

METHODS

ConNIS and labeling instability: New statistical methods for improving the detection of essential genes in TraDIS libraries

Moritz Hanke^{1*}, Theresa Harten², Ronja Foraita¹

1 Department Statistical Methods in Epidemiology, Leibniz Institute for Prevention Research and Epidemiology – BIPS, Bremen, Bremen, Germany, **2** Independent Researcher, Hamburg, Germany

* hanke@leibniz-bips.de



Abstract

The identification of essential genes in *Transposon Directed Insertion Site Sequencing* (*TraDIS*) data relies on the assumption that transposon insertions occur randomly in non-essential regions, leaving essential genes largely insertion-free. While intra-genic insertion-free sequences have been considered as a reliable indicator for gene essentiality, so far, no exact probability distribution for these sequences has been proposed. Further, many methods require setting thresholds or parameter values *a priori* without providing any statistical basis, limiting the comparability of results. Here, we introduce *Consecutive Non-Insertion Sites* (*ConNIS*), a novel method for gene essentiality determination. *ConNIS* provides an analytic solution for the probability of observing insertion-free sequences within genes of given length and considers variation in insertion density across the genome. Based on an extensive simulation study and different real-world scenarios, *ConNIS* was found to be superior to prevalent state-of-the-art methods, particularly when libraries had only a low or medium insertion density. In addition, our results showed that the precision of existing methods can be improved by incorporating a simple weighting factor for the genome-wide insertion density. To set methodically embedded parameter and threshold values of *TraDIS* methods a subsample-based instability criterion was developed. Application of this criterion in real and synthetic data settings demonstrated its effectiveness in selecting well-suited parameter/threshold values across methods. An R package and an interactive web application are provided to facilitate application and reproducibility.

OPEN ACCESS

Citation: Hanke M, Harten T, Foraita R (2026) ConNIS and labeling instability: New statistical methods for improving the detection of essential genes in TraDIS libraries. *PLoS Comput Biol* 22(3): e1013428. <https://doi.org/10.1371/journal.pcbi.1013428>

Editor: Jinyan Li, Shenzhen University of Advanced Technology, CHINA

Received: August 12, 2025

Accepted: February 12, 2026

Published: March 6, 2026

Peer Review History: PLOS recognizes the benefits of transparency in the peer review process; therefore, we enable the publication of all of the content of peer review and author responses alongside final, published articles. The editorial history of this article is available here: <https://doi.org/10.1371/journal.pcbi.1013428>

Copyright: © 2026 Hanke et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution,

Author summary

Identifying essential genes in bacteria is key to understanding their ability to survive, which can, for example, be applied to the development of new treatments. One way to do identify these genes is by creating libraries where small DNA fragments (“insertions”) are randomly placed in the genome: essential genes

and reproduction in any medium, provided the original author and source are credited.

Data availability statement: All data and code used for running experiments, model fitting, and plotting is available on a GitHub repository at https://github.com/bips-hb/ConNIS_results. We have also used Zenodo to assign a DOI to the repository in a zip format: <https://doi.org/10.5281/zenodo.16790977>. Additional real-world results are available under <https://zenodo.org/records/18538450>. All results can be interactively explored under <https://connis.bips.eu>. The new methods are made available as R package under <https://github.com/bips-hb/ConNIS>.

Funding: The author(s) received no specific funding for this work.

Competing interests: The authors have declared that no competing interests exist.

tend to remain insertion-free because insertions disrupt their function. The challenge is to determine whether a (long) uninterrupted sequence is due to chance or because the gene is truly essential. Here, we present *Consecutive Non-Insertion Sites (ConNIS)*, a statistical method that calculates the probability of such insertion-free sequences. Extensive comparisons on simulated and real datasets show that *ConNIS* outperforms existing methods, especially when a library is rather sparse in terms of the total number of insertion sites. Since many analysis methods rely on parameter values that have to be set before the analysis and can heavily influence the final results, we also propose a data-driven approach to set these values, making results more comparable across studies. Our methods are freely available as an R package and all results are presented in a web app.

Introduction

Determination of genes essential for the growth and survival of bacteria has been of major interest in genetic research as it provides a deeper understanding of lifestyle and adaptation [1–3]. While site-directed mutagenesis approaches determine essential genes accurately, such methods are laborious and time-consuming when performed globally. Consequently, whole genome analyses have only been conducted for well-known model organisms such as *Escherichia coli*, i.e., the Keio library [4]. In the last decade, wider availability of high-throughput sequencing methods initiated a shift from single-gene to whole-genome analysis, resulting in the development of *transposon insertion sequencing (TIS)* methods. Employing transposons that are randomly inserted into the genome enables researchers to generate large mutant libraries and to characterize them by the location of insertion sites (IS) via high-throughput sequencing. *Transposon directed insertion site sequencing (TraDIS)* is a widely applied *TIS* method [5–8] and has been established for the determination of essential genes in various scientific set-ups [9–15].

A key challenge in *TIS* studies resides in the statistical analysis, which typically aims to maximize the detection of true positives (essential genes) while minimizing the number of false positives (non-essential genes incorrectly identified as ‘essential’). Although there are multiple software suites and packages, not every method embedded therein will be equally suitable for the analysis of the obtained data set. For example, sliding window approaches [16,17] and Hidden Markov Models [18,19] have been proposed for the analysis of high-density libraries which regularly originate from *mariner*-transposon-based mutagenesis [20,21]. However, many *TIS* studies utilize *Tn5*-based transposons, which have different underlying assumptions and constraints in terms of the data generating process: Unlike *mariner* transposons, *Tn5* insertions do not depend on the presence of specific motifs and can theoretically occur in any non-essential region of the genome [6,22]. Nevertheless, *Tn5*-based libraries reported so far tend to be less dense than *mariner*-based libraries. Consequently, observing larger genomic regions lacking IS just by chance becomes more likely in *Tn5*-based libraries. Furthermore, the set of detected IS across the genome

rarely displays a uniform distribution of gene-wise insertion densities. Reasons may be transposon-driven preferences for GC- or AT-rich regions [23–26] and genomic hot- or coldspots, i.e., genomic regions of notably higher or lower insertion densities, respectively [24,27–29].

So far, a couple of *Tn5*-based statistical methods for identifying essential genes have been proposed. Burger et al. [30] suggest estimating the probability of observing several IS within a gene of a given length based on a binomial distribution using the genome-wide insertion density as success probability. The *Tn5Gaps* method of the Transit package [31] uses a Gumbel distribution to approximate the probabilities of observed IS-free gaps along the genome. Essentiality is determined by the largest gap within or partially overlapping a gene. Since both methods rely on *p*-values derived from thousands of genes, the authors recommend correcting for multiple testing. However, they did not evaluate how different correction approaches might affect the identification of essential genes. Alternatively, the Bio-TraDIS software package [32] avoids the multiple testing problem by heuristically leveraging an often observed bimodal distribution of gene-wise insertion densities. Combining an exponential distribution (for essential genes) with a gamma distribution (for non-essential genes), genes are labeled as ‘essential’, ‘non-essential’, or ‘ambiguous’ based on an *a priori* set \log_2 likelihood ratio threshold. In practice, a clear distinction between the two distributions is not always guaranteed [6], and the threshold values are usually set arbitrarily. The recently proposed Bayesian method *InsDens* calculates the posterior probability of a gene being essential [33] and the authors suggest to use Bayesian decision theory to set a posterior probability threshold. Although this method offers a clear interpretation, it requires choosing an *a priori* probability distribution parameter, too, which can influence the outcome.

Some but rather limited comparative studies of *TIS* methods are available. Based on high-density library data, Larivière et al. [20] draw a comparison between the bimodal approach of Bio-TraDIS, the *Tn5Gaps* method and a custom modification of the bimodal approach [11]. However, only two threshold values were applied and no performance analysis under different controlled data-generating processes was described. While Nlebedim et al. [33] used four different parameter combinations to generate synthetic data, the only method applied in the analysis was their method *InsDens*. The additional analysis of three real-world datasets using Bio-TraDIS suffers from the application of only one and rather low threshold value. Similarly, Ghomi et al. [34] proposed to use the non-parametric clustering algorithm embedded in the DBSCAN R package for the sole reason that it allows omitting the heuristic setting of threshold values required by Bio-TraDIS. Again, only a single and rather low threshold value for comparison using Bio-TraDIS was applied. However, the DBSCAN clustering algorithm itself requires setting two *a priori* parameter values but the authors did not report how different values might affect labeling performance, nor did they provide guidance on how to choose appropriate values.

The widespread use of *TIS*, particularly *Tn5* statistical analysis methods, contrasts with the lack of systematic reviews of these methods, especially when considering different data-generating processes. Furthermore, a transparent, comprehensive statistical method for setting threshold or parameter values in *TIS* methods is missing. As a consequence, publications often only justify the choice of methods and parameters by citing prior studies that used similar approaches and parameters. In addition, most studies lack sensitivity analyses for their threshold and parameter values. A common practice is truncating the 5'- and/or 3' ends of genes by several base pairs or up to 20%, to align with the assumption that the gene ends are generally non-essential [6,13,35–38], yet this approach is never investigated in sensitivity analyses.

At this scientific stage, we provide the following contributions to the statistical detection of essential genes. First, we introduce *Consecutive Non-Insertion Sites (ConNIS)*, a novel method that determines gene essentiality based on insertion-free sequences within genes. *ConNIS* provides an analytical solution for the probability of observing the longest insertion-free sequence within a gene, based on its length and the number of IS under the assumption of being non-essential. Second, we performed an extensive simulation study with 160 parameter combinations mimicking different data-generating processes. Using these synthetic datasets, four additional semi-synthetic datasets and three real datasets, *ConNIS* demonstrated its superiority over five state-of-the-art *Tn5* analysis methods, especially in settings with

low- and medium-dense libraries. In this context, we propose to use a weighting factor when applying genome-wide insertion density values to better represent genomic regions with low insertion densities. This modification also improved three competing methods by reducing the number of false positives without losing too many true positives in many settings. Third, we provide for the first time a data-driven instability criterion for selecting thresholds and parameter values in *TIS* methods, thereby making the results from different studies and methods comparable and more transparent. Applications of this approach to real and synthetic datasets clearly demonstrate its suitability for setting parameter values for all methods considered. An in-depth analysis of biological functions for selected genes illustrates the ability of *ConNIS* to reliably detect essential genes even among short genes that are often excluded from statistical analyses because competing methods cannot distinguish signal from noise in this length range. *ConNIS* and the instability criterion for all competing methods have been made available as an R package. To further explore our results, we provide an interactive web app and publicly available code.

