1

2

# Identifying widespread and recurrent variants of genetic parts to improve annotation of engineered DNA sequences

6

7

8  Matthew J. McGuffie[1], Jeffrey E. Barrick[1]*

9

10  [1] Department of Molecular Biosciences, Center for Systems and Synthetic Biology,

11  The University of Texas at Austin, Austin, Texas, United States of America

12

13  * Corresponding author

14  E-mail: jbarrick@cm.utexas.edu (JEB)

# Abstract

Engineered plasmids have been workhorses of recombinant DNA technology for nearly half a century. Plasmids are used to clone DNA sequences encoding new genetic parts and to reprogram cells by combining these parts in new ways. Historically, many genetic parts on plasmids were copied and reused without routinely checking their DNA sequences. With the widespread use of high-throughput DNA sequencing technologies, we now know that plasmids often contain variants of common genetic parts that differ slightly from their canonical sequences. Because the exact provenance of a genetic part on a particular plasmid is usually unknown, it is difficult to determine whether these differences arose due to mutations during plasmid construction and propagation or due to intentional editing by researchers. In either case, it is important to understand how the sequence changes alter the properties of the genetic part. We analyzed the sequences of over 50,000 engineered plasmids using depositor metadata and a metric inspired by the natural language processing field. We detected 217 uncatalogued genetic part variants that were especially widespread or were likely the result of convergent evolution or engineering. Several of these uncatalogued variants are known mutants of plasmid origins of replication or antibiotic resistance genes that are missing from current annotation databases. However, most are uncharacterized, and 3/5 of the plasmids we analyzed contained at least one of the uncatalogued variants. Our results include a list of genetic parts to prioritize for refining engineered plasmid annotation pipelines, highlight widespread variants of parts that warrant further investigation to see whether they have altered characteristics, and suggest cases where

37    unintentional evolution of plasmid parts may be affecting the reliability and

38    reproducibility of science.

## Author Summary

40          Plasmids are used in molecular biology and biotechnology for a wide variety of

41    tasks such as cloning DNA, expressing recombinant proteins, and creating vaccines.

42    One challenge in working with plasmids is that there has been a long, and often lost

43    history of pieces of plasmids being copied and remixed by researchers to create new

44    plasmids. Current databases used for annotating key genetic parts in plasmids are

45    incomplete, especially with respect to cataloguing closely related versions of parts that

46    can have very different characteristics. Some genetic part variants have arisen due to

47    purposeful editing while others are the result of unplanned mutations and evolution.

48    When a researcher finds differences between a database sequence and a genetic part

49    in their newly constructed plasmid, it is often unclear how and when it arose and

50    whether it will affect their experiments. We identified 217 genetic part variants that are

51    either widespread or have likely arisen independently more than once on plasmids due

52    to convergent evolution or engineering. We propose that these variants should be

53    prioritized for inclusion in curated databases of engineered DNA sequences and for

54    functional characterization to improve the reliability and reproducibility of science.

## Introduction

56          Engineered plasmids are ubiquitous tools in the biological sciences. They are

57    used for a wide variety of tasks, ranging from routine cloning of recombinant DNA and

58    protein overexpression to reprogramming cells with new enzymes, sensors, and genetic

59    circuits [1–3]. Engineering plasmids by assembling DNA from different natural sources

60    began in 1973 with the construction of plasmid pSC101 [4]. Chemically synthesizing

61    DNA sequences and introducing them into plasmids has now been commonplace for

62    decades [5]. Many plasmids have been passed from researcher to researcher, and their

63    genetic parts have been copied and remixed, practices facilitated by plasmid

64    repositories [6–8]. The net result is that the genetic components on any plasmid used in

65    a laboratory today often have long, circuitous, and usually incompletely known histories.

66    It has only been standard practice to check the sequences of certain pieces of plasmids,

67    such as by Sanger sequencing a gene of interest inserted by a researcher into a vector

68    backbone, to validate that they are present exactly as designed. Large portions of these

69    plasmids, including origins of replication and antibiotic resistance genes that are critical

70    for plasmid maintenance, are typically assumed to be immutable or to have only

71    sustained mutations with no effect on their performance.

72         Recently, DNA sequencing has become much more affordable and high-

73    throughput [9,10]. Computational pipelines have been developed for assembling

74    accurate and complete plasmid sequences [11–13], and researchers now have

75    complete information about pieces of plasmids that were rarely sequenced in the past.

76    These full plasmid sequences reveal that there are often discrepancies, usually of one

77    to a few nucleotides, between the actual parts on a plasmid and their expected,

78    canonical sequences. Plasmid DNA sequences need to be annotated with information

79    about the genetic parts they contain so that their contents can be checked. Annotation

80    programs, such as PlasMapper [14], and commercial software, like SnapGene, tolerate

81    some variation in the matches they report to the consensus sequence for a genetic part

82    in a database. However, they do not alert a user when they encounter these imperfect

83    matches, which may obscure changes in the sequence of a part that have functional

84    consequences. We recently developed a plasmid annotation tool, pLannotate [15], that

85    reports the nucleotide identity of imperfect matches so users can evaluate parts that are

86    not in agreement with the reference sequences.

87         When a researcher encounters a change from the consensus sequence for a

88    critical genetic part, they are confronted with questions and choices. Should they use

89    the plasmid "as is" or spend time trying to correct the change? Does the change matter

90    for the function of the genetic part? Was the change an edit that was introduced by a

91    prior researcher for some forgotten purpose or was it due to a random mutation?

92         Unfortunately, there is no comprehensive central repository of genetic part

93    sequences that a researcher can consult to answer these questions. Databases like

94    iGEM's Registry of Standard Biological Parts [16], the Joint BioEnergy Institute's

95    Inventory of Composable Elements (JBEI ICE) [17], and SynBioHub [18] contain many

96    plasmid and genetic part sequences. However, they are not fully curated and are known

97    to also contain spurious and incorrect information [19]. GenoLIB [20] and the related

98    SnapGene database are computationally and manually compiled databases of a

99    fundamental set of 293 common plasmid parts. They include multiple, curated entries

100   for major families of related parts (e.g., different aminoglycoside resistance genes), but

101   do not attempt to capture the functional implications of more subtle sequence variation.

102   Only specialized databases reach this level of precision (e.g., FPbase for fluorescent

103   proteins) [21]. These resources do not exist for most categories of critical genetic parts.

104    How do new variants of genetic parts found on engineered plasmids originate?

105    Often these changes are due to researchers finding ways to improve or modify part

106    performance. For example, the *lacI^q* promoter has a single base change that increases

107    its transcription initiation rate by 10-fold relative to the wild-type *lacI* promoter found in

108    the *E. coli* genome [22]. Hundreds of fluorescent proteins have been engineered by

109    introducing changes into natural sequences to alter their spectra, stability, maturation

110    rates, and other properties for imaging applications [21]. CRISPR interference

111    (CRISPRi) uses a catalytically dead Cas9 (dCas9) for the purposes of knocking down

112    gene expression [23]. This variant has two mutations that inactivate the nuclease

113    domain of Cas9, and these mutations have been engineered independently by different

114    groups in Cas9 proteins encoded by different plasmid lineages [24,25]. Other changes

115    may have purposes that are more difficult to ascertain, such as when researchers

116    introduce silent changes in protein-coding sequences to add or avoid restriction enzyme

117    cut sites to make parts compatible with certain DNA assembly methods.

