

Combining multiple functional annotation tools increases completeness of metabolic annotation

Marc Griesemer^{1,§,*}, Jeffrey Kimbrel^{1,*}, Carol Zhou², Ali Navid¹, Patrik D'haeseleer^{1,2}

¹ Biosciences and Biotechnology Division, Lawrence Livermore National Laboratory, Livermore, CA 94551, USA

² Global Security Computing Applications Division, Lawrence Livermore National Laboratory, Livermore, CA 94551, USA

§ Corresponding author (griesemer1@llnl.gov)

*Co-first authors

Abstract

The dirty little secret behind genome-scale systems biology modeling efforts is that they are invariably based on very incomplete functional annotations. Annotated genomes typically contain 30-50% of genes with little or no functional annotation, severely limiting our knowledge of the "parts lists" that the organisms have at their disposal. In metabolic modeling, these incomplete annotations are often sufficient to derive a reasonably complete model of the core metabolism at least, typically consisting of well-studied (and thus well-annotated) metabolic pathways that are sufficient for growth in pure culture. However secondary metabolic pathways or pathways that are important for growth on unusual metabolites exchanged in complex microbial communities are often much less well understood, resulting in missing or lower confidence functional annotations in newly sequenced genomes. Here, we present preliminary results on a comprehensive reannotation of 27 bacterial Tier 1 and Tier 2 reference genomes from BioCyc, focusing on enzymes with EC numbers annotated by KEGG, RAST, EFICAZ, and the Brenda enzyme database, and on membrane transport annotations by TransportDB, KEGG and RAST.

Introduction

Genome annotation has existed since the very first sequenced genomes. Initially, bacterial genomes were annotated primarily through manual curation of different groups of genes by experts. Today, automated gene annotation tools employing different methodologies with minimal manual curation are widely used, functionally annotating by homology to existing annotations, or by identification of conserved domains/motifs within a coding sequence. Many draft genomes and metagenome bins are often run through a single annotation pipeline where genome annotations are inherited from previous genome annotations. These annotation tools each have their own distinct advantages and limitations, often focusing on particular aspects of annotation or particular organisms. There is often a trade-off between general annotation tools which typically excel at

correctly identifying genes in core metabolic processes, and specialist annotation tools focusing on a particular set of genes or enzyme substrates (such as transporters). The array of different tools can be overwhelming, and the burden of learning the intricacies of multiple annotation pipelines can lead the typical researcher to choose only a single to annotate an entire genome. Additionally, even when multiple annotation tools are used, incorporating the different pieces of information in a cohesive manner remains a major barrier. Individual metabolic annotation tools often return annotations for different subsets of genes, offering the potential to greatly increase the completeness of metabolic annotations by combining the outputs of multiple tools.

“Genome-scale” metabolic models implicitly assume complete and accurate genome annotation, however, only 50-70% of genes in a typical prokaryotic genome are annotated (1), with the most reliably annotated genes often those involved in core metabolic processes. Modern metabolic modeling efforts, however, are moving beyond studying core metabolic pathways in a single organism towards multi-species models, real-world communities and ecosystems. Additionally, incorporation of complex ‘omics and metabolite data is becoming increasingly common, and robust genome annotations are therefore necessary to aggregate these different data streams. There also exist steps in metabolic pathways converting molecules for which the responsible enzymes are unknown. For example, one third of the EC database consists of “orphan enzymes” that have been described in the literature but for which no sequence data is available (1).

In Flux Balance Analysis, the issue of missing metabolic annotations is dealt with by “gap filling” - that is, adding a set of metabolic reactions beyond those that were derived directly from the genome annotation. A variety of gapfilling algorithms have been developed to predict a set of reactions to be added to make the metabolic network model sufficiently complete to be able to produce biomass (2-5). Henry et al. showed that in a broad collection of 130 genome-scale metabolic models added to the ModelSEED database, on average 56 additional gap filled reactions needed to be added to each model to produce biomass from simple defined nutrient media (6). But even after those additions, around a third of the reactions in each model were still inactive, meaning that there were enough holes left in the network to preclude metabolic flux through those reactions (6,7). In addition, the number of reactions that can partake in a gap-filling solution is vast (3,270 in the case of *E. coli* (8)), and the sets of reactions generated by different gap filling algorithms may have little or no overlap with each other (9). Clearly, if we could start with a more complete annotation of metabolic reactions to begin with, that would be preferable over having to add dozens of poorly supported reactions afterwards just to patch the holes in the network.

Recent genome-scale modeling of *Clostridium beijerinckii* NCIMB 8052 (9) demonstrated that the total number of genes and reactions included in the final curated model could be almost doubled by incorporating multiple annotation tools. The reconstruction of the *C. beijerinckii* metabolic network used 3 different database sources (SEED (10), KEGG (11), and RefSeq annotations captured in BioCyc (12)) to evaluate annotation coverage and produce a more robust model. Reactions annotated by all three databases were considered most reliable, but only 34 percent of reactions fit that criteria. Furthermore, for reactions appearing in more than two databases, gene-protein-reaction annotations were compared, finding that one database was missing a gene-reaction relationship rather than suggesting an alternative one. The overlap between annotations from different sources corresponded with an active set of reactions that

represented the core metabolic network of *C. beijerinckii* and illustrated the importance of annotation agreement. However, while this convergence of annotations is true for primary metabolic pathways, a comparison of annotations for secondary metabolic pathways is less clear.

