

# **PanDaTox: a tool for accelerated metabolic engineering**

Gil Amitai<sup>1</sup> and Rotem Sorek<sup>1,\*</sup>

<sup>1</sup>Department of Molecular Genetics, Weizmann Institute of Science, Rehovot 76100, Israel

\* Email for correspondence: rotem.sorek@weizmann.ac.il

## **Abstract**

Metabolic engineering is often facilitated by cloning of genes encoding enzymes from various heterologous organisms into *E. coli*. Such engineering efforts are frequently hampered by foreign genes that are toxic to the *E. coli* host. We have developed PanDaTox ([www.weizmann.ac.il/pandatox](http://www.weizmann.ac.il/pandatox)), a web-based resource that provides experimental toxicity information for more than 1.5 million genes from hundreds of different microbial genomes. The toxicity predictions, which were extensively experimentally verified, are based on serial cloning of genes into *E. coli* as part of the Sanger whole genome shotgun sequencing process. PanDaTox can accelerate metabolic engineering projects by allowing researchers to exclude toxic genes from the engineering plan and verify the clonability of selected genes before the actual metabolic engineering experiments are conducted.

Keywords: Metabolic engineering, toxic genes, pandatox, gene cloning, synthetic biology

## Metabolic engineering and the problem of toxic genes

Metabolic engineering is a rapidly growing field where enzymatic pathways for the biosynthesis of desired molecules are genetically engineered into model microorganisms in order to harness bacterial productivity into industrial use<sup>1</sup>. Metabolic engineering is often practiced in close connection with synthetic biology and systems biology, as significant theoretical modeling is conducted in order to model pathways to be engineered into a given organism<sup>2-5</sup>.

The bacterial species *Escherichia coli* is one of the most widely used microorganisms in metabolic engineering, and had been utilized for numerous bio-production applications. For example, massive production of precursors for the antimalarial drug artemisinin was facilitated by inserting genes from several microorganisms into *E. coli*<sup>6,7</sup>; polyketides, which are the precursors for many antibiotics, have been produced within *E. coli* by introducing a combination of genes from three different bacteria into this organism<sup>8</sup>; various nutritional molecules such as acetate, pyruvate, succinate and an array of amino acids are also among the products produced within *E. coli* through metabolic engineering<sup>9</sup>. Recently, metabolic engineering has taken center stage in the global efforts for the design and generation of biofuel-producing organisms, in attempts to generate biological alternatives to fossil fuels<sup>1,10</sup>.

Thousands of microbial genomes have been sequenced to date<sup>5,11</sup>, and these genomes cumulatively harbor millions of enzymes that can be used as "building blocks" in the toolbox of the metabolic engineer. Following detailed planning of the pathway of new enzymes to engineer into *E. coli*, the metabolic engineer will usually search for these enzymes in the set of available genomes, and will select genes encoding the desired enzymes from organisms in which these genes exist. The selected genes will then be cloned into *E. coli* with or without optimization for expression through codon and promoter alterations<sup>3,12</sup>.

One of the major hurdles in such metabolic engineering efforts is genes that are toxic to the receiving organism<sup>4,6</sup>. Many enzymes are toxic to *E. coli* when cloned into this organism heterologously, and their toxicity may stem from generation of toxic intermediates or other incompatibility issues<sup>13</sup>. Attempts to engineer a toxic gene into *E. coli* will usually fail to produce viable clones, and will significantly slow down the engineering project. Furthermore, since the experimental procedures in metabolic engineering projects are significantly longer and more expensive as compared to the theoretical design, months of effort and considerable funds expenditure could be wasted due to toxic genes.

## PanDaTox: the Pan Genomic Database of Genes Toxic to Bacteria

We have recently introduced PanDaTox, the Pan genomic Database of genes and genomic elements that are Toxic to bacteria ([www.weizmann.ac.il/pandatox](http://www.weizmann.ac.il/pandatox))<sup>13</sup>. PanDaTox contains experimental toxicity information for over 1.5 million genes from hundreds of microorganisms, and reports the results of cloning attempts for each of these genes on single-copy and multiple-copy vectors within *E. coli* (Figure 1). This resource is expected to accelerate metabolic engineering efforts by allowing researchers to perform a-priori exclusion of toxic genes from the engineering plan, and, on the other hand, a-priori verification of clonability of selected genes before the actual experiments.

