# Lawrence Berkeley National Laboratory

Title
MAGI: A method for metabolite, annotation, and gene integration

Authors
Erbilgin, Onur
Rübel, Oliver
Louie, Katherine B
et al.

Peer reviewed

1 MAGI: A method for metabolite, annotation, and gene integration
2

3 Onur Erbilgin[1], Oliver Rübel[2], Katherine B. Louie[3], Matthew Trinh[1], Markus de
4 Raad[1], Tony Wildish[3,4], Daniel Udwary[3,4], Cindi Hoover[3], Samuel Deutsch[1,3],
5 Trent R. Northen[1,3,*], Benjamin P. Bowen[1,3,*]
7    1. Environmental Genomics and Systems Biology Division, Lawrence
8       Berkeley National Laboratory
9    2. Data Analytics and Visualization Group, Computational Research
10       Division, Lawrence Berkeley National Laboratory
11    3. Joint Genome Institute, Lawrence Berkeley National Laboratory
12    4. National Energy Research Scientific Computing Center, Lawrence
13       Berkeley National Laboratory
14

**Author Contributions**
16 OE, BPB, TRN conceived and designed the method
17 OE, BPB, OR, MT, TW, DWU developed the method
18 KBL, MdR, CAH conducted the experiments
19 OE, BPB, SD, TRN, KBL wrote the manuscript.
20

21 **\*Correspondence** should be addressed to bpbowen@lbl.gov and
22 trnorthen@lbl.gov

**Abstract**

Metabolomics is a widely used technology for obtaining direct measures of metabolic activities from diverse biological systems. However, ambiguous metabolite identifications are a common challenge and biochemical interpretation is often limited by incomplete and inaccurate genome-based predictions of enzyme activities (*i.e.* gene annotations). Metabolite, Annotation, and Gene Integration (MAGI) generates a metabolite-gene association score using a biochemical reaction network. This is calculated by a method that emphasizes consensus between metabolites and genes via biochemical reactions. To demonstrate the potential of this method, we applied MAGI to integrate sequence data and metabolomics data collected from *Streptomyces coelicolor* A3(2), an extensively characterized bacterium that produces diverse secondary metabolites. Our findings suggest that coupling metabolomics and genomics data by scoring consensus between the two increases the quality of both metabolite identifications and gene annotations in this organism. MAGI also made biochemical predictions for poorly annotated genes that were consistent with the extensive literature on this important organism. This limited analysis suggests that using metabolomics data has the potential to improve annotations in sequenced organisms and also provides testable hypotheses for specific biochemical functions. MAGI is freely available for academic use both as an online tool at https://magi.nersc.gov and with source code available at https://github.com/biorack/magi.

**Introduction**

Metabolomics approaches now enable global profiling, comparison, and discovery of diverse metabolites present in complex biological samples[1]. Liquid chromatography coupled with electrospray ionization mass spectrometry (LC-MS) is one of the leading methods in metabolomics[1]. A critical measure in metabolomics datasets is known as a "feature," which is a unique combination of mass-to-charge (*m/z*) and chromatographic retention time[1]. Each distinct feature may match to hundreds of unique chemical structures. This makes metabolite identification (the accurate assignment of the correct chemical structure to each feature) one of the fundamental challenges in metabolomics[2-4]. To aid metabolite identification efforts, ions (each with a unique *m/z* and retention time) are typically fragmented, and the resulting fragments are compared against either experimental[5, 6] or computationally predicted[5, 7-11] reference libraries. While this method is highly effective at reducing the search space for metabolite identification, misidentifications are inevitable, especially for metabolites lacking authentic standards.

One strategy for addressing the large search space of compound identifications is to assess identifications in the context of the predicted metabolism of the organism(s) being studied. Several tools do this with varying degrees of complexity with strategies including directly mapping metabolites onto reactions[12] or scoring the likelihood of metabolite identities using reaction networks and predictive pathway mapping[13]. However, many

73  metabolites cannot be included in these approaches. This is due to a number
74  of factors, including the low coverage in reaction databases[14, 15] (especially for
75  secondary metabolites[16-19]), incomplete or inaccurate set of reactions for an
76  organism, and enzyme promiscuity not being taken into account when
77  formulating the potential metabolism of an organism. To help address these
78  challenges computational strategies have been developed including
79  MyCompoundID [20, 21], IIMDB [22], MINES [23] and the ATLAS of biochemistry [24] to
80  enzymatically enlarge compound space similar to the retrosynthesis tools
81  such as Retrorules [25] and rePrime[26]. These approaches can be complimented
82  by chemical networking to help address the limited number of metabolites
83  represented in reactions, by expanding reaction space based on chemical or
84  spectral similarity between metabolites. Effectively, even when a metabolite
85  is not directly involved in a reaction, a linkage can still be made with a
86  reaction based on similarity to another well-studied metabolite[16-19, 27]. In this
87  way, chemical networking is a viable solution that expands reaction
88  databases to integrate with already expansive metabolite databases. This
89  allows more putative metabolite identifications to be assessed using the
90  predicted metabolism of the organism(s).
91
92  Recently, approaches have been developed that span the gap between
93  metabolomics and genomics and allow for some enzyme promiscuity. GNP,
94  developed specifically for discovering new nonribosomal peptides (NRPs) and
95  polyketides, uses a gene-forward strategy that predicts possible chemical
96  structures of NRP and polyketide synthases and generates a set of predicted
97  MS/MS spectra based on those predictions; these predictions are then used to
98  mine MS data [28]. Pep2Path, also developed exclusively for NRPs and post-
99  translationally modified peptides (RiPPs), takes a Bayesian approach to
100 scoring putative NRPs and RiPPs based on the gene sequences present in the
101 assayed organism [29]. Finally, a more general approach has been developed
102 where a mutant library of an organism is assayed for major differences in the
103 mass spectrometry profile, and the major differences are manually annotated
104 with human intuition [30].
105
106 Due to the vast amount of knowledge about Streptomyces species, they are
107 an excellent target for developing new tools for metabolite and genome
108 exploration. Representatives from this genus produce many antibacterials,
109 anticancer compounds, immunosuppresents, antifungals, cardiovascular
110 agents, and veterinary products including erythromycin, tetracycline,
111 doxorubicin, enediyenes, FK-506, rapamycin, avermectin, nemadectin,
112 amphotericin, griseofulvin, nystatin, lovastatin, compactin, monensin, and
113 tylosin [31]. Thus making them a highly relevant group for in depth studies to
114 link natural products with associated genes. In particular, *Streptomyces*
115 *coelicolor* is a model actinomycete secondary metabolite producer [32]; studies
116 from over three thousand papers and over 60 years of work [33] have produced,
117 among other things, a detailed understanding of the secondary metabolites
118 this organism produces, where two are the pigmented antibiotics:
119 actinorhodin and undecylprodiosin. These experiments have identified the
120 biochemical pathways, genes, and regulatory processes that are necessary
121 for producing the associated secondary metabolites [34].
122

Here we report Metabolite, Annotation, and Gene Integration (MAGI), an approach to generate metabolite-gene associations (Figure 1) by scoring consensus between metabolite identifications and gene annotations. MAGI is guided by the principles that the probability of a metabolite identity increases if there is genetic evidence to support that metabolite and that the probability of a gene function increases if there is metabolomic evidence for that function. Inputs to MAGI are typically a metabolite identification file of LCMS features and a protein or gene sequence FASTA file. For each LCMS feature, there are often many plausible metabolite identifications that can be given a probability based on accurate mass error and/or mass fragmentation comparisons. MAGI links these putative compound identifications to reactions both directly and indirectly by a biochemically relevant chemical similarity network. Likewise, MAGI associates input sequences to biochemical reactions by assessing sequence homology to reference sequences in the MAGI reaction database. For each sequence, there are often several plausible reactions with equal or similar probability. While annotation services would typically reduce specificity in these cases (*e.g., by* simply annotating as oxidoreductase), MAGI maintains all specific reactions as possibilities. Since MAGI links both metabolites and sequences to reactions with numerical scores that are proxies for probabilities, a final integrative MAGI score is calculated that magnifies consensus between a gene annotation and a metabolite identification. We applied this approach to one of the best characterized secondary metabolite producing bacteria, *Streptomyces coelicolor* A3(2)[35], by integrating its genome sequence with untargeted metabolomics data. MAGI successfully reduced the metabolite identity search space by scoring metabolite identities based on the predicted metabolism of an organism. Additionally, further investigation of the metabolite-gene associations led to identification of unannotated and misannotated genes that were subsequently validated using literature searches. This simple example illustrates the key aspects of MAGI.