Transporter annotations are rarely used in genome scale metabolic modeling, because of the difficulty in annotating which exact substrate is being transported by that transporter. Because of this, many metabolic modeling methods simply assume that there exist a transporter for any metabolite that needs to be imported into the cell or exported out of the cell. We know that this assumption is incorrect in some cases. For example, the yogurt bacteria *Streptococcus thermophilus* is known to have a highly unusual growth phenotype, in that it grows much better on lactose than on glucose (13), even though it does so by importing lactose, hydrolyzing it to glucose and galactose, and then metabolizing only the glucose, with the galactose secreted back out of the cell. The cause for its growth deficit when fed on glucose directly is that it lacks the usual glucose phosphotransferase system used by many bacteria, and instead it has a very efficient lactose import mechanism that makes it well adapted to grow in milk (13). Better prediction tools such as TransportDB's Transporter Automatic Annotation Pipeline (TransAAP (14)) now allow us to generate substrate predictions that are sufficiently detailed to be included in metabolic pathways, and could give insights into growth or metabolite exchange phenotypes that are not readily apparent from the enzymes present in the genome.

We undertook an investigation into the effectiveness of individual tools in genome annotation and their overlap with each other. Using 27 bacterial reference genomes from BioCyc (12), we evaluated how many genes, EC numbers, and gene-reaction links were unique or shared with the findings of other tools. We also undertook a study of how transporter annotations were handled between RAST (10), KEGG (11), and TransportDB (14).

Methods

Reference Genomes

We utilized a total of 27 prokaryotic genomes of Metacyc Tier 1 & Tier 2 organisms (12). These genomes are from a range of phyla, including 15 Proteobacteria, 6 Firmicutes, 3 Actinobacteria, 2 Bacteroidetes and 1 Cyanobacteria. In addition to a range of phyla, these 27 organisms also display a range of lifestyles, including a human gut symbiont (*Bacteroides thetaiotaomicron* VPI-5482) and pathogens (e.g. *Mycobacterium tuberculosis*), an obligate insect endosymbiont (*Candidatus Evansia muelleri*), a bacterium with interesting metabolism (*Aurantimonas manganoxydans* SI85-9A1), and an important marine primary producer (*Synechococcus elongatus* PCC 7942). These organisms present various types of annotation challenges, such as identifying incomplete pathways in obligate intracellular organisms with small genomes, as well as the challenges of annotating incomplete, "draft" contigs compared to "finished" genomes.

Genbank files were downloaded from NCBI (accessions list in Table 1). Genbank files are designed to hold both standard genomic information, as well as fields with user-specified or project specific information. The specification includes genomic features (e.g. "CDS" or "tRNA") with embedded qualifiers (e.g. "product" or "locus_tag"). The 27 genomes here were all in various

states of completeness, used custom qualifiers, and on occasion used similar qualifier fields for different information. Therefore, our first priority was to “standardize” each genbank file to a similar, reduced state. As this work pertains only to protein-coding genes, only “CDS” features were retained. Additionally, CDS features identified as a pseudogene were removed. Qualifiers within CDS features that were retained included locus_tag, protein_id and translation. Thus, removed qualifiers include any that would bias our downstream analysis, such as included EC numbers or gene products. Our standardized genbank files use the same CDS genome coordinates as the original, and no attempt was made to correct these features.

We sought to match the locus tag prefix for a genome with the one present in BioCyc. Once this was accomplished for all of the 27 organisms, we began the process of running each genome through the various tools.

Table 1. Reference genomes used in this study

Genome Name	Biocyc ID	Phylum	NCBI Accessions	Proteins
<i>Mycobacterium tuberculosis</i> CDC1551	MTBCDC1551	Actinobacteria	AE000516	4189
<i>Mycobacterium tuberculosis</i> H37Rv	MTBH37RV	Actinobacteria	AL123456	4018
<i>Streptomyces coelicolor</i> A3(2)	SCO	Actinobacteria	NC_003888, NC_003903, NC_003904	8152
<i>Bacteroides thetaiotaomicron</i> VPI-5482	BTHE	Bacteroidetes	AE015928, AY171301	4825
<i>Candidatus Cardinium hertigii</i> ^b	CBTQ1	Bacteroidetes	HG422566, CBQZ01000001- CBQZ010000011	739
<i>Synechococcus elongatus</i> PCC 7942	SYNEL	Cyanobacteria	CP000100, CP000101	2661
<i>Listeria monocytogenes</i> 10403S	10403S_RAST	Firmicutes	CP002002	2814
<i>Bacillus anthracis</i> Ames	ANTHRA	Firmicutes	NC_003997, AE017335, AE017336	5602
<i>Bacillus subtilis</i> 168	BSUB	Firmicutes	AL009126	4185
<i>Clostridium saccharoperbutylacetonicum</i> ATCC 27021	CLOSSAC	Firmicutes	CP004121, CP004122	5821
<i>Eubacterium rectale</i> ATCC 33656	EREC	Firmicutes	CP001107	3626
<i>Peptoclostridium difficile</i> 630	PDIF272563	Firmicutes	AM180355, AM180356	3809
<i>Agrobacterium fabrum</i> C58	AGRO	Proteobacteria	AE008687, AE008688, AE008689, AE008690	5402
<i>Aurantimonas manganoxydans</i> SI85-9A1	AURANTIMONAS	Proteobacteria	AAPJ01000001- AAPJ01000035	3650
<i>Caulobacter crescentus</i> CB15	CAULO	Proteobacteria	AE005673	3737
<i>Caulobacter crescentus</i> NA1000	CAULONA1000	Proteobacteria	CP001340	3885
<i>Escherichia coli</i> CFT073	ECOL199310	Proteobacteria	AE014075	5379
<i>Escherichia coli</i> K-12 substr. W3110	ECOL316407	Proteobacteria	NC_007779	4410
<i>Escherichia coli</i> B str. REL606	ECOL413997	Proteobacteria	CP000819	4209

