

1 **Causal genetic variation underlying metabolome differences**

2 Devjane Swain-Lenz^{*,†}, Igor Nikolisky[‡], Jiye Cheng^{*,§}, Priya Sudarsanam^{*,†},

3 Darcy Naylor[†], Max V. Staller^{*,†}, and Barak A. Cohen^{*,†,1}

4 ^{*}Center for Genome Sciences and Systems Biology, Washington University in St.

5 Louis School of Medicine, St. Louis, Missouri 63110, USA

6 [†] Department of Genetics, Washington University in St. Louis School of Medicine,

7 St. Louis, Missouri 63110, USA

8 [‡] Department of Chemistry, Washington University in St. Louis, St. Louis,

9 Missouri, 63130, USA

10 [§] Center for Gut Microbiome and Nutrition Research, Washington University

11 School of Medicine, St. Louis, MO 63110, USA

12 ¹ Corresponding author

13

14 **An ongoing challenge in biology is to predict the phenotypes of individuals**

15 **from their genotypes. Genetic variants that cause disease often change an**

16 **individual's total metabolite profile, or metabolome. In light of our**

17 **extensive knowledge of metabolic pathways, genetic variants that alter the**

18 **metabolome may help predict novel phenotypes. To link genetic variants to**

19 **changes in the metabolome, we studied natural variation in the yeast**

20 ***Saccharomyces cerevisiae*. We used an untargeted mass spectrometry**

21 **method to identify dozens of metabolite Quantitative Trait Loci (mQTL),**

22 **genomic regions containing genetic variation that control differences in**

23 **metabolite levels between individuals. We mapped differences in urea cycle**

24 **metabolites to genetic variation in specific genes known to regulate amino**
25 **acid biosynthesis. Our functional assays reveal that genetic variation in**
26 **two genes, *AUA1* and *ARG81*, cause the differences in the abundance of**
27 **several urea cycle metabolites. Based on knowledge of the urea cycle we**
28 **predicted, and then validated a new phenotype, sensitivity to a particular**
29 **class of amino acid isomers. Our results are a proof of concept that**
30 **untargeted mass spectrometry can reveal links between natural genetic**
31 **variants and metabolome diversity. The interpretability of our results**
32 **demonstrates the promise of using genetic variants underlying natural**
33 **differences in the metabolome to predict novel phenotypes from genotype.**

34 **INTRODUCTION**

35 A fundamental goal in biology is to understand the properties of genetic
36 variants that underlie phenotypic differences between individuals. Because
37 causal genetic variants often change an individual's metabolome (Gauguier
38 2016; Suhrer and Geiger 2012), metabolomics, the systematic study of
39 metabolites, offers an avenue to identify genetic variants that contribute to
40 phenotypic differences through their effects on metabolism. Understanding how
41 genetic variation in specific genes affects metabolic phenotypes is an important
42 step towards the goal of predicting phenotype from genetic variation. For
43 example, therapeutic outcomes are better when stroke patients receive a dose of
44 warfarin that depends on their genotypes at two metabolic genes rather than a
45 fixed dose (Consortium, T. I. W. P. 2009; Pirmohamed *et al.* 2013). To further
46 explore causal genetic variation in the metabolome, we combined improvements

47 in untargeted mass spectrometry with a genetically tractable yeast system to
48 uncover gene variants that underlie metabolome differences. Then using
49 knowledge of known metabolic pathways, we predicted novel drug sensitivity
50 phenotypes from genotype.

51 A large number of metabolite levels can be measured simultaneously
52 either through targeted methods, in which the identities of metabolites are known,
53 or untargeted methods, in which the identities of metabolites are unknown.
54 Targeted methods are typically more quantitative, while untargeted methods can
55 be used to screen a broader range of metabolic phenotypes. Previous
56 metabolomics studies in plants, humans, and lab strains of yeast have either
57 identified genetic variants that effect metabolite levels using targeted methods
58 (Breunig *et al.* 2014; Chen *et al.* 2014; Dong *et al.* 2015), or have identified new
59 metabolic phenotypes between individuals with known causal genetic variation
60 using untargeted methods (Broyart *et al.* 2009; Hu *et al.* 2014, Keurentjes *et al.*
61 2006). However, few investigators have attempted to map genetic variation using
62 untargeted mass spectrometry, which can lead to discovery of both unknown
63 metabolic and genetic variation (Lewis *et al.* 2014). For instance untargeted
64 methods led to the detection of variation in the chloroquine resistant gene in
65 *Plasmodium* that confers different levels of hemoglobin-derived proteins (Lewis
66 *et al.* 2014). Armed with an extensive knowledge of metabolic pathways, we can
67 further such studies to interpret genotypes to predict novel phenotypes.

68 Wild strains of the yeast *Saccharomyces cerevisiae* have proven to be
69 useful models of natural genetic variation. Strains of *S. cerevisiae* are

70 approximately 1% divergent at the nucleotide level, and are phenotypically
71 diverse. For instance, in low nitrogen conditions domesticated wine strains
72 preferentially ferment glucose while natural oak isolates respire and then
73 sporulate (Schacherer *et al.* 2009; Fay and Benavides 2005; Liti *et al.* 2009).
74 Previous studies identified genetic variants that contribute both to the expression
75 and sporulation differences between these natural isolates (Gerke *et al.* 2009),
76 but the genetic variation that causes metabolome differences between wild
77 strains is unknown.

78 We studied natural variation in the metabolome of *S. cerevisiae*. We used
79 an untargeted mass spectrometry method (Fuhrer *et al.* 2011) to identify dozens
80 of metabolite Quantitative Trait Loci (mQTL), genomic regions containing
81 variation that control differences in levels of unknown metabolites between
82 individuals. We mapped variation in urea cycle metabolites to genetic variation in
83 specific genes known to regulate amino acid biosynthesis (Sophianopoloulou and
84 Diallinas 2005; Dubois and Messenguy 1985). Our functional assays reveal that
85 genetic variation in two genes, *AUA1* and *ARG81*, underlie the differences
86 between individuals' abundance of several urea cycle metabolites. Drawing from
87 knowledge of the urea cycle, we predicted and validated a novel phenotypic
88 difference between strains. The interpretability of our results demonstrates the
89 promise of mapping causal genetic variants underlying complex metabolic
90 phenotypes and further using these variants to predict an individual's phenotype.

