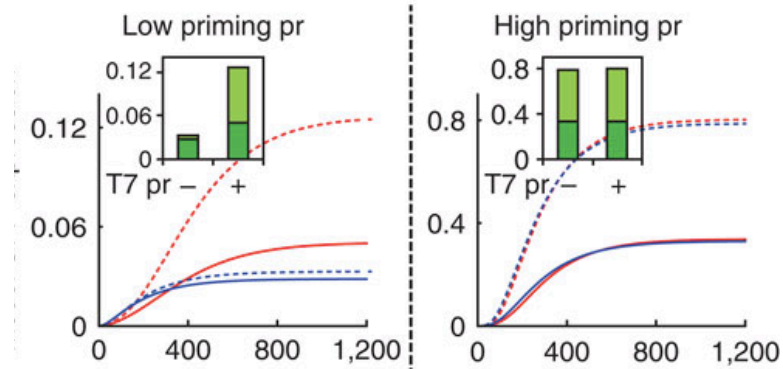
Regulatory elements

Molecular interactions

Strain/Host/Chassis

Biological Context

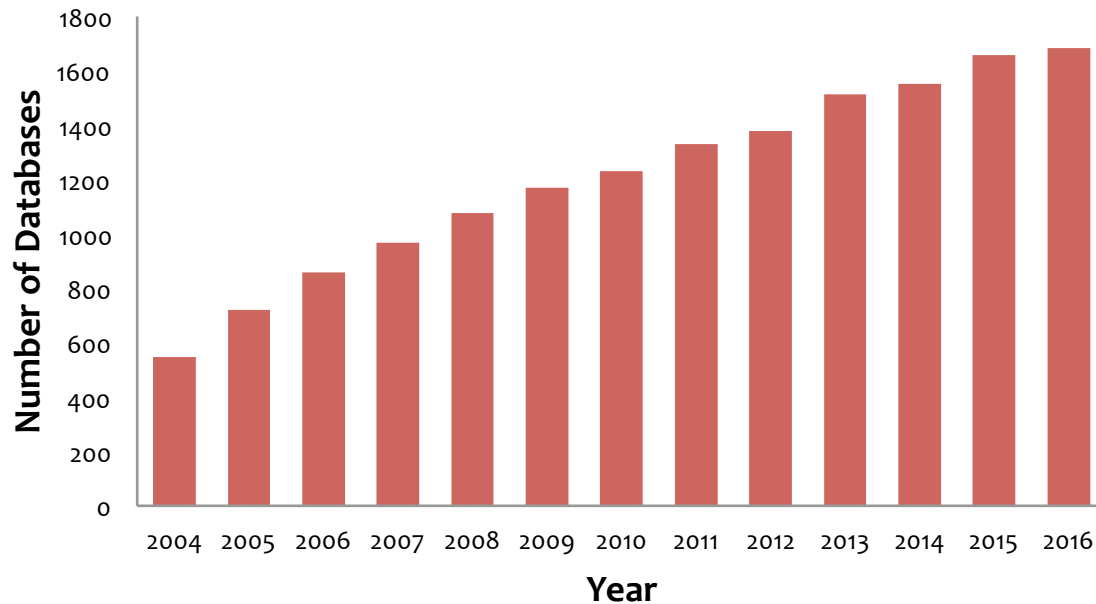
Experimental Conditions



Kushwaha and Salis, 2015

Data integration

- * There is a large amount of information about model organisms such as *B. subtilis* and *E. coli*
- * This information may be spread in
 - * Different databases
 - * Different formats
 - * Different semantics
- * This information can be integrated and used to inform & constrain biological designs



“Knowledge is power”