**Methods**

**Media and culture conditions.** A 20 µL volume of glycerol stock of wild-type *S. coelicolor* spores was cultured in 40 mL R5 medium in a 250-mL flask. One liter of R5 medium base included 103 g sucrose, 0.25 g $K_2SO_4$, 10.12 g $MgCl_2 \cdot 6H_2O$, 10 g glucose, 0.1 g cas-amino acids, 2 mL trace element solution, 5 g yeast extract, and 5.73 g TES buffer to 1 L distilled water. After autoclave sterilization, 1 mL 0.5% $KH_2PO_4$, 0.4 mL 5M $CaCl_2 \cdot 2H_2O$, 1.5 mL 20% L-proline, 0.7 ml 1N NaOH were added as per the following protocol: https://www.elabprotocols.com/protocols/#!protocol=486. Each flask contained a stainless steel spring (McMaster-Carr Supply, part 9663K77), cut to fit in a circle in the bottom of the flask. The spring was used to prevent clumping of *S. coelicolor* during incubation. A foam stopper was used to close each flask (Jaece Industries Inc., Fisher part 14-127-40D). Four replicates of each sample were grown in a 28°C incubator with shaking at 150 rpm. On day six, 1 mL from each replicate were collected in 2 mL Eppendorf tubes in a sterile hood. Samples were centrifuged at 3,200 x g for 8 minutes at 4 °C to pellet the cells. Supernatants were decanted into fresh 2 mL tubes and frozen at -80 °C. Pellets were flash frozen on dry ice and then stored at -80 °C.

173
174 **LCMS sample preparation and data acquisition.** In preparation for LCMS,
175 medium samples were lyophilized. Dried medium was then extracted with
176 150 µL methanol containing an internal standard (2-Amino-3-bromo-5-
177 methylbenzoic acid, 1 µg/mL, Sigma, #631531), vortexed, sonicated in a
178 water bath for 10 minutes, centrifuged at 5,000 rpm for 5 min, and
179 supernatant finally centrifuge-filtered through a 0.22 µm PVDF membrane
180 (UFC40GV0S, Millipore). LC-MS/MS was performed in negative ion mode on a
181 2 µL injection, with UHPLC reverse phase chromatography performed using
182 an Agilent 1290 LC stack and Agilent C18 column (ZORBAX Eclipse Plus C18,
183 Rapid Resolution HD, 2.1 x 50 mm, 1.8 µm) at 60 °C and with MS and MS/MS
184 data collected using a QExactive Orbitrap mass spectrometer (Thermo
185 Scientific, San Jose, CA). Chromatography used a flow rate of 0.4 mL/min, first
186 equilibrating the column with 100% buffer A (LC-MS water with 0.1% formic
187 acid) for 1.5 min, then diluting over 7 minutes to 0% buffer A with buffer B
188 (100% acetonitrile with 0.1% formic acid). Full MS spectra were collected at
189 70,000 resolution from $m/z$ 80-1,200, and MS/MS fragmentation data
190 collected at 17,500 resolution using an average of 10, 20 and 30 eV collision
191 energies.
192
193 **Feature detection**. MZmine (version 2.23) [36] was used to deconvolute mass
194 spectrometry features. The methods and parameters used were as follows (in
195 the order that the methods were applied). MS/MS peaklist builder: retention
196 time between 0.5-13.0 minutes, m/z window of 0.01, time window of 1.00.
197 Peak extender: $m/z$ tolerance 0.01 $m/z$ or 50.0 ppm, min height of 1.0E0.
198 Chromatogram deconvolution: local minimum search algorithm where
199 chromatographic threshold was 1.0%, search minimum in RT range was 0.05
200 minutes, minimum relative height of 1.0%, minimum absolute height of
201 1.0E5, minimum ratio of peak top/edge of 1.2, peak duration between 0.01
202 and 30 minutes. Duplicate peak filter: m/z tolerance of 0.01 m/z or 50.0 ppm,
203 RT tolerance of 0.15 minutes. Isotopic peaks grouper: $m/z$ tolerance of 1.0E-6
204 m/z or 20.0 ppm, retention time tolerance of 0.01, maximum charge of 2,
205 representative isotope was lowest $m/z$. Adduct search: RT tolerance of 0.01
206 minutes, searching for adducts M+Hac-H, M+Cl, with an $m/z$ tolerance of
207 1.0E-5 $m/z$ or 20.0 ppm and max relative adduct peak height of 1.0%. Join
208 aligner: $m/z$ tolerance of 1.0E-6 $m/z$ or 50.0 ppm, weight for $m/z$ of 5,
209 retention time tolerance of 0.15 minutes, weight for RT of 3. Same RT and $m/z$
210 $z$ range gap filler: $m/z$ tolerance of 1.0E-6 $m/z$ or 20.0 ppm.
211
212 **Metabolite identification.** During the LCMS acquisition, two MS/MS
213 spectra were acquired for every MS spectrum. These MS/MS spectra are
214 acquired using data-dependent criteria in which the 2 most intense ions are
215 pursued for fragmentation, and then the next 2 most intense ions such that
216 no ion is fragmented more frequently than every 10 seconds. To assign
217 probable metabolite identities to a spectrum a modified version of the
218 previously described MIDAS approach was used[7]. Our metabolite database is
219 the merger of HMDB, MetaCyc, ChEBI, WikiData, GNPS, and LipidMaps
220 resulting in approximately 180,000 unique chemical structures. For each of
221 these structures, a comprehensive fragmentation tree was pre-calculated to
222 a depth of 5 bond-breakages; these trees were used to accelerate the MIDAS

223 scoring process.   The source code to generate trees and score spectra
224 against trees is available on GitHub (https://github.com/biorack/pactolus).
225 The following procedure was used in the MIDAS scoring.  Precursor *m/z* values
226 were neutralized by 1.007276 Da.  For each metabolite within 10 ppm of the
227 neutralized precursor mass, MS/MS ions were associated with nodes of the
228 fragmentation tree using a window of 0.01 Da using MS/MS neutralizations of
229 1.00727, 2.01510, and -0.00055, as described [7]. For metabolite-features of
230 interest discussed in the text, retention time, m/z, adduct, and fragmentation
231 pattern were used to define a Metabolite Atlas [37] library (Supplementary Data
232 1). For each metabolite, raw data was inspected manually using MZmine [36] to
233 rule out peak misidentifications due to adduct formation and in-source
234 degradation.

235 **MAGI biochemical reaction and reference sequence database.** The
236 MAGI biochemical reaction database was constructed by aggregating all
237 publicly available biochemical reactions in MetaCyc and RHEA biochemical
238 reaction databases [14, 15]. This reaction database currently includes 12,293
239 unique metabolite structures. Identical reactions were collapsed together by
240 calculating a "reaction InChI key," where the SMILES strings of all members
241 of a reaction were joined together, separated by a "." and converted to a
242 single InChI string through an RDkit (https://github.com/rdkit/rdkit) Mol
243 object, and then the InChI key was calculated also using RDKit. Biochemical
244 reactions with identical reaction InChI keys have identical chemical
245 metabolites, indicating they are duplicates, and were collapsed into one
246 database entry, retaining reference sequences. Reference sequences for
247 each biochemical reaction from each database were combined to create a set
248 of curated reference sequences for each biochemical reaction in the
249 database.
250
251 **Chemical Network.** In order to expand the chemical space beyond what is
252 in the biochemical reaction database, a chemical network was constructed to
253 relate all metabolites in the database to metabolites in biochemical reactions
254 by biochemical similarity.   In each molecule, 70 chemical features
255 (Supplementary Table 1) were located.   These features were defined
256 previously as being biochemically relevant [38].  The count of each feature was
257 stored as a vector for each molecule.  The Euclidean distance between two
258 vectors was used to determine similarity between two molecules and
259 construct a similarity network where every molecule is connected to every
260 molecule by the difference in their vectors.  This network was trimmed by
261 calculating a minimum-spanning tree based on frequency of biochemical
262 differences where more frequent differences would be preserved when
263 possible (Supplementary Data 2).
264
265 **Gene Annotations of *Streptomyces coelicolor*.** KEGG annotations were
266 obtained by submitting the S. coelicolor protein FASTA obtained from IMG to
267 the KEGG Automatic Annotation Server version 2.1 [39] and downloading the
268 gene-KO results table. KO numbers were associated with reactions by
269 assessing if there was a link to one or more KEGG reaction entries directly
270 from the webpage of that KO. For BioCyc annotations and reactions, the
271 BioCyc *S. coelicolor* database was downloaded. For the reactions in Table 1,