Materials and methods

Consecutive Non-Insertion sites (ConNIS)

The transcription of an essential gene is, by definition, vital for an organism's survival and its hindrance due to transposon insertion will result in the mutant's removal from the population. Consequently, IS are expected to be exclusively detected in the non-essential genome, which comprises genes, intergenic regions and smaller fractions of essential genes [7]. Based on the assumption that a *Tn5* transposon can occur at any position in the *non-essential* genome, we propose *ConNIS*, a method that classifies a gene as 'essential' or 'non-essential' by analyzing its largest insertion-free sequence in terms of base pairs. Therefore, we derive a novel probability distribution (see [S1 File](#)) that we used to determine the probability of observing an insertion-free sequence in a gene of given length and number of IS.

Let b denote the length of a genome in terms of base pairs that contains p genes with corresponding lengths b_1, b_2, \dots, b_p . We define $\theta = h/b$ as the genome-wide insertion density, where h is the genome-wide number of observed IS. For a gene $j = 1, \dots, p$ let $\hat{h}_j = \lfloor b_j \cdot \theta \rfloor$ be the rounded expected number of IS within gene j under the assumption that gene j is not essential. Furthermore, let L_j be the length of an observed consecutive sequence of non-insertions of gene j under the assumption of uniformly distributed IS ([Fig 1A](#)).

Based on Theorem 1 (see Section 1 in [S1 File](#)), the probability mass function of L_j is given by

$$\mathbb{P}(L_j = i) = f(i, b_j, h_j) = \frac{\binom{b_j - i - 1}{b_j - \hat{h}_j - i}}{\binom{b_j - 1}{b_j - \hat{h}_j - 1}}.$$

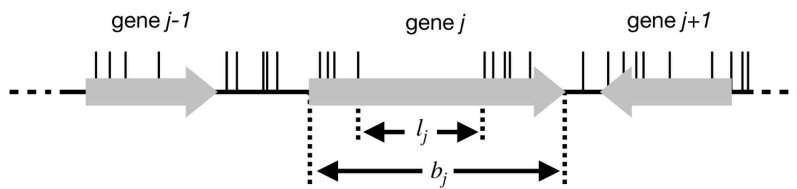
We next consider that IS are often distributed non-uniformly across the genome (see [Fig B](#) in [S2 File](#)). Other approaches use the genome-wide insertion density θ to estimate expected insertions per gene. However, under the assumption of non-essentiality this can inflate false positives in regions with a lower-than-average IS density, where missing insertions are likely due to chance. *ConNIS* corrects for this by introducing a weight factor w ($0 < w \leq 1$) to θ that adjusts for low-density regions ([Fig 1B](#)), a bias not handled by normalization methods focused only on insertion counts.

Observing the longest gap of size l_j *ConNIS* is then defined as the probability of observing an insertion-free consecutive sequence of at least length l_j in gene j :

$$C(l_j, b_j, \hat{h}_j, w) := \mathbb{P}(L_j \geq l_j) = 1 - \mathbb{P}(L_j < l_j) = 1 - \sum_{i=1}^{l_j-1} \frac{\binom{b_j - i - 1}{b_j - \lfloor \hat{h}_j w \rfloor - i}}{\binom{b_j - 1}{b_j - \lfloor \hat{h}_j w \rfloor - 1}}. \quad (1)$$

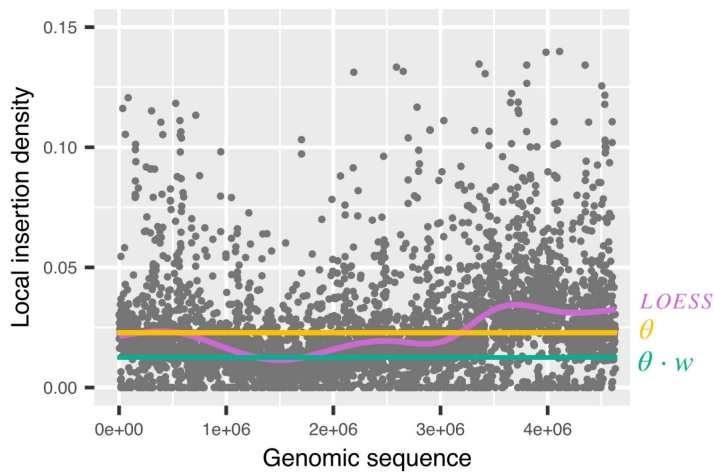
A

Determination of longest insertion-free sequence for gene j



B

Determination of the insertion density θ and weight value w



C

ConNIS with estimated number of insertion sites and weight

$$C(l_j, b_j, \hat{h}_j, w) := 1 - \sum_{i=1}^{l_j-1} \frac{\binom{b_j-i-1}{b_j-\hat{h}_j w-i}}{\binom{b_j-1}{b_j-\hat{h}_j w-1}}$$

$$\text{with } \hat{h}_j = b_j \cdot \theta \text{ and } 0 < w \leq 1$$

Fig 1. Overview of ConNIS. **A** For a gene j , determine its length b_j and its longest insertion-free sequence l_j . **B** Set a weight w for the genome-wide insertion density θ reflecting rather low density regions of the genome. **C** ConNIS: The probability of observing l_j within gene j due to random chance.

<https://doi.org/10.1371/journal.pcbi.1013428.g001>

Given a significance level of $0 < \alpha < 1$, we declare a gene j to be essential if $C(l_j, b_j, \hat{h}_j, w) \leq \alpha$. To control the global type I error when applying ConNIS, we suggest using either the Bonferroni(-Holm) method [39,40] to control the family-wise error rate (FWER) or the Benjamini-Hochberg method [41] to control the false discovery rate (FDR).

Competing state-of-the-art methods with proposed weighting strategy. We compared *ConNIS* with five popular state-of-the-art *Tn5* analysis methods for determining essential genes:

1. the *Binomial* distribution approach in the TSAS 2.0 package [30],
2. the approach of fitting a bimodal distribution based on gene-wise insertion densities included in the Bio-TraDIS package [32] (referred to as *Exp. vs. Gamma* method throughout this paper),
3. the *InsDens* method [33],
4. the *Tn5Gaps* method of the TRANSIT package [31] and
5. the *Geometric* distribution.

Although the geometric distribution has not been published as a stand-alone method, we included it as a competitor because it is the limiting distribution of *ConNIS* (see Theorem 2 in [S1 File](#)) and has been part of an analysis pipeline for determining the probability of insertion-free regions in the genome [42].

The *Binomial* and *Geometric* methods use the genome-wide insertion density θ as success probability, and the *Tn5Gaps* method uses it as a location parameter. However, for the reason outlined above, this naive use of θ can increase the number of false positives within genomic regions with a relatively low insertion density compared to the rest of the genome. To address this potential pitfall, we introduce a weight w to adjust θ when applying these methods, as we do it in *ConNIS*. We then use these modified methods, as well as the original versions, for comparison. See Sect 2 in [S1 File](#) for methodological details. Further, *InsDens* requires several prior hyperparameters. In line with the authors' claim, our tests on selected simulation settings showed minimal impact from these choices [33]. Thus, we used the default settings of the R package *insdens* for all simulations and real data analyses (see <https://github.com/Kevin-walters/insdens>, commit 286f114).

A labeling instability criterion for tuning parameter selection. TIS methods often require *a priori* set parameter or threshold values that will influence the final number of genes labeled 'essential' and therefore the methods' performances in terms of correct classification. The setting of an 'appropriate' parameter/threshold value in a given data scenario can be interpreted as a *tuning* problem. In this context our data-driven tuning approach selects a parameter of threshold value for a TIS method from a set of candidate values.

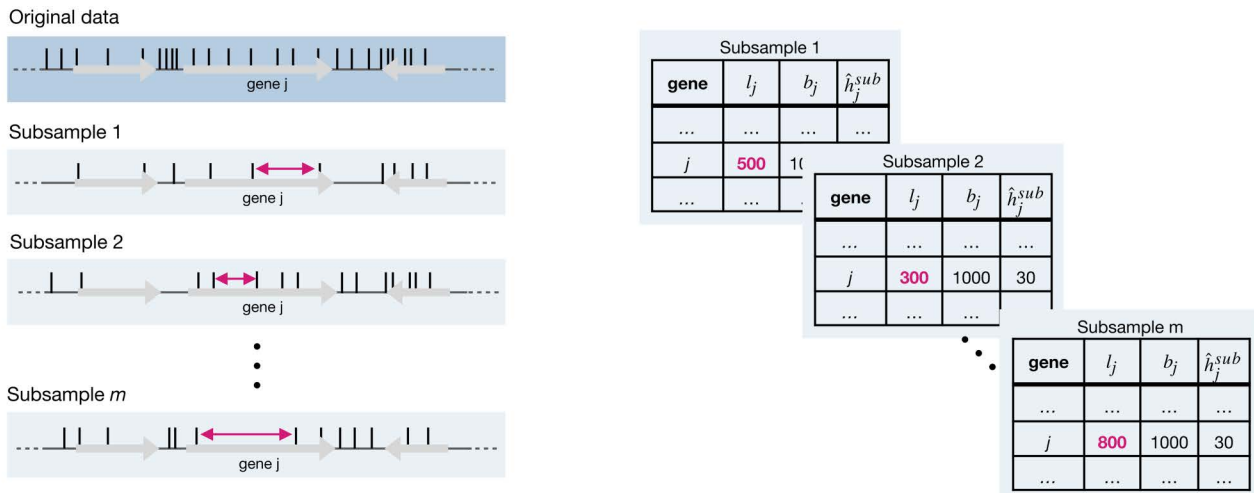
We consider observed IS as realizations of an unknown probability distribution across the genome. This is comparable with a repeated TIS experiment which yields different IS positions in each realization, particularly in non-essential genomic regions. As a result, gene-wise IS metrics, such as the longest insertion-free sequence l_j or gene-wise IS density θ_j , would vary between experiments, potentially altering the set of genes classified as essential. The main idea of our selection criterion is to leverage these variations by quantifying the average variation of gene labeling based on m subsample for a given tuning value. Inspired by stability selection approaches in linear regression and graph estimation problems [43–45], a 'good' tuning value should give rather stable results in terms of gene classification, i.e., the gene labeling should be less sensitive to the random occurrence of IS across the genome. In the following, we detail the procedure for selecting a suitable weight value w using *ConNIS* as an example. A transfer to other TIS methods or to threshold based filters in the data pre-processing steps is straight forward.