91 **Material and Methods**

92 ***Strains, growth conditions and metabolite extractions***

93 We used 147 genotyped segregants derived from a previously described
94 oak (YPS606) and wine strain (UCD2120) hybrid (Gerke *et al.* 2006). We
95 engineered reciprocal hemizygotes by transforming strains with kanMX4 targeted
96 to the gene of interest. For preparing extracts for mass spectrometry we grew
97 strains overnight in synthetic dextrose media (0.145% yeast nitrogen base minus
98 amino acids/ammonium sulfate, 0.5% ammonium sulfate, 2% dextrose) at 30°C.
99 We diluted overnight cultures into 25 ml of SD media to an OD₆₀₀ of 0.20.
100 Cultures were grown in flasks at 30°C and 300 rpm until mid-log phase. We
101 harvested cells by vacuum filter, and extracted hydrophilic metabolites from 0.2
102 um filters using 40:40:20 (v/v/v) methanol/acetonitrile/water (Lu *et al.* 2010). We
103 froze and thawed extracts at -80°C and -20°C, respectively, three times. We
104 pelleted cells, and stored the supernatant at -80°C until we performed mass
105 spectrometry. We replicated growth for the parents fifteen times, segregants
106 three times, and reciprocal hemizygotes nine to fifteen times (ARG81-oak::KAN,
107 9; ARG81-wine::Kan, 12; AUA1-oak::Kan, 15; AUA1-wine::Kan, 15). We
108 randomized samples using a partial block design and extracted biological
109 replicates at different times. To negatively control for non-biological ions specific
110 to the extraction process, we extracted ions from seven samples containing
111 neither media nor cell culture. To test Gap1p activity, parent strains were grown
112 in SD media overnight at 30°C, washed with water, serially diluted and grown at
113 30°C for three days on SD media plus either 1% L-Proline or 0.1% ammonium
114 sulfate, and with or without 0.16% D-Histidine (Sophianopoloulou and Diallinas
115 2005; Regenbergs and Hansen 2000).

116 ***Flow-injection OrbiTrap MS and Data Processing***

117 We directly injected metabolite extracts into a LTQ-Orbitrap Discovery
118 Mass Spectrometer (Thermo Fisher Corporation) without a chromatography
119 phase. The mobile phase for negative mode was isopropanol/water (60:40, v/v)
120 buffered with 5 mM ammonium carbonate at pH 9, and the flow rate was 150
121 $\mu\text{L}/\text{min}$ (fuhrer2011high). We used the R package MALDIquant (Gibb and
122 Strimmer 2012) to process profile mode data. We used a square root
123 transformation on each sample's spectra. For each sample, we removed the
124 baseline and used total ion current to normalize the intensity. We used a signal to
125 noise ratio of 5 and a half-window size of 3 to detect peaks in each sample. To
126 compare peaks across samples, we aligned peaks using a warping function
127 determined by MALDIquant. There were 478 detectable peaks. To eliminate non-
128 biological ions, we filtered out ions that were at least half as abundant in the
129 negative controls as the mean of the segregant samples. Additionally we ensured
130 that our technical replication was within the standard coefficient of variance of
131 10% by creating a standard of the wine and oak parents mixed at equal amounts.
132 We ran the same standard sample at least once for every fifty samples we ran on
133 the mass spectrometer. We ran all of our samples over the course of four days.
134 To confirm metabolite identity, we compared the candidate peaks of our standard
135 to the profile of known metabolites using high performance liquid
136 chromatography coupled with mass spectrometry.

137 ***Statistical and QTL analyses***

138 To identify metabolites that are significantly different in abundance
139 between the two parents, we used a mixed linear model (Bates and Maechler
140 2010) to describe metabolite abundance with batch as a random effect and
141 genotype as a fixed effect (abundance \sim genotype + batch). To determine the
142 significance of the genotypic effect, we compared our full model to a null model
143 (abundance \sim batch) using a two-way ANOVA with a Benjamini-Hochberg
144 adjustment (FDR = 0.1) (Benjamini and Hochberg 1985). We calculated
145 transgression and epistasis for all ions as previously described in Brem and
146 Kruglyak (Gerke *et al.* 2006; Brem and Kruglyak 2005). We chose a conservative
147 cutoff of three standard deviations for transgression to ensure we were not
148 overestimating transgressive effects. We calculated epistasis T as a modified t-
149 test: $T = \Delta / \sigma$, where Δ is the difference between means and σ is the variance. We
150 calculated broad-sense heritability as $H^2 = 1 - \sigma_e / \sigma_o$, where σ_e is the expected
151 variance from the parents and σ_o is the observed variance in the segregants.

152 We used the R package *qtl* (Broman *et al.* 2003) to map quantitative trait
153 loci to the abundance of metabolites. We permuted the data 1000 times to create
154 a null distribution, and used an automated Haley-Knott regression to identify
155 mQTLs with a 5% significance threshold. For ions with significant QTLs, we
156 again permuted the data 1000 times to create a null distribution, and used
157 composite interval mapping to identify weaker QTL peaks ($P < 0.05$), and used
158 linear models to explain the variance in ion abundance due to candidate QTLs. If
159 an mQTL mapped to multiple metabolites, we took the overlap of the mQTLs
160 ranges for each metabolite, and mapped the overlapping mQTL regions to the S.

161 *cerevisiae* reference genome to identify candidate genes (Engel *et al.* 2014). We
162 used linear models to calculate the variance due to specific mQTL.

163 We performed PCA on the segregants and parents using ornithine,
164 glutamine, glutamate, citrulline and arginine as variables using **princomp()**
165 function in R (R Core Team 2014). We used the eigenvectors to calculate broad-
166 sense heritability (Gerke *et al.* 2006). We mapped QTL to the eigenvectors using
167 composite interval mapping as described above. For negative controls we
168 performed the same QTL analyses from all 99 metabolites and the twenty
169 metabolites with individual mQTLs. Additionally we performed the analyses on
170 randomly selected five metabolites from the twenty metabolites with mQTLs, and
171 performed this analysis ten times. We used MANOVA in R to analyze the
172 difference in the urea cycle in reciprocal hemizygotes, and one-way ANOVA in R
173 to analyze the differences in metabolite abundance between reciprocal
174 hemizygotes (R Core Team 2014).

175 ***Data Availability***

176 Strains are available upon request. The raw mass spectrometry data obtained in
177 this study will be accessible at the NIH Common Fund's Data Repository and
178 Coordinating Center (supported by NIH grant, U01-DK097430) website, the
179 Metabolomics Workbench, <http://www.metabolomicsworkbench.org>.

180 Supplemental Table 1 contains a metabolite reporting checklist. Supplemental

181 Table 2 contains processed averages for segregant metabolite abundance.

182 Supplemental Table 3 contains data for transgression, epistasis and heritability.

183 Supplemental Table 4 contains data for linear models of mQTL. Supplemental

184 Table 5 contains processed averages for parents metabolite abundance.
185 Supplemental Table 6 contains processed data for mixed linear models of
186 parents metabolite abundance. Supplemental Table 7 contains processed
187 averages of reciprocal hemizygotes metabolite abundance.

188 **Results and Discussion**

189 ***High-throughput measurement of untargeted metabolites***

190 We employed untargeted mass spectrometry to rapidly and systematically
191 quantify abundances of unknown metabolites in natural isolates of the yeast
192 *Saccharomyces cerevisiae*. Previous studies successfully identified causal
193 genetic variation by targeting specific metabolites (Dubois and Messenguy 1985;
194 Chen *et al.* 2014; Dong *et al.* 2015) or by untargeted metabolic analyses in
195 individuals with known genetic variants (Broyart *et al.* 2009; Hu *et al.* 2014). As a
196 complement to these approaches, we instead quantified unknown metabolites in
197 minimally processed extracts by direct injection into a mass spectrometer (Fuhrer
198 *et al.* 2011; Lu *et al.* 2010) (Supplemental Table 1). We chose to sacrifice the
199 resolution of liquid chromatography coupled MS for the speed of the direct
200 injection method, which allowed us to avoid the analytical challenges of mass
201 spectrometer measurement drift over time, and more accurately measure
202 metabolite abundances.