118    Further complicating the picture, genetic part variants can also arise due to

119    evolution. Mutations occur when DNA sequences are copied and assembled into new

120    plasmids *in vitro*. When a single-cell transformant of a plasmid is picked, any mutations

121    it harbors become fixed in all of that plasmid's progeny. There are further opportunities

122    for mutations to arise due to *in vivo* errors in DNA replication and repair as plasmids are

123    propagated in bacterial cells. If the mutated plasmid functions as expected by a

124    researcher, and they don't detect or reject a mutation when validating the plasmid

125    sequence, it will be retained. In some cases, selection will even favor mutated plasmids.

126    Engineered plasmids can impose a significant fitness burden on the host cell if they

127  divert resources needed for cellular replication or produce toxic products [26–29]. In

128  these cases, there is a strong selection pressure favoring cells with plasmids mutated in

129  ways that alleviate this burden by reducing or eliminating the designed function [30–33].

130  Researchers may also impose other types of selection on part/plasmid function, by

131  picking the most fluorescent or largest colonies after a transformation, for example.

132        Precisely annotating the presence and properties of common genetic part

133  variants—whether they result from undocumented engineering or unintentional

134  evolution—is key to improving reliability and reproducibility in the biological sciences.

135  However, there are many of these variants, and determining which ones to prioritize for

136  time-consuming manual curation and experimental characterization is a challenge.

137  Here, we develop methods for computationally identifying widespread genetic part

138  variants and variants that recurrently arose from convergent engineering or evolution

139  given a large set of plasmid sequences. We use these approaches to create a list of

140  217 currently uncatalogued genetic part variants that should be prioritized for further

141  characterization and inclusion in annotation databases.

142  # Results

143  ## Variants of canonical genetic part sequences are common in

144  ## engineered plasmids

145        We used pLannotate [15] to annotate 983,436 genetic parts in 51,384

146  engineered plasmids in the Addgene repository [6,7] that have been fully sequenced.

147  We found 171,828 examples of parts that did not match their canonical sequences

148  present in the databases used for annotation. These part variants can be broadly

149  classified into 14 different categories (**Fig 1**). As expected, we observed more variants

150  for more common types of parts and for types of parts that generally have longer

151  sequences. The most common non-canonical plasmid parts are protein-coding

152  sequences, with 73,884 total variants observed, which are comprised of 10,406 distinct

153  variant sequences (**Fig 1A**). The part type that had the next greatest number of variants

154  was origins of replication (46,677 observations of 607 distinct variant sequences), and

155  the third most common variant type was promoters (24,319 observations of 905 distinct

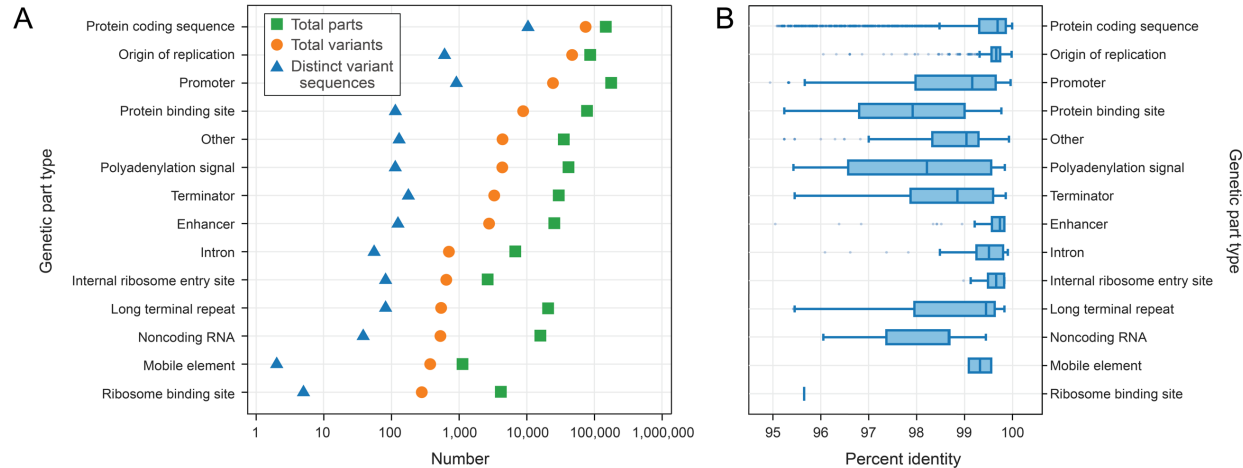156  variant sequences).



157

158  **Fig 1. Many non-canonical genetic parts are found on plasmids.** (A) Overall

159  representation in Addgene plasmids of genetic part variants with sequences that differ

160  slightly from canonical features present in annotation databases. Within each part type,

161  the total number of genetic parts (green squares), total number of genetic parts that are

162  variants (i.e., differ from the canonical sequence) (orange circles), and number of

163  distinct genetic part variant sequences (i.e., counting each unique sequence that differs

164  from the canonical sequence one time) (blue triangles) are plotted. Part types are sorted

165  in descending order by the number of total variants in each category. (B) Distributions of

166    percent identity between distinct genetic part variants in each category and their

167    canonical sequences. Boxes represent lower and upper quartiles (the interquartile

168    range). Vertical lines within each box are medians. The whiskers correspond to 1.5

169    times the interquartile range. Points are outliers outside this range.

170

171        Variants of protein coding sequences and origins of replication are relatively

172    close in sequence to their database counterparts. Variants of smaller parts, such as

173    promoters or protein binding sites, exhibit higher relative levels of sequence divergence

174    (**Fig 1B**). Some of the variants we found are known but not differentiated in current

175    databases used for plasmid annotation. For example, pLannotate and SnapGene

176    currently have a single database entry for the ColE1 plasmid origin of replication, which

177    is the pBR322 variant, the sequence found in a natural plasmid. However, most

178    plasmids contain the engineered pUC19 variant of this origin, which includes a single

179    point mutation that increases plasmid copy number by a factor of about 10-fold [34,35].

## Some widespread genetic part variants are found on plasmids created by many different labs

182        The sheer number of plasmid part variants is a challenge for improving plasmid

183    annotation. Our goal is to determine which variants should be catalogued and prioritized

184    as candidates for further investigation, better documentation, and inclusion in annotation

185    databases. The naïve approach would be to catalog all previously undocumented

186    variants, but this is not practical. Engineered plasmids experience severe population

187    bottlenecks when they are constructed and propagated in the laboratory. When

188    plasmids are transformed into a population of cells, typically only a single plasmid

189    enters a successful transformant. It is also standard practice to re-streak cells and

190    isolate a colony derived from a single cell when obtaining a new plasmid from another

191    researcher or from a repository. Therefore, many part variants may be a result of recent

192    genetic drift (fixation of mutations due to chance) caused by these extreme population

193    bottlenecks. Cataloging these "random" variants is not likely to be particularly

194    informative, especially if they are found in just one or a few plasmids.