Let $\mathbf{w} = w_1, w_2, \dots, w_z$ be a sequence of ordered weights ($0 < w_q < w_r \leq 1$ for $q < r$ and $w_q, w_r \in \mathbf{w}$). Assume further that m subsamples, each of size $h^{sub} < h$, are drawn without replacement from the set of h observed IS and the expected number of insertion sites per gene to be $\hat{h}_j^{sub} = b_j \frac{h^{sub}}{b}$ (Fig 2A).

For all $w_y \in \mathbf{w}$, genes are labeled as 'essential' or 'non-essential' for a given significance level within each of the m subsamples using *ConNIS*. By modeling the gene labeling as a Bernoulli process ('essential' or 'non-essential'), we estimate the probability of labeling a gene j as 'essential' (see Fig 2B) using

A

Drawing of h^{sub} insertion site subsamples



B

Instability values based on subsamples

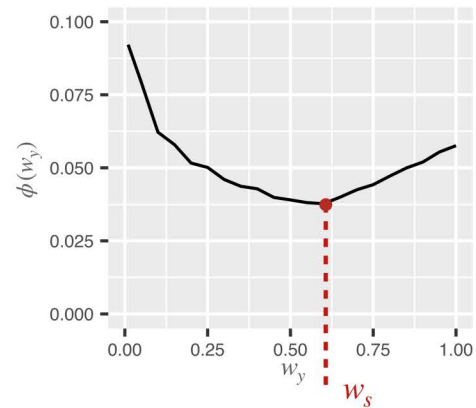
$$\hat{\pi}_j(w_y) = \frac{1}{m} \sum_{d=1}^m \mathbf{1} \left[C(l_j^{(d)}, b_j, \hat{h}_j^{sub}, w_y) \leq \alpha \right]$$

$$\phi(w_y) := \sum_{j=1}^p \frac{\hat{\pi}_j(w_y) \cdot (1 - \hat{\pi}_j(w_y))}{q(w_y)}$$

for $w_y \in \mathbf{w}$ with $\mathbf{w} = \{w_1, w_2, \dots, w_z\}$

C

Determination of $w_s \in \mathbf{w}$



D

Application of ConNIS with selected weight

$$C(l_j, b_j, \hat{h}_j, w_s) \quad \text{for } j = 1, 2, \dots, p$$

| gene | probability gap I | adjusted α | putative essential |
|-------|-------------------|-------------------|--------------------|
| .. | ... | ... | ... |
| $j-1$ | 0.0045 | 0.00005 | FALSE |
| j | 0.000007 | 0.00005 | TRUE |
| ... | ... | ... | ... |

Fig 2. Overview of the labeling instability criterion to select the weight parameter for ConNIS. **A** Drawing m subsamples of the h original observed IS. **B** Calculation of the instability values for all weights $w_y \in \mathbf{w}$ based on the estimated variances of a Bernoulli variable. **C** Selecting the weight w_s with the lowest instability $\phi(w_y)$. **D** Application of ConNIS using w_s followed by a multiple testing correction to identify putative essential genes.

<https://doi.org/10.1371/journal.pcbi.1013428.g002>

$$\hat{\pi}_j(w_y) = \frac{1}{m} \sum_{d=1}^m \mathbb{1} \left[C \left(l_j, b_j, \hat{h}_j^{sub}, w_y \right) \leq \alpha \right].$$

We can then define the *instability criterion* over all genes for a given weight w_y as

$$\phi(w_y) := \sum_{j=1}^p \frac{\hat{\pi}_j(w_y) \cdot (1 - \hat{\pi}_j(w_y))}{q(w_y)}, \quad (2)$$

where $\hat{\pi}_j(w_y) \cdot (1 - \hat{\pi}_j(w_y))$ is the Bernoulli variance and $q(w_y) = \sum_{j=1}^p \mathbb{1}(\hat{\pi}_j(w_y) > 0)$ is the total number of genes that have been labeled at least once as ‘essential’ in the m subsamples. This normalization factor ensures $0 \leq \phi_j(w_y) \leq 0.25$. A value of $\phi(w_y) = 0$ indicates complete consistency in gene labeling for each gene across subsamples, whereas $\phi(w_y, d) = 0.25$ reflects total instability, equivalent to randomly assigning labels by flipping a fair coin for each gene in each subsample.

After calculating the instability for all weights, we have a sequence of instability values $\phi = \{\phi(w_y)\}_{y=1}^z$ and select then the weight w_y that minimizes the instability of labeling (see Fig 2C):

$$w_s = \arg \min_{w_y \in \mathbf{w}} \phi(w_y).$$

Finally, *ConNIS* is applied to the original data with w_s (Fig 2D).

Depending on the range of \mathbf{w} and the number of observed IS, it is possible that very small weight values may lead to instability values approaching or even reaching zero with (nearly) all genes being labeled as ‘non-essential’. Following other stability approaches [43–45], these values are excluded from the set of candidate tuning values because they provide no useful information. For our instability criterion, we propose to omit all weights smaller than the smallest weight that maximizes the function $\phi(w)$ to ensure that only the most informative weights are considered, i.e.,

$$w_{max} = \min \left\{ w_y \in \mathbf{w} \mid \arg \max_{w_y \in \mathbf{w}} \phi(w_y) \right\}$$

$$w_s = \arg \min_{\tilde{w}_y \in \mathbf{w}: \tilde{w}_y \geq w_{max}} \phi(\tilde{w}_y).$$

Results

Comparison of ConNIS with state of the art methods

To evaluate the performance of *ConNIS* and its competitors we applied all methods to synthetic, semi-synthetic and real-world data. A detailed explanation of our simulation schemes for generating synthetic data that covers different assumptions about the data generating process can be found in S3 File. For the application to real-world data we used three publicly available *Tn5* libraries of different organisms and insertion densities. The semi-synthetic datasets were generated from a high-density *Tn5* library by randomly deleting IS. Note that the performance of all methods depends on the chosen values of their respective parameters and thresholds.

Since the number of essential genes is small compared to the number of non-essential genes, we used the Mathew’s Correlation Coefficient (MCC), a metric suitable for imbalanced data [46–48], as main performance measure. The MCC equals 1 when the method perfectly labels all genes as ‘essential’ or ‘non-essential’ (perfect agreement), 0 when the

labeling is completely random, and -1 when there is perfect disagreement between the true and predicted labels. For comparison the MCC is plotted given the number of genes labeled 'essential' which can be controlled by the methods' different parameter and threshold values. In addition, the precision-recall-curve (PRC) is shown to investigate two desirable, yet occasionally, conflicting objectives: selecting as many true positives as possible (recall) while avoiding an inflated number of false positives (precision).

In all applications the weight value of *Binomial*, *ConNIS*, *Geometric* and *Tn5Gaps* was set to $w = 0.1, 0.2, \dots, 1$ with $w = 1$ being the original, unweighted version. For *Exp. vs. Gamma*, we set the \log_2 -likelihood ratio threshold $t \in \{2, 3, \dots, 12\}$, covering the range of values commonly reported in the literature. The posterior probability threshold of *InsDens* was set at $r \in \{0.01, 0.1, 0.2, \dots, 0.9, 0.99\}$. Furthermore, we truncated genes by excluding the distal ends by either 0%, 5% or 10% and applied Bonferroni(-Holm) and Benjamini-Hochberg procedures for multiple testing correction.

Synthetic data settings. We present three illustrative synthetic data settings offering an overview of performance the methods. The results of all 160 different settings and additional classification metrics can be interactively explored at <https://connis.bips.eu>, supporting our findings.

Synthetic data example 1 (SDE1) had 200,000 IS randomly distributed along the genome in a sinusoidal shape. 'Essential' genes were defined to contain insertion-free sequences of at least 75% of the gene's length. This represented scenarios where essential genes can contain IS relatively far from their distal ends. Fig 3A shows *ConNIS* clearly outperforming the other methods with regard to the MCC. The PRC demonstrates that *ConNIS* effectively enhanced the identification of true essential genes without substantially inflating false positives by maintaining high precision. Notably, all methods showed their peak performance when the number of genes labeled 'essential' was about the number of true essential genes (dashed orange line). The plot also highlights the beneficial effect of applying a weight $w < 1$ to *Binomial* and *Geometric* (the points indicate the average performance if no weight is applied, i.e., $w = 1$).

Synthetic data example 2 (SDE2) mimicked scenarios where only a low insertion density can be achieved, e.g., due to bottleneck effects by environmental pressure. Therefore, 50,000 IS were randomly distributed along the genome in a sinusoidal shape, and the essential genes had insertion-free sequences of at least 85% of their length. All methods suffered from the sparse library (Fig 3C). *ConNIS* achieved the highest MCC value if the number of genes labeled 'essential' was close to the number of true essential genes. *Binomial* and *Tn5Gaps* could only achieve mediocre values at best. For *Exp. vs. Gamma*, *InsDens* and *Tn5Gaps*, the range of the number of genes labeled 'essential' never contained the number of true genes. However, the first two could achieve MCC values that were slightly worse than *ConNIS*. With respect to the PRC, all methods induced false positives due to larger non-insertion sequences occurring by chance compared to denser libraries. *ConNIS* tended to have a rather high precision while *Exp. vs. Gamma* and *InsDens* achieved rather high recall values, but at the price of an inflation of false positives.

Synthetic data example 3 (SDE3) covered scenarios with so-called 'cold-spots' along the genome, which have a much lower chance of containing IS. 'Essential' genes were defined to contain insertion-free sequences of at least 80%. In non-essential sections of the genome, each base pair had the same probability to contain one of the 200,000 IS, yet, in 25 randomly placed sections of size 10,000bp these probabilities were lowered by factor 10. *ConNIS* achieved clearly the best MCC values and PRC performance (Fig 3D). However, all methods tended to overestimate the number of essential genes (indicated by the rather low precision values) due to the higher chance of false positives in cold spots.