272  KEGG and BioCyc reactions were manually inspected and compared to MAGI
273  reactions.
274
275  **MAGI workflow.** An input metabolite structure is expanded to similar
276  metabolite structures as suggested by the chemical network and all
277  tautomers of those metabolites. Searching all tautomeric forms of a
278  metabolite structure is a known method to enhance metabolite database
279  searches [40]. Tautomers were generated by using the MolVS package. The
280  reaction database is then queried to find reactions containing these
281  metabolites or their tautomers. Direct matches are stereospecific, but
282  tautomer matches are not. This is due to limitations in the tautomer
283  generating method and in how the chemical network was constructed. The
284  metabolite score, $C$, is inherited from the MS/MS scoring algorithm and is a
285  proxy for the probability that a metabolite structure is correctly assigned. In
286  our case, it is the MIDAS score, but could be any score due to the use of the
287  geometric mean to calculate the MAGI score. The metabolite score is set to 1
288  as a default.
289
290  If the reaction has a reference sequence associated with it, this reference
291  sequence is used as a BLAST query against a sequence database of the input
292  gene sequences to find genes that may encode that reaction. The reciprocal
293  BLAST is also performed, where genes in the input gene sequences are
294  queries against the reaction reference sequence database; this finds the
295  reactions that a gene may encode for. The BLAST results are joined by their
296  common gene sequence and are used to calculate a homology score:
297  $H = F + R - |F - R|$ where $F$ and $R$ are log-transformed e-values of the BLAST
298  results (a proxy for the probability that two gene sequences are homologs),
299  with $F$ representing the reaction-to-gene BLAST score, and $R$ the gene-to-
300  reaction BLAST score. The homology score is set to 1 if no sequence is
301  matched.
302
303  The reciprocal agreement between both BLAST searches is also assessed,
304  namely whether they both agreed on the same reaction or not, formulating a
305  reciprocal agreement score: $\alpha$. $\alpha$ is equal to 2 for reciprocal agreements, 1 for
306  disagreements that had BLAST score within 75% of the larger score, 0.01 for
307  disagreements with very different BLAST scores, and 0.1 for situations where
308  one of the BLAST searches did not yield any results. For cases where
309  metabolites are linked to reactions but there is not a reference protein
310  sequence available, a weight factor, $X$, is needed. We chose, $X$, such that: i)
311  X=0.01 when a metabolite is not in any reaction; ii) X=1.01 when a
312  metabolite is in reaction missing a reference sequence; and iii) X=2.01 when
313  a metabolite is in a reaction with a sequence. These arbitrary scores were
314  selected solely to distinguish between different "agreement" states during
315  the reciprocal BLAST. We did not observe much difference in the plurality of
316  compound annotations depending on these weights (data not shown),
317  however, they did have an impact on the number of annotations that agreed
318  with KEGG and MetaCyc (Supplementary Figure 4). The most impactful
319  weight appeared to be the "close reciprocal disagreement," meaning that
320  there was not an exact match in the bidirectional BLAST, but the e-scores

321 were within the given threshold. If this weight was low (0.01) or high (2.0),
322 there were fewer annotations that agreed with KEGG and MetaCyc.
323
324 The final MAGI-score $M = \mathrm{GM}([C, H, a, X])/n^{L}$ is a proxy for the probability that
325 a gene and metabolite are associated. $M$ is generated by calculating the
326 geometric mean ($GM$) of the metabolite score ($C$), homology score ($H$),
327 reciprocal agreement score ($\alpha$) and weight factor ($X$), and whether or not the
328 metabolite is present in a reaction ($n^{L}$) where $L$ is the network level
329 connecting the metabolite to a reaction (a proxy for the probability that a
330 compound is involved in a reaction) and $n$ is a penalty factor for the network
331 level. Currently, n is equal to 4, but this parameter may change as the
332 scoring function is optimized and more training data is acquired. We did not
333 observe this penalty factor to greatly affect the number of gene annotations
334 that agreed with KEGG or MetaCyc, though this did have a large impact on
335 the number of features with multiple suggestions for compound identities;
336 the higher the penalty factor, the lower the number of compound
337 suggestions. MAGI often gives a high score to multiple metabolites, which is
338 not surprising given the relatedness of many metabolites (e.g. isomers).
339 Therefore, we recommend carefully considering the top scoring molecules
340 and not assuming that the top ranked one is correct (Supplementary Figure
341 5). Additional benchmarking analysis shows agreement between KEGG,
342 BioCyc and MAGI annotations for high MAGI homology scores (Supplementary
343 Discussion and SI Figures 2 and 3). The geometric mean was used  to account
344 for the different scales of the individual scores, but weights may be applied to
345 each individual score during the geometric mean calculation to further fine-
346 tune the MAGI scoring process. We expect the weights to become further
347 optimized as more results are processed through MAGI.
348
349 The final output is a table representing all unique metabolite-reaction-gene
350 associations, their individual scores, and their integrated MAGI score
351 (Supplementary Table 2). For scoring metabolite identities, a slice of this final
352 output is created by retaining the top scoring metabolite-reaction-gene
353 association for each unique metabolite structure; these can be mapped back
354 onto the mass spectrometry results table to aid the identification of each
355 mass spectrometry feature. For assessing gene functions, another slice of
356 this final output is created by retaining the top scoring metabolite-reaction-
357 gene association for each unique gene-reaction pair. For a typical bacterial
358 genome of ~ 6000 genes and a metabolites file of ~ 6000 compounds, the
359 MAGI calculation performed via the web service at https://magi.nersc.gov/
360 should take about thirty minutes to complete. While MAGI can provide
361 valuable insights into primary metabolism, these reactions tend to be better
362 characterized and therefore a particularly important application of MAGI is for
363 secondary metabolite pathways.
364
365 **Data Availability**
366 All source code is available at https://github.com/biorack/magi, and the *S.*
367 *coelicolor* mass spectrometry data (.mzML files) and MIDAS results
368 (metabolite_0ae82b08.csv) can be found here: https://magi.nersc.gov/jobs/?
369 id=0ae82b08-b2a3-40d8-bb9a-e64b567cacd2.
370

**Application Availability and Usage**
Potential MAGI users may use the application on their personal computers by downloading the source code from the GitHub repository, or may upload their data files to the web service. In order to use MAGI, users must provide at least one of the following: a FASTA file of genes they wish to be associated to reactions, and/or a metabolites file they wish to associate to reactions. The metabolites file should be in a table file format (*e.g.* .csv, .tsv, Excel), and must have a column named "original_compound" that describes the InChI Key for each metabolite of interest. If both FASTA and metabolite files are provided, then  associations between genes and metabolites will be made as well.

**Results and Discussion**

**Improved metabolite identification for metabolomics.** To examine how MAGI uses genomic information to filter and score possible metabolite identities from a metabolomics experiment, sequencing and metabolomics data were obtained for *S. coelicolor*. After processing the raw LCMS data to find chromatograms and peaks, 878 features with a unique *m/z* and retention time were found in the dataset. After neutralizing the *m/z* values, accurate mass searching, and conducting MS/MS fragmentation pattern analysis, 6,604 unique metabolite structures were tentatively associated with these features (Supplementary Table 3). This means on average there were almost 8 candidate structures for each feature. For a candidate structure to be associated with a feature, it must have at least one matching fragmentation spectrum. As this is often the method for identifying metabolites, it highlights the problem in deconvolution of a signal to a specific chemical structure. 2,786 of these structures were then linked to a total of 10,265 reactions either directly or via the chemical similarity network, and the reactions were associated with 3,181 (out of 8,210) *S. coelicolor* genes by homology. Finally, a MAGI score was calculated for each metabolite-reaction-gene association (Supplementary Table 4).

An example that illustrates MAGI's utility in determining the most likely correct metabolite identification is the feature putatively identified as 1,4-dihydroxy-6-naphthoic acid.  Here, a feature with an *m/z* of 203.0345 was observed. This feature was associated with the chemical formula $C_{11}H_8O_4$, which could be derived from 16 unique chemical structures in the metabolite database (Supplementary Table 5). Mass fragmentation spectra were collected for this feature and analyzed using MIDAS[7], a tool that scores the observed fragmentation spectrum against its database of *in-silico* fragmentation trees for the 16 potential structures.  Based only on the MIDAS metabolite score, the top scoring structure was 5,6-dihydroxy-2-methylnaphthalene-1,4-dione. However, after calculating the MAGI scores, a different metabolite received the highest score. Of the 16 potential metabolites, only 1,4-dihydroxy-6-naphthoic acid was in a reaction that had a perfect match to genes in *S. coelicolor* (an E-value of 0.0 to *SCO4326*; Table 1). This metabolite is a known intermediate in an alternative menaquinone biosynthesis pathway discovered in *S. coelicolor*[41, 42], making it much more

