
COBREX.jl: constraint-based reconstruction and exascale analysis

Miroslav Kratochvíl^{1,†}, Laurent Heirendt^{1,4,†}, St. Elmo Wilken², Taneli Pusa^{1,3}, Sylvain Arreckx¹, Alberto Noronha³, Marvin van Aalst², Venkata P Satagopam^{1,4}, Oliver Ebenhöf², Reinhard Schneider^{1,4}, Christophe Trefois^{1,4,*}, and Wei Gu^{1,4,*}

¹Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Campus Belval, Belvaux, Luxembourg

²Institute of Quantitative and Theoretical Biology, Heinrich-Heine-Universität Düsseldorf, Universitätsstraße 1, 40225 Düsseldorf, Germany.

³Nium S.à.r.l., 6 Avenue des hauts fourneaux, L-4362 Esch-sur-Alzette, Luxembourg

⁴ELIXIR Luxembourg, 6, avenue du Swing, Campus Belval, L-4367 Belvaux, Luxembourg.

[†]Equal contributors.

*To whom correspondence should be addressed.

Summary: COBREX.jl is a Julia package for scalable, high-performance constraint-based reconstruction and analysis of very large-scale biological models. Its primary purpose is to facilitate the integration of modern high performance computing environments with the processing and analysis of large-scale metabolic models of challenging complexity. We report the architecture of the package, and demonstrate how the design promotes analysis scalability on several use-cases with multi-organism community models.

Availability and implementation: <https://doi.org/10.17881/ZKCR-BT30>.

Contact: christophe.trefois@uni.lu, wei.gu@uni.lu

1 Introduction

Understanding metabolic interactions in cells is a crucial step to investigate disease mechanisms and to discover new therapeutics (Cook and Nielsen, 2017; Apaolaza et al., 2018; Brunk et al., 2018). Constraint-Based Reconstruction and Analysis (COBRA) is a promising methodology for analyzing various metabolic processes at the organism- and community- levels (Fang, Lloyd, and Palsson, 2020). The main idea behind COBRA is to represent an organism as a constrained set of interconnected reactions and metabolites based on genomic sequencing data. This leads to a straightforward interpretation of metabolism as a constrained linear system, which enables the utilization of a wide range of well-developed analysis methods (Orth, Thiele, and Palsson, 2010).

12 The increasing ubiquity of genomic sequencing has led to a rapid expansion
13 in the number and complexity of genome-scale metabolic models, e.g. the human
14 metabolic model that has more than 80,000 reactions (Thiele et al., 2020). Recent
15 automated reconstruction tools can generate models spanning the entire primary
16 metabolism of both pro- and eukaryotes (Machado et al., 2018). Consequently,
17 metabolic models are becoming considerably larger in scale than their predeces-
18 sors, which is further compounded by the construction of multi-member commu-
19 nity models. This growth implies increasing analysis complexity (see Figure S1),
20 which in turn drives the need to develop analysis software that can accommodate
21 this complexity. While computing the solutions to the underlying constrained op-
22 timization problems is hard to accelerate and parallelize, many analysis types can
23 be decomposed into individual invocations of the optimizer, which may be paral-
24 leled. However, despite continued efforts (Heirendt, Thiele, and Fleming, 2017),
25 this remains challenging due to the scalability limits of existing software imple-
26 mentations.

27 Here, we present COBREXA.jl, a package for implementing and running dis-
28 tributed COBRA workflows. The package is implemented in the Julia program-
29 ming language (Bezanson et al., 2017), enabling facile extension with user-defined
30 numeric-computing routines, and interoperability with many high-performance
31 computing packages. It provides a ‘batteries-included’ solution for scaling analy-
32 ses to make efficient use of high-performance computing (HPC) facilities, giving
33 researchers a powerful toolkit for executing complicated high-volume workflows,
34 such as the creation and exploration of digital metabolic twins in personalized
35 medicine (Björnsson et al., 2020), and analysis of extensive microbial commu-
36 nities in ecology and biotechnology. We report the implementation architecture,
37 and substantiate how the design accommodates future extensions and scaling of
38 common analysis tasks.