Real-world data. In the first example, we applied all methods to an *E. coli* BW25113 strain library comprising approximately 102,000 IS at time point T0 [49]. As ground truth, we used the results of the single-knockout study by Baba et al. [4], which is often considered as the gold standard. Fig 4A shows *ConNIS* outperforming the other methods by reaching MCC values up to 0.65. *InsDens* and *Exp. vs. Gamma* labeled too many genes as 'essential' (at least 575) even for their strictest thresholds ($t = 12$ and $r = 0.99$). However, the thresholds of *Exp. vs. Gamma* had only a marginal influence on the number of genes labeled 'essential', which resembles in parts the results of the simulation study. The

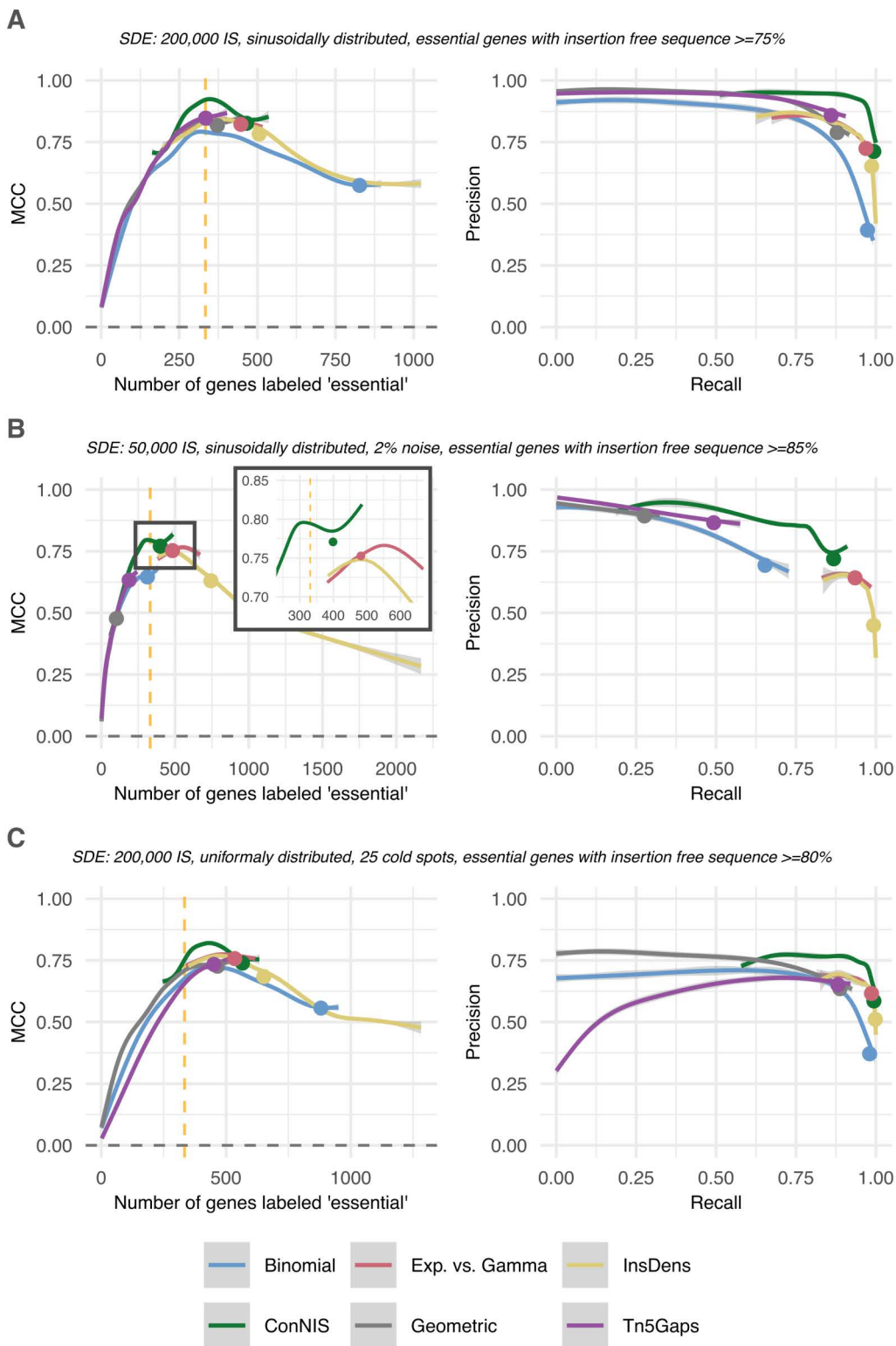


Fig 3. MCC and PRC performance based on synthetic data. The plots show the LOESS-smoothed curves with 95% confidence intervals for the MCC and PRC in three synthetic data settings. At the vertical dotted line the number genes labeled 'essential' matches the number of true essential genes. Dots on the curves indicate the average performance without weights or the least stringent threshold ($t = 2$ for *Exp. vs. Gamma* and $r = 0.1$ for *InsDens*). In all settings, 5 % of both ends of each gene were trimmed.

<https://doi.org/10.1371/journal.pcbi.1013428.g003>

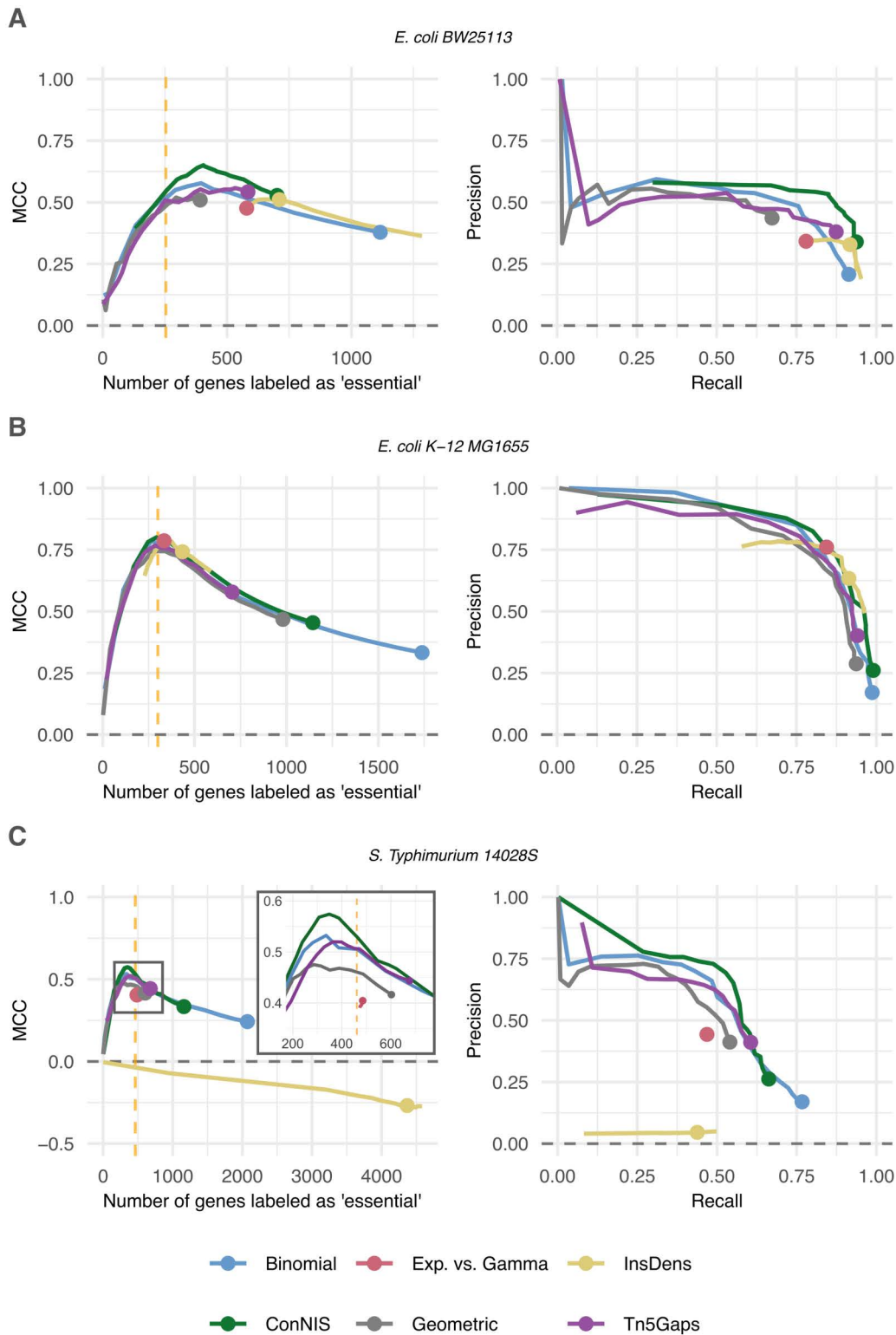


Fig 4. MCC and PRC performances based on real-world data. The vertical dotted line shows the true number of genes. Dots indicate the performance of the original methods ($w = 1$). **A** *E. coli* BW25113 strain with $\approx 102,000$ IS [49]. **B** *E. coli* K-12 MG1655 strain with $\approx 390,000$ IS [50]. **C** *Salmonella enterica* serovar Typhimurium 14028S strain with $\approx 186,000$ IS [37]. Note, in applications A and B, most of the results of *Exp. vs. Gamma* are covered by those of *InsDens*.

<https://doi.org/10.1371/journal.pcbi.1013428.g004>

PRCs reveal that all methods achieved mediocre precision values at best with *ConNIS* having relatively stable precision values for rising recall values.

A high-density *E. coli* K-12 MG1655 library characterized by approximately 390,000 IS [50] was used as second example. As truth, we used the gene essentiality classification from the Profiling of *E. coli* Chromosome (PEC) database [51]. Here, results were consistent with the simulation study, i.e., all methods benefited from the high number of IS and performed best when the number of selected genes approached the number of true essential genes (Fig 4B). Weighting ($w < 1$) was highly beneficial for all methods, as the original versions ($w = 1$, indicated by the points) had high recall values at the cost of low precision values.