203 Using a stringent filter for reproducibility, we reliably measured the relative
204 abundance of ninety-nine distinct ions (Supplemental Table 2). To control for
205 non-biological ions, we ensured the ninety-nine ions were more than twice as
206 abundant in the biological samples than the negative controls (Supplemental

207 Figure 1). To determine the reproducibility of the direct injection approach, we
208 created a reference standard by extracting and pooling metabolites from two
209 independently grown strains. We ran this standard eleven times over the course
210 of the four-day run and determined the median coefficient of variance across
211 biological metabolites was 10%, well within the range of acceptable experimental
212 variation (Supplemental Figure 2) (Lu *et al.* 2010). This conservative analysis
213 revealed that we can use untargeted methods to consistently measure the
214 relative abundance of unknown, biological metabolites (Supplemental Table 2).

215 ***Complex genetic architecture underlying natural variation in metabolite***
216 ***differences.***

217 Metabolite abundances are genetically complex traits with alleles that
218 have both small additive and non-additive effects. To define the genetic
219 architecture of metabolite levels we quantified the abundances of metabolites in
220 147 diploid segregants derived from a cross between a yeast strain isolated from
221 the bark of an oak tree and a yeast strain isolated from a commercial wine barrel
222 (Gerke *et al.* 2006). The continuous distribution of metabolite abundances in the
223 segregants indicates that metabolite levels are controlled by many alleles of
224 small effect (Figure 1A and 1B). We also found statistical evidence from the
225 shape of the phenotype distributions for genetic interactions among alleles that
226 influence metabolite levels, especially for metabolites that displayed
227 transgressive segregation patterns (Figure 1C, Supplemental Table 3) (Brem and
228 Kruglyak 2005). Thus alleles with small additive effects and alleles that display
229 epistatic interactions contribute to natural variation in metabolite levels.

230 For more than half of all metabolites, abundance in some segregants was
231 more than three standard deviations away from both parents' abundances, which
232 is evidence for pervasive transgression (Brem and Kruglyak 2005) (Figure 1B,
233 Supplemental Table 3). In the most striking examples sixteen metabolites had
234 very low or undetectable levels in both parents, while 75% or more of segregants
235 had high levels of the metabolite (Figure 1B and 1C, e.g. 159.08 m/z). This
236 transgressive segregation pattern is consistent with the hypothesis that the wild
237 parental strains contain compensatory alleles with both positive and negative
238 effects on metabolite levels, which together maintain low levels of certain
239 intermediate metabolites. Recombination of compensatory alleles during meiosis
240 leads to the accumulation of high levels of metabolites in the segregants.

241 ***mQTL influence metabolites in the urea cycle***

242 We next identified mQTL by correlating segregating polymorphisms with
243 metabolite levels of unknown metabolites in the panel of segregants (Figure 2).
244 We previously genotyped 225 markers in our 147 segregants (Gerke *et al.* 2009).
245 We detected a genetically complex network of mQTLs with several mQTLs
246 influencing the same metabolite, and several metabolites with multiple mQTL. In
247 total we detected sixteen significant mQTL (Figure 2B) that contribute to the
248 variation of twenty metabolites ($P < 0.05$). Seven mQTL are shared among two or
249 more metabolites. Most metabolites have either one or two detectable mQTL,
250 and four metabolites have either four or five detectable mQTL. To determine the
251 fraction of the variance in metabolite abundance explained by mQTL, we used
252 linear models (Supplemental Table 4). On average individual mQTL explain

253 11.0% of the variance in metabolite levels, with a range of 6.0% - 22.6%. As
254 expected from our general transgression analysis, we found that of the seven
255 metabolites with multiple mQTLs, four metabolites had mQTLs with effects in
256 opposite directions. This finding further supports the hypothesis that parental
257 strains contain compensatory alleles that maintain optimal metabolite
258 abundances that are similar to each other. To determine whether contributions to
259 mQTL are additive or non-additive, we analyzed the interactions of alleles. Our
260 analysis of epistasis from the phenotype distribution of segregants suggested
261 that interactions between QTL contribute to variation in metabolite levels (Figure
262 1C). Typically the additive contribution to QTL are larger than non-additive
263 interaction. As the additive contribution of mQTL are small, we expected the
264 effects of interactions to be even smaller. Consistent with this idea, linear models
265 revealed one small but significant interaction term (Supplemental Table 4). Thus,
266 metabolite abundances are largely shaped by many loci with small additive
267 effects, and while interactions do play some role in shaping the distributions of
268 ion abundances, most interaction effects are likely quite small.

269 We identified several segregating loci that impact urea cycle metabolism.
270 We organized metabolites into pathways by determining the identities of
271 metabolites with the most mQTL. After searching yeast mass spectrometry
272 databases (Jewison *et al.* 2012) for candidate metabolites, we used traditional
273 liquid chromatography coupled with targeted mass spectrometry to compare our
274 samples to standards of these candidate metabolites. In this way we identified
275 glutamine and citrulline as mQTL targets. As citrulline is produced during the

276 urea cycle, and glutamine biosynthesis is closely connected to the urea cycle, we
277 searched for other possible candidates in the urea cycle (Jewison *et al.* 2012).
278 We confirmed the identity of five metabolites in or adjacent to the urea cycle:
279 citrulline, ornithine, arginine, glutamine and glutamate (Figure 2A). Our results
280 demonstrate that several segregating genetic variants impact urea cycle
281 metabolism and that our rapid untargeted method identified mQTLs that affect an
282 important biochemical pathway.

283 Because segregating variation in the recombinant progeny influenced
284 metabolites in the urea cycle, we predicted that the parental strains would harbor
285 differences in urea cycle metabolism. Consistent with this prediction we found
286 significantly different levels of citrulline, ornithine, glutamine and glutamate
287 between the parents (Figure 2A, Figure 3A, Supplemental Tables 5 and 6).
288 Notably, humans domesticated wine strains in low nitrogen environments, which
289 may have provided selective pressure on the urea cycle, a nitrogen reclamation
290 pathway (Marsit and Dequin 2015). Our genetic data reveal significant natural
291 variation in urea cycle metabolism between strains from different ecological
292 niches.

293 Our initial mQTL analysis assumed that each metabolite was independent,
294 but metabolite abundances in the urea cycle are intrinsically linked to one
295 another. Given that multiple metabolites in the urea cycle map to overlapping
296 mQTL, we reasoned that combining metabolite measurements from the urea
297 cycle would improve our power and allow us to narrow the linkage region. In
298 other words, an mQTL could have effects spread across correlated metabolites,

299 and may have stronger effects on pooled measurements from correlated
300 metabolites. We performed a principal component analysis on the segregants
301 using the five amino acids in the urea cycle, and then remapped mQTL to these
302 principal components (PCs). Using PCs as phenotypes increases the statistical
303 power to detect QTL for correlated and variable data (Chase *et al.* 2002; Mangin
304 *et al.* 1998). Four PCs explain 99.0% of the variance (Table 1). We calculated the
305 broad sense heritabilities (H^2) of each PC, which measures the proportion of
306 phenotypic variability due to genetic variation (Gerke *et al.* 2006; Brem and
307 Kruglyak 2005). The first two PCs have low H^2 , which indicates that the majority
308 of phenotypic variability of intracellular metabolites is due to environmental
309 effects. In contrast, PC3 and PC4 have higher H^2 , supporting a genetic
310 component to phenotypic variability in the urea cycle. As a negative control, we
311 attempted to map mQTL to PCs derived from all ninety-nine metabolites, the
312 twenty metabolites with mQTL, as well as five random metabolites with mQTL,
313 but found no significant peaks. This suggests that the mQTL with the strongest
314 genetic signal are specific to the urea cycle.