<i>Escherichia coli</i> K-12 substr. MG1655 ^a	ECOLI	Proteobacteria	U00096	4140
<i>Escherichia coli</i> O157:H7 str. EDL933	ECO0157	Proteobacteria	AE005174, AF074613	5449
<i>Candidatus</i> Evansia muelleri ^b	EVA	Proteobacteria	LM655252	330
<i>Helicobacter pylori</i> 26695	HPY	Proteobacteria	CP003904	1594
<i>Methylosinus trichosporium</i> OB3b	MOB3B	Proteobacteria	NZ_ADVE02000001- NZ_ADVE02000003	4344
<i>Candidatus</i> Portiera aleyrodidarum BT-QVLC ^b	PABTQVLC	Proteobacteria	CP003867	280
<i>Shigella flexneri</i> 2a str. 2457T	SHIGELLA	Proteobacteria	AE014073	4068
<i>Vibrio cholerae</i> O1 biovar El Tor str. N16961	VCHO	Proteobacteria	AE003852, AE003853	3828

a: Tier 1 Pathway Genome Database (EcoCyc)

b: Endosymbiont with reduced genome

Annotation Tools

RAST (Rapid Annotation Subsystem Technology, (10)) is an open-source web server for genome annotation, using a "Highest Confidence First" assignment propagation strategy based on manually curated subsystems and subsystem-based protein families that automatically guarantees a high degree of assignment consistency. RAST returns an analysis of the genes and subsystems in each genome, as supported by comparative and other forms of evidence. We used the NMPDR website (rast.nmpdr.org) to generate genome-wide annotations for our 27 reference genomes.

KEGG (Kyoto Encyclopedia of Genes and Genomes, (11)) is a collection of genome and pathway databases for systems biology. We used the KAAS (KEGG Automatic Annotated Server, www.genome.jp/tools/kaas, (15)) to generate genome-wide annotations for our 27 reference genomes. KAAS assigns KEGG Orthology (KO) numbers using the bi-directional best hit method (BBH) against a set of default prokaryotic genomes in the KEGG database.

EFICAz (16) uses large scale inference to classify enzymes into functional families, combining 4 methods into a single approach without the need for structural information. Recognition of functionally discriminating residues (FDR) allows EFICAz to use a method called evolutionary footprinting. The latest EFICAz has high precision and recall ability: under test using sequence similarity of >40%, precision and recall were 0.88 and 0.85, while at sequence similarity of >60%, they were near unity. We used a local install of EFICAz2.5 to generate EC number predictions for our 27 reference genomes.

Brenda (17) is a publicly available enzyme database (containing 82,568 enzymes and 7.2 million enzyme sequences as of 2017) based on the literature and contains functional and molecular information such as nomenclature, enzyme structure, and reactions and specificity. We annotated our 27 reference genomes based on a BLAST search at >60% sequence identity against a local copy of the Brenda database of enzyme reference sequences.

Transporter annotations

Where available, pregenerated transporter annotations were downloaded from the TransportDB 2.0 database (14). For those reference genomes that were not already present in the database (*M. trichosporium* OB3b, *Candidatus C. hertigii*, *Candidatus E. muelleri*, and *A. manganoxydans* SI85-9A1), we submitted the genome to the TransAAP web-based transporter annotation tool (http://www.membranetransport.org/transportDB2/TransAAP_login.html). We also retrieved transporter annotations from RAST and KEGG. Substrate names were standardized to allow comparison and ranked 1-5 from most to least specific. Substrates that consist of a single metabolite that can be incorporated as a transport reaction reaction in a metabolic model are ranked 1. Substrate predictions that map to a small number of possible reactions are ranked 2. Broader substrate classes that are not directly usable to construct a metabolic network but could be used for gapfilling or interpretation of transcriptomics data are ranked 3 and 4. Finally, annotated transporters without substrate prediction are ranked 5.