Marx, Nature, 2013

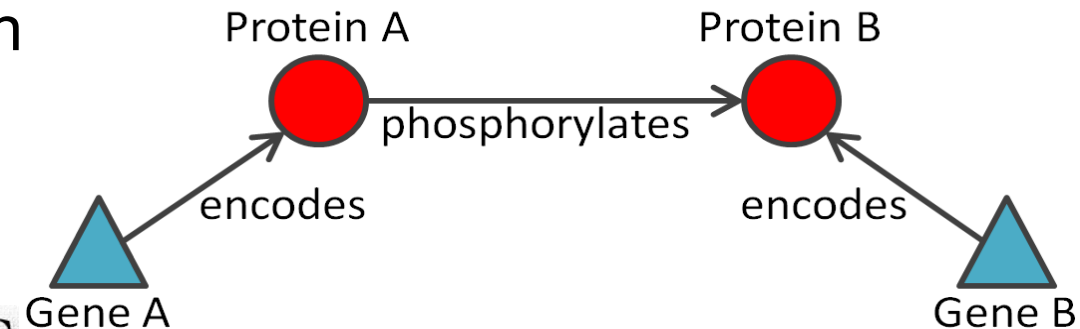
What is an ontology?

The representation of the entities of a domain is consistent and unambiguous

An abstract and simplified view of a domain being modelled

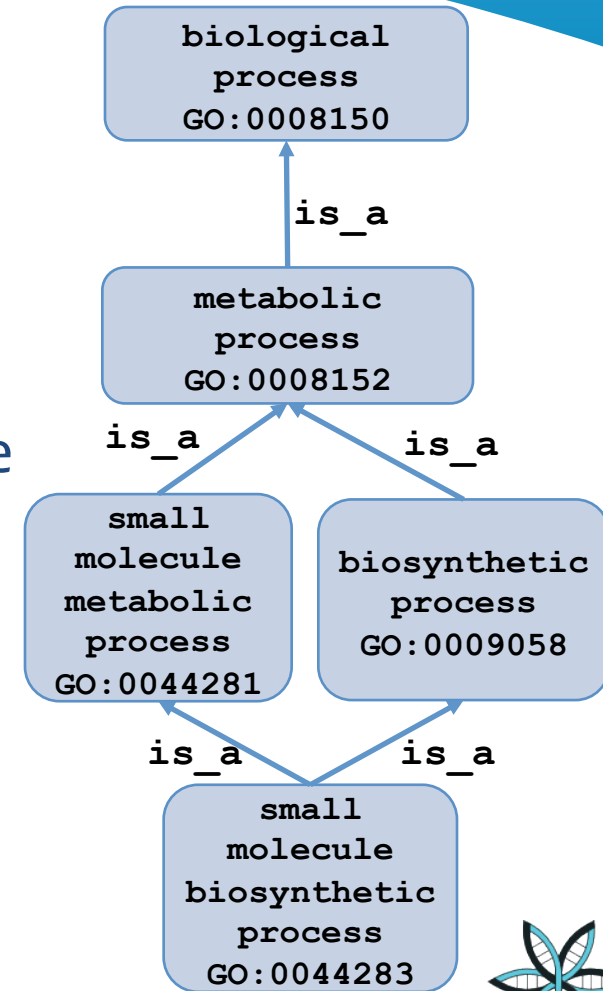
“An explicit and formal specification of a conceptualisation”
Gruber

Shared understanding of a domain



Semantic Web resources for synthetic biology

- * Gene Ontology
- * Sequence Ontology
- * Systems Biology Ontology
- * Synthetic Biology Open Language
- * Standard Biological Parts Knowledgebase
- * SBOL Stack
- * Ontologies are needed to
 - * Capture different relationships between biological parts
 - * Facilitate data mining



Synthetic Biology Ontology (SyBiOnt)