195          One might, therefore, propose documenting part variants with the most overall

196    observations. However, this strategy still encounters the same issue. Most variants are

197    found on sets of plasmids deposited by just one or two labs (**Fig 2A**), and some of these

198    variants have become prevalent due to chance (**Fig 2B**). These cases typically occur

199    when a single lab deposits a collection of hundreds of related plasmids that all share the

200    same unique variant of a genetic part. For example, one lab deposited 597 highly

201    similar plasmids, which includes their general lab plasmids as well as a subset used for

202    expressing human SH3 domains [36]. These plasmids all share a single base change in

203    the ColE1 origin of replication. This mutation was almost certainly present in the

204    backbone of an ancestral plasmid they inherited, and its propagation does not seem to

205    be intentional. Even though this variant is the most common origin of replication variant

206    measured in terms of the gross number of observations (besides the canonical pUC19

207    variant), we would assign it a relatively low priority for characterization since it appears

208    to be a one-off mutation that was unintentionally cloned into one set of related plasmids.
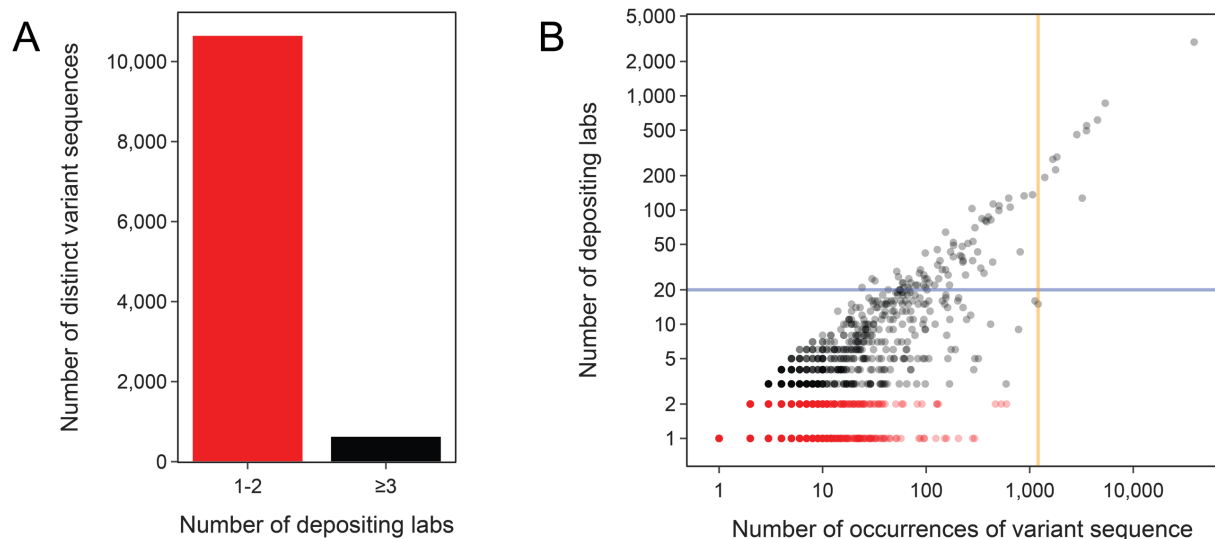
**Fig 2. Most genetic part variants are found in plasmids from one or two labs, but some are more widespread.** (A) Total number of distinct variant sequences found in plasmids from one or two depositing labs (1-2) versus found in plasmids from three or more depositing labs (≥3). (B) All genetic part variants plotted by how many times they were observed versus the number of labs that deposited a plasmid with that variant. The blue horizontal line at 20 labs is the minimum threshold we used for selecting variants that were widespread. The orange vertical line at 1205 variant observations is the cutoff above which we did not perform the authorship analysis to find cases of convergent evolution or engineering.

While deciding which variants to prioritize based on their raw frequency may not be particularly useful, we believe that cataloging variants found in plasmids deposited by many independent labs does have value. In this case, these variants may also have arisen due to chance in a single progenitor plasmid, but this event likely occurred years or decades in the past, so the potential impact has spread such that it could be affecting

225      many more researchers and experiments. Therefore, we flagged all 75 genetic part

226      variants found in plasmids from least 20 labs (**Fig 2B**, above the blue horizontal line) for

227      inclusion in our set of high-priority variants of interest.

228      **Recurrent engineering or evolution of unannotated genetic**

229      **part variants can be predicted using a design similarity score**

230           Variants that are from a few or a middling number of labs are harder to classify. If

231      a variant appears in unrelated plasmids, it could be an engineered variant that is

232      missing from current annotation databases or an evolved variant that arose more than

233      once in unrelated plasmid lineages. Whether designed or evolved, these recurrent

234      mutations are especially likely to affect the function of a part, so it is a high priority to

235      document these cases even if they are in fewer total plasmids. To identify likely

236      examples of convergent engineering and evolution, we analyzed plasmids as authored

237      works. In the natural language processing and information retrieval fields, inverse

238      document frequency (IDF) [37,38] is a metric employed to predict shared authorship

239      [39–41]. IDF scores the rarity of a word or phrase by counting the observations within a

240      document and compares that to its relative frequency in an entire corpus of documents.

241      We created an IDF-inspired metric for use with biological sequences, calculating a

242      quantity that we term the design similarity (DS) score and using it to group plasmids.

243           The procedure we developed to analyze sets of plasmids containing the same

244      part variant (shared unique word) for signs of shared authorship is shown in **Fig 3**. We

245      began by identifying all other contiguous sequence segments shared by these plasmids

246      (shared phrases between documents) and tabulating the frequencies of each of these

247      segments in the entire database of all plasmids (how rare the phrases are). We

248    calculated a DS score for each pair of plasmids from these frequencies. Then, we

249    grouped plasmids by constructing a network graph from an adjacency matrix of these

250    DS scores. This step used a score cutoff determined by examining the distribution of DS

251    scores between random plasmids from different labs (**Fig 4**, top). Finally, we divided the

252    resulting network graph into connected clusters that represent groups of plasmids that

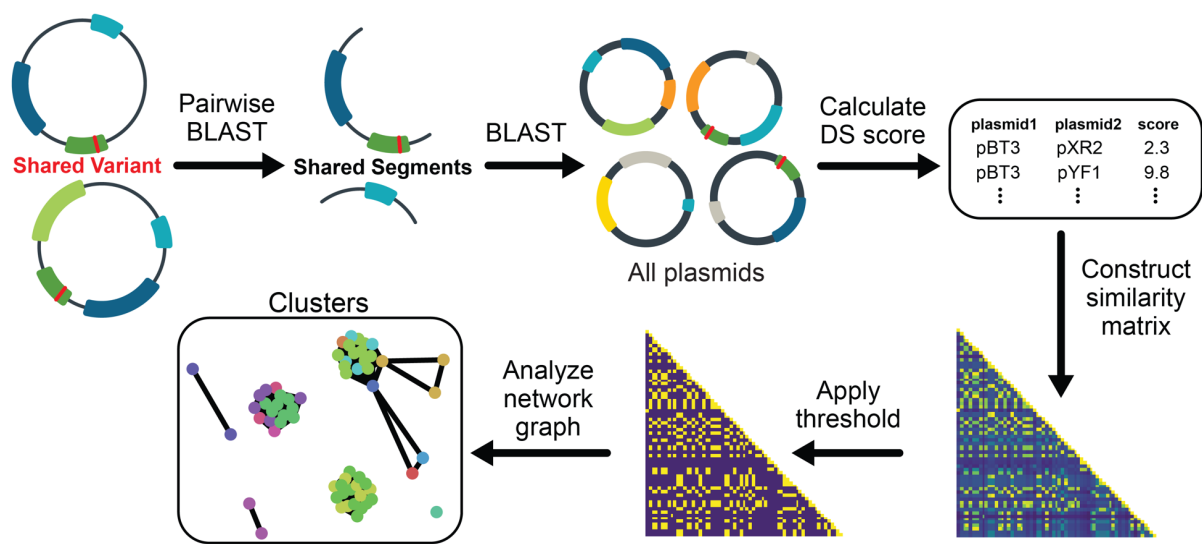253    are unlikely to share the part variant due to common descent or copying of the part.