In the third example, all methods were applied to a *Salmonella enterica* serovar Typhimurium 14028S library comprising approximately 186,000 IS [37]. Following Nlebedim et al. [33], we used as truth the combined set of essential genes provided by Baba et al. [4] and Porwollik et al. [52]. In this scenario all methods showed at best mediocre performance. The removal of IS with low read counts improved the performances of the methods slightly and might be a sign of the presence of spurious IS [31,53] (we tried minimum read count thresholds with value 1, 2, 3, 5 and 10). *ConNIS* achieved the best MCC and precision values. *Exp. vs. Gamma* and especially *InsDens* performed very poorly (Fig 4C), with the latter labeling far too many genes as 'essential', resulting in negative MCC values. *Binomial*, *ConNIS*, and *Tn5Gaps* benefited from a fairly low weight w , which seemed to reduce the number of false positives compared to the original versions ($w = 1$).

Semi-synthetic data settings. To investigate the influence of the number of observed IS on the methods' performances, we generated semi-synthetic data by drawing IS subsamples of sizes 50,000, 100,000, 200,000 and 400,000 from a very high density *Tn5* library [11]. The Kaio library [4] was used as a reference for true gene essentiality. In low and medium density libraries (subsamples of 50,000 to 200,000 IS) *ConNIS* outperformed the other methods, clearly (see Fig 5A; for medium-sized libraries see Fig A in S2 File). Similar to the synthetic and real-world data settings, *ConNIS* showed its best performance in terms of MCC when the number of selected genes was about the number of true essential genes. In case of the rather high-density library (subsample size of 400,000 IS), all methods were on par, (Fig 5B).

Application of gene labeling instability criterion for tuning parameter selection

The performance of the gene labeling instability criterion for tuning parameter selection was investigated using the three previously described real-world datasets and three randomly chosen dataset examples from the simulation study. We applied the instability criterion to select the tuning parameter of each method. For each setting, $m = 500$ subsamples were drawn without replacement, with each subsample containing 50% of randomly picked IS from the original data. Genes were truncated by 5 % at each distal end. We used the MCC to evaluate the performance of the labeling instability criterion and compared it to MCC values for the *optimal* parameters (those that produce the highest MCC) as well as for parameters used in earlier studies, such as unweighted versions or heuristic choices.

The results in Table 1 demonstrate that the application of the gene labeling instability criterion for determining a weight for *ConNIS* is highly beneficial. For the real-world datasets *E. coli* BW25113 and *Salmonella enterica* serovar Typhimurium 14028S, as well as the synthetic dataset 2, the application of instability criterion successfully identified the 'optimal' weights based on the corresponding MCC values (Table 1). In these examples, *ConNIS* also achieved the highest MCC with its selected w compared to all other methods. Furthermore, for synthetic datasets 1 and 3, the MCCs obtained by the instability criterion were close to its highest possible MCCs in these settings (0.90 vs. 0.94 and 0.76 vs. 0.77). Only for the MG1655 strain data, our tuning approach was less successful for *ConNIS* (0.67 vs. 0.79).

The instability criterion also successfully tuned the other methods, resulting in many cases where the best possible MCC was achieved. For *Exp. vs. Gamma* in all six settings \log_2 thresholds were selected that resulted in MCC values close or even to values obtained by applying an optimal threshold value. In comparison to \log_2 thresholds used in recent studies

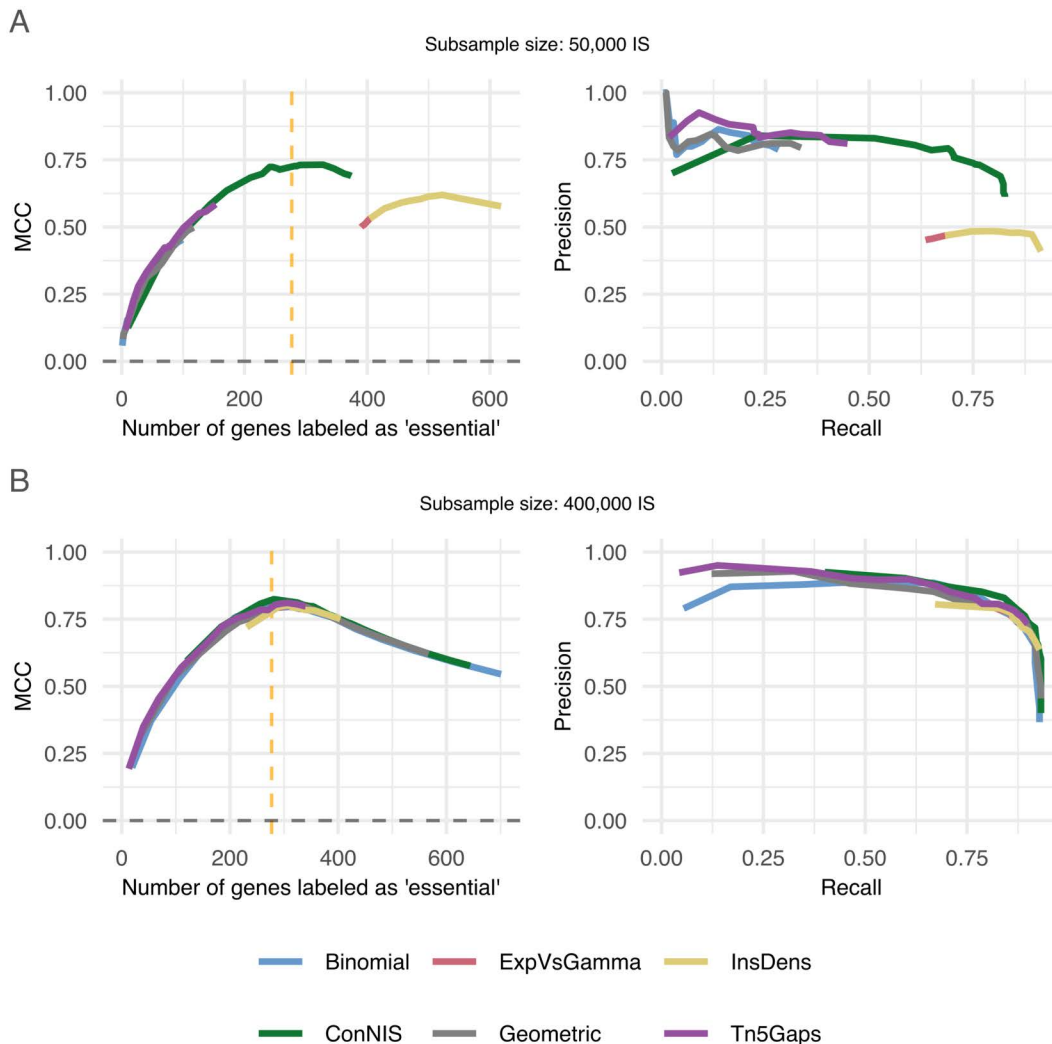


Fig 5. MCC and PRC performances for semi-synthetic data. Subsamples were generated by randomly drawing IS from a very high-density library to generate a low- (**A**) and high-density (**B**) library of *E. coli* BW25113 [11]. The Kaio library [4] was used as reference for ‘true’ gene essentiality. The vertical dotted line indicates where the number of genes labeled “essential” corresponded with number of true essential genes.

<https://doi.org/10.1371/journal.pcbi.1013428.g005>

[11,20] we found our instability approach to give similar or better MCC values, yet, the range of possible MCC was rather small. Applying labeling instability to the posterior probability threshold r in *InsDens* yielded favorable results in five settings and achieved the highest possible MCC in three of them. It was also able to select a good choice between a very strict ($r = 0.99$) and a very relaxed value ($r = 0.05$) which have been used before [33]. Only in the *Salmonella enterica* serovar Typhimurium 14028S real-world dataset, the selected r had a bad MCC value. However, since *InsDens* generally showed a weak performance in this setting before (Fig 4C), the result is not surprising. For *Binomial*, *Geometric* and *Tn5Gaps* the application of the instability approach was also beneficial compared to the unweighted (i.e., original) version of the methods.

Biological relevance

To highlight biological relevance beyond global performance metrics, we analyzed genes where *ConNIS* systematically disagrees with the other methods. Here, we focus on ‘major discrepancies’, i.e., genes called ‘essential’ by *ConNIS* but

Table 1. Tuning performance of the gene labeling instability criterion.

| | | BW25113 | MG1655 | 14028S | Syn. Data 1 | Syn. Data 3 | Syn. Data 3 |
|----------------|-----------------------|-------------|--------|-------------|-------------|-------------|-------------|
| Binomial | instability | 0.55 | 0.71 | 0.51 | 0.83 | 0.73 | 0.61 |
| | optimal | 0.58 | 0.79 | 0.53 | 0.89 | 0.73 | 0.69 |
| | unweighted | 0.38 | 0.33 | 0.24 | 0.51 | 0.64 | 0.52 |
| ConNIS | instability | 0.64 | 0.67 | 0.57 | 0.9 | 0.88 | 0.76 |
| | optimal | 0.64 | 0.79 | 0.57 | 0.94 | 0.88 | 0.77 |
| | unweighted | 0.17 | 0.08 | 0.33 | 0.81 | 0.34 | 0.21 |
| Exp. vs. Gamma | instability | 0.47 | 0.78 | 0.39 | 0.91 | 0.81 | 0.72 |
| | optimal | 0.49 | 0.79 | 0.41 | 0.92 | 0.81 | 0.74 |
| | log ₂ (4) | 0.47 | 0.77 | 0.39 | 0.88 | 0.81 | 0.73 |
| | log ₂ (12) | 0.47 | 0.71 | 0.39 | 0.91 | 0.74 | 0.73 |
| Geometric | instability | 0.52 | 0.58 | 0.46 | 0.86 | 0.59 | 0.63 |
| | optimal | 0.52 | 0.74 | 0.48 | 0.91 | 0.75 | 0.70 |
| | unweighted | 0.22 | 0.08 | 0.42 | 0.83 | 0.42 | 0.22 |
| | instability | 0.51 | 0.74 | -0.26 | 0.91 | 0.81 | 0.74 |
| InsDens | optimal | 0.51 | 0.78 | 0.00 | 0.92 | 0.81 | 0.74 |
| | <i>r</i> = 0.05 | 0.48 | 0.72 | -0.27 | 0.81 | 0.74 | 0.63 |
| | <i>r</i> = 0.3 | 0.51 | 0.78 | -0.25 | 0.86 | 0.79 | 0.70 |
| | <i>r</i> = 0.6 | 0.49 | 0.76 | -0.22 | 0.90 | 0.81 | 0.73 |
| | <i>r</i> = 0.99 | 0.47 | 0.64 | - | 0.90 | 0.74 | 0.72 |
| Tn5Gaps | instability | 0.55 | 0.63 | 0.52 | 0.91 | 0.77 | 0.70 |
| | optimal | 0.56 | 0.77 | 0.52 | 0.91 | 0.77 | 0.70 |
| | unweighted | 0.54 | 0.58 | 0.44 | 0.91 | 0.72 | 0.60 |

