

# The accuracy of absolute differential abundance analysis from relative count data

Kimberly E. Roche<sup>1\*</sup>, Sayan Mukherjee<sup>123</sup>

1. Program in Computational Biology and Bioinformatics, Duke University, Durham, NC 27708, United States
2. Departments of Statistical Science, Mathematics, Computer Science, Biostatistics & Bioinformatics, Duke University, Durham, NC 27708, United States
3. Center for Genomic and Computational Biology, Duke University, Durham, NC 27708, United States

\* kimberly.roche@duke.edu

## Abstract

Concerns have been raised about the use of relative abundance data derived from next generation sequencing as a proxy for absolute abundances. For example, in the differential abundance setting, compositional effects in relative abundance data may give rise to spurious differences (false positives) when considered from the absolute perspective. In practice however, relative abundances are often transformed by renormalization strategies intended to compensate for these effects and the scope of the practical problem remains unclear. We used simulated data to explore the consistency of differential abundance calling on renormalized relative abundances versus absolute abundances and find that, while overall consistency is high, with median sensitivities (true positive rates) and specificities (1 - false positive rates) each of around 0.90, consistency can be much lower where there is widespread change in the abundance of features across conditions. We confirm these findings on a large number of real data sets drawn from 16S metabarcoding, expression array, bulk RNA-seq, and single-cell RNA-seq experiments, where data sets with the greatest change between experimental conditions are also those with the highest false positive rates. Finally, we evaluate the predictive utility of summary features of relative abundance data themselves. Estimates of sparsity and the prevalence of feature-level change in relative abundance data give accurate predictions of discrepancy in differential abundance calling in simulated data and provide useful bounds for worst-case outcomes in real data.

## Introduction

Warnings about the consequences of compositional effects in sequence count data have been published repeatedly in the decades since the technology's advent and its application to a host of biological problems. The issue relates to a loss of scale information during sample processing, which renders counts of genes, transcripts, or bacterial species as relative abundances. No consensus solution for this problem exists. In this work, we use simulated and real data on differential abundance calling to quantify the discrepancy between differential abundance estimates made on relative versus "absolute" abundances. Our simulations show that methods which heuristically rescale sample abundances are often highly consistent across relative and absolute count data and we confirm that the low complexity case, roughly corresponding to bacterial community profiling, is the most problematic. Further, we show that data sets which are especially susceptible to distortion by compositional effects can often be predicted on the basis of "signatures" of this distortion.

## Compositionality in sequence count data

Compositionality refers to the nature of sequence count data as containing relative abundance information only. In the differential abundance setting, several authors [1, 2, 3] have described the problem this poses: whereas researchers would like to interpret change in absolute abundances, compositional effects mean using change in relative abundances as a proxy can lead to false discoveries. A few authors have cited instances of these false discoveries in real data. Coate and Doyle [4, 5] discussed the issue of transcriptome size variation in plants and other systems and the impact of this on accurate transcriptome profiling. Nie *et al* and Lin *et al* [6, 7] documented the phenomenon of widespread "transcription amplification" by the transcription factor *c-Myc* and Lovén *et al* [8] used *c-Myc* data and parallel RNA quantification assays to show that substantial differences in total abundance between control and elevated *c-Myc* conditions resulted in very different interpretations of apparent differential expression.

Common to these studies of transcriptomes is a recommendation that, where feasible, researchers leverage RNA spike-ins as controls against which changes in observed abundance can be scaled [9, 10]. But this practice has fallen short of widespread adoption. While several papers have expressed confidence in the utility of spike-ins [11, 12, 13, 14, 15], the doubt cast by reports of widespread batch effects [16] and technical noise [17] have had the effect of reducing researcher confidence in their use. Further, the introduction of spike-ins is not practical on all platforms.

### **Box 1: Measuring *relative* abundances**

Sequence counting has become widespread as a means of census-taking in microscopic biological systems. Genomic material, typically RNA, is captured and quantified at the component level. Sampled cells are lysed, messenger RNA is captured and fragmented, transcribed into cDNA, sequenced, classified, and quantified. The results are relative abundances of gene products in the cell (in the case of single-cell RNA-seq) or tissue (in bulk RNA-seq). In another instance, whole bacterial communities are profiled by barcoding of the 16S subunit of the ribosome. Ribosomal RNA associated with this piece of translation machinery is ubiquitously present across the bacterial kingdom but variations in the genetic sequence of this component can uniquely identify bacteria to the species or strain level in well-characterized systems, allowing a researcher to profile bacterial community composition. Absent measurements of microbial load or transcriptome size, however, the observed sequence counts in all these cases represent relative abundances.

Sequence count data is compositional due to steps in sample processing. Across domains, samples are typically normalized to some optimal total amount of genetic material prior to sequencing in accordance with manufacturer recommendations for best performance. This step removes variation in total abundance across samples. Saturation of sequencing has been cited [18] as another mechanism by which abundances are rendered relative: a finite amount of reagent means there is an upper limit on biological material which can be captured; rare components can be forced out by a "competition" to be sampled. These factors withstanding, observed total abundances would likely still be noisy. Repeated subsampling of small amounts of material and variation in the efficiency of library preparation steps can distort observed totals.

In transcriptomics and in microbial community profiling, residual variation in observed total abundances across samples is generally taken to be technical noise and most analytics pipelines involve steps to rescale observed abundances. The simplest of these is the counts per million (CPM) transformation which converts observed counts to relative abundances, then scales by 1 million.