315 ***Causal variation in two genes underlies natural variation in urea cycle***
316 ***metabolites***

317 Our pathway level analysis of metabolite abundances narrowed mQTL
318 and revealed promising candidate genes. When we mapped mQTL to PC3 and
319 PC4, we detected multiple QTL peaks (Table 2), two of which overlap with peaks
320 mapped to individual metabolites and contain excellent candidate genes (Engel
321 *et al.* 2014). One peak covers the gene *AUA1*. *AUA1* regulates amino acid

322 transport in the presence of ammonia, which is removed from the cell via the
323 urea cycle. The wine variant of *AUA1* contains a premature stop codon, which
324 truncates the eighty-four amino acid peptide, to just thirteen amino acids. The
325 mutation rate (dN/dS) between strains is not higher than expected, which
326 suggests the wine strain mutation is relatively new. Another QTL contains
327 *ARG81*, a zinc-finger transcription factor that represses arginine biosynthesis.
328 The number of nonsynonymous mutations between strains is not higher than
329 expected, but when we analyzed our previously published expression data
330 (Gerke *et al.* 2006), we indeed see differential expression of nine out of twenty-
331 six *ARG81* targets, all of which show reduced expression in the wine strain.

332 We found that *ARG81* and *AUA1* contain causal variants for differences in
333 the urea cycle. We tested our hypothesis that *ARG81* and *AUA1* contain causal
334 genetic variation modulating urea cycle activity using reciprocal hemizyosity
335 assays (Steinmetz *et al.* 2002). We used a multivariate ANOVA (MANOVA) to
336 test whether the genotype of *ARG81* or *AUA1* has an effect across the whole
337 urea cycle. We find that the genotype of *ARG81* has a significant effect on the
338 abundance across all urea cycle metabolites ($P=0.03$), while the genotype of
339 *AUA1* does not ($P=0.23$). When we split the MANOVA into separate components,
340 we find the genotypes of *ARG81* and *AUA1* have significant effects on different
341 individual metabolites. We found that the wine *ARG81* allele produces a higher
342 abundance of ornithine than the oak allele, which matches the direction of the
343 effect between the parents but not the QTL model (one-way ANOVA, $P=0.03$,
344 Figure 4A, Supplemental Table 7). Although our QTL mapping did not detect an

345 effect of the *ARG81* peak on citrulline, the wine *ARG81* allele also produces a
346 higher abundance of citrulline than the oak allele (one-way ANOVA, $P=0.005$).
347 Arginine can passively turn into citrulline, which can explain the discrepancies of
348 the mQTL and metabolite data. Additionally, the wine allele of *AUA1* produces a
349 higher abundance of glutamine than the oak allele in the hybrid, which matches
350 the direction of effect between the parents and the mQTL model (one-way
351 ANOVA, $P=0.02$, Figure 4B, Supplemental Table 7). Both the difference in
352 directionality between the hybrid and parental backgrounds and original mQTL
353 mapping results suggest that there are other alleles that influence glutamine
354 abundance. We conclude that *ARG81* and *AUA1* are novel mQTGs (metabolite
355 Quantitative Trait Genes).

356 ***Predicting phenotype from genotype: A novel phenotype deduced from***
357 ***variation in the urea cycle***

358 In principle, genetic variation in metabolism can predict new phenotypes.
359 We hypothesized that both mQTGs control nitrogen metabolism by regulating the
360 gene General Amino acid Permease (*GAP1*) (Figure 5A). *AUA1* post-
361 translationally controls Gap1p by deactivating transport activity in the presence of
362 a strong nitrogen source, such as ammonia (Sophianopoloulou and Diallinas
363 2005). In poor nitrogen sources such as proline, both *ARG81* and Gap1p are
364 active (Sophianopoloulou and Diallinas 2005). We predict that the small 13 a.a.
365 truncated version of the wine *AUA1* gene is effectively a null allele, which allows
366 the wine strain to upregulate Gap1p to increase amino acid uptake. This would
367 give the wine strain a selective advantage to continue fermenting instead of

368 sporulating in low nitrogen environments, such as a wine barrel. According to this
369 model, under nitrogen poor conditions, Gap1p should be deactivated in the oak
370 parent relative to the wine strain. To test this prediction we leveraged the fact that
371 stereoisomers of L-amino acids, D-amino acids, are toxic to yeast and only enter
372 the cell through Gap1p. If Gap1p activity is higher in the wine parent than the oak
373 parent, then the wine parent will be more sensitive to the toxic amino acid D-
374 Histidine (Sophianopoloulou and Diallinas 2005; Regenberg and Hansen 2000).
375 Consistent with this prediction we found that the wine strain does not grow as
376 well as the oak strain in the presence of proline and D-Histidine, indicating higher
377 Gap1p activity in the wine parent (Figure 5B). Additionally, in the presence of a
378 strong nitrogen source in which Gap1p is not induced, both parents grow
379 similarly regardless of toxin, indicating that the phenotype is Gap1p dependent.
380 This example demonstrates how linking metabolic pathways to mQTGs can lead
381 to prediction of novel organismal phenotypes.

382 This work demonstrates the value of identifying genetic variation that
383 underlies natural differences in the metabolome. We have presented a rapid
384 approach for measuring metabolites in an untargeted fashion to systematically
385 identify causal alleles controlling variation in a core metabolic pathway. Most
386 importantly, by leveraging decades of biochemistry to interpret our results, we
387 predicted and then validated a novel cellular phenotype from measured
388 genotypes. The genes we identify as containing causal variation in the urea cycle
389 are coherent with our existing knowledge of natural selection and metabolic
390 pathways. With the current growth in metabolomics and genetics in human

391 studies (Wishart *et al.* 2013; Dharuri *et al.* 2014; Shin *et al.* 2014), similar
392 predictive methods can be used and tested in cell culture to further understand
393 how causal loci of one metabolic phenotype can affect other phenotypes, ranging
394 from metabolic biomarkers to drug sensitivity.

395 **Acknowledgements**

396 We thank Jeffrey Gordon for use of the mass spectrometer, Amy Caudy,
397 Heather Lawson and Gary Patti for advice and discussions, and members of the
398 Cohen lab for valuable feedback. This work was supported by a grant from the
399 National Institutes of Health (R01 GM092910-5).

400 **Author contributions**

401 DSL, PS, and BAC designed experiments. IN and JC developed and performed
402 mass spectrometry protocols. DSL, PS, DN and MVS ran experiments. DSL
403 analyzed data. DSL and BAC wrote the manuscript.

404 **Corresponding author**

405 The authors declare no competing financial interests. Correspondence should be
406 addressed to BAC (cohen@genetics.wustl.edu).