Tuning was applied to three real-world and three synthetic datasets. For each method, we report the MCC obtained when using the gene labeling *instability* criterion. For comparison, we also show the MCC achieved using (i) the *optimal* tuning value (i.e., the “oracle” value yielding the highest possible MCC) and (ii) heuristic values used in previous studies. The entries in bold indicate cases where the MCC based on the instability criterion was identical to the MCC value of the optimal tuning value. **Syn. Data 1:** 400,000 sinusoidal distributed IS, essential genes contained an insertion free sequence of $\geq 80\%$. 8000 IS were added as noise. **Syn. Data 2:** 100,000 sinusoidal distributed IS and essential genes contained an insertion-free sequence of $\geq 75\%$. **Syn. Data 3:** 200,000 uniformly distributed IS with 25 cold spots, 4000 noise IS and essential genes contained an insertion-free sequence of $\geq 75\%$.

<https://doi.org/10.1371/journal.pcbi.1013428.t001>

“non-essential” by four to five comparator methods, or vice versa. Across the three libraries (*E. coli* BW25113, *E. coli* MG1655 and *S. Typhimurium* 14028s), this affects 59 genes in total (15, 26 and 18 genes, respectively). Of these, 44 genes are called ‘essential’ by *ConNIS* but ‘non-essential’ by the comparator methods, whereas 15 genes show the opposite case (see [S4-S6 Files](#)). Overall, the analysis shows that *ConNIS* agrees well with experimental gold-standard essentiality sets while providing specific gains in low-insertion regimes and for short genes that have often been excluded *a priori* from the analysis due to lack of detection power of established methods. For all methods the threshold/parameter values were set by our instability criterion.

In the first group, *ConNIS*-specific essential calls have a median length of 328 bp (interquartile range 131 to 477bp; minimum gene length 74bp). For example, *ftsL* (365bp), *ffs* (113 bp), *argU* (76bp), and *folK* (479bp) were correctly identified as being essential by *ConNIS*. *FtsL* encodes a cytoplasmic membrane protein which essentiality manifests in rapid cell division blockade upon mutation [54]. *Ffs* together with *Ffh* builds up the well-known signal recognition particle (SRP) in *E. coli*, a multifunctional ribonucleoprotein complex fundamental for membrane protein targeting. As described by

Peterson et al. [55], both the *Ffh* protein and the *ffs* encoded 4.5S RNA are essential for cell viability and correct localization of proteins to the cytoplasmic membrane. Another RNA-gene correctly identified by *ConNIS* is *argU*. Lack of function mutations have been reported to cause DNA replication defects [56] manifesting in inhibition of cell growth [57]. Essentiality of *folK* has been critically analyzed by Goodall et al. [11] who, in contrast to the Keio library [4] and PEC database [51], classified the gene as ‘conditionally essential’ and not as ‘essential’. However, *ConNIS* also named *folK* essential. Interestingly, Goodall et al. [11] as well as the Keio library [4], the PEC database [51] and *ConNIS* are correct and the explanation highlights the importance of the chosen growth conditions on which basis essentiality is defined. While Goodall et al. [11] determine essentiality by using a library obtained directly from LB-agar plates and conditional essentiality by a library obtained after successive growth to 5–6 generations in liquid LB medium, Wetmore et al. [49], which data were used in this study, grew their mutant library on LB-plates followed by growth in liquid LB medium (as Goodall et al. [11] to an OD of 1.0). Consequently, *ConNIS* identifies *folK* correctly as being essential for the Wetmore et al. [49] mutant library. The disagreement pattern between *ConNIS* and its competitors is consistent with the known limitation of density/count-based methods on short or low-insertion genes, which are often excluded or down-weighted. Thus, the results highlights that *ConNIS* retains statistical power even for short genes. Yet, in the case of *nusB*, one of the four Nus factor encoding genes of *E. coli*, *ConNIS* wrongly assigns ‘essentiality’, while Bubunenko et al. [58] have shown that *NusB*, despite of being important for cell growth, is not essential. A closer inspection revealed that the incorrect assignment resides in the fact that the analyzed library carries only one insertion at the far 3’-end of *nusB*. Consequently, the relatively long insertion-free gap yields a relatively low *p*-value.

The second group comprises genes that *ConNIS* called ‘non-essential’ but that ≥ 4 the other methods classified as ‘essential’. These genes are typically much longer (median length 1,440bp, interquartile range 1132 to 1659bp). Two examples of genes correctly identified as being non-essential by *ConNIS* are *ptsI* and *ybcK*. As shown by Wu et al. [59] and Wu et al. [60] viable loss-of-function mutants of *ptsI* and *ybcK*, respectively, can be recovered. On the other hand, *ConNIS* classified *pssA* as non-essential, whereas the encoded phosphatidylserine synthase (*PssA*) is known to be essential for vitality in various pathogenic bacteria including *E. coli* [61]. The wrong assignment can be explained by the combination of three factors: first, only one insertion was observed close to the middle of the gene making the observed gap nearly as small as possible for a single insertion site. Second, a rather low number of expected insertion sites was used due to the low weighting factor of $w = 0.15$, increasing the probability to observe bigger insertion free gaps under the null model. Third, the applied Bonferroni-Holm correction method is relatively conservative and can label even small *p*-values non-significant when thousands of genes are examined.

Discussion and conclusion

In this work, we addressed three main challenges inherent in statistical analysis in *TraDIS* studies. The first challenge arises from the fact that in *Tn5* datasets every base pair of the genome serves a potential insertion site, while reported insertion densities often remain far below saturation levels. Considering this, *ConNIS* gives an analytic solution for the probability of observing an insertion-free sequence within a gene of a given length and number of insertion sites. The second challenge is the often observed non-uniform distribution of IS across the genome. Neglecting this factor can lead to an increased number of (nearly) insertion-free genes being incorrectly labeled as essential in regions with relatively low insertion densities. Addressing non-uniformity, *ConNIS* contains a weighting parameter that increases the precision by making it more difficult to label genes as ‘essential’ in low-density regions. We extended this idea to three state-of-the-art methods to improve their precision. The third challenge lies in the fact that many *TIS* methods rely on *a priori* set threshold or parameter values, which can substantially influence labeling performance. However, an ‘objective’ criterion for setting these values has been lacking, often resulting in arbitrarily chosen values. By introducing the concept of gene labeling instability based on subsamples of observed IS, we proposed a data-driven approach to select appropriate parameters and threshold values.

An extensive simulation study and application to three real-world datasets and four semi-synthetic datasets was conducted to compare the performance of *ConNIS* to multiple state-of-the-art *Tn5* analysis methods. In most settings, *ConNIS* outperformed these methods or was at least on par with the best of them. Unlike its competitors, *ConNIS* showed usually robust performances for (arbitrarily) chosen truncation and filter values. The results also confirm our idea of weighting the genome-wide insertion density when applied to existing methods: it could reduce the number of false positives without sacrificing too many true positives. Applying our proposed gene labeling instability criterion for tuning parameter selection in various real-world and multiple synthetic data scenarios demonstrated its potential to select favorable weight and threshold values for all methods. By inspecting major discrepancies between the classification results of *ConNIS* and its competitors, we showed that *ConNIS* was able to correctly classify even very short genes, thereby avoiding the standard practice of dismissing such genes from the analysis *a priori*.

Given that *ConNIS* demonstrated superior performance, especially in low and medium insertion density settings, its application is expected to improve the precision of results in experimental settings characterized by high selective pressure or observation of bottleneck effects. While we have investigated *ConNIS*' ability to identify essential genes, we anticipate that its application might be similarly beneficial for the determination of *conditionally* essential genes, for example by comparing gene-wise *ConNIS* scores between conditions and by defining quasi-essential genes based on differences in these scores between time points or conditions. This would broaden the scope of *ConNIS* to settings where a non-binary characterization such as relative gene fitness is desirable. As a first step, we provide a proof-of-concept (see [S7 File](#)), where we illustrate how *ConNIS* in combination with the instability approach yields a continuous gene-wise essentiality evidence score and how this score can be used to classify quasi-essential genes and fitness-like effects. This framework could be extended in future work by explicitly modeling the loss of insertion sites over time or across conditions [27]. Further, the weighting approach could be improved by incorporating multiple weighting values to target different genomic regions more effectively.

Our gene labeling instability criterion was originally developed for selecting threshold and parameter values of *Tn5* analysis methods, but it may also be applicable to other *TIS* methods that employ alternative transposons, such as the popular *mariner* transposon. It might also serve as a criterion in pre-processing steps like quality filters or trimming of distal gene ends. Last but not least, our work showed the crucial role of the underlying data-generating process on the performance of all methods. Future work could expand the range of scenarios considered, helping researchers choose the most appropriate method for analyzing their data. In this context, a systematic re-analysis of publicly available *TraDIS* datasets could raise the confidence in essential gene prediction and allow for further hypothesis generation. As a first step, we provide a curated multi-study resource comprising eight publicly available *Tn5*-based *TraDIS/Tn-Seq* datasets, including transparent per-study processing scripts and the resulting *ConNIS* essential-gene predictions, via Zenodo (DOI: <https://doi.org/10.5281/zenodo.18538449>).

Supporting information

S1 File. Proofs and methodological extensions. PDF file with proofs for *ConNIS* and formal definition of the extension of existing methods.

(PDF)

S2 File. Additional plots. PDF file with plots of additional analysis results.

(PDF)

S3 File. Detailed description of the generation of synthetic data.

(PDF)

S4 File. Gene classifications for *E. coli* BW25113.

(CSV)

S5 File. Gene classifications for *E. coli* MG1655.