Table 2. Examples of substrate annotation ranking, from most specific (rank 1) to least specific (no substrate, rank 5)

Rank	Substrate examples
1	<ul style="list-style-type: none">• Fe• lysine
2	<ul style="list-style-type: none">• Mg/Co/Ni• aromatic amino acid
3	<ul style="list-style-type: none">• ferric siderophore• sugar
4	<ul style="list-style-type: none">• multidrug efflux protein
5	No substrate annotated

Results

EC number annotations

In total, the four annotation tools produces 47,447 unique Gene-EC annotations (“gene X codes for an enzyme with EC number Y”) across the 27 reference genomes, or an average of 1757 annotations per genome. The Venn diagram in Figure 1 illustrates the degree of overlap - and non-overlap - between the sets of annotations produced by the four tools.

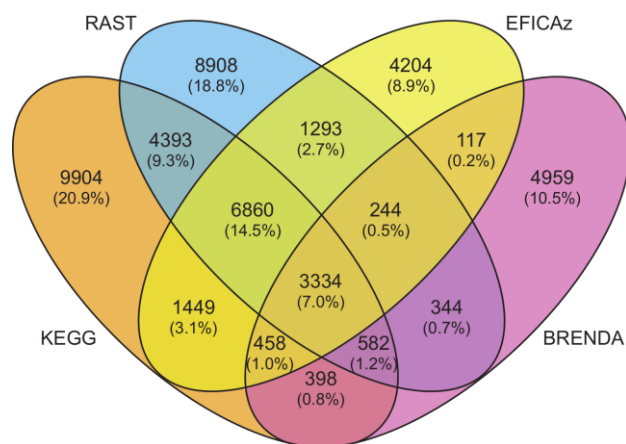


Figure 1. Large differences exist between the sets of Gene-EC annotations generated by the four annotation tools across the 27 reference genomes.

It is clear that the metabolic annotations produced by these automated genome-wide annotation tools can differ drastically. Each tool produced on average between 23% (EFICAz) and 48% (Brenda) unique gene-EC annotations that were not predicted by any of the other tools. Overall, fewer than a quarter of all gene-EC annotations are agreed on by at least 3 tools.

Table 3. Number of gene-EC annotation disagreements that exist across pairs of tools.

Tool Combination	Gene-EC Disagreements
KEGG-RAST	4218/20915 (20.2%)
KEGG-EFICAz	2264/16677 (13.6%)
KEGG-BRENDA	2971/6748 (44.0%)
RAST-EFICAz	2717/15694 (17.3%)
RAST-BRENDA	2381/6288 (37.9%)
EFICAz-BRENDA	1699/5601 (30.3%)

The EC numbers on which the different tools most often agree across the 27 reference genomes tend to belong to well studied core metabolic pathways, such as glycolysis, amino acid and nucleotide biosynthesis, etc.

The EcoCyc database (18) is an extensively hand-curated and continuously updated database summarizing all the experimentally determined enzymatic functions in *Escherichia coli* K-12 substr. MG1655, the single best studied model organism in the history of modern biology. We can use the EC numbers annotated in EcoCyc as a gold standard to evaluate how well the

automated annotation tools are able to annotate the enzymes in *E. coli* K-12. Figure 2 shows how the set of gene-EC annotations generated by KEGG and RAST compare against EcoCyc.

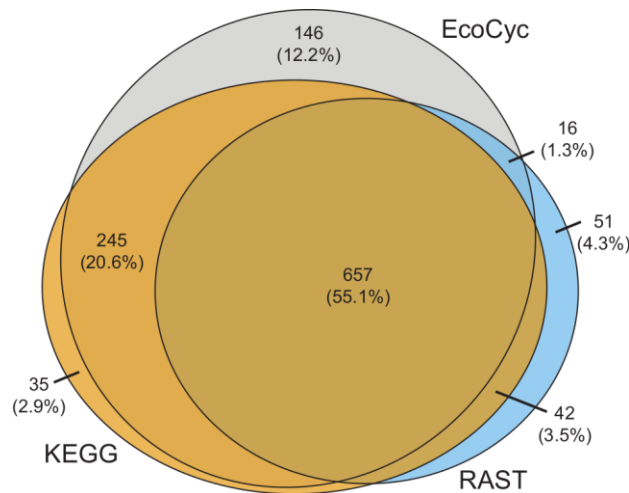


Figure 2. Gene-EC annotations produced by KEGG and RAST for *E. coli* K12, compared to the EcoCyc gold standard. The sets and intersections are drawn proportionally to the number of annotations in each.

Using EcoCyc as the Gold Standard for *E. coli* K12, the different tools achieve a Precision of 78% (BRENDA) to 92% (KEGG), and a Recall of 33% (BRENDA) to 85% (KEGG). Note that even on *E. coli* K12, which we expect to be a best-case situation, there are significant differences in the annotations produced by the different tools, and each tool is only able to cover a subset of the known enzymes in EcoCyc.

The annotation tools show much more agreement on *E. coli* than on more remote lineages such as Actinomycetes, Bacteroidetes, or Clostridia. For *E. coli* K12, 60% of EC numbers are agreed on by 3 or more tools, while 28% EC numbers come from only a single tool. In contrast, for *P. difficile* 630, only 33% of EC numbers are agreed on by 3 or more tools, and 48% of EC numbers come from only a single tool.