254    

255    **Fig 3. Method for identifying recurrent genetic part variants that likely arose from**

256    **convergent evolution or engineering.** All plasmids containing the same genetic part

257    variant are analyzed as a set. Segments shared by each pair of these plasmids are

258    identified and queried against the full plasmid database. The results are used to

259    calculate a design similarity (DS) score between the two plasmids. DS scores for all

260    comparisons are used to construct a network graph of plasmid relatedness. Each

261    separate cluster in the final graph is predicted to represent a set of plasmids in which

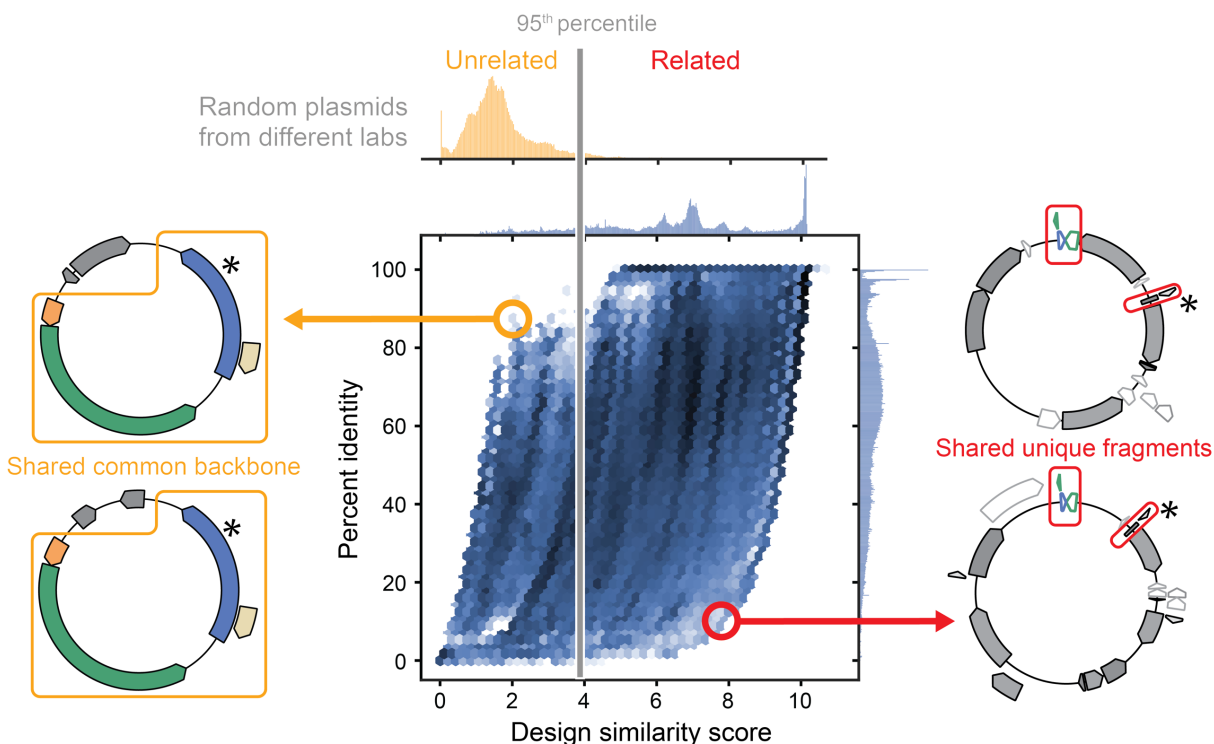262    the variant arose independently.

263

264

**Fig 4. Design similarity scores reliably identify plasmids that are likely to be related while percent identity does not.** The distributions of DS scores and percent identities for pairwise comparisons of plasmids that share undocumented part variants are plotted. Every plasmid containing a given genetic part variant that was observed 1205 or fewer total times was compared to every other plasmid with that part variant for a total of 7,508,114 comparisons. High pairwise percent identity is not compelling evidence that plasmids are related when they share a commonly used backbone, as illustrated by the plasmid pair shown to the left. The DS score of these two plasmids is low in this instance. Low pairwise percent identity also does not necessarily indicate that plasmids are unrelated, as illustrated by the plasmid pair shown to the right. In this case, a high DS score highlights small, but unique sequences present in both plasmids, which is evidence of shared authorship. Asterisks indicate the location of the shared mutation in the associated genetic part variant that differentiates it from the canonical

278    sequence in the annotation database. The distribution of DS scores between 100,000

279    randomly selected pairs of plasmids from different labs is shown above the plot. The

280    grey line indicates the 95th percentile of the distribution, which was used as the score

281    cutoff for shared plasmid authorship.

282

283        If multiple distinct authorship clusters are predicted for a variant, it likely had

284    more than one independent origin due to recurrent engineering or evolution. In this

285    case, it should be a priority to document the variant and further characterize whether its

286    function differs from that of the canonical sequence. Because the DS scoring algorithm

287    involves making pairwise comparisons of all plasmids containing a given genetic part

288    variant, it was only computationally feasible for us to apply it to variants with 1205 or

289    fewer observations (**Fig. 2B**, left of orange vertical line), which included all variants

290    found on plasmids deposited by fewer than 20 labs that we had not already flagged as

291    being of interest simply because they were widespread. As expected, plasmids sharing

292    a variant that were deposited by the same lab are almost always found within a single

293    cluster at the end of this procedure. This tracks with the intuition that a depositing lab

294    likely recycles their plasmid backbones and pieces of those plasmids for various

295    purposes. In total, 149 of the variants tested using the DS clustering procedure were

296    predicted to occur in two or more author groups. This total includes 7 of the 64 variants

297    tested in this way that were found in plasmids deposited by 20 or more labs.

298        Using the DS score as a metric has advantages over using a percent identity-

299    cutoff to determine if instances of the same genetic part variant on two plasmids are

300    related (**Fig 4**). Any two plasmids often share extensive stretches of DNA, but this may

301    not actually indicate anything about how related the plasmids are to each other. For

302    example, the ColE1 origin of replication is used in nearly 95% of the plasmids in our

303    dataset, and 62% of plasmids contain β-lactamase as an antibiotic resistance marker.

304    Since these features are widely used, their co-occurrence is not convincing evidence

305    that a pair of plasmids is related, even if they constitute a majority of the shared

306    sequence identity between them (**Fig. 4**, left). The DS metric weights features based on

307    their overall rarity rather than their length or context, so that even a small part or cloning

308    scar can be a strong signal of shared authorship (**Fig. 4**, right).