(CSV)

S6 File. Gene classifications for *S. Typhimurium* 14028S.

(CSV)

S7 File Proof-of-concept for gene fitness and quasi-essentiality. An R Cookbook.

(PDF)

Acknowledgments

The authors would like to thank Ian Henderson, Emily Goodall and Ash Robinson for providing the list of insertion sites of their high-density library of the *E. coli* BW25113 strain.

Author contributions

Conceptualization: Moritz Hanke.

Data curation: Moritz Hanke.

Formal analysis: Moritz Hanke, Ronja Foraita.

Investigation: Moritz Hanke, Theresa Harten.

Methodology: Moritz Hanke.

Software: Moritz Hanke.

Validation: Moritz Hanke, Theresa Harten, Ronja Foraita.

Visualization: Moritz Hanke, Theresa Harten.

Writing – original draft: Moritz Hanke, Theresa Harten, Ronja Foraita.

Writing – review & editing: Moritz Hanke, Theresa Harten, Ronja Foraita.

References

- Gerdes S, Edwards R, Kubal M, Fonstein M, Stevens R, Osterman A. Essential genes on metabolic maps. *Curr Opin Biotechnol.* 2006;17(5):448–56. <https://doi.org/10.1016/j.copbio.2006.08.006> PMID: [16978855](https://pubmed.ncbi.nlm.nih.gov/16978855/)
- Zhang Z, Ren Q. Why are essential genes essential? - The essentiality of *Saccharomyces* genes. *Microb Cell.* 2015;2(8):280–7. <https://doi.org/10.15698/mic2015.08.218> PMID: [28357303](https://pubmed.ncbi.nlm.nih.gov/28357303/)
- Shang W, Wang F, Fan G, Wang H. Key elements for designing and performing a CRISPR/Cas9-based genetic screen. *J Genet Genomics.* 2017;44(9):439–49. <https://doi.org/10.1016/j.jgg.2017.09.005> PMID: [28967615](https://pubmed.ncbi.nlm.nih.gov/28967615/)
- Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol Syst Biol.* 2006;2:2006.0008. <https://doi.org/10.1038/msb4100050> PMID: [16738554](https://pubmed.ncbi.nlm.nih.gov/16738554/)
- Langridge GC, Phan M-D, Turner DJ, Perkins TT, Parts L, Haase J, et al. Simultaneous assay of every *Salmonella* Typhi gene using one million transposon mutants. *Genome Res.* 2009;19(12):2308–16. <https://doi.org/10.1101/gr.097097.109> PMID: [19826075](https://pubmed.ncbi.nlm.nih.gov/19826075/)
- Chao MC, Abel S, Davis BM, Waldor MK. The design and analysis of transposon insertion sequencing experiments. *Nat Rev Microbiol.* 2016;14(2):119–28. <https://doi.org/10.1038/nrmicro.2015.7> PMID: [26775926](https://pubmed.ncbi.nlm.nih.gov/26775926/)
- Cain AK, Barquist L, Goodman AL, Paulsen IT, Parkhill J, van Opijnen T. A decade of advances in transposon-insertion sequencing. *Nat Rev Genet.* 2020;21(9):526–40. <https://doi.org/10.1038/s41576-020-0244-x> PMID: [32533119](https://pubmed.ncbi.nlm.nih.gov/32533119/)
- Liang Y-T, Luo H, Lin Y, Gao F. Recent advances in the characterization of essential genes and development of a database of essential genes. *Imeta.* 2024;3(1):e157. <https://doi.org/10.1002/imt2.157> PMID: [38868518](https://pubmed.ncbi.nlm.nih.gov/38868518/)
- Christen B, Abeliuk E, Collier JM, Kalogeraki VS, Passarelli B, Collier JA, et al. The essential genome of a bacterium. *Mol Syst Biol.* 2011;7:528. <https://doi.org/10.1038/msb.2011.58> PMID: [21878915](https://pubmed.ncbi.nlm.nih.gov/21878915/)
- Rubin BE, Wetmore KM, Price MN, Diamond S, Shultzaberger RK, Lowe LC, et al. The essential gene set of a photosynthetic organism. *Proc Natl Acad Sci U S A.* 2015;112(48):E6634–43. <https://doi.org/10.1073/pnas.1519220112> PMID: [26508635](https://pubmed.ncbi.nlm.nih.gov/26508635/)

11. Goodall ECA, Robinson A, Johnston IG, Jabbari S, Turner KA, Cunningham AF. The Essential Genome of *Escherichia coli* K-12. *mBio*. 2018;9(1).
12. Sternon J-F, Godessart P, Gonçalves de Freitas R, Van der Henst M, Poncin K, Francis N, et al. Transposon Sequencing of *Brucella abortus* Uncovers Essential Genes for Growth In Vitro and Inside Macrophages. *Infect Immun*. 2018;86(8):e00312-18. <https://doi.org/10.1128/IAI.00312-18> PMID: [29844240](https://pubmed.ncbi.nlm.nih.gov/29844240/)
13. Poulsen BE, Yang R, Clatworthy AE, White T, Osmulski SJ, Li L, et al. Defining the core essential genome of *Pseudomonas aeruginosa*. *Proc Natl Acad Sci U S A*. 2019;116(20):10072–80. <https://doi.org/10.1073/pnas.1900570116> PMID: [31036669](https://pubmed.ncbi.nlm.nih.gov/31036669/)
14. Luo H, Lin Y, Liu T, Lai F-L, Zhang C-T, Gao F, et al. DEG 15, an update of the Database of Essential Genes that includes built-in analysis tools. *Nucleic Acids Res*. 2021;49(D1):D677–86. <https://doi.org/10.1093/nar/gkaa917> PMID: [33095861](https://pubmed.ncbi.nlm.nih.gov/33095861/)
15. Rivas-Marin E, Moyano-Palazuelo D, Henriques V, Merino E, Devos DP. Essential gene complement of *Planctopirus limnophila* from the bacterial phylum Planctomycetes. *Nat Commun*. 2023;14(1):7224. <https://doi.org/10.1038/s41467-023-43096-3> PMID: [37940686](https://pubmed.ncbi.nlm.nih.gov/37940686/)
16. Pritchard JR, Chao MC, Abel S, Davis BM, Baranowski C, Zhang YJ, et al. ARTIST: high-resolution genome-wide assessment of fitness using transposon-insertion sequencing. *PLoS Genet*. 2014;10(11):e1004782. <https://doi.org/10.1371/journal.pgen.1004782> PMID: [25375795](https://pubmed.ncbi.nlm.nih.gov/25375795/)
17. Zhang YJ, Ioerger TR, Huttenhower C, Long JE, Sasseti CM, Sacchetti JC, et al. Global assessment of genomic regions required for growth in *Mycobacterium tuberculosis*. *PLoS Pathog*. 2012;8(9):e1002946. <https://doi.org/10.1371/journal.ppat.1002946> PMID: [23028335](https://pubmed.ncbi.nlm.nih.gov/23028335/)
18. Chao MC, Pritchard JR, Zhang YJ, Rubin EJ, Livny J, Davis BM, et al. High-resolution definition of the *Vibrio cholerae* essential gene set with hidden Markov model-based analyses of transposon-insertion sequencing data. *Nucleic Acids Res*. 2013;41(19):9033–48. <https://doi.org/10.1093/nar/gkt654> PMID: [23901011](https://pubmed.ncbi.nlm.nih.gov/23901011/)
19. DeJesus MA, Ioerger TR. A Hidden Markov Model for identifying essential and growth-defect regions in bacterial genomes from transposon insertion sequencing data. *BMC Bioinformatics*. 2013;14:303. <https://doi.org/10.1186/1471-2105-14-303> PMID: [24103077](https://pubmed.ncbi.nlm.nih.gov/24103077/)
20. Larivière D, Wickham L, Keiler K, Nekrutenko A, Galaxy Team. Reproducible and accessible analysis of transposon insertion sequencing in Galaxy for qualitative essentiality analyses. *BMC Microbiol*. 2021;21(1):168. <https://doi.org/10.1186/s12866-021-02184-4> PMID: [34090324](https://pubmed.ncbi.nlm.nih.gov/34090324/)
21. Ioerger TR. Analysis of gene essentiality from TnSeq data using Transit. *Essential genes and genomes*. Springer US; 2021. p. 391–421.
22. Kwon YM, Ricke SC, Mandal RK. Transposon sequencing: methods and expanding applications. *Appl Microbiol Biotechnol*. 2016;100(1):31–43. <https://doi.org/10.1007/s00253-015-7037-8> PMID: [26476650](https://pubmed.ncbi.nlm.nih.gov/26476650/)
23. Zhang H, Lu T, Liu S, Yang J, Sun G, Cheng T, et al. Comprehensive understanding of Tn5 insertion preference improves transcription regulatory element identification. *NAR Genom Bioinform*. 2021;3(4):lqab094. <https://doi.org/10.1093/nargab/lqab094> PMID: [34729473](https://pubmed.ncbi.nlm.nih.gov/34729473/)
24. van Opijnen T, Levin HL. Transposon Insertion Sequencing, a Global Measure of Gene Function. *Annual Review of Genetics*. 2020;54(1):337–65.
25. Lluch-Senar M, Delgado J, Chen W-H, Lloréns-Rico V, O'Reilly FJ, Wodke JA, et al. Defining a minimal cell: essentiality of small ORFs and ncRNAs in a genome-reduced bacterium. *Mol Syst Biol*. 2015;11(1):780. <https://doi.org/10.15252/msb.20145558> PMID: [25609650](https://pubmed.ncbi.nlm.nih.gov/25609650/)
26. Green B, Bouchier C, Fairhead C, Craig NL, Cormack BP. Insertion site preference of Mu, Tn5, and Tn7 transposons. *Mob DNA*. 2012;3(1):3. <https://doi.org/10.1186/1759-8753-3-3> PMID: [22313799](https://pubmed.ncbi.nlm.nih.gov/22313799/)
27. Mahmutovic A, Abel Zur Wiesch P, Abel S. Selection or drift: the population biology underlying transposon insertion sequencing experiments. *Comput Struct Biotechnol J*. 2020;18:791–804. <https://doi.org/10.1016/j.csbj.2020.03.021> PMID: [32280434](https://pubmed.ncbi.nlm.nih.gov/32280434/)
28. Kimura S, Hubbard TP, Davis BM, Waldor MK. The nucleoid binding protein H-NS biases genome-wide transposon insertion landscapes. *mBio*. 2016;7(4).
29. Manna D, Porwollik S, McClelland M, Tan R, Higgins NP. Microarray analysis of Mu transposition in *Salmonella enterica*, serovar Typhimurium: transposon exclusion by high-density DNA binding proteins. *Mol Microbiol*. 2007;66(2):315–28. <https://doi.org/10.1111/j.1365-2958.2007.05915.x> PMID: [17850262](https://pubmed.ncbi.nlm.nih.gov/17850262/)
30. Burger BT, Imam S, Scarborough MJ, Noguera DR, Donohue TJ. Combining genome-scale experimental and computational methods to identify essential genes in *Rhodobacter sphaeroides*. *mSystems*. 2017;2(3).
31. DeJesus MA, Ambadipudi C, Baker R, Sasseti C, Ioerger TR. TRANSIT - A Software Tool for Himar1 TnSeq Analysis. *PLOS Computational Biology*. 2015;11(10):e1004401. <https://doi.org/10.1371/journal.pcbi.1004401>
32. Barquist L, Mayho M, Cummins C, Cain AK, Boinett CJ, Page AJ, et al. The TraDIS toolkit: sequencing and analysis for dense transposon mutant libraries. *Bioinformatics*. 2016;32(7):1109–11. <https://doi.org/10.1093/bioinformatics/btw022> PMID: [26794317](https://pubmed.ncbi.nlm.nih.gov/26794317/)
33. Nlebedim VU, Chaudhuri RR, Walters K. Probabilistic identification of bacterial essential genes via insertion density using TraDIS data with Tn5 libraries. *Bioinformatics*. 2021;37(23):4343–9. <https://doi.org/10.1093/bioinformatics/btab508> PMID: [34255819](https://pubmed.ncbi.nlm.nih.gov/34255819/)
34. Ghomi A, Jung JJ, Langridge GC, Cain AK, Boinett CJ, Abd El Ghany M. High-throughput transposon mutagenesis in the family Enterobacteriaceae reveals core essential genes and rapid turnover of essentiality. *mBio*. 2024;15(10).
35. Zhang C, Phillips APR, Wipfler RL, Olsen GJ, Whitaker RJ. The essential genome of the crenarchaeal model *Sulfolobus islandicus*. *Nat Commun*. 2018;9(1):4908. <https://doi.org/10.1038/s41467-018-07379-4> PMID: [30464174](https://pubmed.ncbi.nlm.nih.gov/30464174/)
36. Jana B, Cain AK, Doerler WT, Boinett CJ, Fookes MC, Parkhill J, et al. The secondary resistome of multidrug-resistant *Klebsiella pneumoniae*. *Sci Rep*. 2017;7:42483. <https://doi.org/10.1038/srep42483> PMID: [28198411](https://pubmed.ncbi.nlm.nih.gov/28198411/)
37. Mandal RK, Kwon YM. Global screening of *Salmonella enterica* serovar Typhimurium genes for desiccation survival. *Frontiers in Microbiology*. 2017;8.