ACS
SyntheticBiology

Research Article

pubs.acs.org/synthbio

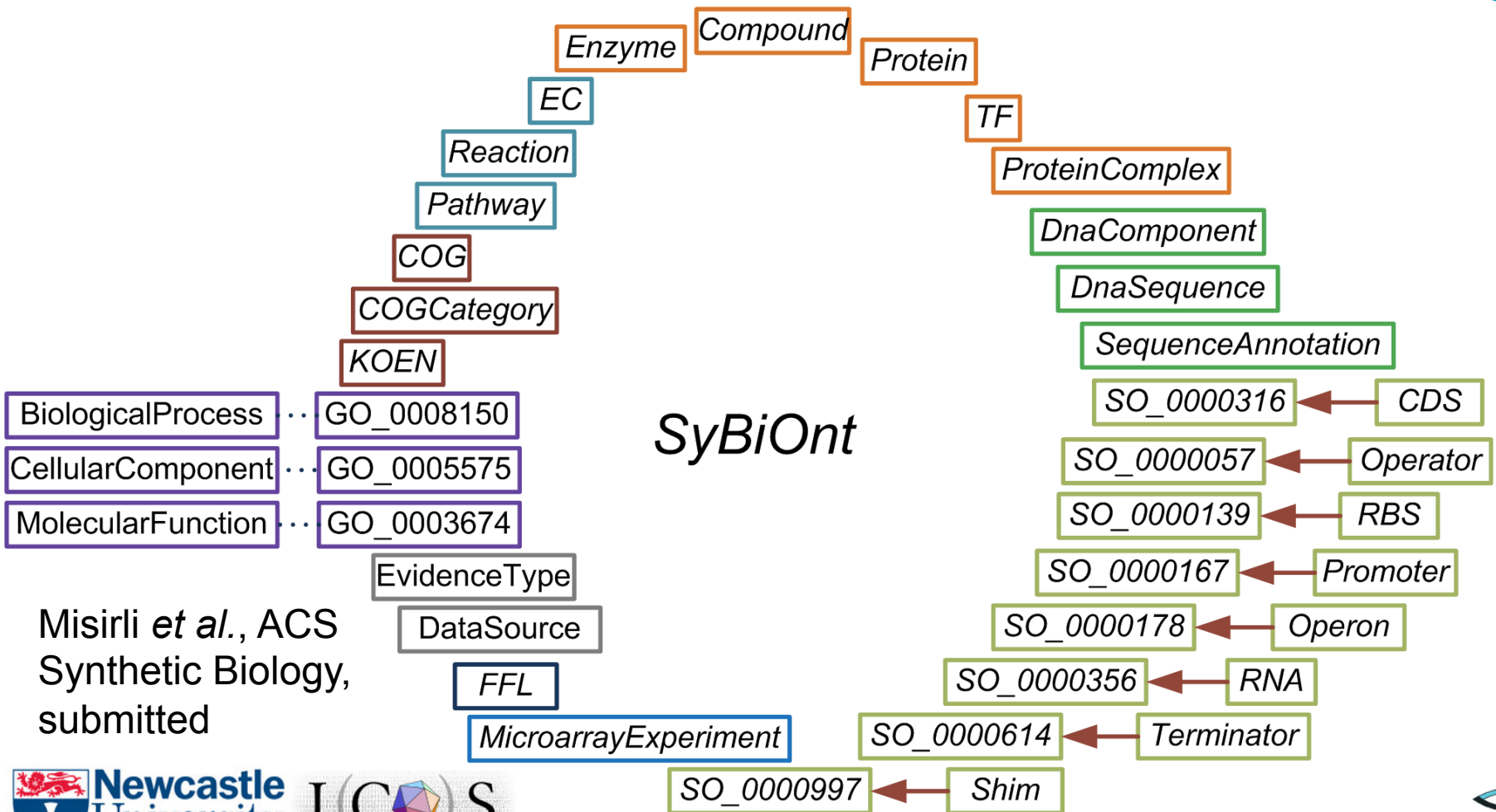
Data Integration and Mining for Synthetic Biology Design

Göksel Mısırlı,[†] Jennifer Hallinan,^{†,||} Matthew Pocock,^{†,‡} Phillip Lord,[†] James Alastair McLaughlin,[†]
Herbert Sauro,[§] and Anil Wipat^{*,†}