Where approaches that rely on spike-ins are undesirable or infeasible, sample rescaling procedures have proliferated. These methods typically assume the existence of a stable set of features and attempt to normalize compositions in such a way as to recover this stable set across samples. In fact, in transcriptomics, these methods predominate.

In the setting of microbial community profiling, the prevailing assumption is that typical compositions are

too simple for rescaling methods to work well (although results in benchmarking studies have been mixed [19, 20]). Competing approaches have been developed for dealing with compositionality in microbial sequence count data. Quantitative microbiome profiling [21] and similar approaches combine relative abundances with complementary measurements of microbial load to reconstruct absolute abundances. In contrast, so-called compositional methods are also utilized. These involve log relative representations which can give approximate log-normality, such that workhorse statistical methods for continuous data may be applied. However, interpretation of these quantities can be challenging (e.g. as with the isometric logratio [22]).

Though there is evidence from simulated and real data that *scale* - i.e. increasing complexity of composition in terms of numbers of genes, transcripts, or bacterial sequence variants - mitigates the problem of compositionality [2, 20], it remains unclear whether there are instances where it is reasonable to substitute relative abundances for absolute abundances and several fields could benefit from clarity on this issue. In this work, we quantify the discrepancy in differential abundance calling on simulated and real data sets representative of 16S metabarcoding, bulk RNA-seq, and single-cell RNA-seq experiments.

In simulations exploring a broad range of differential abundance scenarios, we find median false positive rates of differential abundance calls made on absolute versus relative counts are around 0.10 and that false positive rates increase with the proportion of differentially abundant features. The complexity of composition (number of features) plays little role in observed outcomes. An exploration of sequencing data collected from twelve external studies reveals similar trends. Further, we show that summaries of sparsity and the prevalence of apparent feature-level change can provide excellent predictions of the amount of discrepancy in differential abundance calls in simulated data and can bound expectations of discrepancy in real data sets.

## Results

We simulated differentially abundant count data in paired sets of absolute and relative abundances. In the absolute abundances, differentially abundant features experienced either an increase or decrease in abundance - often large - between each of two simulated conditions. Large numbers of differentially abundant features frequently had the effect of changing the overall total abundance, potentially resulting in several-fold changes in scale between conditions. These changes in scale were purposely removed from a paired relative abundances by a simple resampling procedure. Compositional effects would present as the *relative* change experienced by non-differentially abundant features.

We explored ranges of complexity in composition and in the amount of differential abundance, grouping simulations into three partially overlapping settings: a **Microbial** setting, characterized by low complexity and high differential abundance; a **Bulk transcriptomic** setting having high complexity and low differential abundance; and an intermediate **Cell transcriptomic** setting. The full results from almost 6000 simulations are shown in Fig. 1, where agreement between differential abundance calls on absolute and relative abundances for each data set are summarized by means of sensitivity and specificity statistics, as described below. The same results are presented in simpler terms, as a function increasing complexity of composition (expressed as increasing feature number), in Supplemental Figs. ?? and ??. We evaluated the consistency of a small set of widely utilized differential abundance testing methods: ALDEx2 [23], DESeq2 [24], and scran [25]. Each of these rescales the observed counts against a reference quantity - typically, a subset of putatively stable features. Where such a stable reference exists, differences in these rescaled counts should approximate differences in absolute counts. In all cases, baseline or "true" differential abundance for each feature in the absolute count data was determined by the use of a simple generalized linear model. This yielded a common standard against which to compare differential abundance calls made by ALDEx2, DESeq2, and scran. Details on the simulation procedure and these differential abundance calling algorithms are given in Methods.

We report outcomes in terms of sensitivity (true positive rate) and specificity (1 - false positive rate). Perfect concordance of differential abundance calls made on observed versus absolute counts would yield a sensitivity of 1.0 and a specificity of 1.0. Sensitivity drops as more features deemed differentially abundant from the perspective of the absolute counts fail to appear significantly different in the observed data and specificity drops as an increasing number of features appear differentially abundant in the observed counts alone. In simulated data, the primary mechanism of low specificity was the compositional effect: relative changes in non-differentially abundant features. The sensitivities and specificities we report summarize agreement and disagreement at the level of whole data sets. We highlight key observations made on simulated data below.

## **Simulated data**

### **Moderate specificity was typical but cases of low specificity were observed in all settings for all methods**

Median specificity was moderately high at 0.90 - lowest for DESeq2 in the Microbial setting at 0.86 and highest for ALDEx2 in transcriptomic settings at 0.92 (see Table 1). However, a minority of simulated

data sets yielding very large false positive rates were observed in every setting. Data sets with very low specificity ( $< 0.5$ ) made up less than 10% of simulated cases in the Transcriptomic settings and more than 14% in the Microbial setting.

### **Increasing feature number improves specificity but the effect is modest**

For all methods, median specificity was lowest in the Microbial setting and highest in either the Bulk or Cell transcriptomic setting, as reported in Table 1. Simply increasing the simulated feature number while fixing all other parameters also resulted in higher overall specificity for all methods (see Table 2) indicating that feature number mediates this reduced false positive rate. The improvement in specificity is small, however. In essence larger feature number merely eliminates the very worst outcomes.