420   likely to be a metabolite detected from the metabolome of *S. coelicolor* as
421   opposed to the metabolite found by looking at mass fragmentation alone.
422
423   **Metabolomics-driven gene annotations.** MAGI keeps the biochemical
424   potential of an organism unconstrained by considering a plurality of probable
425   gene product functions. One effect of this is that more reactions are
426   associated with genes than other services (Figure 2A). Because reactions are
427   the pivotal link between metabolites and genes, this allows integration of a
428   larger fraction of a metabolomics dataset with genes. Furthermore, MAGI
429   associates many genes that are not annotated using traditional approaches
430   with at least one reaction (Figure 2B). Out of a total of 8,210 predicted coding
431   sequences in *S. coelicolor*, KEGG and BioCyc have one or more reactions
432   associated with 1,106 and 1,294 genes, respectively. On the other hand,
433   MAGI associated 5,209 genes with one or more reactions, out of which 3,719
434   genes had no reaction associated with them in either KEGG or BioCyc (Figure
435   2B). Of these 3,719 genes, 1,883 were linked to at least one metabolite in the
436   metabolomics data (Supplementary Table 4). Certainly, not all MAGI gene-
437   reaction associations are correct, however, this does provide many testable
438   hypotheses that give footholds to discover new biochemistry As can be seen
439   in Figure 2C, many of these new gene-reaction associations have high scores,
440   indicating a likely connection.
441
442   **Validation of gene-metabolite integration in pathways.** One of the
443   most well-known biosynthetic pathways in *S. coelicolor* is the pathway to
444   synthesize the pigmented antibiotic actinorhodin[35]. We examined the MAGI
445   results involving the metabolites and genes of actinorhodin biosynthesis as a
446   proof-of-principle that MAGI successfully integrates metabolites and genes,
447   and that these results can be mapped onto a reaction network. Actinorhodin
448   and all of its detected intermediates were correctly identified and accurately
449   mapped to the correct genes (Figure 3A), despite some intermediates having
450   several plausible metabolite identities (Supplementary Table 6). Notably,
451   KEGG did not annotate the majority of actinorhodin biosynthesis genes, and
452   the one gene that it did annotate was incorrect (Table 1).
453
454   In another example, we examined the menaquinone biosynthesis pathway,
455   which is essential for respiration in bacteria[43] and thus should be included in
456   every metabolic reconstruction for organisms that produce menaquinone. An
457   alternative menaquinone biosynthesis pathway was recently discovered and
458   validated in *S. coelicolor*[41, 42], serving as another proof-of-principle exercise for
459   assessing the MAGI platform. MAGI linked 4 of 7 intermediate metabolites of
460   the pathway to the appropriate genes (Figure 3B, Supplementary Table 7).
461   Interestingly, while KEGG accurately assigned reactions to all but one of the
462   genes in this biosynthetic pathway, BioCyC had vague textual annotations
463   and no reactions (Table 1). Therefore, a metabolomics tool that relies on
464   BioCyc model for *S. coelicolor* would be unable to integrate any of these
465   metabolites with genes for the purpose of either improved metabolite
466   identifications or gene annotations.
467
468   **Correction of annotation errors.** Gene annotation pipelines are
469   notoriously error-prone[44] and yield inconsistent results based on the

470 bioinformatic analyses used: the database used for homology searches, and
471 what kind of additional data (*e.g.* PFams, genetic neighborhoods, and
472 literature mining) are incorporated into the annotation algorithm or not (see
473 Table 1 for some examples). For example, the undecylprodigiosin synthase
474 gene is known[45], yet was incorrectly annotated in the KEGG genome
475 annotation for *S. coelicolor*. KEGG annotated this gene as "PEP utilizing
476 enzyme" with an EC number of 2.7.9.2 (pyruvate, water phosphotransferase
477 with paired electron acceptors). This error is notable because the
478 undecylprodigiosin synthase reaction has an EC number of 6.4.1.-: ligases
479 that form carbon-carbon bonds. On the other hand, BioCyc correctly
480 annotates *SCO5896* as undecylprodigiosin synthase, presumably using
481 manual curation or a thorough literature-searching algorithm.
482
483 MAGI used metabolomics data to score the possible gene annotations for
484 *SCO5896* in addition to homology scoring (*i.e.* E-value). In the absence of
485 metabolomics data, MAGI initially associated the *SCO5896* gene sequence
486 with the prodigiosin synthase and norprodigiosin synthase reactions via
487 BLAST searches against the MAGI reaction reference sequence database
488 (Figure 4). Metabolomics analysis revealed that the feature with an *m/z* of
489 392.2720 could potentially be undecylprodigiosin, which MAGI associated
490 with only the undecylprodigiosin synthase reaction (Figure 4). Because this
491 reaction does not have a reference sequence in our database, it could not be
492 queried against the *S. coelicolor* genome. However, the chemical network
493 revealed that prodigiosin is a similar metabolite that is in a reaction that does
494 have a reference sequence (Figure 4). When the prodigiosin synthase
495 reaction's reference sequence was queried against the *S. coelicolor* genome,
496 the top hit was *SCO5896*, thus making a reciprocal connection between the
497 mass spectrometry feature and gene via the prodigiosin synthase reaction
498 (Figure 4).
499
500 **Making nonexistent or vague annotations specific.** The vast majority of
501 sequenced genes have no discrete functional predictions, preventing the in-
502 depth understanding of metabolic processes of most organisms. *S. coelicolor*
503 is well known to produce several polyketides and is known to have the
504 genetic potential to produce many more. The *SCO5315* gene product is WhiE,
505 a known polyketide aromatase involved in the biosynthesis of a white
506 pigment characteristic of *S. coelicolor*[46, 47]. KEGG and BioCyC textually
507 annotated the gene as "aromatase" or "polyketide aromatase," but neither
508 links the gene to a discrete reaction. Although the text annotations are
509 correct, the lack of a biochemical reaction prohibits the association of this
510 gene with metabolites. On the other hand, MAGI was successfully able to
511 associate *SCO5315* with an observed metabolite (20-carbon polyketide
512 intermediate with an *m/z* of 401.0887) via a polyketide cyclization reaction
513 with a MAGI consensus score of 4.59 (Table 1). While the physiological
514 function of WhiE is to cyclize a 24-carbon polyketide intermediate, the
515 enzyme has been shown to also catalyze the cyclization of similar polyketides
516 with varying chain length, including the 20-carbon species observed in the
517 metabolomics data presented here[48-50].
518

519    In another example where other annotation services were unable to assign
520    any reactions to a gene product, MAGI associated *SCO7595* with the anhydro-
521    NAM kinase reaction via the detected metabolite anhydro-N-acetylmuramic
522    acid (anhydro-NAM) (*m/z* 274.0941) (Table 1). Anhydro-NAM is an
523    intermediate in bacterial cell wall recycling, a critically important and
524    significant metabolic process in actively growing bacterial cells; *E. coli* and
525    other bacteria were observed to recycle roughly half of cell wall components
526    per generation[51, 52]. MAGI also associated anhydro-NAM to *SCO6300* via an
527    acetylhexosaminidase reaction (Table 1) that produces the metabolite. KEGG
528    and RAST both annotate this gene to be acetylhexosaminidase with a total of
529    5 possible reactions, but none involve anhydro-NAM (Table 1). The detection
530    of anhydro-NAM may be considered orthogonal experimental evidence to
531    indicate that *SCO6300* can act on N-acetyl-β-D-glucosamine-anhydro-NAM
532    along with the other acetylhexosamines predicted by KEGG and RAST,
533    forming an early stage in anhydromurpoeptide recycling. In the absence of
534    MAGI, a researcher may have been able to manually curate a metabolic
535    model by manually assessing the text annotations and adding reactions to
536    the model, but the MAGI framework not only makes this process easier, it
537    also connects an experimental observation that supports the predicted
538    function of the gene.
539
540    **Potential for making novel annotations.** In addition to these few
541    examples, there are hundreds more gene-reaction-metabolite associations
542    that could be used to strengthen, validate, or correct existing annotations
543    from KEGG or BioCyc, as well as discover new annotations through
544    experimentation. These MAGI associations can be sorted by their MAGI score
545    to generate a ranked list of candidate genes and gene functions, with
546    optional hierarchical grouping and filtering of the list by homology,
547    metabolite, chemical network, and/or reciprocal score. For example, of the
548    1,883 *S. coelicolor* genes that were uniquely linked to a metabolite via a
549    reaction by MAGI, roughly one-third were connected directly to a metabolite;
550    that is, the chemical similarity network was not used to expand reaction
551    space (Figure 5A and Figure 2C teal markers). Furthermore, one-third of
552    these genes had perfect reciprocal agreement between the metabolite-to-
553    gene and gene-to-metabolite search directions (Figure 5B and Figure 2C teal
554    circles). These 190 genes can be further separated or binned based on their
555    homology score or MAGI score (Figure 5C), resulting in an actionable number
556    of high-priority and high-strength novel gene function hypotheses to test in
557    future studies.
558
559    **Limitations of this study.** In this study, we show that MAGI produces
560    plausible associations between genes and metabolites from *Streptomyces*
561    *coelicolor*. Since the associations shown in this paper are judged by manual
562    inspection, there are not enough validated links to compute a reliable false
563    discovery rate or applicability to other systems. Therefore an important
564    future work will be to broadly apply MAGI across many organisms and
565    evaluate the generality of this approach. This will ensure that the parameters
566    used are not over fit specifically to *Streptomyces coelicolor*. In addition, given
567    the paucity of direct biochemical validations of gene functions, it will likely be
568    necessary to integrate MAGI with high throughput mutagenesis studies to

569 accurately determine false discovery rates. Lastly, more unique metabolites
570 can be observed by combining data collected from polar and lipid fractions of
571 metabolites along with combining positive and negative ionization modes.
572 The results here are based on measured signals from a small subset of the
573 Streptomyces coelicolor metabolome.
574
575

576 **Conclusion**

577
578 In this work we describe MAGI, a method for integrating metabolomics
579 observations with genomic predictions to help overcome the limitations of
580 each and strengthen the biological conclusions made by both. Using
581 *Streptomyces coelicolor* as a test case, we find that this method can help
582 strengthen metabolite identifications, suggests specific biochemical
583 predictions about genes that may otherwise be ambiguous, and suggests
584 new biochemistry via the chemical network. It will be important to also
585 evaluate this approach for diverse organisms to determine the generality of
586 the method. In order to facilitate broad usage by the academic community,
587 we provide MAGI through the National Energy Research Scientific Computing
588 Center (NERSC) at https://magi.nersc.gov, where users can upload their own
589 metabolite and FASTA files for analysis through MAGI.
590