39 **2 Implementation and results**

40 COBREXA.jl is an open architecture solution, providing interchangeable build-
41 ing blocks for implementing complicated COBRA workflows. Common analysis
42 methods, such as flux balance, flux variability, and gene knockout analyses (Gud-
43 mundsson and Thiele, 2010), are implemented as ready-to-use functions that may
44 be easily composed and customized. Most importantly, the building blocks are de-
45 signed so that the constructed workflows can be easily separated into parallelizable
46 analysis steps and executed on multiple computation nodes in HPC environments

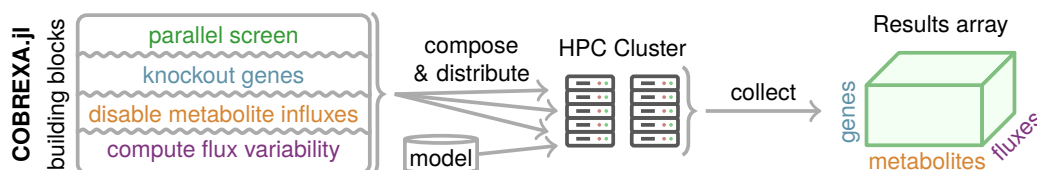


Figure 1: Schema of an example custom analysis construction that examines flux variability in many variants of a model, its distributed execution with COBREXA.jl, and collection of many results in a multi-dimensional array.

47 (as illustrated in Fig. 1). The concurrent execution of such workflows results in
48 significant computational speedups, without requiring user expertise in parallel
49 programming.

50 The design of COBREXA.jl distinguishes it from other COBRA implemen-
51 tations, which typically provide parallelization support for only a few selected
52 methods, and no current support for parallelization of custom method variants.
53 For example, parallel single-gene deletion analysis is commonly supported, but a
54 variant that explores the flux variability in knockouts must be reimplemented and
55 parallelized by the user.

56 A variety of model exchange and representation formats are supported, includ-
57 ing MATLAB format (Heirendt, Arreckx, et al., 2019); object-oriented JSON for-
58 mat (Ebrahim et al., 2013), and SBML (Keating et al., 2020). Additionally, im-
59 plementation of the workflows in Julia results in highly optimized execution of the
60 code at the cost of minor pre-compilation overhead, which benefits large, data-
61 heavy use cases. A detailed architecture overview is provided in Supplementary
62 Section S1.

63 To evaluate the effect of the new architecture and optimizations on the per-
64 formance and scalability of COBRA analyses, we benchmarked COBREXA.jl on
65 use-cases that benefit from parallelization. We compared its performance to that
66 obtained with COBRAPy (Ebrahim et al., 2013) and COBRA Toolbox (Heirendt,
67 Arreckx, et al., 2019), which are the widely adopted tools for running COBRA
68 workflows. Running on a 256-CPU multi-node cluster, COBREXA.jl was able
69 to fully utilize the available distributed computing resources and outperform the
70 implementation of flux variability analysis in other packages by a factor of be-
71 tween 2× and 10×, even on relatively small models (Supplementary Table S2).
72 We further demonstrated that COBREXA.jl is able to parallelize and distribute
73 custom workloads by re-implementing the production envelope functionality of
74 COBRAPy; leading to speedups of over 10×, even on a single 16-core computa-

75 tion node (Supplementary Table S3). Consequently, we expect that the COBRA
76 methods implemented in COBREXA.jl will enable reliable acceleration of many
77 current and future workloads by simply adding more computing resources. The
78 results are further discussed in Supplementary Section S3.4.