## 309 Final list of widespread and recurrent genetic part variants

## 310 includes known but uncatalogued mutants

311       We combined the widespread and recurrent part variants we identified into a final

312    list of 217 currently uncatalogued genetic part variants (**S1 Table**). This list includes

313    diverse genetic parts with a wide range of functions that are used for engineering all

314    kinds of organisms (**Fig 5**). For parts designed to function in bacteria, most of the newly

315    identified variants of interest were plasmid origins of replication or antibiotic resistance

316    markers. For eukaryotic parts, promoter variants were most common. Many fluorescent

317    proteins, which function in both types of organisms, were also present in this list of

318    uncatalogued variants not found in current annotation databases.

319

**Fig 5**. **Uncatalogued genetic part variants to prioritize for characterization and inclusion in annotation databases. (**A) The final 217 variants of interest categorized by part type and by the kind of organism in which the part is typically used. Bars are shaded according to the method by which each variant was judged to be a priority for characterization and annotation: either it occurred in plasmids from ≥20 depositing labs (widespread, orange) or it was in plasmids from fewer labs but there was evidence that it was engineered or evolved multiple times from the authorship analysis (convergent, blue). (B) Names of the canonical parts to which the 217 variants are most closely related. Parts are categorized and sorted by function.

329

330     To validate our inclusion criteria, we looked for cases of known variants that were

331     uncatalogued in the initial annotation databases but were identified by our analysis. The

332     top two variants with 38,693 and 25,995 total observations are the pUC19 variant of the

333     ColE1 origin of replication and TEM-116 β-lactamase antibiotic resistance marker,

334     respectively (Fig 5B). These are both engineered variants that differ from their parent

335     sequences, pBR322 and TEM-1, by one or two bases, respectively [35,42]. These

336     variants were included in our list because they occurred in ≥20 labs. We also identified

337     one other canonical variant, TEM-171, which was both a frequent and recurrent variant.

338     TEM-171 has one of the two mutations that TEM-116 has relative to TEM-1 [42].

339     As an example of how these predictions can aid in directing efforts to refine

340     annotations of engineered DNA, one fluorescent protein variant in our list had a clear

341     signal of a recurrent origin due to convergent engineering. Seventeen plasmids with the

342     variant that were deposited by five different labs were from four authorship clusters.

343     This variant is a derivative of enhanced GFP (eGFP) originally described in 1996 by

344     Cormack et al. [43] with additional A164V and G176S amino acid substitutions. This

345     derivative of eGFP is not currently listed in FPbase, and none of the five publications

346     associated with the plasmids containing this derivative mention its provenance or the

347     mutations it harbors [44–48], so their effects on its function are unknown.

## Discussion

349     It is becoming standard practice for researchers to fully sequence plasmids and

350     other engineered DNA constructs they use in their experiments [11,49]. These

351     sequences need to be validated by precisely annotating the genetic parts they contain

352     and recognizing unexpected sequence variation in these parts in order to ensure the

353    reliability and reproducibility of science. In the work reported here, we created a list of

354    217 currently uncatalogued variants of common genetic parts that can be added to

355    databases used by annotation pipelines. These variants are a priority because they are

356    either already widespread in plasmids being exchanged by researchers or they appear

357    to have originated multiple times due to convergent engineering or evolution.

358        Many of the variants in our final list are in high-copy ColE1-family origins of

359    replication or in antibiotic resistance cassettes that are commonly paired with these

360    origins in *E. coli* vectors used for cloning and replicating DNA. These are by far the most

361    common genetic parts in Addgene plasmids because pUC vectors are used to

362    manufacture high-quality DNA for many applications, ranging from *in vitro* transcription

363    of RNA for biochemical studies to transfection into mammalian cells. Sequence variation

364    in these backbone components might affect cloning success or DNA yields, if a

365    mutation alters plasmid copy number, for example. But, these differences would be

366    unlikely to affect the results of downstream experiments after DNA is isolated from

367    bacterial cells. On the other hand, variants in other origins of replication that we

368    identified, such as the medium-copy p15A origin that is commonly used in plasmids

369    encoding synthetic biology devices meant to function in *E. coli* and the broad-host-range

370    pBBR1 origin that is used for engineering diverse bacteria, are more likely to affect

371    research outcomes. Overall, this logic argues for prioritizing characterization of part

372    variants that are important in the ultimate context in which the DNA will be used, which

373    includes many variants in our final list related to eukaryotic gene expression.

374        To detect recurrent variants that likely arose multiple times, we developed an

375    approach for grouping plasmids based on signals of shared authorship. Previously,

authorship of plasmid sequences has been analyzed from a biosecurity standpoint, with the aim of attributing an unknown plasmid to a specific lab [50,51]. All of these prior studies analyzed the Addgene plasmid corpus. The first used deep learning to train a convolutional neural network to predict the lab of origin of a plasmid from its DNA sequence [52]. It correctly identified the source lab 48% of the time and the source lab appeared in the top 10 predicted labs 70% of the time. A comparable method, deteRNNt, used recurrent neural networks trained on plasmid sequences and associated phenotype data to identify DNA motifs indicative of different genetic designers [53]. It demonstrated an improvement in accuracy to 70% correct attribution to one lab among 1,300 in the dataset. An alternative approach, PlasmidHawk [54], opted to not use deep learning, citing the higher accuracy and higher interpretability of sequence alignment-based techniques compared to machine learning approaches. Their approach had 76% accuracy in identifying the lab that deposited an unknown plasmid and could precisely single out the signature sub-sequences responsible for a prediction. Notably, this study used an approach similar to our own where they down-weighted observations of sequence segments that are frequent in the overall dataset, though their metrics differ from our IDF-inspired design similarity score.

We had to infer shared authorship of plasmids to predict when a variant had arisen multiple times because the cloning history of most plasmids is not fully known. Ideally, one would be able to track the provenance of plasmids and their parts using the scientific literature and/or metadata in plasmid repositories to understand which changes to the sequence of a genetic part were intentional and when and how many times they were introduced or arose due to mutations. QUEEN is a recent framework

399  which proposes to record traceable linages of engineered plasmids by having

400  researchers meticulously document their construction process and store this information

401  as metadata in GenBank flat files [55]. Addgene is now encouraging researchers to use

402  QUEEN when submitting new constructs. If this or a similar metadata format for tracking

403  how engineered DNA sequences have been copied, remixed, and modified is widely

404  adopted, it will be very useful for tracking the engineering and evolution of plasmids in

405  the future. Many scientists who performed foundational research creating key plasmid

406  backbones and genetic parts in the early days of recombinant DNA technology are

407  retired or will be soon. It would be extremely valuable if the community could also

408  capture or reconstruct their knowledge of earlier plasmid construction efforts.

409       pLannotate and other plasmid annotation pipelines use BLAST to find matches to

410  genetic part sequences in a database. This simple approach has some potential

411  shortcomings with respect to variant detection and prediction. One is that BLAST

412  matches may not detect instances of a part or properly delineate their extent when there

413  are mutations at or near its ends. For example, if a bacterial promoter variant has a

414  mismatch in the −35 box at the end of the canonical promoter core sequence and this is

415  also where the part sequence in a database ends, the BLAST hit may only match the

416  downstream part of the promoter. This could result in reporting an incomplete match

417  that is not recognized as a variant or potentially no match at all. Compounding this

418  problem is the issue that some types of genetic parts and important functional variants

419  of these parts can be defined on multiple, overlapping scales. For a bacterial promoter,

420  the database sequence could be just the core element containing the −10  and −35

421  boxes, or it could be an extended element that includes upstream sequences such as

422   UP-elements [56] or adjacent cis-regulatory elements. Computational matching

423   methods that force extending alignments to the boundaries of part sequences and

424   expert curation of how a core part and elaborated variants of that part are related could

425   help annotation programs deal with these difficult cases.