### **Though scran was the top performer, outcomes were similar across methods**

scran had the highest overall accuracy with a median sensitivity of 0.93 and specificity of 0.90. Sensitivity was lower in ALDEx2 (median = 0.77) and specificity was lower in DESeq2 (median = 0.86). All in all, methods exhibited similar performance in terms of their distributions of sensitivity and specificity across settings. See Figs. 1 and 2.

Next, we explored characteristics of relative abundances associated - either positively or negatively - with observed sensitivities and specificities in simulated data. It is possible to imagine characteristics which might indicate the presence of distortion by compositional effects, for example, an increase in the percent of rare features from one condition to the next (in effect, *dropouts*). While we might not expect any single such characteristic to predict the sensitivity or specificity of calls made on relative versus absolute abundances, composites of such characteristics might be informative. We describe the strongest associations we observed below.

### **Sensitivity is anti-correlated with estimates of sparsity**

Large proportions of zero- and one-counts in the simulated relative abundances correlated with low sensitivity (Spearman's correlation between proportion zeros and sensitivity,  $\rho = -0.60$ ). In effect, this echoes similar findings that decreasing sequencing depth decreases power in genomics studies [14] and reflects a lesser expected statistical confidence in the observed change of low-count features. Interestingly, in our simulations, evidence of extremes in terms of the apparent correlation of relative abundances (in fact, the skew of that distribution of correlations) was also inversely associated with sensitivity ( $\rho = -0.51$ ).

## Specificity is strongly anti-correlated with the estimated proportion of differential features

The proportions of features undergoing large fold decreases or large fold increases in abundance relative to the mean were highly informative with respect to specificity ( $\rho = -0.67$  and  $\rho = -0.45$  respectively). The standard deviation of the change in log counts between simulated conditions was also anti-correlated with specificity ( $\rho = -0.59$ ). In each case, these characteristics supplied evidence of the existence (or lack) of widespread change in composition. Methods which rescale observed abundances rely on a pool of stable reference features against which to estimate per-sample "scaling factors." The larger the number of components in the composition apparently changing, the greater the extent to which this rescaling is impaired.

We utilized these features derived from relative abundance data as well as several dozens more to train per-method models of both sensitivity and specificity, with the aim of predicting the discrepancy in relative versus absolute differential abundance calls from features of observed data alone, reasoning that these "signatures" might be highly informative. All model features are outlined in Supplemental Table ???. Predictive models were trained on 80% of our simulated data and their performance was evaluated on the held-out 20% of simulated data sets. Predictions of specificity were more accurate than predictions of sensitivity (mean specificity  $R^2 = 0.91$ ; mean sensitivity  $R^2 = 0.82$ ). Prediction of outcomes was easiest for DESeq2, where the  $R^2$  values for sensitivities and specificities on held-out data were 0.83 and 0.94, respectively. Full results are given in Table 3. These results suggested it might be straightforward to predict discrepancy between differential abundance calls made on absolute versus relative abundances using characteristics of the relative abundance data alone. However, the study of some representative real data sets will suggest prediction may be much more difficult in this setting.

## Real data

Next, we examined a variety of real data sets across many experimental settings in order to explore outcomes in real data. We collected publicly available data from twelve studies [26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37] and attempted to reconstruct absolute abundances by normalizing observed total sample abundances against reference quantities provided in the same published materials. In most cases, these reference quantities were external RNA spike-in sequences. In others, reconstructed absolute abundances had already been estimated, as in [26] through quantitative microbiome profiling (QMP [21]). In one case [36], we normalized against *Gapdh*, a stable, highly expressed housekeeping gene [38] and in another [31], against paired measurements of cell mass. Visual and textual summaries of these data

sets are available in Fig. 4 and Table 4. While not exhaustive, these studies illustrate a wide range of experimental conditions. We acknowledge the difficulty in reconstructing absolute abundances and caution that these estimates are likely noisy. Rather than call them absolute abundances, we will adopt the helpful terminology of [9] and call these reconstructions "nominal abundances" to distinguish them from some theoretical ground truth. While these quantities are an abstraction, since they derive from real experiments, we believe they capture some of the real variation in composition and scale we could expect to see in data from typical experiments and should serve to ground expectations. As with simulated data, we proceed by comparing differential abundances called on the relative count data (via ALDEx2, DESeq2, or scran) with differential abundances called on the nominal abundances - our proxy for "true" differential abundance.

We have endeavored to include among the real data sets some of the most challenging possible cases for differential abundance calling from relative count data. Many of the data sets we have selected involve a large amount of absolute and/or compositional change across experimental or biological conditions. These data sets help us answer the questions: What is the scale of discrepancy in the worst cases? How useful are summaries derived from relative abundances in these cases? Could a researcher reasonably predict error from relative abundance data alone?

Observed sensitivities and specificities for real data sets are given in Table 5.

### **Low sensitivity, high specificity cases were common**

Results of differential abundance calling on the data sets of Vieira-Silva et al., Muraro et al., Hagai et al., Hashimshony et al., Grün et al., and Kimmerling et al. are shown in Fig. 5 and were characterized by a combination of high specificity and very low sensitivity. These data consist of five transcriptomic studies and a single gut microbial data set. In general there was less detectable differential abundance in these six data sets than in the others we examined (see Table 4) but perhaps more importantly, this differential abundance was often subtle - small in absolute terms - and therefore difficult to detect. Features which were significantly differentially abundant in the absolute counts typically fell below the threshold of significance in the observed data. An example of this from the data of Kimmerling et al. is shown in Fig. 7.