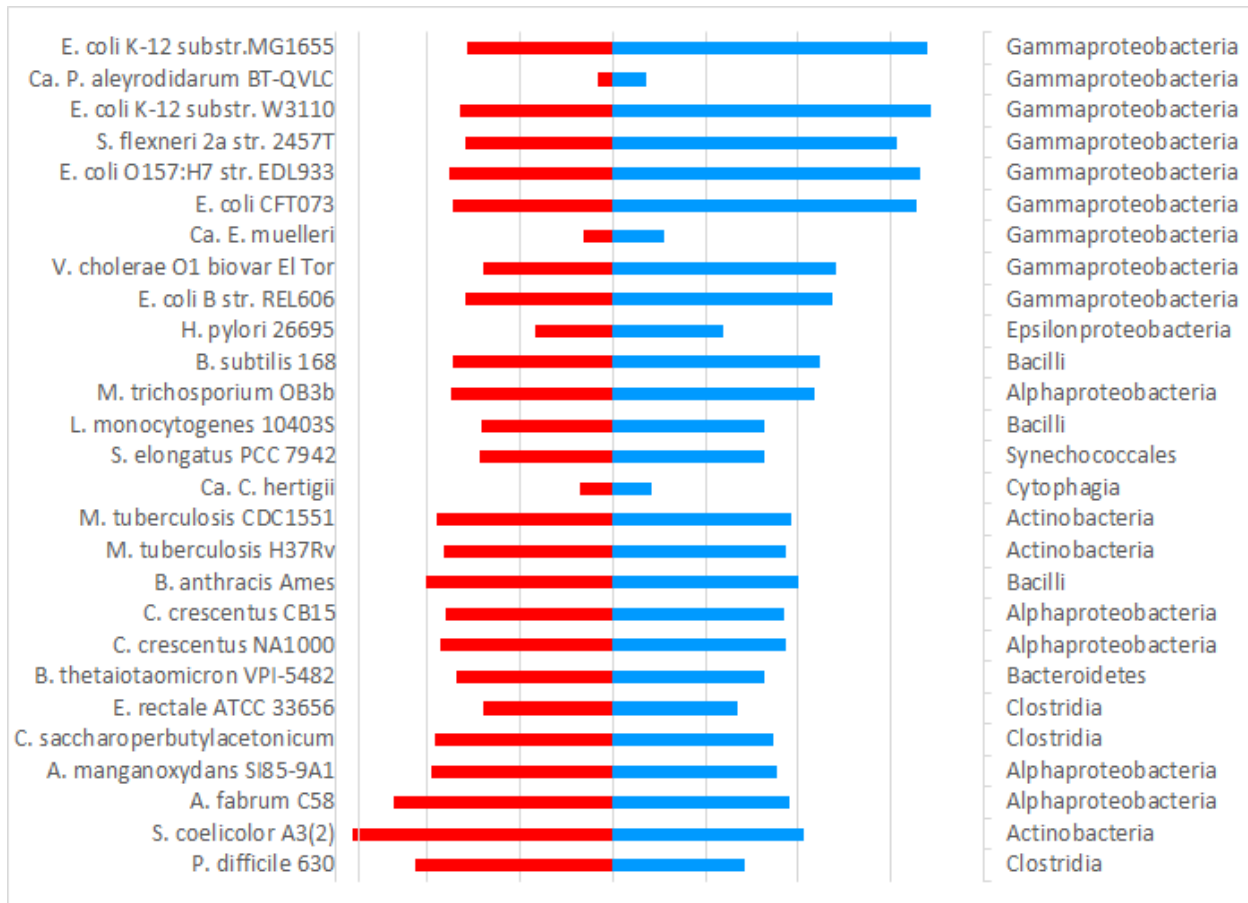


Figure 3. The 27 reference genomes were sorted with respect to the fraction of EC numbers that were predicted by 3 or more tools (blue bars). The top of the list is dominated by model organisms such as *E. coli*, *B. subtilis*, and closely related organisms. As we move farther away from such well-studied model organisms, the fraction of unique EC numbers predicted only by a single tool (red bars) increases, at the expense of those predicted by multiple tools.

The disjoint sets of annotations produced by the different tools provide us with an opportunity to trade off confidence in the annotations versus coverage. If higher confidence is required, we can focus solely on the subset of annotations that is agreed upon by multiple tools. Conversely, if the lack of genome coverage or metabolic network coverage is considered a problem, we can use the union of multiple tools to achieve a wider annotation.

Figure 4A shows the number of unique EC numbers annotated by each of the tools across the 27 reference genomes, which reflects the size of the metabolic network reconstruction (average 868 EC numbers per genome). Figure 4B shows the number of genes in all of the reference genomes annotated with one or more EC numbers by each of the tools, which reflects the overall genome annotation coverage (average 1361 genes per genome, or 34% of the genes).