426        Ideally, we would be able to provide annotation programs with detailed

427   information to accompany the sequences of the 217 high-priority variants we identified,

428   including their provenance and functional characteristics. It may be possible to trace

429   more of our variants of interest to existing publications in which a researcher engineered

430   mutations on purpose. However, this will require analyzing hundreds or thousands of

431   publications. Since some variants are bound to be the result of *de novo* mutations in the

432   laboratory, these searches will sometimes come up empty. In these cases, one needs

433   to test whether and how the performance of the part variant differs from the canonical

434   sequence and associate that information with the database sequence. Such efforts will

435   take years of expert curation and laboratory experiments by a community of scientists.

436   A framework is needed to centrally collect and organize this information and encourage

437   community participation. FBbase is an outstanding example of continuous and expert

438   curation of a specific type of engineered part [21]. This type of resource needs to be

439   extended to more types of genetic parts. Integrating work on documenting part variants

440   using a micropublication [57,58] or wiki model [59] could be ways to recognize the

441   contributions of curators and researchers to this kind of resource, hopefully including

442   those with first-hand knowledge of the histories of important genetic parts. In the end, a

443   combination of computational and community-based curation efforts will likely be the

444   most effective path forward for improving plasmid annotation.

# Conclusions

As fully sequencing engineered plasmids becomes commonplace, researchers are encountering an overwhelming number of uncatalogued variants of canonical genetic parts and being forced to reckon with whether these differences are important or not. We developed a procedure for predicting variants that are likely to have arisen due to convergent evolution or engineering. We combined these predictions with genetic part variants that are found in plasmids from many labs, under the premise that both widespread and recurrent variants are more likely to affect the function of a genetic part and the reproducibility of research than random one-off changes. Genetic part variants in our final list of 217 predictions warrant further investigation and should be integrated into tools that annotate engineered DNA. This work is a promising step towards automating better plasmid annotation, but there is still a need for integrating this information with expert curation to create comprehensive databases of genetic parts.

# Materials and Methods

## Identification of genetic part variants in engineered plasmids

We downloaded 51,359 complete plasmid sequences from Addgene, a non-profit plasmid repository based in Cambridge, Massachusetts, on August 9$^{th}$, 2021. Plasmid sequences were annotated using pLannotate v1.2.0, which identifies matches to the Swissprot [60] (release 2021_03), Snapgene (2021-07-23), FPbase [21] (2020-09-02), and Rfam [61] (release 14.5) databases. We extracted all annotated features from every plasmid, keeping matches that pLannotate identified as covering ≥ 95% of the length of the feature in the database. Matches that were 100% identical at either the nucleotide or

467    amino acid level to annotation database entries were removed. Protein-coding

468    sequence features with 3′ or 5′ deletions were also removed. The remaining non-

469    consensus features were considered genetic part variants and further analyzed.

470    **Grouping genetic part variants on related plasmids**

471         The design similarity (DS) score is calculated based on a formula that is similar

472    to that for the Inverse Document Frequency (IDF) of the most common segment shared

473    by two plasmids, except extra terms are added when there are multiple segments

474    shared by the two plasmids. For each genetic part variant found in plasmids from two or

475    more depositing labs, we first performed a pairwise BLASTN search (BLAST 2.10.1+)

476    [62] between all plasmids that contained that variant to identify shared plasmid

477    segments. Each of these segments was then queried against the entire database using

478    BLASTN to find the number of plasmids that contained the segment. The following

479    BLASTN parameters were used in both cases: mismatch penalty −8, match reward 2,

480    gap open penalty 4, gap extend penalty 6, and word size 28. These parameters were

481    chosen to maximize the reporting of matches consisting of contiguous segments with

482    few point mutations. A segment match was defined as having ≥98% identity, an E-value

483    ≤ $10^{-5}$, and a length difference of at most 10 bp. The DS score was then calculated

484    using the following equation:

485
$$\text{Design Similarity} \; = \; \log\left(\frac{p}{x_1} + \frac{\sum_{i=2}^{n}\frac{p}{x_i}}{n}\right)$$

486         Where, $x$ is a vector of length $n$ containing the number of plasmids matching

487    each segment query, sorted from the smallest to the largest value. $p$ is the number of

488    reference plasmids in the database. The right term of the equation is an extra score

489    heuristic that is applied when there is more than one matching segment.

490         We also cataloged all variants that were found in plasmids from ≥ 20 depositing

491    labs, irrespective of DS. It was not computationally feasible to calculate pairwise DS

492    scores for variants with > 1,205 observations, but all 11 of these variants were

493    catalogued because they were found on plasmids originating in ≥ 20 labs.

494    **Determining a threshold for plasmid relatedness**

495         To determine a DS score threshold that indicates two examples of a genetic part

496    variant on different plasmids likely shared an ancestor, we examined the distribution of

497    DS scores for 100,000 random plasmid pairs. We picked only plasmid pairs that did not

498    share a common depositing lab to increase the likelihood that we did not include pairs

499    that did share a construction history in this set. We picked a DS cutoff for plasmid

500    relatedness that gave a 5% false-positive rate on this dataset as the metric for calling

501    two plasmids as related.

502         After calculating the pairwise DS scores for each group of plasmids that shared

503    the same genetic part variant, we binarized the results based on the DS score cutoff

504    threshold. The binary adjacency matrices were then analyzed as a network, and we

505    counted the number of unlinked subgraphs within each plasmid network to estimate the

506    number of times the variant had independently appeared.

# Acknowledgments

# References

1. Itakura K, Hirose T, Crea R, Riggs AD, Heyneker HL, Bolivar F, et al. Expression in *Escherichia coli* of a chemically synthesized gene for the hormone somatostatin. Science. 1977;198: 1056–1063. doi:10.1126/science.412251

2. Goeddel DV, Kleid DG, Bolivar F, Heyneker HL, Yansura DG, Crea R, et al. Expression in *Escherichia coli* of chemically synthesized genes for human insulin. Proc Natl Acad Sci U S A. 1979;76: 106–110. doi:10.1073/pnas.76.1.106

3. Van Gaal EVB, Hennink WE, Crommelin DJA, Mastrobattista E. Plasmid engineering for controlled and sustained gene expression for nonviral gene therapy. Pharm Res. 2006;23: 1053–1074. doi:10.1007/s11095-006-0164-2

4. Cohen SN, Chang AC, Boyer HW, Helling RB. Construction of biologically functional bacterial plasmids *in vitro*. Proc Natl Acad Sci U S A. 1973;70: 3240–3244. doi:10.1073/pnas.70.11.3240

527    5.    Itakura K, Rossi JJ, Wallace RB. Synthesis and use of synthetic

528           oligonucleotides. Annu Rev Biochem. 1984;53: 323–356.