We attempted to predict sensitivity and specificity for all real data sets using the same features defined for simulated data and the high specificity associated with these six data sets represented one of the easier-to-predict outcomes. Predictive intervals for specificity were generally small and outcomes fell within or above the predicted specificity 0.83 of the time (15 of 18 cases). However, corresponding sensitivities were



not as well predicted, especially for the single-cell experiments of Hashimshony et al. and Kimmerling et al. where the models overpredicted the observed sensitivity by as much as 50%.

### **Methods vary in their ability to handle borderline cases**

Results for the data sets of Song et al., Barlow et al., Monaco et al., and Yu et al. are summarized in Fig. 6. Accuracy of outcomes was mixed, both in terms of sensitivity and specificity, with results often varying by method. One method tended to outperform the rest in these intermediate cases - more frequently than not, DESeq2, for example on Yu et al. (Fig. 6j). These data sets generally had a larger proportion of differentially abundant features, approaching or exceeding 50%, making the task of accurate renormalization of sample abundances in these cases exceedingly difficult for all methods.

### **Extreme cases with large absolute change had very low specificity**

The data sets of Yu et al., Owens et al., and Klein et al. each had a (sometimes overwhelming) majority of differentially abundant features in their nominal counts. Experimentally, these studies involved either a genome-wide perturbation of expression or differences in transcriptome size across conditions. Together, this group presented the most difficult challenges for accurate renormalization. As a result, though sensitivity was at its highest in this data, specificity was overall very low, seemingly due to compositional effects, as seen in a representative example from Klein et al. in Fig. 7.

Two data sets stand out in their superficial similarity and divergent outcomes: the high-specificity count data of Kimmerling et al. and low-specificity Owens et al. Both experiments involved large biological perturbations resulting in more than a doubling of transcriptome volume between conditions. However compositional effects are only apparent in the relative abundances of Owens et al. The differentiating factor is ostensibly noise: the dispersion of counts in Kimmerling et al. is much greater - that is, the data are noisier - and confident calls of differential abundance on the observed data are few.

Sensitivity was greatly reduced in the real data sets we studied relative to our simulated data, possibly due to two factors. We observed a trend of diminishing sensitivity with increasing feature number in our simulations (see Fig. 2) and the largest real data sets evaluated here have 20,000 or more features. These data were also noisier than our simulations, due to our spike-in rescaling strategy. In both cases, reduced sensitivity might have been expected.

Cases where false positive rates should have been low generally were: low average fold change across conditions and a strict minority of features in Muraro et al., Hashimshony et al., and Song et al. resulted in high specificity for all methods (often at a cost of low sensitivity). In these data sets there is little

change in scale available to drive compositional effects and large pools of relatively stable features should make effective sample renormalization possible. In increasingly challenging cases like Grün et al. and Hagai et al., which exhibited larger fold change across conditions and larger proportions of differential features, most methods still performed well with respect to Type 1 error control.

Direct prediction of these outcomes from summary statistics was more successful for specificity than for sensitivity but accuracy of point predictions was reduced relative to simulated data ( $R^2 = 0.260$  for predicted vs. observed sensitivity, all methods;  $R^2 = 0.643$  for specificity). However, the same models can be induced to yield predictive intervals and when utilized as bounds on plausible outcomes for a given data set, these may still be highly useful. Consider the lower bound of the 90% predictive interval: observed sensitivity was at least as high as that lower bound in 83% of cases (29 of 35) and observed specificity was at least as high as that lower bound in over 91% of cases (32 of 35). This implies that this and similar models, trained on features of the observed data alone, may be able to bound expectations about the "worst case" for a given data set.

## Discussion

While the potential for compositional effects to drive differential abundance has repeatedly been described in the literature, uncertainty remains about the scope of this problem. Previous studies have shown that rates of false positive differential abundance calls can be high in certain settings [19, 20, 2]. Our results indicate that the problem is at least partially a function of the amount of change in the system under study and that differential abundance estimates from experiments characterizing extreme change across observed conditions are likely to be distorted. Both simulated and published experimental data contributed to this picture. Data sets with low fold change across conditions and a strict minority of differentially abundant features had high specificities. Compositional effects and high rates of false positive differential abundance calls were observed in simulated and real data where the scale of change between conditions was roughly two-fold or greater, in line with similar evidence of from 16S data [39]. We found prediction of outcomes was possible in a limited sense. In simulated data, especially discrepant outcomes were usually predictable by a few summary statistics derived from their relative abundances. Low sensitivity was predictable from features which captured information about sequencing depth. Low specificity was best predicted by estimates of the number of differentially abundant features and (implicitly) their effect on the efficacy of renormalization. Though accuracy was not exceptionally high for direct prediction of true and false positive rates in real data, the prediction of a lower bound on these

same quantities was reasonably reliable, suggesting the possibility of identifying worst case scenarios for a given data set based on observed counts alone. These estimates might, for example, motivate a researcher to pursue more or less stringent false discovery control.

Our analysis of real data admits a shortcoming: in attempting to reconstruct absolute abundances in real data sets, we generally utilized spike-ins as control quantities against which to scale total sample abundances. The resulting reconstructions are undoubtedly noisy and simply a best approximation of real change in the system. A further important investigation in the spirit of recent work by Lloréns-Rico et al. [40] might assess the usefulness of partial reconstructions of abundance (e.g. by quantitative PCR or spike-in renormalization) by quantifying the effect of this restored information on the accuracy of downstream analyses.