407 **Literature Cited**

- 408 Bates, D., and M. Maechler, 2010 Package 'lme4.' [http://lme4.r-forge.r-](http://lme4.r-forge.r-project.org/)
409 [project.org/](http://lme4.r-forge.r-project.org/)
410
411 Benjamini, Y. and Y. Hochberg, 1995 Controlling the False Discovery Rate: A
412 Practical and Powerful Approach to Multiple Testing. *Journal of the Royal*
413 *Statistical Society* **57**: 289–300.
414
415 Brem, R. B. and L. Kruglyak, 2005 The landscape of genetic complexity across
416 5,700 gene expression traits in yeast. *PNAS* **102**: 1572–1577.
417
418 Broman, K. W., H. Wu, S. Sen, and G.A. Churchill, 2003 R/qtl: QTL mapping in
419 experimental crosses. *Bioinformatics* **19**: 889–890.
420
421 Breunig, J. S., S.R. Hackett, J.D. Rabinowitz, and L. Kruglyak, 2014 Genetic
422 Basis of Metabolome Variation in Yeast. *PLoS Genetics* **10**: 1–15.
423
424 Broyart, C., J. Fontaine, R. Moliné, D. Callieu, T. Tercé-Laforgue, *et al.*, 2009
425 Metabolic profiling of maize mutants deficient for two glutamine synthetase
426 isoenzymes using ¹H-NMR-based metabolomics. *Phytochem. Anal.* **21**:102–109.
427
428 Chase, K., D.R. Carrier, F.R. Adler, T. Jarvik, E.A. Ostrander, T.D Lorentzen, and
429 K.G. Lark, 2002 Genetic basis for systems of skeletal quantitative traits: Principal
430 component analysis of the canid skeleton. *PNAS* **99**: 9930–9935.
431
432 Chen, W., Y. Gao, W. Xie, L. Gong, K. Lu, *et al.*, 2014 Genome-wide association
433 analyses provide genetic and biochemical insights into natural variation in rice
434 metabolism. *Nature Genetics* **46**: 714-721.
435
436 Consortium, T. I. W. P. 2009 Estimation of the Warfarin Dose with Clinical and
437 Pharmacogenetic Data. *New England Journal of Medicine* **360**: 753-764.
438
439 Dharuri H., A. Demirkan, J.B. van Klinken, D.O. Mook-Kanamori, C.M. van Dujhn,
440 *et al.*, 2014 Genetics of the human metabolome, what is next? *Biochim Biophys*
441 *Acta* **1842**:1921-1931.
442
443 Dong, X., Y. Gao, W. Chen, W. Wang, L. Gong, *et al.*, 2015 Spatiotemporal
444 Distribution of Phenolamides and the Genetics of Natural Variation of
445 Hydroxycinnamoyl Spermidine in Rice. *Molecular Plant* **8**: 111-121.
446
447 Dubois E. and F. Messenguy, 1985 Isolation and characterization of the yeast
448 ARGR11 gene involved in regulating both anabolism and catabolism of arginine.
449 *Mol. Gen. Genet.* **198**: 283-9.
450

451 Engel, S. R., F.S. Dietrich, D.G. Fisk, G. Binkley, R. Balakrishnan, *et al.*, 2014
452 The Reference Genome Sequence of *Saccharomyces cerevisiae*: Then and Now.
453 *G3* **4**: 389–398.

454
455 Fay J.C. and J.A. Benavides, 2005 Evidence for domesticated and wild
456 populations of *Saccharomyces cerevisiae*. *PLoS Gen.* **1**: 66-71. Fernie, A.R., A.
457 Aharoni, L. Willmitzer,
458
459 Fernie A.R., A. Aharoni, L. Willmitzer, M. Stitt, T. Tohge, *et al.* 2011.
460 Recommendations for Reporting Metabolite Data. *Plant Cell* **23**: 2477-2482.

461
462 Fuhrer, T., D. Heer, B. Begemann, and N. Zamboni, N., 2011 High-Throughput,
463 Accurate Mass Metabolome Profiling of Cellular Extracts by Flow Injection–Time-
464 of-Flight Mass Spectrometry. *Anal. Chem.* **83**: 7074–7080.

465
466 Gauguier, D., 2016 Application of quantitative metabolomics in systems genetics
467 in rodent models of complex phenotypes. *Archives of Biochemistry and*
468 *Biophysics* **589**: 158–167.

469
470 Gerke, J. P., C.T. L. Chen, and B.A. Cohen, 2006 Natural Isolates of
471 *Saccharomyces cerevisiae* Display Complex Genetic Variation in Sporulation
472 Efficiency. *Genetics* **174**: 985–997.

473
474 Gerke, J. P., K. Lorenz, and B.A. Cohen, 2009 Genetic interactions between
475 transcription factors cause natural variation in yeast. *Science* **323**: 498-501.

476
477 Gibb, S. and K. Strimmer, 2012 MALDIquant: a versatile R package for the
478 analysis of mass spectrometry data. *Bioinformatics* **28**: 2270–2271.

479
480 Hu, C., J. Shi, S. Quan, B. Cui, S. Kleessen, Z. Nikoloski, *et al.*, 2014 Metabolic
481 variation between japonica and indica rice cultivars as revealed by non-targeted
482 metabolomics. *Scientific Reports* **4**.

483
484 Jewison T., V. Neveu, J. Lee, C. Knox, P. Liu, *et al.*, 2012 YMDB: The Yeast
485 Metabolome Database. *Nucleic Acids Res.* **40**.

486
487 Keurentjes J.J.B., J. Fu, C.H.R. de Vos, A. Lommen A, R.D.Hall, *et al.*, 2006 The
488 genetics of plant metabolism. *Nat. Genet.* **38**: 842–849.

489
490 Lewis I.A., M. Wacker, K.L. Olszewski, S.A. Cobbold, K.S. Baska, *et al.*, 2014
491 Metabolic QTL analysis links chloroquine resistance in *Plasmodium falciparum* to
492 impaired hemoglobin catabolism. *PLoS Genet.* **10**: e1004085. doi:
493 10.1371/journal.pgen.1004085 pmid:24391526.

494
495 Liti G., D.M. Carter, A.M. Moses, J. Warringer, L. Parts, S.A. James SA, *et al.*,
496 2009 Population genomics of domestic and wild yeasts. *Nature* **458**: 337-341.