529           doi:10.1146/annurev.bi.53.070184.001543

530    6.    Herscovitch M, Perkins E, Baltus A, Fan M. Addgene provides an open forum

531           for plasmid sharing. Nat Biotechnol. 2012;30: 316–317. doi:10.1038/nbt.2177

532    7.    Kamens J. The Addgene repository: an international nonprofit plasmid and data

533           resource. Nucleic Acids Res. 2015;43: D1152–D1157. doi:10.1093/nar/gku893

534    8.    Seiler CY, Park JG, Sharma A, Hunter P, Surapaneni P, Sedillo C, et al.

535           DNASU plasmid and PSI:Biology-Materials repositories: resources to

536           accelerate biological research. Nucleic Acids Res. 2014;42: D1253–D1260.

537           doi:10.1093/nar/gkt1060

538    9.    Kumar KR, Cowley MJ, Davis RL. Next-generation sequencing and emerging

539           technologies. Semin Thromb Hemost. 2019;45: 661–673. doi:10.1055/s-0039-

540           1688446

541    10.   Marx V. Method of the year: long-read sequencing. Nat Methods. 2023;20: 6–

542           11. doi:10.1038/s41592-022-01730-w

543    11.   Gallegos JE, Rogers MF, Cialek CA, Peccoud J. Rapid, robust plasmid

544           verification by de novo assembly of short sequencing reads. Nucleic Acids Res.

545           2020;48: e106. doi:10.1093/nar/gkaa727

546    12.   Emiliani FE, Hsu I, McKenna A. Multiplexed assembly and annotation of

547           synthetic biology constructs using long-read nanopore sequencing. ACS Synth

548           Biol. 2022;11: 2238–2246. doi:10.1021/acssynbio.2c00126

549    13.  Brown SD, Dreolini L, Wilson JF, Balasundaram M, Holt RA. Complete

550        sequence verification of plasmid DNA using the Oxford Nanopore

551        Technologies' MinION device. BMC Bioinformatics. 2023;24: 116.

552        doi:10.1186/s12859-023-05226-y

553    14.  Dong X, Stothard P, Forsythe IJ, Wishart DS. PlasMapper: a web server for

554        drawing and auto-annotating plasmid maps. Nucleic Acids Res. 2004;32:

555        W660–W664. doi:10.1093/nar/gkh410

556    15.  McGuffie MJ, Barrick JE. pLannotate: engineered plasmid annotation. Nucleic

557        Acids Res. 2021;49: W516–W522. doi:10.1093/nar/gkab374

558    16.  Peccoud J, Blauvelt MF, Cai Y, Cooper KL, Crasta O, DeLalla EC, et al.

559        Targeted development of registries of biological parts. PloS One. 2008;3:

560        e2671. doi:10.1371/journal.pone.0002671

561    17.  Ham TS, Dmytriv Z, Plahar H, Chen J, Hillson NJ, Keasling JD. Design,

562        implementation and practice of JBEI-ICE: an open source biological part

563        registry platform and tools. Nucleic Acids Res. 2012;40: e141–e141.

564        doi:10.1093/nar/gks531

565    18.  McLaughlin JA, Myers CJ, Zundel Z, Mısırlı G, Zhang M, Ofiteru ID, et al.

566        SynBioHub: a standards-enabled design repository for synthetic biology. ACS

567        Synth Biol. 2018;7: 682–688. doi:10.1021/acssynbio.7b00403

568    19.  Mante J, Roehner N, Keating K, McLaughlin JA, Young E, Beal J, et al.

569        Curation principles derived from the analysis of the SBOL iGEM data set. ACS

570        Synth Biol. 2021;10: 2592–2606. doi:10.1021/acssynbio.1c00225

571    20. Adames NR, Wilson ML, Fang G, Lux MW, Glick BS, Peccoud J. GenoLIB: a

572          database of biological parts derived from a library of common plasmid features.

573          Nucleic Acids Res. 2015;43: 4823–4832. doi:10.1093/nar/gkv272

574    21. Lambert TJ. FPbase: a community-editable fluorescent protein database. Nat

575          Methods. 2019;16: 277. doi:10.1038/s41592-019-0352-8

576    22. Calos MP. DNA sequence for a low-level promoter of the *lac* repressor gene

577          and an "up" promoter mutation. Nature. 1978;274: 762–765.

578          doi:10.1038/274762a0

579    23. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, et al.

580          Repurposing CRISPR as an RNA-guided platform for sequence-specific control

581          of gene expression. Cell. 2013;152: 1173–83. doi:10.1016/j.cell.2013.02.022

582    24. Bikard D, Jiang W, Samai P, Hochschild A, Zhang F, Marraffini LA.

583          Programmable repression and activation of bacterial gene expression using an

584          engineered CRISPR-Cas system. Nucleic Acids Res. 2013;41: 7429–7437.

585          doi:10.1093/nar/gkt520

586    25. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A

587          programmable dual-RNA-guided DNA endonuclease in adaptive bacterial

588          immunity. Science. 2012;337: 816–821. doi:10.1126/science.1225829

589    26. Sandoval CM, Ayson M, Moss N, Lieu B, Jackson P, Gaucher SP, et al. Use of

590          pantothenate as a metabolic switch increases the genetic stability of farnesene

591          producing *Saccharomyces cerevisiae*. Metab Eng. 2014;25: 215–226.

592          doi:10.1016/j.ymben.2014.07.006

593    27.    Ceroni F, Algar R, Stan G-B, Ellis T. Quantifying cellular capacity identifies

594           gene expression designs with reduced burden. Nat Methods. 2015;12: 415–

595           418. doi:10.1038/nmeth.3339

596    28.    Bentley WE, Mirjalili N, Andersen DC, Davis RH, Kompala DS. Plasmid-

597           encoded protein: the principal factor in the "metabolic burden" associated with

598           recombinant bacteria. Biotechnol Bioeng. 1990;35: 668–681.

599           doi:10.1002/bit.260350704

600    29.    Oliveira PH, Prather KJ, Prazeres DMF, Monteiro GA. Structural instability of

601           plasmid biopharmaceuticals: challenges and implications. Trends Biotechnol.

602           2009;27: 503–511. doi:10.1016/j.tibtech.2009.06.004

603    30.    Sleight SC, Bartley BA, Lieviant JA, Sauro HM. Designing and engineering

604           evolutionary robust genetic circuits. J Biol Eng. 2010;4: 12. doi:10.1186/1754-

605           1611-4-12

606    31.    Rugbjerg P, Myling-Petersen N, Porse A, Sarup-Lytzen K, Sommer MOA.

607           Diverse genetic error modes constrain large-scale bio-based production. Nat

608           Commun. 2018;9. doi:10.1038/s41467-018-03232-w

609    32.    Renda BA, Hammerling MJ, Barrick JE. Engineering reduced evolutionary

610           potential for synthetic biology. Mol Biosyst. 2014;10: 1668–1678.

611           doi:10.1039/C3MB70606K

612    33.    Ellis T. Predicting how evolution will beat us. Microb Biotechnol. 2019;12: 41–

613           43. doi:10.1111/1751-7915.13327

614    34. Yanisch-Perron C, Vieira J, Messing J. Improved M13 phage cloning vectors

615        and host strains: nucleotide sequences of the M13mpl8 and pUC19 vectors.