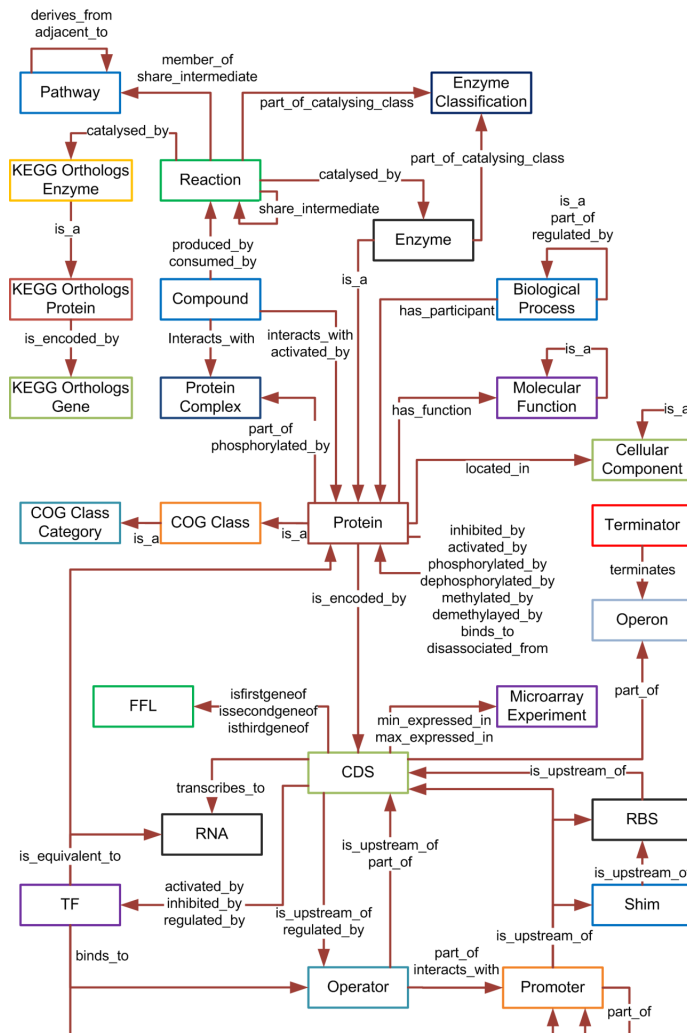
Open Access, available from ACS Synthetic Biology

w3id.org/synbio/ont

Synthetic Biology Ontology (SyBiOnt)