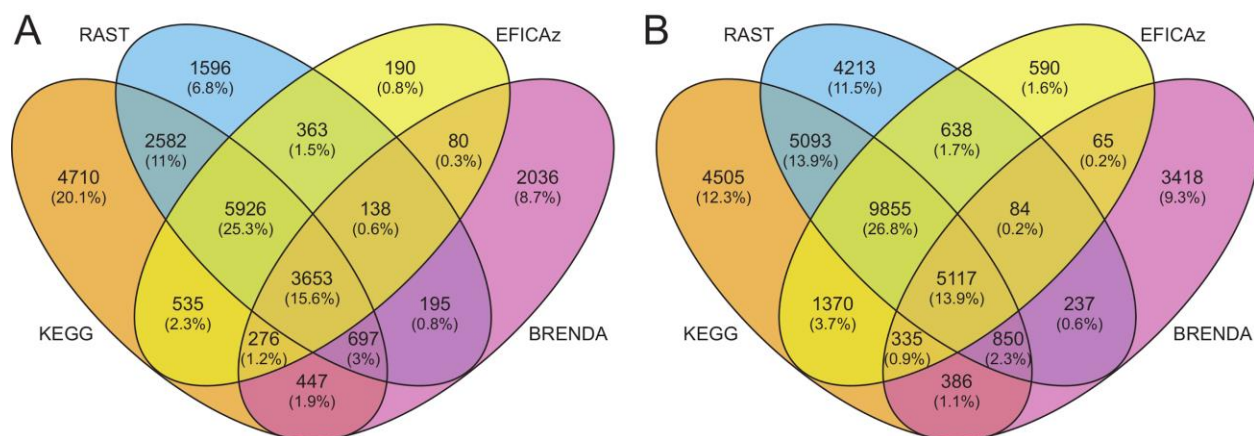


Figure 4. Combining Annotation Tools Improves Genome annotation and Metabolic Network Coverage. A: Total unique EC numbers annotated. B: Total genes annotated with EC numbers

Note that each tool adds a significant number of reactions to the metabolic network model, and each tool significantly contributes to the number of genes covered with metabolic annotations.

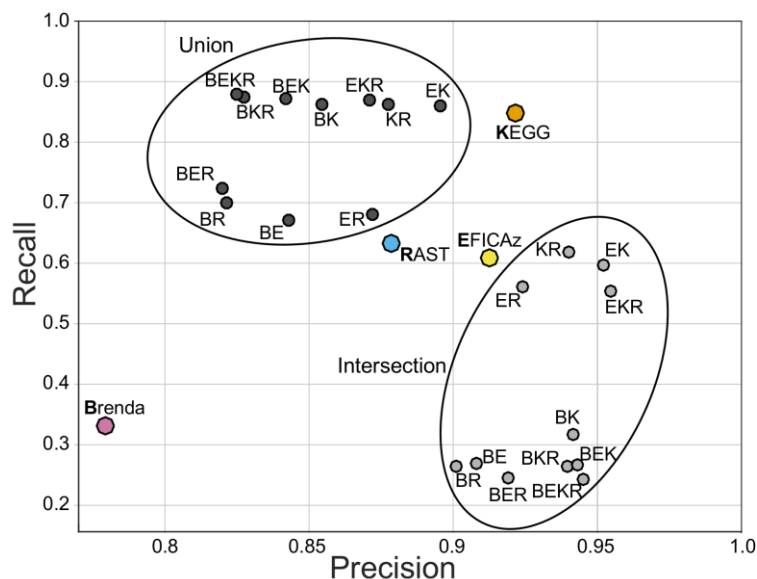


Figure 5. Precision vs Recall for different combinations of tools on EcoCyc. Individual tools are denoted by B,E,K, or R for Brenda, EFICAz, KEGG, and RAST, respectively.

Choosing *E. coli* K12 as the “gold standard”, we took the EC numbers annotated by each tool and performed the union and intersection of each of the combinations. These combinations included all pairs, triplets as well as the list of ECs from the union and intersection of all 4 tools combined. We then compared each of these EC lists to the 1064 EC numbers from EcoCyc. The number of true positives, false positives, and false negatives were calculated for each combination. True positives (TP) correspond to EC

numbers predicted by the annotation tools, and present in EcoCyc. False positives (FP) are EC numbers annotated by the tools, but not found in EcoCyc. False negatives (FN) are those EC number that occur in EcoCyc, but were not predicted by our annotation tools.. Precision is defined as the $(TP)/(TP+FP)$ and recall as $(TP)/(TP+FN)$.

Figure 5 shows a plot of precision versus recall for all the different combinations of tools. Three groupings of points stand out from this graph. First, all the combinations that contain some union of the tools have high precision and high recall. Secondly, the combinations that are a result of intersections of EC lists are bunched in an area of very high precision but much lower recall. In between these two groups lie the single tool results. One can see that precision and recall are lowest for Brenda alone; KEGG in this case has the best precision and recall of the four alone. When any set of tools intersects with Brenda, precision remains about the same but recall drops dramatically (e.g. comparing ER and BER). On the other hand, a union with KEGG increases recall. In some cases, the combination of tools makes the precision and/or recall worse.