616        Gene. 1985;33: 103–119. doi:10.1016/0378-1119(85)90120-9

617    35. Lin-Chao S, Chen W-T, Wong T-T. High copy number of the pUC plasmid

618        results from a Rom/Rop-suppressible point mutation in RNA II. Mol Microbiol.

619        1992;6: 3385–3393. doi:10.1111/j.1365-2958.1992.tb02206.x

620    36. Teyra J, Huang H, Jain S, Guan X, Dong A, Liu Y, et al. Comprehensive

621        analysis of the human SH3 domain family reveals a wide variety of non-

622        canonical specificities. Struct Lond Engl 1993. 2017;25: 1598-1610.e3.

623        doi:10.1016/j.str.2017.07.017

624    37. Sparck Jones K. A statistical interpretation of term specificity and its application

625        in retrieval. J Doc. 1972;28: 11–21. doi:10.1108/eb026526

626    38. Fung BCM, Wang K, Ester M. Hierarchical document clustering using frequent

627        itemsets. Proceedings of the 2003 SIAM International Conference on Data

628        Mining (SDM). Society for Industrial and Applied Mathematics; 2003. pp. 59–70.

629        doi:10.1137/1.9781611972733.6

630    39. Cota RG, Gonçalves MA, Laender AHF. A heuristic-based hierarchical

631        clustering method for author name disambiguation in digital libraries. XXII

632        Simpósio Brasileiro de Banco de Dados. 2007. pp. 20–34.

633    40. Layton R, McCombie S, Watters P. Authorship attribution of IRC messages

634        using inverse author frequency. 2012 Third Cybercrime and Trustworthy

635        Computing Workshop. 2012. pp. 7–13. doi:10.1109/CTC.2012.11

636     41.   Nizamani S, Memon N. CEAI: CCM-based email authorship identification

637           model. Egypt Inform J. 2013;14: 239–249. doi:10.1016/j.eij.2013.10.001

638     42.   Jacoby GA, Bush K. The curious case of TEM-116. Antimicrob Agents

639           Chemother. 2016;60: 7000–7000. doi:10.1128/AAC.01777-16

640     43.   Cormack BP, Valdivia RH, Falkow S. FACS-optimized mutants of the green

641           fluorescent protein (GFP). Gene. 1996;173: 33–38. doi:10.1016/0378-

642           1119(95)00685-0

643     44.   Schlüter OM, Xu W, Malenka RC. Alternative N-terminal domains of PSD-95

644           and SAP97 govern activity-dependent regulation of synaptic AMPA receptor

645           function. Neuron. 2006;51: 99–111. doi:10.1016/j.neuron.2006.05.016

646     45.   Lin R, Wang R, Yuan J, Feng Q, Zhou Y, Zeng S, et al. Cell-type-specific and

647           projection-specific brain-wide reconstruction of single neurons. Nat Methods.

648           2018;15: 1033–1036. doi:10.1038/s41592-018-0184-y

649     46.   Santos TE, Schaffran B, Broguière N, Meyn L, Zenobi-Wong M, Bradke F. Axon

650           growth of CNS neurons in three dimensions is amoeboid and independent of

651           adhesions. Cell Rep. 2020;32: 107907. doi:10.1016/j.celrep.2020.107907

652     47.   Wrobel CN, Mutch CA, Swaminathan S, Taketo MM, Chenn A. Persistent

653           expression of stabilized beta-catenin delays maturation of radial glial cells into

654           intermediate progenitors. Dev Biol. 2007;309: 285–297.

655           doi:10.1016/j.ydbio.2007.07.013

656     48.   Beier KT, Kim CK, Hoerbelt P, Hung LW, Heifets BD, DeLoach KE, et al.

657           Rabies screen reveals GPe control of cocaine-triggered plasticity. Nature.

658           2017;549: 345–350. doi:10.1038/nature23888

659   49.   Thuronyi BW, DeBenedictis EA, Barrick JE. No assembly required: Time for

660         stronger, simpler publishing standards for DNA sequences. PLoS Biol. 2023;21:

661         e3002376. doi:10.1371/journal.pbio.3002376

662   50.   Lewis G, Jordan JL, Relman DA, Koblentz GD, Leung J, Dafoe A, et al. The

663         biosecurity benefits of genetic engineering attribution. Nat Commun. 2020;11:

664         6294. doi:10.1038/s41467-020-19149-2

665   51.   Crook OM, Warmbrod KL, Lipstein G, Chung C, Bakerlee CW, McKelvey TG, et

666         al. Analysis of the first genetic engineering attribution challenge. Nat Commun.

667         2022;13: 7374. doi:10.1038/s41467-022-35032-8

668   52.   Nielsen AAK, Voigt CA. Deep learning to predict the lab-of-origin of engineered

669         DNA. Nat Commun. 2018;9. doi:10.1038/s41467-018-05378-z

670   53.   Alley EC, Turpin M, Liu AB, Kulp-McDowall T, Swett J, Edison R, et al. A

671         machine learning toolkit for genetic engineering attribution to facilitate

672         biosecurity. Nat Commun. 2020;11: 6293. doi:10.1038/s41467-020-19612-0

673   54.   Wang Q, Kille B, Liu TR, Elworth RAL, Treangen TJ. PlasmidHawk improves

674         lab of origin prediction of engineered plasmids using sequence alignment. Nat

675         Commun. 2021;12: 1167. doi:10.1038/s41467-021-21180-w

676   55.   Mori H, Yachie N. A framework to efficiently describe and share reproducible

677         DNA materials and construction protocols. Nat Commun. 2022;13: 2894.

678         doi:10.1038/s41467-022-30588-x

679   56.   Ross W, Gosink KK, Salomon J, Igarashi K, Zou C, Ishihama A, et al. A third

680         recognition element in bacterial promoters: DNA binding by the alpha subunit of

681    RNA polymerase. Science. 1993;262: 1407–1413.

682    doi:10.1126/science.8248780

683    57. Clark T, Ciccarese PN, Goble CA. Micropublications: a semantic model for

684    claims, evidence, arguments and annotations in biomedical communications. J

685    Biomed Semant. 2014;5: 28. doi:10.1186/2041-1480-5-28

686    58. Raciti D, Yook K, Harris TW, Schedl T, Sternberg PW. Micropublication:

687    incentivizing community curation and placing unpublished data into the public

688    domain. Database. 2018;2018: bay013. doi:10.1093/database/bay013

689    59. Burge SW, Daub J, Eberhardt R, Tate J, Barquist L, Nawrocki EP, et al. Rfam

690    11.0: 10 years of RNA families. Nucleic Acids Res. 2013;41: D226–D232.

691    doi:10.1093/nar/gks1005

692    60. Bairoch A, Boeckmann B. The SWISS-PROT protein sequence data bank.

693    Nucleic Acids Res. 1991;19: 2247–2249.

694    61. Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K,

695    Marz M, et al. Rfam 14: expanded coverage of metagenomic, viral and

696    microRNA families. Nucleic Acids Res. 2021;49: D192–D200.

697    doi:10.1093/nar/gkaa1047

698    62. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment

699    search tool. J Mol Biol. 1990;215: 403–410. doi:10.1016/s0022-2836(05)80360-

700    2

701

## Supporting Information

703    **S1 Table. Final list of 217 widespread and/or recurrent genetic part variants.**