497
498 Lu, W., M.F. Clasquin, E. Melamund, D. Amador-Noguez, A.A. Caudy, and J.D.
499
500
501 Mangin, B., P. Thoquet and N. Grimsley, 1998 Pleiotropic QTL analysis.
502 *Biometrics* 54: 88–99.
503
504 Marsit, S. and S. Dequin, 2015 Diversity and adaptive evolution of
505 *Saccharomyces wine yeast*: a review. *FEMS Yeast Research* 15: fov067.
506
507 Pirmohamed, M., G. Burnside, N. Eriksson, A.L. Jorgensen, C. Hok Toh, *et al.*,
508 2013 A Randomized Trial of Genotype-Guided Dosing of Warfarin. *New England*
509 *Journal of Medicine* 369: 2294-2303.
510
511 Rabinowitz 2010 Metabolomic Analysis via Reversed-Phase Ion-Pairing Liquid
512 Chromatography Coupled to a Stand Alone Orbitrap Mass Spectrometer. *Anal.*
513 *Chem.* 82: 3212–3221.
514
515 R Core Team, 2014 R: A language and environment for statistical computing. R
516 Foundation for Statistical Computing, Vienna, Austria. URL [http://www.R-](http://www.R-project.org/)
517 [project.org/](http://www.R-project.org/).
518
519 Regenberg, B. and J. Hansen, 2000 GAP1, a novel selection and counter-
520 selection marker for multiple gene disruptions in *Saccharomyces cerevisiae*.
521 *Yeast* 16: 1111-9.
522
523 Schacherer J., J.A. Shapiro, D.M. Ruderfer, and L. Kruglyak, 2009
524 Comprehensive polymorphism survey elucidates population structure of
525 *Saccharomyces cerevisiae*. *Nature* 458: 342-5.
526
527 Shin S.Y., E.B. Fauman, A.K. Petersen, J. Krumsiek, R. Santos, *et al.*, 2014 An
528 atlas of genetic influences on human blood metabolites. *Nat Genet* 46: 543-550.
529
530 Sophianopoloulou, V. and G. Diallinas, 2005 AUA1, a gene involved in ammonia
531 regulation of amino acid transport in *Saccharomyces cerevisiae*. *Molecular*
532 *Microbiology* 8: 167–178.
533
534 Steinmetz, L. M., H. Sinha, D.R. Richards, J.I. Spiegelman, P.J. Oefner PJ, *et al.*,
535 2002 Dissecting the architecture of a quantitative trait locus in yeast. *Nature* 416:
536 326–330.
537
538 Suhre, K. and C. Geiger, 2012 Genetic variation in metabolic phenotypes: study
539 designs and applications. *Nature Reviews Genetics* 13: 759–769.
540

541 Wishart D.S., T. Jewison, A.C. Guo, M. Wilson, C. Knox *et al.*, 2013 HMDB 3.0
542 — The Human Metabolome Database in 2013. *Nucleic Acids Res.* **41(D1)**:D801-
543 7.
544
545

546 Table 1. Principal components for urea cycle

PC	Variance explained	Standard deviation	orn	cit	gln	glu	arg	H ²
1	47.5%	1.5%	-0.572	- 0.540	- 0.471	- 0.385	- 0.106	< 0
2	22.0%	1.5%	--	- 0.377	--	- 0.109	- 0.916	0.008
3	19.4%	0.99%	0.430	0.359	- 0.427	- 0.678	0.212	0.28
4	9.9%	0.71%	-0.160	--	0.769	- .0614	--	0.25

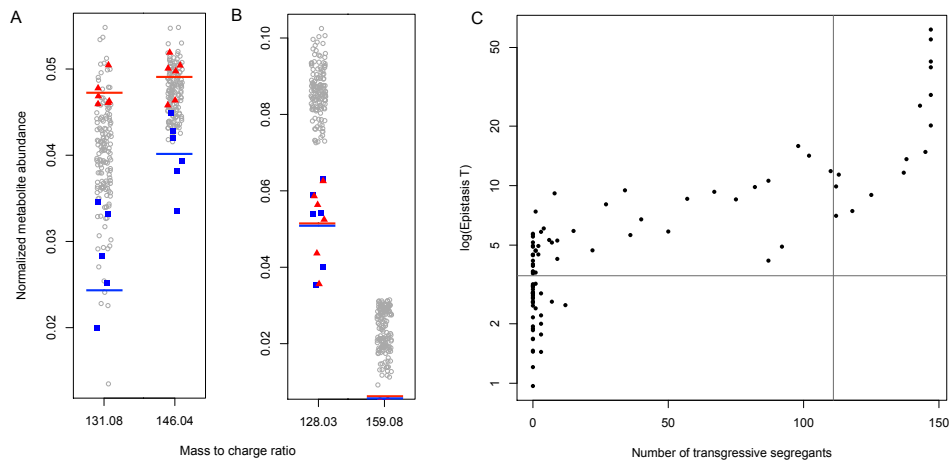
547

548 Table 2. mQTL for urea cycle principal components

PC	Chromosome	cM	Nearest marker	LOD	LOD P	Variance Explained	Additive Effect	Model P
3	6	25	L63	3.26	0.033	3.4%	19.3%	8.4x10 ⁻³
	10	24 7	L1016	3.31	0.027	4.6%	-21.9%	2.1x10 ⁻³
	13	5	L132	3.27	0.031	6.7%	26.2%	2.5x10 ⁻⁴
	16	65	L165	5.43	<0.00 1	9.4%	-31.8%	1.6x10 ⁻⁵
4	6	6	L61	6.62	<0.00 1	19.0%	-31.2%	1.7x10 ⁻⁹
	11	11	L113	5.31	<0.00 1	15.8%	-28.0%	2.9x10 ⁻⁸