Lastly, it should be noted that the concerns motivating this and similar studies may be moot for some types of sequence count data. In particular, some single-cell platforms generate library sizes (i.e. total per-sample observed abundances) which are already roughly proportional to absolute abundances. This is in line a view that abundances in deeply sequenced UMI-barcoded single cells are likely to be a good proxy for absolute abundances [12, 14]. The effect of compositionality may be a minor concern under these circumstances.

## Materials & Methods

We simulated general purpose molecular count data. These counts are interpretable as a variety of biological quantities, for example, transcript abundance in a cell or bacterial strain abundance in a microbial community. The simulated abundances undergo a sampling step intended to loosely replicate the process of measurement itself and, crucially, the normalization of total abundance across samples, giving a second set of count data. We refer to the first set of count data as "absolute" counts and the second, resampled set as relative or observed counts and explored the degree to which this loss of information about changes in total abundance alters the results of a mock differential abundance analysis by simulating a wide range of settings in our data, where key characteristics like complexity of composition (e.g. gene number) and fold change across simulated conditions varied widely.

### Simulation model

We designed a simulation framework to generate count data corresponding to two arbitrarily different conditions, denoted by superscripts in the equations below. First, for  $p = 1, \dots, P$  features in the first

condition, a set of log mean abundances was drawn as

$$\theta_p^{(1)} \sim N(m, S^2)$$

where hyperparameters  $m$  and  $S$  tune the mean and standard deviation of baseline log abundances. A correlation matrix was drawn as

$$\Omega \sim \text{Inverse-Wishart}(n, Q)$$

where scale matrix  $Q$  was supplied as either the identity matrix (for a minority of simulations) or a dense correlation matrix with net positive elements. The matrix  $\Omega$  is subsequently re-scaled to a correlation matrix and used to draw correlated feature perturbations in a second condition as

$$\theta_p^{(2)} \sim \text{MVN}(\theta_p^{(1)}, \Omega \cdot a)$$

where the hyperparameter  $a$  exists in order to tune the overall scale of the correlated log perturbations. Mean abundances on the scale of sequence counts for each condition are calculated as

$$\gamma_p^{(1)} = \exp(\theta_p^{(1)}), \quad \gamma_p^{(2)} = \exp(\theta_p^{(2)})$$

Differentially abundant features in some desired proportion,  $c$ , are obtained as follows: features are selected as differentially abundant with probability  $c$ . For those selected features only, the perturbed  $\gamma_p^{(2)}$  serves as the mean abundance in the second condition; for all other features, the mean abundance in both the first and second conditions is given by  $\gamma_p^{(1)}$ . Let these new vectors be  $\mu_p^{(1)}, \mu_p^{(2)}$ . These represent mean the abundances of  $P$  features in two conditions, some of which differ across conditions, others of which are identical. Replicates  $i = 1, \dots, 10$  are then generated for each condition as follows. A fixed dispersion parameter for absolute counts is defined as  $d_{\text{abs}} = 1000$  and those counts are drawn as

$$y_{i,p}^{(1)} \sim \text{NegBinom}(\mu_p^{(1)} \cdot \delta, 1000)$$

where

$$\delta \sim \max(0.1, N(1, g))$$

(Note that the dispersion parameter has been chosen such that the resulting counts are only barely

overdispersed with respect to a Poisson.) The purpose of the truncated, per-sample multiplier  $\delta$  is to re-scale all abundances within a given sample by some factor of approximately 1 but by increasing the scale of hyperparameter  $g$ , increasing replicate noise can be added within a condition. This process is repeated for the second condition to give a set of absolute counts  $y_p^{(2)}$ .

A new average observed total abundance (or library size) is randomly drawn as

$$u \sim \text{Unif}(5000, 2 \times 10^6)$$

Finally, observed abundances  $z$  are generated through a multinomial resampling procedure similar to that of [10, 20, 23] which gives relative abundances as counts per million. Where  $i$  and  $k$  index different samples prior to resampling,  $i'$  indexes the sample  $i$  after resampling, and total counts for sample  $i$  prior to resampling are given by  $n_i = \sum y_i$ , we have

$$z_{i'} \sim \text{Mult}(\pi_{i'} = y_i/n_i, n_{i'} = n_k)$$

where superscripts have been suppressed as this procedure is identical across simulated "conditions." The resulting  $P$ -length vector of counts for a given sample contains relative but not absolute abundance information. These vectors are collapsed into a  $P \times 20$  count matrix containing 10 replicate samples for each of two simulated conditions. In order to evaluate the discrepancy of differential abundance calling on observed versus absolute counts, we apply differential abundance methods to count matrices  $Z$  and  $Y$  respectively and score the differences.

## Breadth of simulations

In order to generate simulations with a wide variety of characteristics, we swept in a grid over all our hyperparameters. Feature number  $P$  was stepped through values 100, 1000, and 5000. A maximum feature number of 5000 was chosen as these simulations were computationally intensive and major trends had become apparent at that scale. The degree of feature correlation was encoded in five realizations of scale matrices  $Q$ , encoding fully independent features at one extreme and 50% strongly positively correlated features at the other extreme. Log mean abundance ( $m$ ) and the log variance ( $S$ ) were independently incremented through low to high values. Likewise, the average log perturbation size ( $a$ ) was swept from low to high in five steps, as a proportion of log mean abundance.