79 **3 Conclusion**

80 COBREXA.jl is a new package developed for large-scale distributed processing
81 of constraint-based biological models. It differs from the other implementations
82 of COBRA methods (Heirendt, Arreckx, et al., 2019; Ebrahim et al., 2013) by fo-
83 cusing on computational efficiency, and simplifies high-level construction of par-
84 allelized user-defined analysis methods. This is required for performing extensive
85 analyses of large models, future-proof extensibility, and workload distribution that
86 enables effective utilization of the common HPC infrastructure resources. The
87 package thus enables fast analysis of datasets that may pose challenges for the cur-
88 rently available tools, such as the comprehensive human gut microbiome models.

89 **Funding**

90 The research leading to these results has received funding from the European
91 Union’s Horizon 2020 Programme under the PerMedCoE Project (www.per-medcoe.eu), grant agreement n° 951773. This work was also partially funded
92 by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)
93 under Germany’s Excellence Strategy–EXC-2048/1–project ID 390686111 and
94 EU’s Horizon 2020 research and innovation program under the Grant Agreement
95 862087. The presented experiments were carried out using the HPC facilities of
96 the University of Luxembourg (see <https://hpc.uni.lu>).

98 **References**

- 99 Apaolaza, Iñigo et al. (2018). COBRA methods and metabolic drug targets in can-
100 cer. *Mol & Cell Oncol* 5.1, e1389672.
- 101 Bezanson, Jeff et al. (2017). Julia: A fresh approach to numerical computing. *SIAM*
102 *Rev* 59.1, pp. 65–98.
- 103 Björnsson, Bergthor et al. (2020). Digital twins to personalize medicine. *Genome*
104 *Med* 12.1, pp. 1–4.

- 105 Brunk, Elizabeth et al. (2018). Recon3D enables a three-dimensional view of gene
106 variation in human metabolism. *Nat Biotechnol* 36.3, p. 272.
- 107 Cook, Daniel J and Jens Nielsen (2017). Genome-scale metabolic models applied
108 to human health and disease. *Wires Syst Biol Med* 9.6, e1393.
- 109 Ebrahim, Ali et al. (2013). COBRApy: CONstraints-Based Reconstruction and
110 Analysis for Python. *BMC Syst Biol* 7.1, pp. 1–6.
- 111 Fang, Xin, Colton J Lloyd, and Bernhard O Palsson (2020). Reconstructing organ-
112 isms in silico: genome-scale models and their emerging applications. *Nat Rev*
113 *Microbiol* 18.12, pp. 731–743.
- 114 Gudmundsson, Steinn and Ines Thiele (2010). Computationally efficient flux vari-
115 ability analysis. *BMC Bioinformatics* 11.1, pp. 1–3.
- 116 Heirendt, Laurent, Sylvain Arreckx, et al. (2019). Creation and analysis of bio-
117 chemical constraint-based models using the COBRA Toolbox v.3.0. *Nat Protoc*
118 14.3, pp. 639–702.
- 119 Heirendt, Laurent, Ines Thiele, and Ronan MT Fleming (2017). DistributedFBA.jl:
120 high-level, high-performance flux balance analysis in Julia. *Bioinformatics*
121 33.9, pp. 1421–1423.
- 122 Keating, Sarah M et al. (2020). SBML Level 3: an extensible format for the ex-
123 change and reuse of biological models. *Mol Syst Biol* 16.8, e9110.
- 124 Machado, Daniel et al. (2018). Fast automated reconstruction of genome-scale
125 metabolic models for microbial species and communities. *Nucleic Acids Res*
126 46.15, pp. 7542–7553.
- 127 Orth, Jeffrey D, Ines Thiele, and Bernhard O Palsson (2010). What is flux balance
128 analysis? *Nat Biotechnol* 28.3, pp. 245–248.
- 129 Thiele, Ines et al. (2020). Personalized whole-body models integrate metabolism,
130 physiology, and the gut microbiome. *Mol Syst Biol* 16.5, e8982.