Transporter Annotations

Both RAST and KEGG yield surprisingly few transporter annotations, with an average of 114 and 204 transporter predictions per genome respectively. In addition, most of the annotated transporters lack substrate predictions (52% of transporter annotations in RAST, 47% in KEGG) or have ambiguous substrate predictions (ranks 3-4; 20% in RAST, 21% in KEGG), while less than a third have substrate predictions that are sufficiently detailed that they could be incorporated in a metabolic model (ranks 1-2; 28% in RAST, 33% in KEGG). In contrast, TransportDB produces an average of 426 transport annotations per genome, and most of those have specific substrate predictions (59% rank 1-2; 32% rank 3-4, 10% rank 5).

Transporter annotations by the different tools show remarkable little overlap. Out of the more than 15,000 genes annotated as transporters (regardless of substrate prediction), the three tools only agree on 2.8% (423/15161). Out of those, only 69 genes are annotated by all three tools with a specific substrate prediction (ranks 1-2)

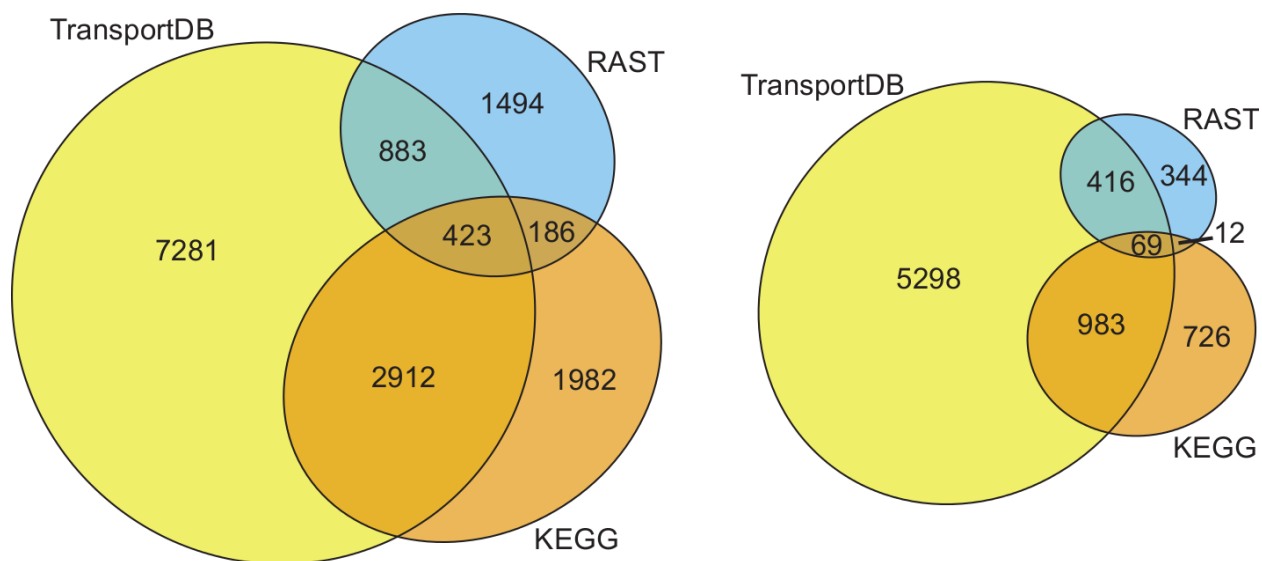


Figure 6. A: Total number of genes annotated as transporters, regardless of substrate. B: Transporter annotations with substrates predictions specific enough to be included in metabolic models (rank 1 or 2). Note that the sizes of the sets in these Venn diagrams are drawn proportional to the number of genes in each.

When two or more tools provide a sufficiently specific substrate annotation, the substrate annotations tend to agree 68% of the time, even if they may not be perfectly identical (for example, one transporter was annotated as “leucine/valine”, “leucine”, and “branched-chain amino acid” by TransportDB, RAST and KEGG respectively)

Discussion

RAST and KEGG are the most widely used tools for metabolic network reconstruction. However, they do not necessarily produce identical annotations. In our analysis, KEGG produced the best Precision and Recall on *E. coli* K12. In general, KEGG produces a larger number of unique EC numbers, which could indicate more over-prediction, or more comprehensive pathway coverage. Note that both also generate many reactions without official EC numbers.

EFICAz produces the least number of unique EC numbers, but can be used in combination with RAST or KEGG to highlight high confidence annotations. EFICAz also produced 3-digit EC number annotations which may be used for hole filling.

BLAST against the Brenda database of reference enzymes produced the smallest number of annotations, but a high fraction of unique EC numbers. Of the top 10 unique EC numbers produced by this method, only one is also covered by RAST and KEGG, two of the EC numbers have been deprecated, and six are EC numbers that have been assigned in 2000 or newer, and may not have been incorporated into the predictions by the other annotation tools yet.

Recommendations

1. Do not just use a single annotation tool unless you are only interested in core metabolism.
2. Trade off confidence versus coverage by looking at the intersection or union of multiple annotation tools.
3. EFICAz can be used to identify higher confidence annotations, or partial EC numbers for hole filling
4. BLASTing against a database of reference sequences such as the Brenda database is generally an inefficient method for annotating enzymes, but may be useful to cover more recently assigned EC numbers not yet included by other tools.
5. More tool development is needed to allow merging of annotations beyond simple EC numbers.