**References**

1. Liu, X. J., and Locasale, J. W. (2017) Metabolomics: A Primer, *Trends in Biochemical Sciences 42*, 274-284.
2. Creek, D. J., Dunn, W. B., Fiehn, O., Griffin, J. L., Hall, R. D., Lei, Z. T., Mistrik, R., Neumann, S., Schymanski, E. L., Sumner, L. W., Trengove, R., and Wolfender, J. L. (2014) Metabolite identification: are you sure? And how do your peers gauge your confidence?, *Metabolomics 10*, 350-353.
3. Wolfender, J. L., Marti, G., Thomas, A., and Bertrand, S. (2015) Current approaches and challenges for the metabolite profiling of complex natural extracts, *J Chromatogr A 1382*, 136-164.
4. Vaniya, A., and Fiehn, O. (2015) Using fragmentation trees and mass spectral trees for identifying unknown compounds in metabolomics, *Trac-Trend Anal Chem 69*, 52-61.
5. Smith, C. A., O'Maille, G., Want, E. J., Qin, C., Trauger, S. A., Brandon, T. R., Custodio, D. E., Abagyan, R., and Siuzdak, G. (2005) METLIN: a metabolite mass spectral database, *Ther Drug Monit 27*, 747-751.
6. Horai, H., Arita, M., Kanaya, S., Nihei, Y., Ikeda, T., Suwa, K., Ojima, Y., Tanaka, K., Tanaka, S., Aoshima, K., Oda, Y., Kakazu, Y., Kusano, M., Tohge, T., Matsuda, F., Sawada, Y., Hirai, M. Y., Nakanishi, H., Ikeda, K., Akimoto, N., Maoka, T., Takahashi, H., Ara, T., Sakurai, N., Suzuki, H., Shibata, D., Neumann, S., Iida, T., Tanaka, K., Funatsu, K., Matsuura, F., Soga, T., Taguchi, R., Saito, K., and Nishioka, T. (2010) MassBank: a public repository for sharing mass spectral data for life sciences, *J Mass Spectrom 45*, 703-714.
7. Wang, Y., Kora, G., Bowen, B. P., and Pan, C. (2014) MIDAS: a database-searching algorithm for metabolite identification in metabolomics, *Anal Chem 86*, 9496-9503.
8. Wolf, S., Schmidt, S., Muller-Hannemann, M., and Neumann, S. (2010) In silico fragmentation for computer assisted identification of metabolite mass spectra, *BMC Bioinformatics 11*, 148.
9. Allen, F., Greiner, R., and Wishart, D. (2015) Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification, *Metabolomics 11*, 98-110.
10. Ridder, L., van der Hooft, J. J. J., Verhoeven, S., de Vos, R. C. H., Bino, R. J., and Vervoort, J. (2013) Automatic Chemical Structure Annotation of an LC-MSn Based Metabolic Profile from Green Tea, *Analytical Chemistry 85*, 6033-6040.
11. Duhrkop, K., Shen, H. B., Meusel, M., Rousu, J., and Bocker, S. (2015) Searching molecular structure databases with tandem mass spectra using CSI:FingerID, *Proceedings of the National Academy of Sciences of the United States of America 112*, 12580-12585.
12. Dhanasekaran, A. R., Pearson, J. L., Ganesan, B., and Weimer, B. C. (2015) Metabolome searcher: a high throughput tool for metabolite identification and metabolic pathway mapping directly from mass spectrometry and using genome restriction, *Bmc Bioinformatics 16*.

644    13. Li, S. Z., Park, Y., Duraisingham, S., Strobel, F. H., Khan, N., Soltow, Q. A., Jones,
645            D. P., and Pulendran, B. (2013) Predicting Network Activity from High
646            Throughput Metabolomics, *Plos Computational Biology 9*.
647    14. Caspi, R., Billington, R., Ferrer, L., Foerster, H., Fulcher, C. A., Keseler, I. M.,
648            Kothari, A., Krummenacker, M., Latendresse, M., Mueller, L. A., Ong, Q., Paley,
649            S., Subhraveti, P., Weaver, D. S., and Karp, P. D. (2016) The MetaCyc database
650            of metabolic pathways and enzymes and the BioCyc collection of
651            pathway/genome databases, *Nucleic Acids Research 44*, D471-D480.
652    15. Morgat, A., Lombardot, T., Axelsen, K. B., Aimo, L., Niknejad, A., Hyka-Nouspikel,
653            N., Coudert, E., Pozzato, M., Pagni, M., Moretti, S., Rosanoff, S., Onwubiko, J.,
654            Bougueleret, L., Xenarios, I., Redaschi, N., and Bridge, A. (2017) Updates in
655            Rhea - an expert curated resource of biochemical reactions, *Nucleic Acids
656            Research 45*, D415-D418.
657    16. Yang, J. Y., Sanchez, L. M., Rath, C. M., Liu, X. T., Boudreau, P. D., Bruns, N.,
658            Glukhov, E., Wodtke, A., de Felicio, R., Fenner, A., Wong, W. R., Linington, R.
659            G., Zhang, L. X., Debonsi, H. M., Gerwick, W. H., and Dorrestein, P. C. (2013)
660            Molecular Networking as a Dereplication Strategy, *J Nat Prod 76*, 1686-1699.
661    17. Hadadi, N., Hafner, J., Shajkofci, A., Zisaki, A., and Hatzimanikatis, V. (2016)
662            ATLAS of Biochemistry: A Repository of All Possible Biochemical Reactions for
663            Synthetic Biology and Metabolic Engineering Studies, *Acs Synthetic Biology 5*,
664            1155-1166.
665    18. Hatzimanikatis, V., Li, C. H., Ionita, J. A., Henry, C. S., Jankowski, M. D., and
666            Broadbelt, L. J. (2005) Exploring the diversity of complex metabolic networks,
667            *Bioinformatics 21*, 1603-1609.
668    19. Li, C. H., Henry, C. S., Jankowski, M. D., Ionita, J. A., Hatzimanikatis, V., and
669            Broadbelt, L. J. (2004) Computational discovery of biochemical routes to
670            specialty chemicals, *Chem Eng Sci 59*, 5051-5060.
671    20. Li, L., Li, R., Zhou, J., Zuniga, A., Stanislaus, A. E., Wu, Y., Huan, T., Zheng, J., Shi,
672            Y., Wishart, D. S., and Lin, G. (2013) MyCompoundID: using an evidence-based
673            metabolome library for metabolite identification, *Anal Chem 85*, 3401-3408.
674    21. Huan, T., Tang, C., Li, R., Shi, Y., Lin, G., and Li, L. (2015) MyCompoundID
675            MS/MS Search: Metabolite Identification Using a Library of Predicted Fragment-
676            Ion-Spectra of 383,830 Possible Human Metabolites, *Anal Chem 87*, 10619-
677            10626.
678    22. Menikarachchi, L. C., Hill, D. W., Hamdalla, M. A., Mandoiu, II, and Grant, D. F.
679            (2013) In silico enzymatic synthesis of a 400,000 compound biochemical
680            database for nontargeted metabolomics, *J Chem Inf Model 53*, 2483-2492.
681    23. Jeffryes, J. G., Colastani, R. L., Elbadawi-Sidhu, M., Kind, T., Niehaus, T. D.,
682            Broadbelt, L. J., Hanson, A. D., Fiehn, O., Tyo, K. E., and Henry, C. S. (2015)
683            MINEs: open access databases of computationally predicted enzyme promiscuity
684            products for untargeted metabolomics, *J Cheminform 7*, 44.
685    24. Hadadi, N., Hafner, J., Shajkofci, A., Zisaki, A., and Hatzimanikatis, V. (2016)
686            ATLAS of Biochemistry: A Repository of All Possible Biochemical Reactions for
687            Synthetic Biology and Metabolic Engineering Studies, *ACS Synth Biol 5*, 1155-
688            1166.