38. Sarmiento F, Mrázek J, Whitman WB. Genome-scale analysis of gene function in the hydrogenotrophic methanogenic archaeon *Methanococcus maripaludis*. *Proc Natl Acad Sci U S A*. 2013;110(12):4726–31. <https://doi.org/10.1073/pnas.1220225110> PMID: [23487778](https://pubmed.ncbi.nlm.nih.gov/23487778/)
39. Dunn OJ. Multiple Comparisons among Means. *Journal of the American Statistical Association*. 1961;56(293):52–64. <https://doi.org/10.1080/01621459.1961.10482090>
40. Holm S. A Simple Sequentially Rejective Multiple Test Procedure. *Scandinavian Journal of Statistics*. 1979;6(2):65–70.
41. Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B: Statistical Methodology*. 1995;57(1):289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>
42. Goodall ECA, Azevedo Antunes C, Möller J, Sangal V, Torres VVL, Gray J, et al. A multiomic approach to defining the essential genome of the globally important pathogen *Corynebacterium diphtheriae*. *PLoS Genet*. 2023;19(4):e1010737. <https://doi.org/10.1371/journal.pgen.1010737> PMID: [37099600](https://pubmed.ncbi.nlm.nih.gov/37099600/)
43. Meinshausen N, Bühlmann P. Stability selection. *Journal of the Royal Statistical Society Series B: Statistical Methodology*. 2010;72(4):417–73.
44. Liu H, Roeder K, Wasserman L. Stability Approach to Regularization Selection (StARS) for High Dimensional Graphical Models. *Adv Neural Inf Process Syst*. 2010;24(2):1432–40. PMID: [25152607](https://pubmed.ncbi.nlm.nih.gov/25152607/)
45. Müller CL, Bonneau R, Kurtz Z. Generalized stability approach for regularized graphical models. 2016.
46. Chicco D, Jurman G. The Matthews correlation coefficient (MCC) should replace the ROC AUC as the standard metric for assessing binary classification. *BioData Mining*. 2023;16(4).
47. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*. 2020;21(1):6. <https://doi.org/10.1186/s12864-019-6413-7> PMID: [31898477](https://pubmed.ncbi.nlm.nih.gov/31898477/)
48. Matthews BW. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta*. 1975;405(2):442–51. [https://doi.org/10.1016/0005-2795\(75\)90109-9](https://doi.org/10.1016/0005-2795(75)90109-9) PMID: [1180967](https://pubmed.ncbi.nlm.nih.gov/1180967/)
49. Wetmore KM, Price MN, Waters RJ, Lamson JS, He J, Hoover CA. Rapid quantification of mutant fitness in diverse bacteria by sequencing randomly bar-coded transposons. *mBio*. 2015;6(3).
50. Ma Y, Pirolo M, Jana B, Mebus VH, Guardabassi L. The intrinsic macrolide resistome of *Escherichia coli*. *Antimicrob Agents Chemother*. 2024;68(8):e0045224. <https://doi.org/10.1128/aac.00452-24> PMID: [38940570](https://pubmed.ncbi.nlm.nih.gov/38940570/)
51. Yamazaki Y, Niki H, Kato J-i. Profiling of *Escherichia coli* chromosome database. In: Osterman AL, Gerdes SY, editors. *Microbial gene essentiality: protocols and bioinformatics*. Totowa, NJ: Humana Press; 2008. p. 385–9.
52. Porwollik S, Santiviago CA, Cheng P, Long F, Desai P, Fredlund J, et al. Defined single-gene and multi-gene deletion mutant collections in *Salmonella enterica* sv Typhimurium. *PLoS One*. 2014;9(7):e99820. <https://doi.org/10.1371/journal.pone.0099820> PMID: [25007190](https://pubmed.ncbi.nlm.nih.gov/25007190/)
53. Bai J, Dai Y, Farinha A, Tang AY, Syal S, Vargas-Cuevas G, et al. Essential Gene Analysis in *Acinetobacter baumannii* by High-Density Transposon Mutagenesis and CRISPR Interference. *J Bacteriol*. 2021;203(12):e0056520. <https://doi.org/10.1128/JB.00565-20> PMID: [33782056](https://pubmed.ncbi.nlm.nih.gov/33782056/)
54. Guzman LM, Barondess JJ, Beckwith J. FtsL, an essential cytoplasmic membrane protein involved in cell division in *Escherichia coli*. *J Bacteriol*. 1992;174(23):7716–28. <https://doi.org/10.1128/jb.174.23.7717-7728.1992> PMID: [1332942](https://pubmed.ncbi.nlm.nih.gov/1332942/)
55. Peterson JM, Phillips GJ. Characterization of conserved bases in 4.5S RNA of *Escherichia coli* by construction of new F' factors. *J Bacteriol*. 2008;190(23):7709–18. <https://doi.org/10.1128/JB.00995-08> PMID: [18805981](https://pubmed.ncbi.nlm.nih.gov/18805981/)
56. Slagter-Jäger JG, Puzis L, Gutsell NS, Belfort M, Jain C. Functional defects in transfer RNAs lead to the accumulation of ribosomal RNA precursors. *RNA*. 2007;13(4):597–605. <https://doi.org/10.1261/rna.319407> PMID: [17293391](https://pubmed.ncbi.nlm.nih.gov/17293391/)
57. Sakamoto K, Ishimaru S, Kobayashi T, Walker JR, Yokoyama S. The *Escherichia coli* argU10(Ts) phenotype is caused by a reduction in the cellular level of the argU tRNA for the rare codons AGA and AGG. *J Bacteriol*. 2004;186(17):5899–905. <https://doi.org/10.1128/JB.186.17.5899-5905.2004> PMID: [15317795](https://pubmed.ncbi.nlm.nih.gov/15317795/)
58. Bubunenko M, Baker T, Court DL. Essentiality of ribosomal and transcription antitermination proteins analyzed by systematic gene replacement in *Escherichia coli*. *J Bacteriol*. 2007;189(7):2844–53. <https://doi.org/10.1128/JB.01713-06> PMID: [17277072](https://pubmed.ncbi.nlm.nih.gov/17277072/)
59. Wu X, Lv X, Lu J, Yu S, Jin Y, Hu J, et al. The role of the ptsI gene on AI-2 internalization and pathogenesis of avian pathogenic *Escherichia coli*. *Microb Pathog*. 2017;113:321–9. <https://doi.org/10.1016/j.micpath.2017.10.048> PMID: [29111323](https://pubmed.ncbi.nlm.nih.gov/29111323/)
60. Wu T, Liu J, Li M, Zhang G, Liu L, Li X, et al. Improvement of sabinene tolerance of *Escherichia coli* using adaptive laboratory evolution and omics technologies. *Biotechnol Biofuels*. 2020;13:79. <https://doi.org/10.1186/s13068-020-01715-x> PMID: [32346395](https://pubmed.ncbi.nlm.nih.gov/32346395/)
61. Lee E, Cho G, Kim J. Structural basis for membrane association and catalysis by phosphatidylserine synthase in *Escherichia coli*. *Sci Adv*. 2024;10(51):eadq4624. <https://doi.org/10.1126/sciadv.adq4624> PMID: [39693441](https://pubmed.ncbi.nlm.nih.gov/39693441/)