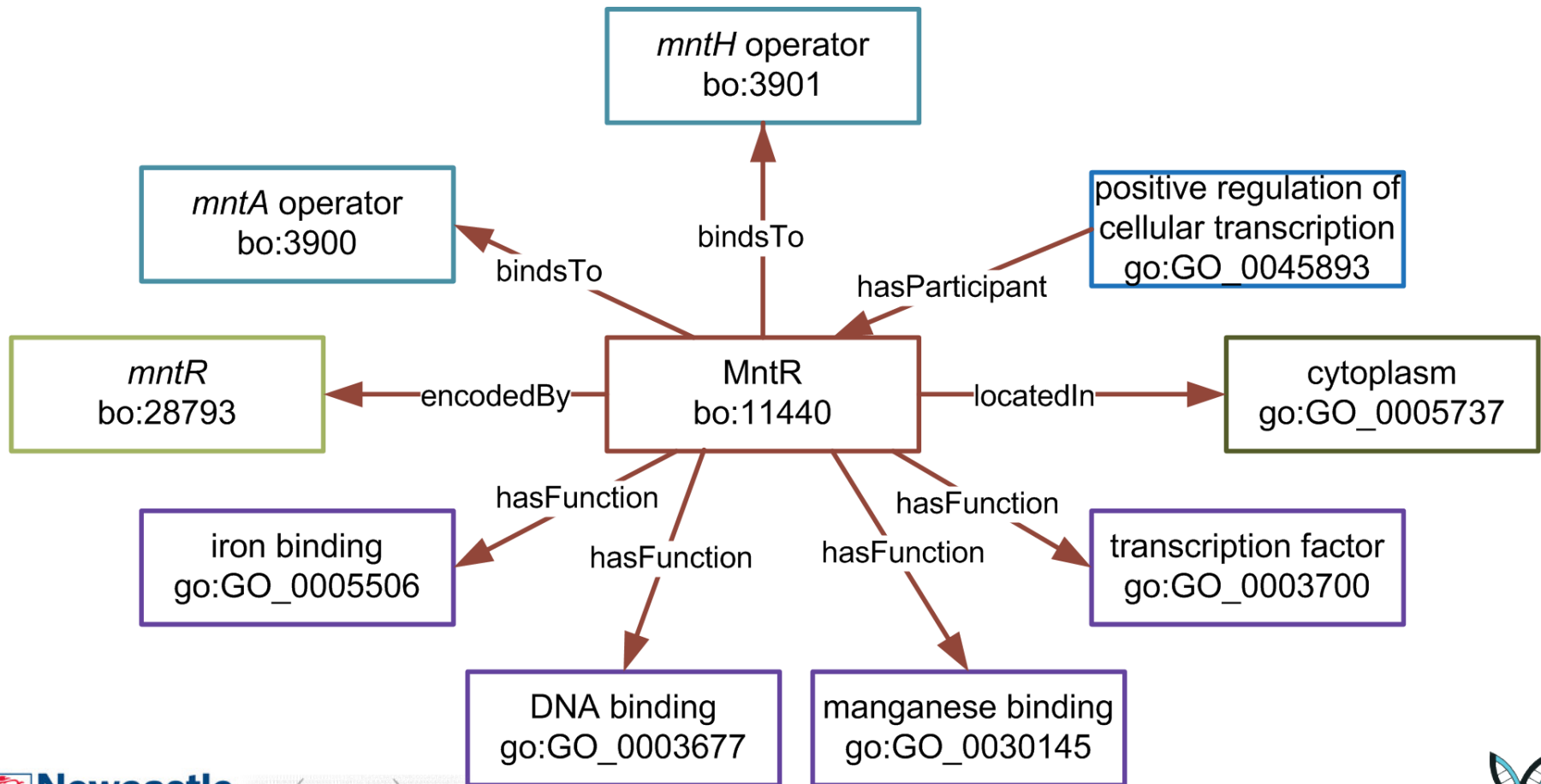
SyBiOntKB knowledge base



- * The knowledge base includes information about
 - * Sequences, annotations
 - * Metabolic pathways
 - * Gene regulatory networks
 - * Protein-protein interactions
 - * Gene expression

Misirli *et al.*, the
Journal of Integrative
Bioinformatics, 2013

An example network from SyBiOntKB



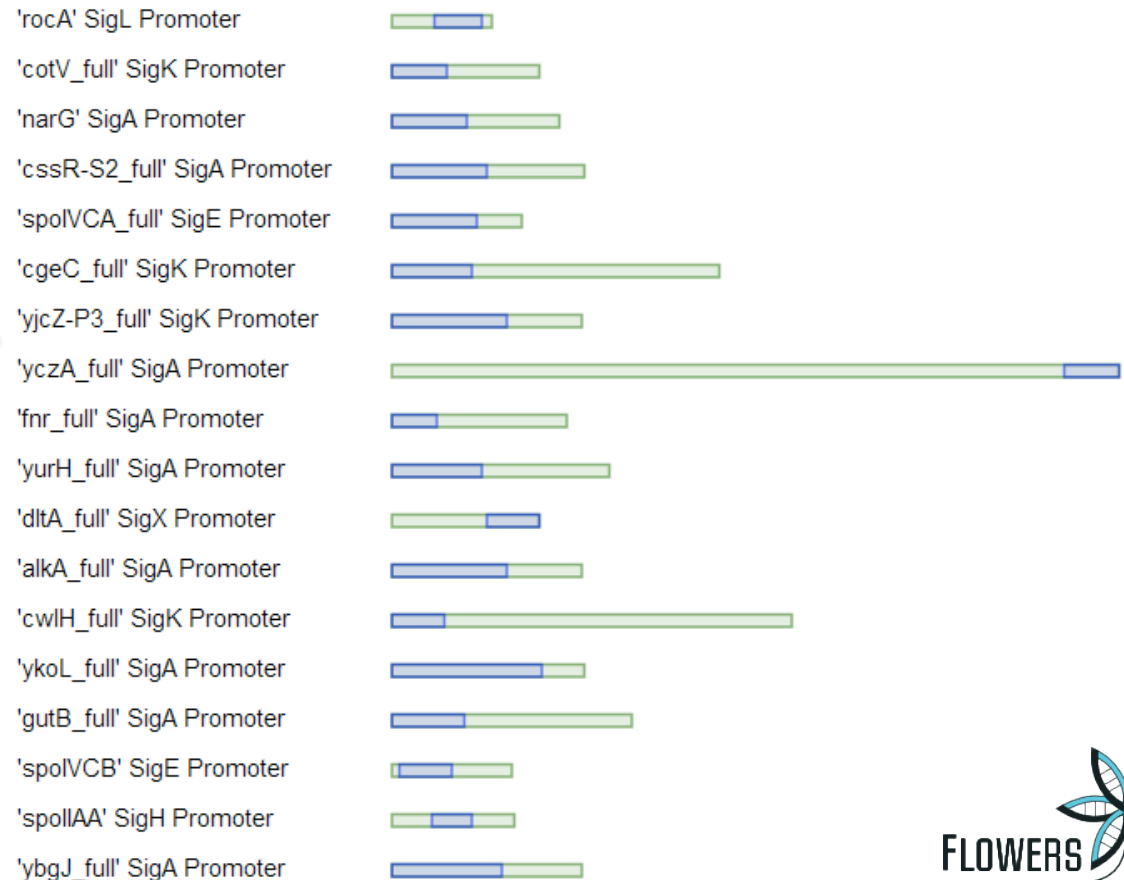
Testing the ontology with competency questions