689 25. Duigou, T., du Lac, M., Carbonell, P., and Faulon, J. L. (2018) RetroRules: a
690         database of reaction rules for engineering biology, *Nucleic Acids Res*.
691 26. Kumar, A., Wang, L., Ng, C. Y., and Maranas, C. D. (2018) Pathway design using de
692         novo steps through uncharted biochemical spaces, *Nat Commun 9*, 184.
693 27. Hattori, M., Tanaka, N., Kanehisa, M., and Goto, S. (2010) SIMCOMP/SUBCOMP:
694         chemical structure search servers for network analyses, *Nucleic Acids Res 38*,
695         W652-656.
696 28. Johnston, C. W., Skinnider, M. A., Wyatt, M. A., Li, X., Ranieri, M. R., Yang, L.,
697         Zechel, D. L., Ma, B., and Magarvey, N. A. (2015) An automated Genomes-to-
698         Natural Products platform (GNP) for the discovery of modular natural products,
699         *Nat Commun 6*, 8421.
700 29. Medema, M. H., Paalvast, Y., Nguyen, D. D., Melnik, A., Dorrestein, P. C., Takano,
701         E., and Breitling, R. (2014) Pep2Path: automated mass spectrometry-guided
702         genome mining of peptidic natural products, *PLoS Comput Biol 10*, e1003822.
703 30. Sevin, D. C., Fuhrer, T., Zamboni, N., and Sauer, U. (2017) Nontargeted in vitro
704         metabolomics for high-throughput identification of novel enzymes in Escherichia
705         coli, *Nat Methods 14*, 187-194.
706 31. Lamb, D. C., Guengerich, F. P., Kelly, S. L., and Waterman, M. R. (2006) Exploiting
707         Streptomyces coelicolor A3(2) P450s as a model for application in drug
708         discovery, *Expert Opin Drug Metab Toxicol 2*, 27-40.
709 32. Chater, K. F. (2016) Recent advances in understanding Streptomyces, *F1000Res 5*,
710         2795.
711 33. Worthen, D. B. (2008) Streptomyces in Nature and Medicine: The Antibiotic Makers,
712         *Journal of the History of Medicine and Allied Sciences 63*, 273-274.
713 34. Craney, A., Ahmed, S., and Nodwell, J. (2013) Towards a new science of secondary
714         metabolism, *J Antibiot (Tokyo) 66*, 387-400.
715 35. Craney, A., Ahmed, S., and Nodwell, J. (2013) Towards a new science of secondary
716         metabolism, *J Antibiot 66*, 387-400.
717 36. Pluskal, T., Castillo, S., Villar-Briones, A., and Oresic, M. (2010) MZmine 2:
718         modular framework for processing, visualizing, and analyzing mass spectrometry-
719         based molecular profile data, *BMC Bioinformatics 11*, 395.
720 37. Bowen, B. P., and Northen, T. R. (2010) Dealing with the unknown: metabolomics
721         and metabolite atlases, *J Am Soc Mass Spectrom 21*, 1471-1476.
722 38. Hattori, M., Okuno, Y., Goto, S., and Kanehisa, M. (2003) Development of a
723         chemical structure comparison method for integrated analysis of chemical and
724         genomic information in the metabolic pathways, *J Am Chem Soc 125*, 11853-
725         11865.
726 39. Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A. C., and Kanehisa, M. (2007) KAAS:
727         an automatic genome annotation and pathway reconstruction server, *Nucleic
728         Acids Res 35*, W182-185.
729 40. Oellien, F., Cramer, J., Beyer, C., Ihlenfeldt, W. D., and Selzer, P. M. (2006) The
730         impact of tautomer forms on pharmacophore-based virtual screening, *J Chem Inf
731         Model 46*, 2342-2354.

732  41. Hiratsuka, T., Furihata, K., Ishikawa, J., Yamashita, H., Itoh, N., Seto, H., and Dairi,
733      T. (2008) An alternative menaquinone biosynthetic pathway operating in
734      microorganisms, *Science 321*, 1670-1673.
735  42. Mahanta, N., Fedoseyenko, D., Dairi, T., and Begley, T. P. (2013) Menaquinone
736      Biosynthesis: Formation of Aminofutalosine Requires a Unique Radical SAM
737      Enzyme, *Journal of the American Chemical Society 135*, 15318-15321.
738  43. Nowicka, B., and Kruk, J. (2010) Occurrence, biosynthesis and function of isoprenoid
739      quinones, *Bba-Bioenergetics 1797*, 1587-1605.
740  44. Schnoes, A. M., Brown, S. D., Dodevski, I., and Babbitt, P. C. (2009) Annotation
741      Error in Public Databases: Misannotation of Molecular Function in Enzyme
742      Superfamilies, *Plos Computational Biology 5*.
743  45. Haynes, S. W., Sydor, P. K., Stanley, A. E., Song, L. J., and Challis, G. L. (2008)
744      Role and substrate specificity of the Streptomyces coelicolor RedH enzyme in
745      undecylprodiginine biosynthesis, *Chem Commun*, 1865-1867.
746  46. Shen, Y. M., Yoon, P., Yu, T. W., Floss, H. G., Hopwood, D., and Moore, B. S.
747      (1999) Ectopic expression of the minimal whiE polyketide synthase generates a
748      library of aromatic polyketides of diverse sizes and shapes, *Proceedings of the
749      National Academy of Sciences of the United States of America 96*, 3622-3627.
750  47. Yu, T. W., Shen, Y. M., McDaniel, R., Floss, H. G., Khosla, C., Hopwood, D. A., and
751      Moore, B. S. (1998) Engineered biosynthesis of novel polyketides from
752      Streptomyces spore pigment polyketide synthases, *Journal of the American
753      Chemical Society 120*, 7749-7759.
754  48. Alvarez, M. A., Fu, H., Khosla, C., Hopwood, D. A., and Bailey, J. E. (1996)
755      Engineered biosynthesis of novel polyketides: Properties of the whiE
756      aromatase/cyclase, *Nature Biotechnology 14*, 335-338.
757  49. Mcdaniel, R., Hutchinson, C. R., and Khosla, C. (1995) Engineered Biosynthesis of
758      Novel Polyketides - Analysis of Tcmn Function in Tetracenomycin Biosynthesis,
759      *Journal of the American Chemical Society 117*, 6805-6810.
760  50. Ames, B. D., Korman, T. P., Zhang, W. J., Smith, P., Vu, T., Tang, Y., and Tsai, S. C.
761      (2008) Crystal structure and functional analysis of tetracenomycin ARO/CYC:
762      Implications for cyclization specificity of aromatic polyketides, *Proceedings of
763      the National Academy of Sciences of the United States of America 105*, 5349-
764      5354.
765  51. Park, J. T., and Uehara, T. (2008) How bacteria consume their own exoskeletons
766      (Turnover and recycling of cell wall peptidoglycan), *Microbiology and Molecular
767      Biology Reviews 72*, 211-227.
768  52. Johnson, J. W., Fisher, J. F., and Mobashery, S. (2013) Bacterial cell-wall recycling,
769      *Ann Ny Acad Sci 1277*, 54-75.
770  53. Cooper, L. E., Fedoseyenko, D., Abdelwahed, S. H., Kim, S. H., Dairi, T., and
771      Begley, T. P. (2013) In Vitro Reconstitution of the Radical S-
772      Adenosylmethionine Enzyme MqnC Involved in the Biosynthesis of Futalosine-
773      Derived Menaquinone, *Biochemistry 52*, 4592-4594.
774  54. Ichinose, K., Surti, C., Taguchi, T., Malpartida, F., Booker-Milburn, K. I.,
775      Stephenson, G. R., Ebizuka, Y., and Hopwood, D. A. (1999) Proof that the actVI
776      genetic region of Streptomyces coelicolor A3(2) is involved in stereospecific