Although our identification of toxic genes was based on massive experimental cloning of over 9 million clones into *E. coli*, this was not an intentionally-designed experiment. In fact, the detection of toxic genes was a by-product of the Sanger whole-genome sequencing procedure through which hundreds of microbial genomes had been sequenced. Within this procedure, multiple copies of the sequenced genome are randomly fragmented into overlapping pieces of DNA that are serially inserted into *E. coli* on plasmids. These fragments are then sequenced and assembled based on sequence overlaps. While usually most fragments can be cloned into the *E. coli* host, a small fraction of the genome fails to be cloned, and this results in sequence gaps that interfere with proper genome assembly (Figure 2).

In an earlier study as well as in our recent study<sup>13,14</sup> we have demonstrated that gaps in microbial genome sequences are frequently caused by genes that are toxic to the *E. coli* host. These genes inhibit *E. coli* growth when cloned into it, and are hence missing from the set of sequences available for genome assembly. Following gap closure (which is performed by cloning-independent methods), these toxic genes can readily be identified. By experimenting with cumulatively more than 100 gap-residing genes from numerous different organisms, we have verified that our gap-based prediction of gene toxicity exceeds 80% of accuracy<sup>13,14</sup>.

The accuracy in our predictions of gene clonability into *E. coli* stems from several factors. First, the high sequencing coverage, needed for proper genome assembly, dictates that on average more than 25 independent clones would contain each gene. This ensures that the absence of a specific gene from the set of sequencing clones is not likely to occur by chance only, and provides high statistical power for identification of toxic genes. Furthermore, the diversity of cloning libraries used for genome sequencing, where some clones are propagated on a high-copy number plasmid while others are found on single-copy plasmid allows us to differentiate between stronger and weaker toxicity, i.e., differentiate between genes that will be toxic only when highly expressed and those which are toxic even in low doses.

PanDaTox contains data on over 40,000 genes that are predicted to be toxic to *E. coli* based on their clonability. It allows text-based searches for enzymes of interest according to multiple keywords, as well as sequence-based searches (blast) with a user-provided gene; homology searches enable checking whether homologs of a selected toxic gene are also toxic when cloned from a different genome; and clone coverage analysis allows the user to evaluate the toxicity inference. PanDaTox presents toxicity information in both graphical and numerical manners, and provides links to the experimental results collected.

With the advent of cheap genome sequencing, sequences of whole microbial genomes accumulate at a rapid pace. These, in turn, enrich the toolbox of metabolic engineers with numerous potential enzymes and pathways, but also increase the challenge in selecting the proper genes from the proper organisms. PanDaTox, which contains toxicity information for over 1.5 million microbial genes, can accelerate metabolic engineering by helping scientists to avoid toxic genes and to select genes that were experimentally proven to be clonable within *E. coli*. With the rapidly growing interest in metabolic engineering both in the academy and in the industry, we expect PanDaTox to form a useful tool for the community.

## **Disclaimer**

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California.

## **Acknowledgements**

The author thanks Hila Sberro and Gil Amitai for comments on the manuscript. R.S. was supported by the NIH R01AI082376-01, ISF FIRST program (grant 1615/09), and ERC-StG grant 260432. This work was supported by the Director, Office of Science, Office of

Basic Energy Sciences, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

## Figure legends

**Figure 1: The PanDaTox online web tool.** PanDaTox holds toxicity information for over 1.5 million genes. It presents detailed information for genes of choice; allows searching by multiple keywords; and provides links to multiple external sources. The home page provides access to the database; the gene page holds multiple clickable details on gene annotations, sequences, toxicity analyses, and presents results of experiments done with these genes (when applicable); the search page enables searching for genes of interest by keywords and various filters; the homologs page presents toxicity of homologs from other genomes. Users can also search for their genes of interest by sequence-based search (using the Blast page).