Replicate noise  $g$  varied from low to high in three steps. And finally, the proportion of differentially abundant features ranged across 20%, 30%, 50%, 70%, and 90%. Note that because many "perturbations"

were very small, detectable differential abundance was generally only a fraction of the parameterized amount and most data sets contain a minority of differentially abundant features. Overall this 5625 simulated data sets were generated with almost continuous variation characteristics of interest.

We suggest that ranges of these parameter settings approximately represent different data modalities. We term the Microbial setting that with low to moderate feature number ( $P \leq 1000$ ) and largest average perturbation, in accordance with a belief that bacterial communities are often simple (in terms of sequence variants with more than negligible abundance) and that they are highly variable even at short time scales [41].

We designate the Bulk transcriptomic setting as that with the largest feature number ( $P = 5000$ ) and having a lower average perturbation, the rationale being that transcriptomes sampled in aggregate over many cells are complex but largely stable compositions. Similarly, we define the intermediate Cell transcriptomic setting, approximately representative of single-cell RNA-seq data, to comprise simulations with moderate to large feature numbers ( $P \geq 1000$ ) and moderate perturbation sizes. These categories are intended as rough outlines and we note that within these settings the realized data varies in terms of 1) degree of feature correlation, 2) overall abundance, 3) (un)evenness of composition, and 4) within-condition variation in totals.

## Calling differential abundance

Three differential abundance calling methods were used in this study; each relies upon the use of a reference quantity to renormalize sample total abundances. For ALDEx2 this reference quantity is a trimmed version of per-sample mean log abundance. DESeq2 rescales samples using a sophisticated Bayesian model. Implicit in DESeq2's procedure is an assumption that a large set of stable features exists against which the observed changes in other features can be adjusted. scran's procedure follows a similar rationale, but employs a local rescaling, within clusters of like samples.

For simplicity, we omit from consideration models which lack a rescaling. We also omit zero-inflation models. Although these are popular in single-cell mRNA sequencing data, the debate continues about whether these models are appropriate for these data [42, 43].

Finally, differential abundance calls made on observed counts must be evaluated against a reference in order to calculate discrepancy. For this reference, we used calls made by a negative binomial generalized linear model on absolute abundances as a pseudo-gold standard or "oracle" in all cases. Features with significantly different means according to the GLM were considered "true" instances of differential abundance. One caveat is that some of the discrepancy in calls between absolute and relative abundances will

be due to differences in sensitivity between the models applied to the reference and "observed" data sets - i.e. between a stock NB GLM and DESeq2. In all settings, true positive, true negative, false positive, and false negative calls were manually spot checked to verify disagreements between methods were generally unambiguous.

Accuracy (in terms of sensitivity and specificity) was calculated between differential abundance calls made on absolute abundances and observed abundances. Three methods were evaluated: ALDEx2, DESeq2, and scran. ALDEx2 was called using the `aldex()` from the associated R package [23] using the interquartile logratio reference. DESeq2 was called using the Seurat wrapper (`FindMarkers(test.use = ...)`) [24, 44]. scran was called using the method described in the vignette associated with its Bioconductor package [25]. Unadjusted p-values were collected from all methods and multiple test correction applied via `p.adjust` in R using Benjamini-Hochberg method.

## Predictive modeling

In total, we trained 6 random forest models over 57 summary features, for each combination differential abundance calling method (ALDEx2, DESeq2, or scran) and accuracy measure (sensitivity or specificity). All such predictive models were fit with the `randomForest` package in R [45]. A random forest is an ensemble of decision trees and this tree-based approach was chosen because, while extensive, our simulations were not exhaustive. We anticipated that learning sets of decision rules might generalize well to unseen conditions, in particular feature numbers larger than those we explored in simulation. Feature importance was measured as "gain," or the relative increase in predictive accuracy achieved by the inclusion in the model of a given feature, as computed by the `caret` package in R [46].

Predictive models built from these summary features attempted to estimate sensitivity and specificity values explicitly. Details on these models are given in the Methods. All models were trained on 80% of the simulated data and their predictive accuracy was assessed on the reserved 20%.

## Real data processing

Publicly available data was downloaded from sources provided in the published materials for the studies cited. In general, relative abundances were converted to counts per million and sequences for external spike-ins extracted from these counts. Simple per-sample scaling factors were calculated from mean spike-in abundances and applied to the relative data to give "nominal" abundances.

In the inhibited and control cells of Klein et al., the expression of housekeeping gene *Gapdh* was used to normalize per-sample total abundance. In the 16S metagenomic data of Barlow et al. and Vieira-Silva et

al., nominal abundances had already been estimated by the authors of those studies using quantitative PCR-based methods. In the single-cell expression data of Yu et al.,  $\log_2$  observed total abundances correlated well with  $\log_2$  spike-in abundances and the observed data themselves were treated as true abundances. Relative abundances were derived by normalizing per-sample library sizes and scaling to give counts per million. Finally, in the coupled cell mass and expression measurements of Kimmerling et al., the estimated cell mass was used as scaling factor for observed expression to give nominal abundances. Nominal and relative abundances were then filtered to exclude features (genes or bacterial sequence variants) present at an average abundance of less than a single count in either the nominal or relative count tables. This both reduced the size of the largest data sets - making them more computationally manageable - and reduced sparsity in the most extreme cases.