549

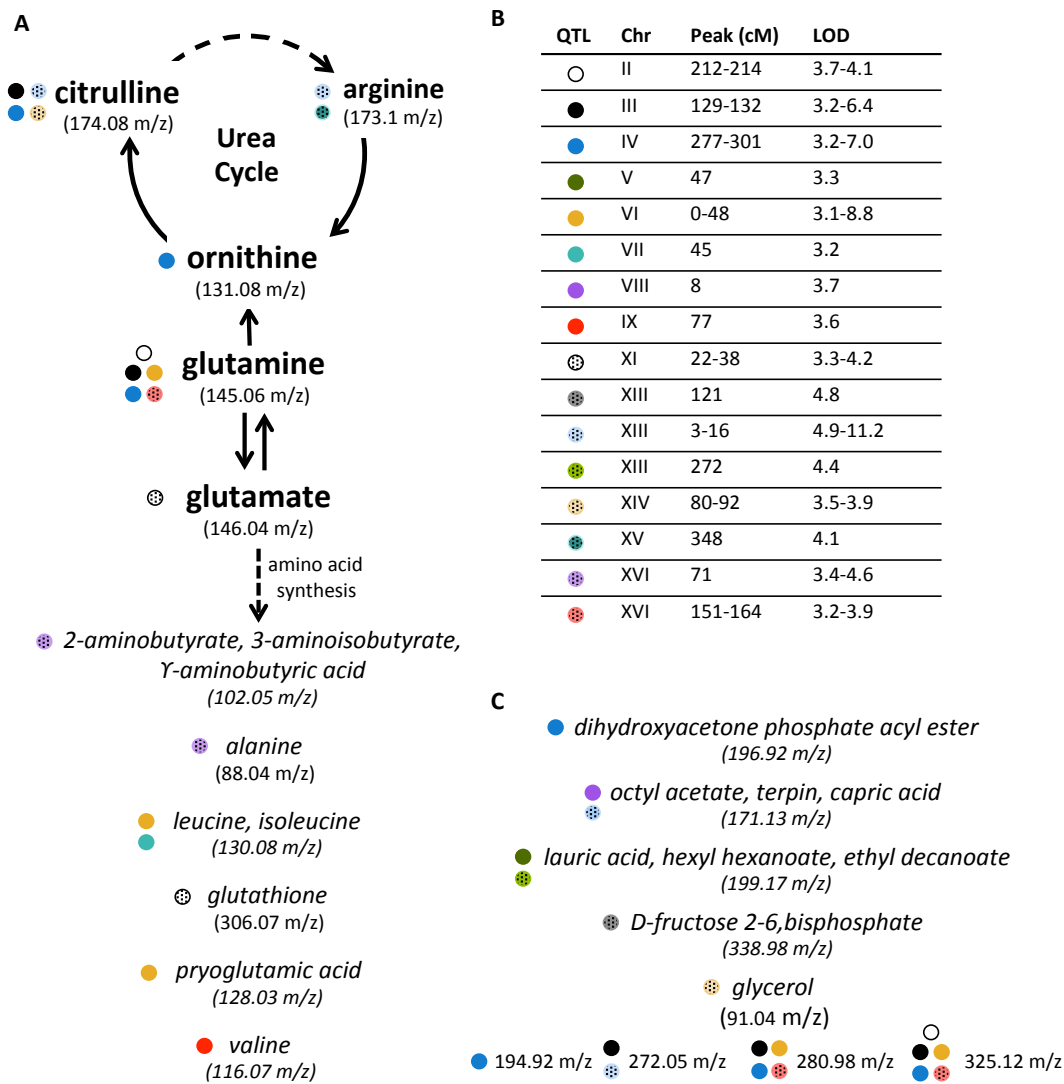
550



551

552 **Figure 1.** Metabolite abundances are complex traits. **(A)** Examples of ions for
 553 which the metabolite levels of the segregants (gray) fall between the sample
 554 means of the oak (blue) and wine (red) parents (45% of all metabolites) **(B)**
 555 Examples of ions for which the metabolites abundance is more than three
 556 standard deviations from the parental mean, indicating transgression. **(C)** For
 557 each metabolite, the number of transgressive segregants is plotted against T, a
 558 score for epistasis (Gerke *et al.* 2006; Brem and Kruglyak 2005). The horizontal
 559 grey line indicates a significant T for epistasis. For sixteen metabolites, at least
 560 75% of segregants are transgressive (vertical gray line).

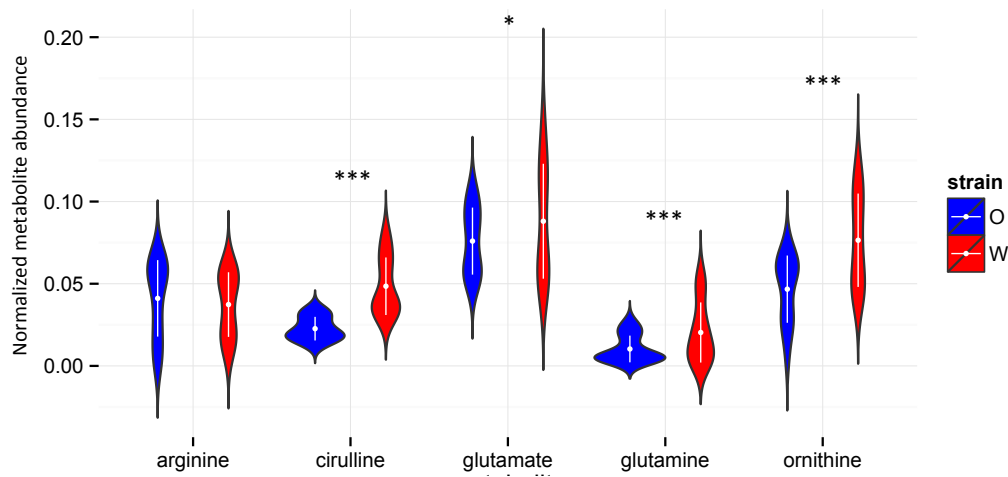
561



562

563 **Figure 2.** Mapping of untargeted metabolites reveals 20 metabolites share
 564 sixteen QTL. Asterisks represent metabolites significantly different in abundance
 565 between parents. **(A)** The confirmed metabolite identity of five amino acids (bold
 566 font), which are involved in the urea cycle and nitrogen utilization. Candidate
 567 metabolites are italicized. Circles represent individual QTL contributing to
 568 metabolite abundance as listed in **(B)**. **(C)** Additional metabolites with at least
 569 one QTL.

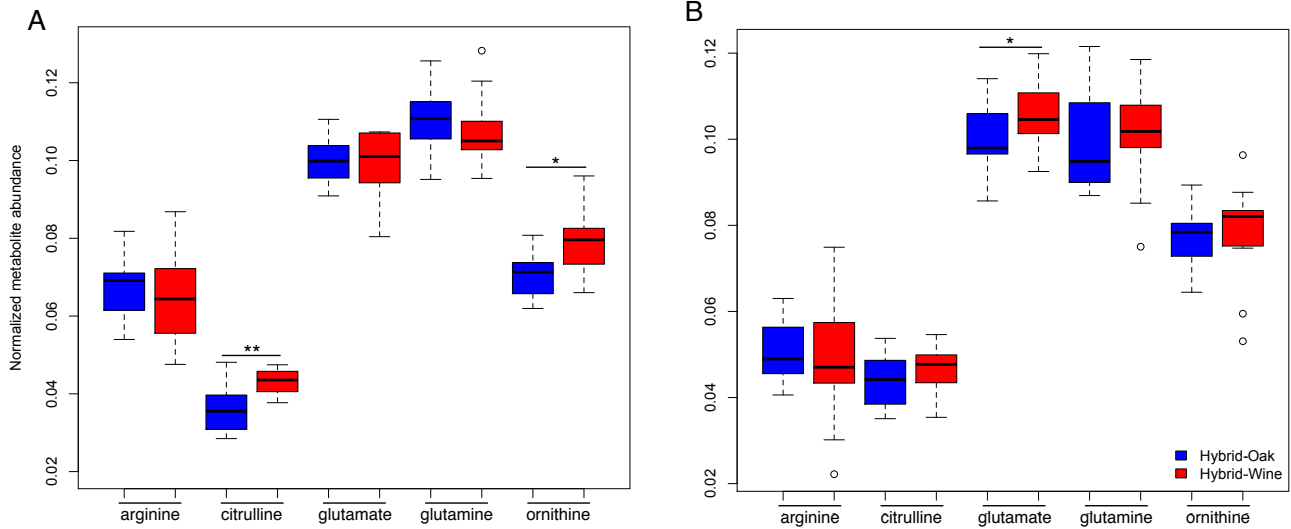
570



571

572 **Figure 3.** The abundance of urea cycle amino acids differs between parent
573 strains. We measured metabolite abundance in fifteen biological replicates of the
574 oak (blue) and wine (red) parents. White dots and bars represent the mean and
575 standard deviation respectively. We used mixed linear models to measure the
576 variance in abundance due to batch and genotype, and measured significance
577 due to genotype. (* $P < 0.05$, *** $P < 0.005$, ANOVA, Benjamini Hochberg
578 correction).