- * *Which parts can be used as inducible promoters?*
- * **Operators** have **regulation type** restrictions to indicate whether they are used positively or negatively in regulating gene expression. A **Promoter** with one **Operator** part that has the 'Positive' **regulation type** restriction is an **Inducible Promoter**.

InduciblePromoter query and the results

QUERY: Promoter
and (has_part exactly 1 Operator)
and (has_part exactly 1 PositivelyRegulatedOperator)

RESULTS: In total, 51 promoters were classified



Examples about the automated identification of biological parts

| Part type | Count | Part type | Count |
|--|-------|----------------------------------|-------|
| Activator sites | 222 | SigA promoters | 465 |
| Repressor sites | 333 | SigB promoters | 67 |
| Inducible promoters | 51 | Constitutive promoters | 311 |
| Repressible promoters | 85 | Repressor encoding CDSs | 55 |
| Inducible promoters with two inputs (AND/OR gates) | 15 | Activator encoding CDSs | 44 |
| Repressible promoters with two inputs (NAND/NOR gates) | 25 | Response regulator encoding CDSs | 40 |
| | | Kinase encoding CDSs | 38 |

More examples

- * Which pathways should be targeted for the over-production of ammonium?
'Ammonium' is **produced by Reactions** that are **member of** 'Arginine and proline metabolism' and 'Purine metabolism' **Pathways**.
- * Which parts can be used to upregulate the production of ammonium?
The Compound 'Ammonia' with the **accession** of 'C00014' is **produced by the Reaction** 'RN:R00131', which **consumes the Compound** 'Carbamide' (C00086). 'Carbamide' is **produced by a Reaction** that is **catalysed by an Enzyme**, which is a subclass of a **Protein encoded by** the *argI* CDS with the **accession** of 'BSU40320'.
- * How can the SpooA protein, the master regulator of sporulation, be phosphorylated to trigger sporulation?
'SpooA' is **phosphorylated by** the 'KinC' and 'SpooB' **Proteins**. The 'SpooB' **Protein** is **phosphorylated by** 'SpooF' **Protein** which is further **phosphorylated by** the 'KinA' and 'KinB' **Proteins**.

Summary

- * SyBiOnt to capture complex biological data for synthetic biology computationally
- * Facilitates semantic reasoning
- * Complementary to existing standards and ontologies
- * Open access and available at w3id.org/synbio/ont
- * Extendible with new terms

Thanks

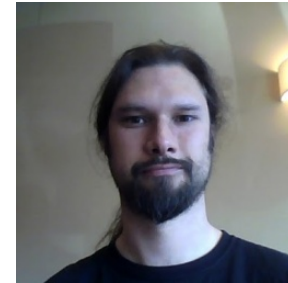
goksel.misirli@ncl.ac.uk
anil.wipat@ncl.ac.uk



Prof. Anil Wipat



Dr Jennifer Hallinan



Dr Matthew Pocock



Dr Phillip Lord



James McLaughlin



Prof. Herbert Sauro