## Acknowledgements

The authors would like to acknowledge conversations with Justin Silverman, Lawrence David, and Tim Reddy for useful comments. The authors would like to acknowledge partial funding from HFSP RGP005, NSF DMS 17-13012, NSF BCS 1552848, NSF DBI 1661386, NSF IIS 15-46331, NSF DMS 16-13261, as well as high-performance computing partially supported by grant 2016-IDG-1013 from the North Carolina Biotechnology Center and a Forge Fellowship through Duke University. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of any of the funders.

## Code availability

All R code related to this study is available on Github at <https://github.com/kimberlyroche/codaDE>.



## References

- [1] Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome Datasets Are Compositional: And This Is Not Optional. *Front Microbiol.* 2017;8:57.
- [2] Evans C, Hardin J, Stoebel DM. Selecting between-sample RNA-Seq normalization methods from the perspective of their assumptions. *Brief Bioinform.* 2018;19(5):776–792.
- [3] Vieth B, Parekh S, Ziegenhain C, Enard W, Hellmann I. A systematic evaluation of single cell RNA-seq analysis pipelines. *Nat Commun.* 2019;10(1):4667.
- [4] Coate JE, Doyle JJ. Quantifying whole transcriptome size, a prerequisite for understanding transcriptome evolution across species: an example from a plant allopolyploid. *Genome Biol Evol.* 2010;2:534–546.
- [5] Coate JE, Doyle JJ. Variation in transcriptome size: are we getting the message? *Chromosoma.* 2015;124(1):27–43.
- [6] Nie Z, Hu G, Wei G, Cui K, Yamane A, Resch W, et al. c-Myc is a universal amplifier of expressed genes in lymphocytes and embryonic stem cells. *Cell.* 2012;151(1):68–79.
- [7] Lin CY, Lovén J, Rahl PB, Paranal RM, Burge CB, Bradner JE, et al. Transcriptional amplification in tumor cells with elevated c-Myc. *Cell.* 2012;151(1):56–67.
- [8] Lovén J, Orlando DA, Sigova AA, Lin CY, Rahl PB, Burge CB, et al. Revisiting global gene expression analysis. *Cell.* 2012;151(3):476–482.
- [9] Athanasiadou R, Neymotin B, Brandt N, Wang W, Christiaen L, Gresham D, et al. A complete statistical model for calibration of RNA-seq counts using external spike-ins and maximum likelihood theory. *PLoS Comput Biol.* 2019;15(3):e1006794.
- [10] McGee WA, Pimentel H, Pachter L, Wu JY. Compositional data analysis is necessary for simulating and analyzing RNA-Seq data; 2019.
- [11] Wu AR, Neff NF, Kalisky T, Dalerba P, Treutlein B, Rothenberg ME, et al. Quantitative assessment of single-cell RNA-sequencing methods. *Nat Methods.* 2014;11(1):41–46.
- [12] Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in single-cell transcriptomics. *Nat Rev Genet.* 2015;16(3):133–145.

- [13] Lun ATL, Calero-Nieto FJ, Haim-Vilmovsky L, Göttgens B, Marioni JC. Assessing the reliability of spike-in normalization for analyses of single-cell RNA sequencing data. *Genome Res.* 2017;27(11):1795–1806.
- [14] Svensson V, Natarajan KN, Ly LH, Miragaia RJ, Labalette C, Macaulay IC, et al. Power analysis of single-cell RNA-sequencing experiments. *Nat Methods.* 2017;14(4):381–387.
- [15] Ziegenhain C, Vieth B, Parekh S, Reinius B, Guillaumet-Adkins A, Smets M, et al. Comparative Analysis of Single-Cell RNA Sequencing Methods. *Mol Cell.* 2017;65(4):631–643.e4.
- [16] Tung PY, Blischak JD, Hsiao CJ, Knowles DA, Burnett JE, Pritchard JK, et al. Batch effects and the effective design of single-cell gene expression studies. *Sci Rep.* 2017;7(1):1–15.
- [17] Risso D, Ngai J, Speed TP, Dudoit S. Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat Biotechnol.* 2014;32(9):896–902.
- [18] Quinn TP, Erb I, Richardson MF, Crowley TM. Understanding sequencing data as compositions: an outlook and review. *Bioinformatics.* 2018;34(16):2870–2878.
- [19] Hawinkel S, Mattiello F, Bijnens L, Thas O. A broken promise: microbiome differential abundance methods do not control the false discovery rate. *Brief Bioinform.* 2019;20(1):210–221.
- [20] Calgaro M, Romualdi C, Waldron L, Risso D, Vitulo N. Assessment of statistical methods from single cell, bulk RNA-seq, and metagenomics applied to microbiome data. *Genome Biol.* 2020;21(1):191.
- [21] Vandeputte D, Kathagen G, D’hoë K, Vieira-Silva S, Valles-Colomer M, Sabino J, et al. Quantitative microbiome profiling links gut community variation to microbial load. *Nature.* 2017;551(7681):507–511.
- [22] Egozcue JJ, Pawlowsky-Glahn V, Mateu-Figueras G, Barceló-Vidal C. Isometric Logratio Transformations for Compositional Data Analysis. *Math Geol.* 2003;35(3):279–300.
- [23] Fernandes AD, Macklaim JM, Linn TG, Reid G, Gloor GB. ANOVA-like differential expression (ALDEx) analysis for mixed population RNA-Seq. *PLoS One.* 2013;8(7):e67019.
- [24] Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 2014;15(12):550.
- [25] Lun ATL, McCarthy DJ, Marioni JC. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Res.* 2016;5:2122.