**Figure 2. Genes that are toxic to *E. coli* are exposed as a byproduct of the microbial genome sequencing process.** The sequenced genome is physically sheared into overlapping fragments of DNA, which are transformed into *E. coli* bacteria. Fragments are sequenced and assembled into contigs. Toxic genes result in *E. coli* growth inhibition and are hence not properly sequenced, creating "gaps" in genome assemblies. After gap closure the gap sequences are retrieved. Post-analysis of gap-residing genes exposes the toxic genes.

## References

- 1 Keasling JD. Manufacturing molecules through metabolic engineering. *Science* 2010; 330: 1355-8.
- 2 Keasling JD. Synthetic biology for synthetic chemistry. *ACS Chem Biol* 2008; 3: 64-76.
- 3 Nielsen J and Keasling JD. Synergies between synthetic biology and metabolic engineering. *Nat Biotechnol* 2011; 29: 693-5.
- 4 Keasling JD. Synthetic biology and the development of tools for metabolic engineering. *Metab Eng* 2012.
- 5 Sorek R and Serrano L. Bacterial genomes: from regulatory complexity to engineering. *Curr Opin Microbiol* 2011; 14: 577-8.
- 6 Martin VJ, Pitera DJ, Withers ST, Newman JD, and Keasling JD. Engineering a mevalonate pathway in *Escherichia coli* for production of terpenoids. *Nat Biotechnol* 2003; 21: 796-802.
- 7 Newman JD, Marshall J, Chang M, Nowroozi F, Paradise E, Pitera D et al. High-level production of amorpha-4,11-diene in a two-phase partitioning bioreactor of metabolically engineered *Escherichia coli*. *Biotechnol Bioeng* 2006; 95: 684-91; Tsuruta H, Paddon CJ, Eng D, Lenihan JR, Horning T, Anthony LC et al. High-level production of amorpha-4,11-diene, a precursor of the antimalarial agent artemisinin, in *Escherichia coli*. *PLoS One* 2009; 4: e4489.
- 8 Pfeifer BA, Admiraal SJ, Gramajo H, Cane DE, and Khosla C. Biosynthesis of complex polyketides in a metabolically engineered strain of *E. coli*. *Science* 2001; 291: 1790-2.
- 9 Wendisch VF, Bott M, and Eikmanns BJ. Metabolic engineering of *Escherichia coli* and *Corynebacterium glutamicum* for biotechnological production of organic acids and amino acids. *Curr Opin Microbiol* 2006; 9: 268-74; Yu C, Cao Y, Zou H, and Xian M. Metabolic engineering of *Escherichia coli* for biotechnological production of high-value organic acids and alcohols. *Appl Microbiol Biotechnol* 2011; 89: 573-83.
- 10 Clomburg JM and Gonzalez R. Biofuel production in *Escherichia coli*: the role of metabolic engineering and synthetic biology. *Appl Microbiol Biotechnol* 2010; 86: 419-34; Lee SK, Chou H, Ham TS, Lee TS, and Keasling JD. Metabolic engineering of microorganisms for biofuels production: from bugs to synthetic biology to fuels. *Curr Opin Biotechnol* 2008; 19: 556-63.
- 11 Pagani I, Liolios K, Jansson J, Chen IM, Smirnova T, Nosrat B et al. The Genomes OnLine Database (GOLD) v.4: status of genomic and metagenomic projects and their associated metadata. *Nucleic Acids Res* 2012; 40: D571-9.
- 12 Boyle NR and Gill RT. Tools for genome-wide strain design and construction. *Curr Opin Biotechnol* 2012.
- 13 Kimelman A, Levy A, Sberro H, Kidron S, Leavitt A, Amitai G et al. A vast collection of microbial genes that are toxic to bacteria. *Genome Res* 2012.
- 14 Sorek R, Zhu Y, Creevey CJ, Francino MP, Bork P, and Rubin EM. Genome-wide experimental determination of barriers to horizontal gene transfer. *Science* 2007; 318: 1449-52.