777     pyran ring formation in the biosynthesis of actinorhodin, *Bioorganic & Medicinal*
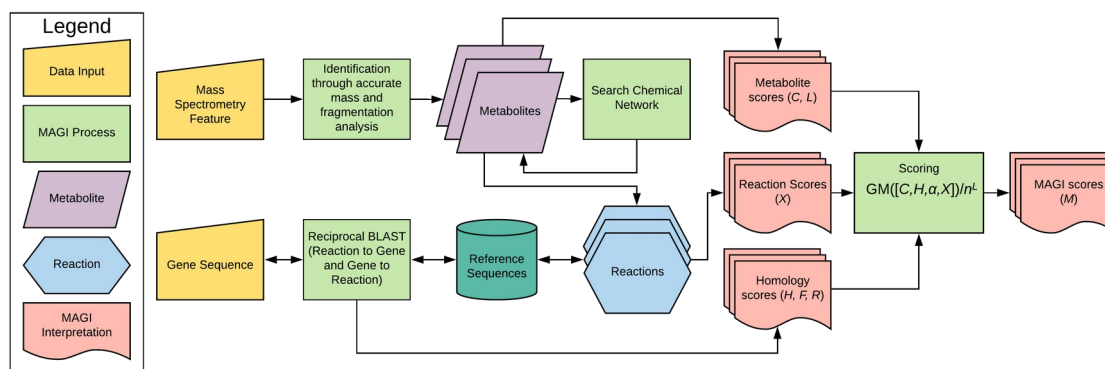778         *Chemistry Letters 9*, 395-400.
779  55. Taguchi, T., Itou, K., Ebizuka, Y., Malpartida, F., Hopwood, D. A., Surti, C. M.,
780         Booker-Milburn, K. I., Stephenson, G. R., and Ichinose, K. (2000) Chemical
781         characterisation of disruptants of the Streptomyces coelicolor A3(2) actVI genes
782         involved in actinorhodin biosynthesis, *J Antibiot 53*, 144-152.
783  56. Valton, J., Filisetti, L., Fontecave, M., and Niviere, V. (2004) A two-component
784         flavin-dependent monooxygenase involved in actinorhodin biosynthesis in
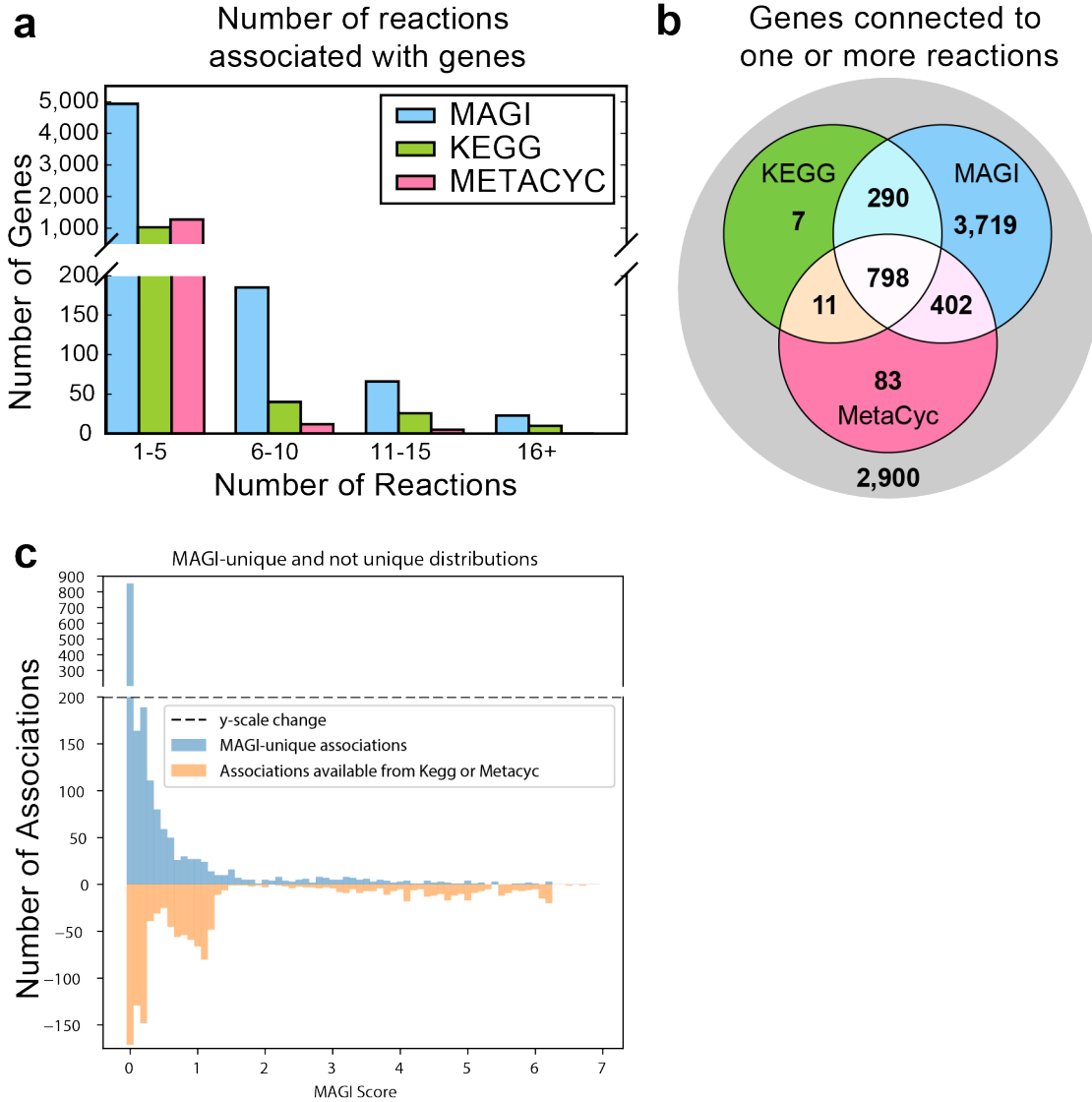785         Streptomyces coelicolor, *Journal of Biological Chemistry 279*, 44362-44369.
786  57. Kendrew, S. G., Hopwood, D. A., and Marsh, E. N. G. (1997) Identification of a
787         monooxygenase from Streptomyces coelicolor A3(2) involved in biosynthesis of
788         actinorhodin: Purification and characterization of the recombinant enzyme,
789         *Journal of Bacteriology 179*, 4305-4310.
790  58. Mcdaniel, R., Ebertkhosla, S., Fu, H., Hopwood, D. A., and Khosla, C. (1994)
791         Engineered Biosynthesis of Novel Polyketides - Influence of a Downstream
792         Enzyme on the Catalytic Specificity of a Minimal Aromatic Polyketide Synthase,
793         *Proceedings of the National Academy of Sciences of the United States of America*
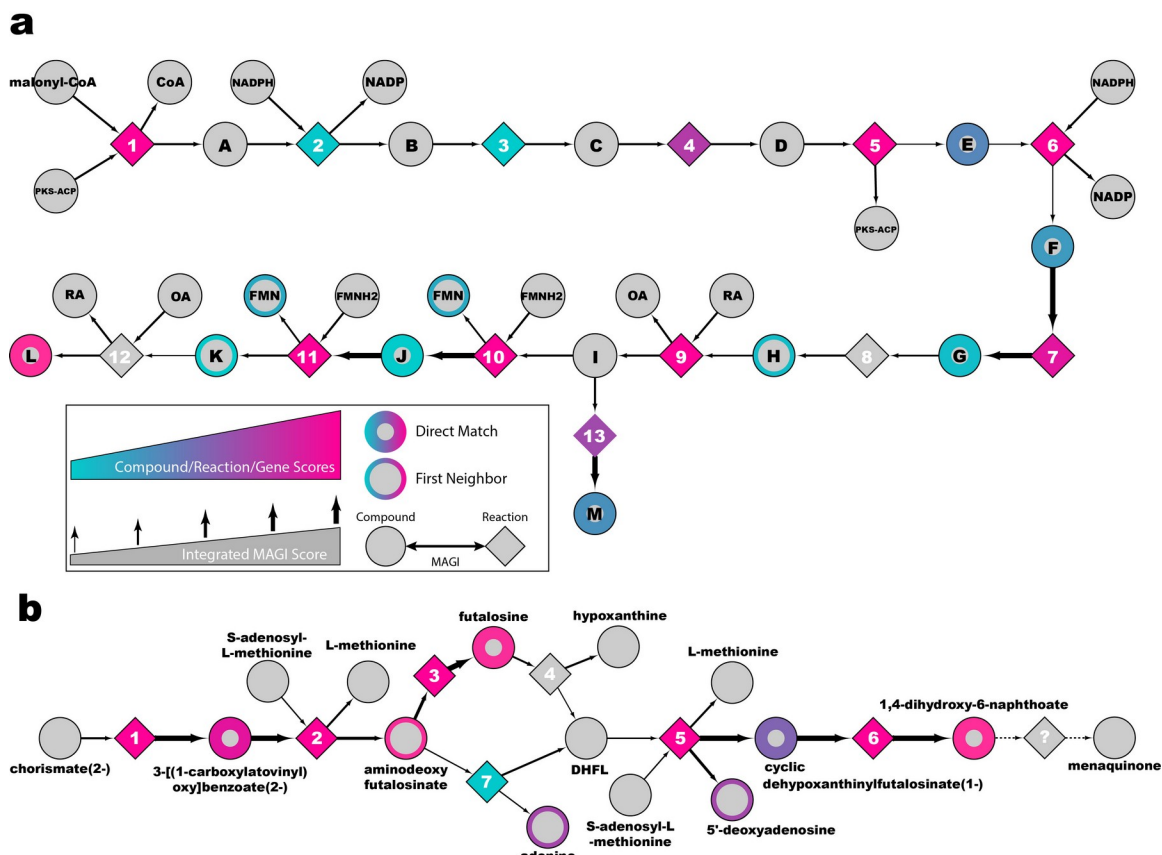794         *91*, 11542-11546.
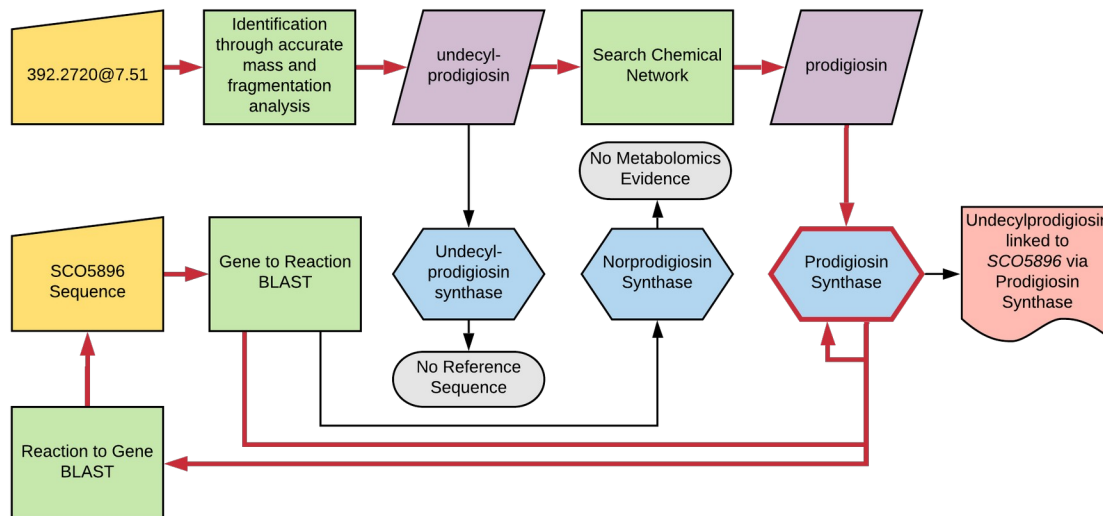795
796