Acknowledgements

This work was supported by the Department of Energy through the Genome Sciences Program as part of the LLNL Biofuels SFA (contract SCW1060). Work at LLNL was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344

References

1. Hanson, A.D., Pribat, A., Waller, J.C. and de Crécy-Lagard, V. (2009) 'Unknown' proteins and 'orphan' enzymes: the missing half of the engineering parts list--and how to find it. *Biochem J.*, 425(1), 1-11.
2. Reed JL, Patel TR, Chen KH, Joyce AR, Applebee MK, Herring CD, Bui OT, Knight EM, Fong SS, Palsson BO. (2006) Systems approach to refining genome annotation. *Proc Natl Acad Sci U S A.* 103(46):17480-4.
3. Satish Kumar V, Dasika MS, Maranas CD. (2007) Optimization based automated curation of metabolic reconstructions. *BMC Bioinformatics.* 8:212.
4. Kumar VS, Maranas CD. (2009) GrowMatch: an automated method for reconciling in silico/in vivo growth predictions. *PLoS Comput Biol.* 5(3):e1000308.
5. Benedict MN, Mundy MB, Henry CS, Chia N, Price ND. (2014) Likelihood-based gene annotations for gap filling and quality assessment in genome-scale metabolic models. *PLoS Comput Biol.* 10(10):e1003882.
6. Henry CS, DeJongh M, Best AA, Frybarger PM, Linsay B, Stevens RL. (2010) High-throughput generation, optimization and analysis of genome-scale metabolic models. *Nat Biotechnol.* 28(9):977-82
7. Ponce-de-Leon M, Calle-Espinosa J, Peretó J, Montero F. (2015) Consistency Analysis of Genome-Scale Models of Bacterial Metabolism: A Metamodel Approach. *PLoS One.* 10(12):e0143626.

8. Krumholz EW, Libourel IG. (2015) Sequence-based Network Completion Reveals the Integrality of Missing Reactions in Metabolic Networks. *J Biol Chem.* 2015 290(31):19197-207.
9. Milne, C.B., Eddy JA, Raju, R., Ardekani, S., Kim, P.J., Senger, R.S., Jin, Y.S., Blaschek, H.P. and Price, N.D. (2011). Metabolic network reconstruction and genome-scale model of butanol-producing strain *Clostridium beijerinckii* NCIMB 8052. *BMC Syst Biol.*, 5, 130.
10. Overbeek, R., Olson, R., Pusch, G.D., Olsen, G.J., Davis, J.J., Disz, T., Edwards, R.A., Gerdes, S., Parrello, B., Shukla, M., Vonstein, V., Wattam, A.R., Xia, F. and Stevens, R. (2014) The SEED and the Rapid Annotation of microbial genomes using Subsystems Technology (RAST). *Nucleic Acids Res.*, 42(Database issue), D206-214.
11. Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. (2016) KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res.* 44(D1):D457-62.
12. Caspi, R., Billington, R., Ferrer, L., Foerster, H., Fulcher, C.A., Keseler, I.M., Kothari, A., Krummenacker, M., Latendresse, M., Mueller, L.A., Ong, Q., Paley, S., Subhraveti, P., Weaver, D.S. and Karp, P.D. (2016) The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases. *Nucleic Acids Res.*, 44(D1), D471-D480.
13. Poolman B. (1993) Energy transduction in lactic acid bacteria. *FEMS Microbiol Rev.* 12(1-3):125-47
14. Elbourne, L.D.H., Tetu, S.G., Hassan, K.A., Paulsen, I.T. (2017) TransportDB 2.0: a database for exploring membrane transporters in sequenced genomes from all domains of life. *Nucleic Acids Res.* 45, D320–D324.
15. Moriya, Y., Itoh, M, Okuda, S., Yoshizawa, A.C. and Kanehisa, M. (2007) KAAS: an automatic genome annotation and pathway reconstruction server. *Nucleic Acids Res.*, 35(Web Server issue), W182-W185.
16. Kumar, N. and Skolnick, J. (2012) EFICAz2.5: application of a high-precision enzyme function predictor to 396 proteomes. *Bioinformatics*, 28(20), 2687-2688.
17. Placzek, S., Schomburg, I., Chang, A., Jeske, L., Ulbrich, M., Tillack, J. and Schomburg, D. (2017) BRENDA in 2017: new perspectives and new tools. *Nucleic Acids Res.*, 35(D1), D380-388.
18. Keseler IM, Mackie A, Santos-Zavaleta A, Billington R, Bonavides-Martínez C, Caspi R, Fulcher C, Gama-Castro S, Kothari A, Krummenacker M, Latendresse M, Muñoz-Rascado L, Ong Q, Paley S, Peralta-Gil M, Subhraveti P, Velázquez-Ramírez DA, Weaver D, Collado-Vides J, Paulsen I, Karp PD. (2017) The EcoCyc database: reflecting new knowledge about *Escherichia coli* K-12. *Nucleic Acids Res.* 45(D1):D543-D550.