579



580

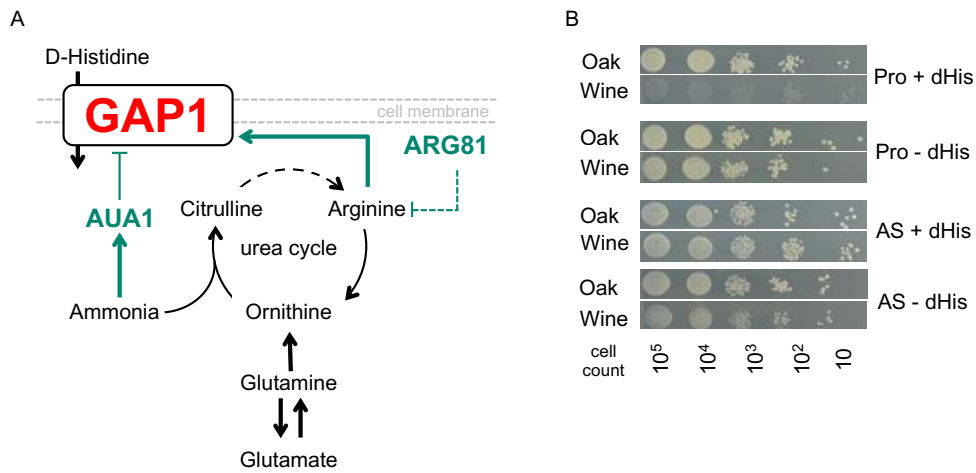
581 **Figure 4.** Reciprocal hemizygosity assays reveal causal variation in *ARG81* and

582 *AUA1*. Hybrid strains that contain only the wine (red) and oak (blue) allele of **(A)**

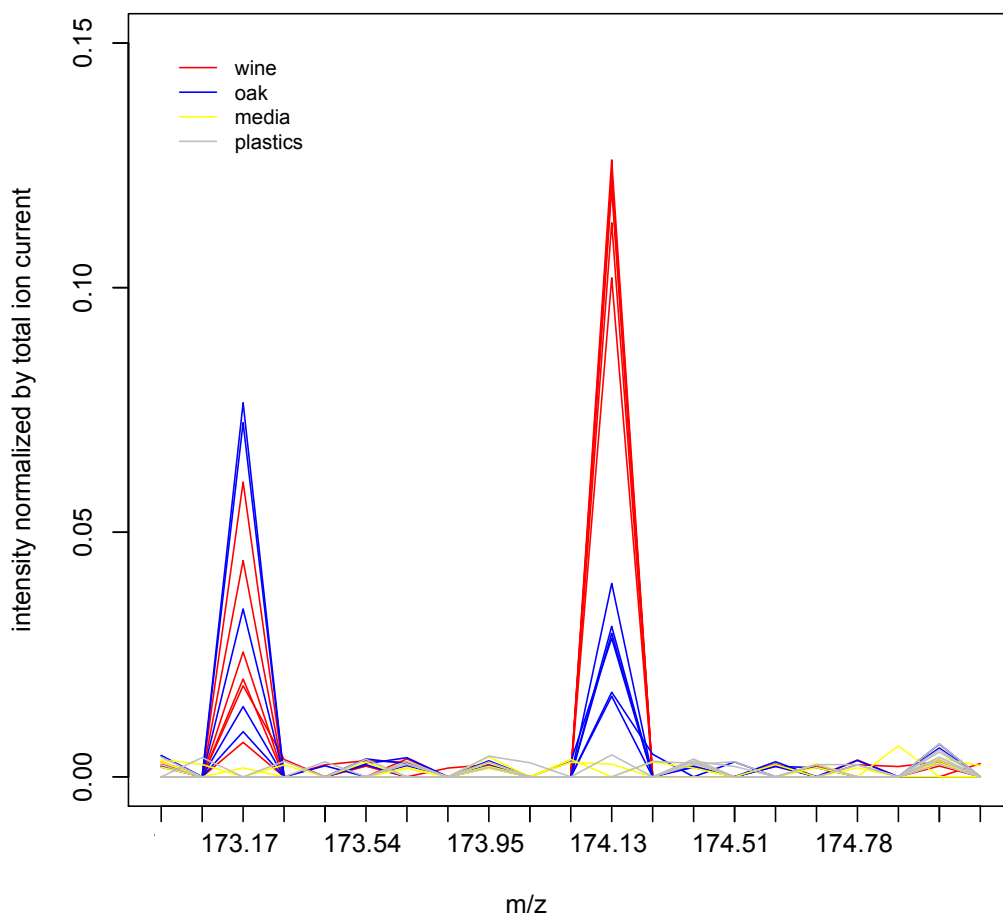
583 *ARG81* and **(B)** *AUA1*. The amino acids in the urea cycle are depicted. * $P < 0.05$,

584 ** $P < 0.005$ (one-way ANOVA).

585

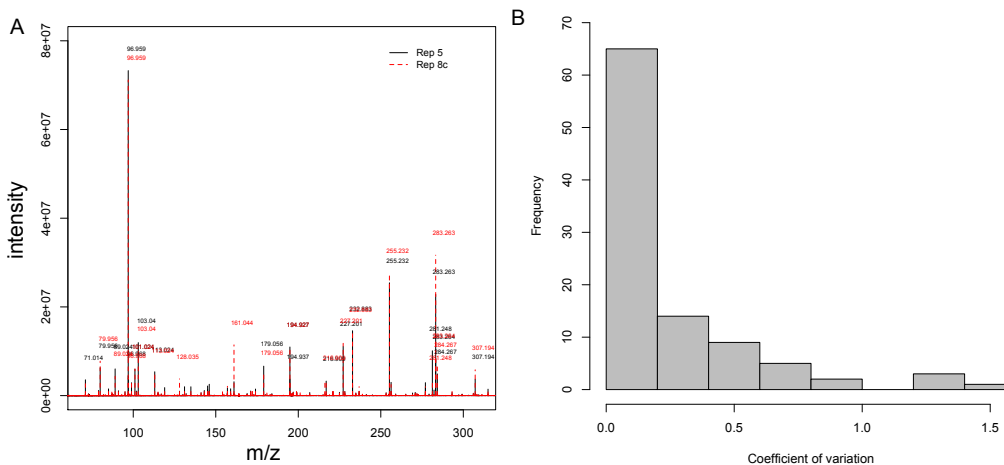


586
 587 **Figure 5.** Genetic variation predicts novel drug sensitivity phenotype. **(A)** Model
 588 for how genetic variation in *AUA1* and *ARG81* impacts *GAP1* activity. The
 589 relative abundances of metabolites (black) induce regulators (green) to modulate
 590 Gap1p activity. We hypothesize that the wine alleles of *AUA1* and *ARG81* lead to
 591 increased activity of *GAP1*. **(B)** The wine strain does not grow as well as the oak
 592 parent in the presence of a poor nitrogen source, proline (Pro) and toxic D-
 593 Histidine (dHis), indicating that the wine parent has higher *GAP1* activity than the
 594 oak parent. *GAP1* is downregulated in the presence of a strong nitrogen source
 595 like ammonium sulfate (AS).



596

597 **Supplemental Figure 1.** Direct injection mass spectrometry captures biological
598 metabolites. We directly injected extractions into the mass spectrometer and
599 normalized reads by total ion current (red and blue lines depict biological
600 samples wine parent and oak parent respectively; yellow and gray depicts
601 negative samples without cells and without either media or cells respectively). To
602 focus on biologically relevant metabolites, we excluded metabolites whose
603 abundance in the negative controls were equal or greater than half of the mean
604 of all biological samples. We excluded the peaks depicted here except the peaks
605 at 173.17 m/z and 174.13 m/z.



606

607 **Supplemental Figure 2.** Direct injection mass spectrometry reproducibility. **(A)**

608 We mixed the oak and wine parents in equal amounts, and extracted the mix to

609 create a standard that was directly injected into the mass spectrometry at least

610 once every 48 runs. We show the raw spectra for two runs of the standard. **(B)**

611 Histogram of the coefficient of variation for each metabolite measured over 11

612 replicates of the standard. Frequency is the number of metabolites for a given

613 coefficient of variation.