797 **Figures and Tables**
798



799
800 **Figure 1. MAGI workflow for consensus scoring**. Mass spectrometry
801 features are connected to metabolites via methods such as accurate mass
802 searching or fragmentation pattern matching. These metabolites are
803 expanded to include similar metabolites by using the Chemical Network.
804 These metabolites are then connected to reactions, which are reciprocally
805 linked to input gene sequences via homology  (Reciprocal BLAST box). The
806 metabolite, reaction, and homology scores generated throughout the MAGI
807 process are integrated to form MAGI scores (Scoring box). For details on MAGI
808 scores, see **Methods**.

**a** Number of reactions associated with genes

**b** Genes connected to one or more reactions

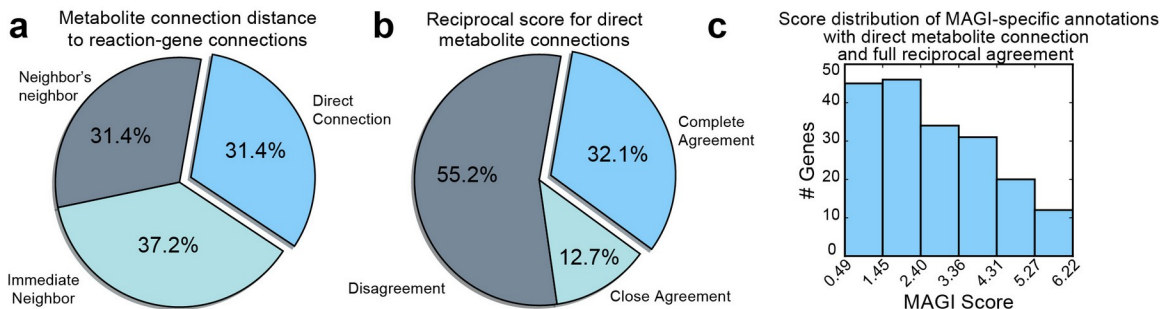**c** MAGI-unique and not unique distributions

809
810 **Figure 2. MAGI associates more genes with reactions that can be**
811 **ranked in *S. coelicolor*.** a) Number of reactions associated with each gene
812 by MAGI, KEGG, and BioCyc. b) Venn diagram showing the genes connected
813 to one or more reactions by MAGI, KEGG, and/or BioCyc. c) Distributions of
814 the associations between a gene and a reaction for genes that have
815 annotations in MetaCyc or Kegg (orange), or are unique to MAGI (blue),
816 highlighting that there are several high-scoring MAGI associations for genes
817 with no annotation.
818

819
820 **Figure 3. Pathway views of MAGI results.** Metabolite, homology, and
821 integrative MAGI scores throughout the (a) actinorhodin and (b) menaquinone
822 biosynthesis pathways guides MAGI interpretations by visualizing results in a
823 broader context. Circular nodes represent metabolites, diamond nodes
824 represent reactions, and edges represent MAGI consensus scores. Border
825 color of circular nodes corresponds to the MIDAS metabolite score, and
826 border width corresponds to the chemical network level searched in MAGI. Fill
827 color of diamond nodes correspond to the homology score. The line width of
828 the edges corresponds to the MAGI score. Abbreviations and legends for
829 metabolites and reactions are in supplementary table 8. The final step(s) in
830 the menaquinone biosynthesis are currently not known and are represented
831 by dashed edges and a "?" as the reaction.
832

833
834 **Figure 4. Flowchart illustrating the key components of the MAGI**
835 **algorithm and process for associating undecylprodigiosin with**
836 **SCO5896.** In the upper half of the flowchart, the mass spectrometry feature
837 with *m/z* 392.2720 at retention time 7.51 minutes was potentially identified
838 to be undecylprodigiosin, which is in the undecylprodigiosin synthase
839 reaction. This reaction has no reference sequence, so could not be directly
840 connected to any *S. coelicolor* genes. Undecylprodigiosin was queried for
841 similar metabolites in the chemical network, finding prodigiosin, which is in
842 the prodigiosin synthase reaction. This reaction does have a reference
843 sequence, which was used in a homology search against the *S. coelicolor*
844 genome (Reaction to Gene BLAST), finding *SCO5896* as the top hit. In the
845 lower half of the flowchart, the *SCO5896* gene sequence was queried against
846 the entire MAGI reaction reference sequence database in a homology search
847 (Gene to Reaction BLAST), finding the prodigiosin synthase and
848 norprodigiosin synthase reactions. Norprodigiosin synthase did not have any
849 metabolomics evidence, The metabolite-to-reaction and gene-to-reaction
850 results were connected via the shared prodigiosin synthase reaction,
851 effectively linking the feature 392.2720 to undecylprodigiosin and to
852 *SCO5896*.
853

**Figure 5. Prioritization of MAGI gene function suggestions.** a) Of the 1,883 MAGI-specific gene-metabolite linkages (Figure 2C), 591 genes were associated with a reaction that was directly connected to an observed metabolite (*i.e.* the chemical similarity network was not used to link a metabolite to the reaction) (light blue). b) Of those, 190 genes had reciprocal agreement in bidirectional BLAST searches (light blue). c) Histogram of the top MAGI scores of the 190 genes from panel (b). Through this process an actionable number of high-priority and high-strength novel gene function hypotheses to test in future studies can be identified.

865 **Table 1. Comparison between MAGI, KEGG, and BioCyC annotations**
866 **for *S. coelicolor* genes discussed in this study.**

| Gene | MAGI annotation (reaction) | MAGI score | Observed Metabolite Evidence | KEGG annotation (name) | KEGG Reaction Agreement with MAGI | BioCyC annotation (name) | BioCyc Reaction Agreement with MAGI |
|---|---|---|---|---|---|---|---|
| SCO4326 | RXN-10622 | 5.68 | Dihydroxy-naphthoate | 1,4-dihydroxy-6-naphthoate synthase | Agree | ORF | None |
| SCO4327 | RHEA:25907 | 5.16 | Futalosine | None | None | ORF | None |
| SCO4494 | RXN-15264 | 5.57 | Carboxy-vinyloxy-benzoic acid | Aminodeoxy-futalosine synthase | Agree | ORF | None |
| SCO4506 | RXN-12345 | 5.57 | Carboxy-vinyloxy-benzoic acid | chorismate dehydratase | Agree | ORF | None |
| SCO4550 | RXN-10620 | 5.03 | Cyclic-DHFL | cyclic dehypoxanthinyl futalosine synthase | Agree | ORF | None |
| SCO5074 | RXN1A0-6312 | 5.37 | Bicyclic intermediate F & (S)-Hemiketal | None | None | ActVI-ORF3 | Agree |
| SCO5075 | RXN1A0-6316 | 1.22 | Dihydro-kalafungin | None | None | ActVI-ORF4 | Agree |
| SCO5080 | RXN-18115 | 4.87 | DHK-red | 3-hydroxy-9,10-secoandrosta-1,3,5(10)-triene-9,17-dione monooxygenase [EC:1.14.14.12] | Disagree: R09819 | ActVA-ORF5 | Agree |
| SCO5081 | RXN1A0-6318 | 4.63 | Dihydro-kalafungin | None | None | ActVA-ORF6 | Agree |
| SCO5091 | RXN1A0-6307 | 5.95 | Bicyclic intermediate E | None | None | ActIV | Agree |
| SCO5315 | RXN-15413 | 4.58 | WhiE_20C_substrate | None | None | Polyketide aromatase | None |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| SCO5896 | RXN-15787* | 1.32 | Undecyl-prodigiosin | pyruvate, water dikinase | Disagree: R00199 | RedH | Agree* |
| SCO6300 | RXN0-5226 | 3.22 | Anhydro-NAM | beta-N-acetyl-hexosaminidase | Disagree: R00022, R05963, R07809, R07810, R10831 | hydrolase | None |
| SCO7595 | RHEA:24952 | 5.23 | Anhydro-NAM | anhydro-N-acetylmuramic acid kinase | None | ORF | None |

867   * Due to chemical network search, this reaction was listed as the prodigiosin
868   synthase reaction but the metabolite connected to it was undecylprodigiosin,
869   requiring manual interpretation to determine the actual reaction connected
870   to the gene was undecylprodigiosin synthase.