- [26] Vieira-Silva S, Sabino J, Valles-Colomer M, Falony G, Kathagen G, Caenepeel C, et al. Quantitative microbiome profiling disentangles inflammation- and bile duct obstruction-associated microbiota alterations across PSC/IBD diagnoses. *Nat Microbiol.* 2019;4(11):1826–1831.
- [27] Muraro MJ, Dharmadhikari G, Grün D, Groen N, Dielen T, Jansen E, et al. A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst.* 2016;3(4):385–394.e3.
- [28] Hagai T, Chen X, Miragaia RJ, Rostom R, Gomes T, Kunowska N, et al. Gene expression variability across cells and species shapes innate immunity. *Nature.* 2018;563(7730):197–202.
- [29] Hashimshony T, Senderovich N, Avital G, Klochendler A, de Leeuw Y, Anavy L, et al. CEL-Seq2: sensitive highly-multiplexed single-cell RNA-Seq. *Genome Biol.* 2016;17:77.
- [30] Grün D, Kester L, van Oudenaarden A. Validation of noise models for single-cell transcriptomics. *Nat Methods.* 2014;11(6):637–640.
- [31] Kimmerling RJ, Prakadan SM, Gupta AJ, Calistri NL, Stevens MM, Olcum S, et al. Linking single-cell measurements of mass, growth rate, and gene expression. *Genome Biol.* 2018;19(1):207.
- [32] Song SG, Kim S, Koh J, Yim J, Han B, Kim YA, et al. Comparative analysis of the tumor immune-microenvironment of primary and brain metastases of non-small-cell lung cancer reveals organ-specific and EGFR mutation-dependent unique immune landscape. *Cancer Immunol Immunother.* 2021;70(7):2035–2048.
- [33] Barlow JT, Bogatyrev SR, Ismagilov RF. A quantitative sequencing framework for absolute abundance measurements of mucosal and lumenal microbial communities. *Nat Commun.* 2020;11(1):2590.
- [34] Monaco G, Lee B, Xu W, Mustafah S, Hwang YY, Carré C, et al. RNA-Seq Signatures Normalized by mRNA Abundance Allow Absolute Deconvolution of Human Immune Cell Types. *Cell Rep.* 2019;26(6):1627–1640.e7.
- [35] Yu Y, Fuscoe JC, Zhao C, Guo C, Jia M, Qing T, et al. A rat RNA-Seq transcriptomic BodyMap across 11 organs and 4 developmental stages. *Nat Commun.* 2014;5(1):1–11.
- [36] Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al. Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells. *Cell.* 2015;161(5):1187–1201.

- [37] Owens NDL, Blitz IL, Lane MA, Patrushev I, Overton JD, Gilchrist MJ, et al. Measuring Absolute RNA Copy Numbers at High Temporal Resolution Reveals Transcriptome Kinetics in Development. *Cell Rep.* 2016;14(3):632–647.
- [38] Padovan-Merhar O, Nair GP, Biaesch AG, Mayer A, Scarfone S, Foley SW, et al. Single mammalian cells compensate for differences in cellular volume and DNA copy number through independent global transcriptional mechanisms. *Mol Cell.* 2015;58(2):339–352.
- [39] Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al. Normalization and microbial differential abundance strategies depend upon data characteristics. *Microbiome.* 2017;5(1):27.
- [40] Lloréns-Rico V, Vieira-Silva S, Gonçalves PJ, Falony G, Raes J. Benchmarking microbiome transformations favors experimental quantitative approaches to address compositionality and sampling depth biases. *Nat Commun.* 2021;12(1):3562.
- [41] Johnson AJ, Vangay P, Al-Ghalith GA, Hillmann BM, Ward TL, Shields-Cutler RR, et al. Daily Sampling Reveals Personalized Diet-Microbiome Associations in Humans. *Cell Host Microbe.* 2019;25(6):789–802.e5.
- [42] Silverman JD, Roche K, Mukherjee S, David LA. Naught all zeros in sequence count data are the same. *Comput Struct Biotechnol J.* 2020;18:2789–2798.
- [43] Sarkar A, Stephens M. Separating measurement and expression models clarifies confusion in single-cell RNA sequencing analysis. *Nat Genet.* 2021;53(6):770–777.
- [44] Hafemeister C, Satija R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* 2019;20(1):296.
- [45] Liaw A, Wiener M. Classification and Regression by randomForest. *R News.* 2002;2(3):18–22.
- [46] Kuhn M. Building predictive models in R using the caret package. *J Stat Softw.* 2008;

## List of Tables

1	Median and thresholded simulated specificity by setting . . . . .	22
2	Median and thresholded simulated specificity by feature number . . . . .	22
3	Predictive model performance summaries . . . . .	23
4	Overview of real data sets . . . . .	24
5	Sensitivity and specificity on real data . . . . .	25

## List of Figures

1	Simulated sensitivity and specificity as a function of setting and proportion of differentially abundant features . . . . .	26
2	Simulated sensitivity and specificity as a function of setting and fold change . . . . .	27
3	Visual summaries of real data sets 1-6 . . . . .	28
4	Visual summaries of real data sets 7-12 . . . . .	29
5	Results for real data sets 1-6 . . . . .	30
6	Results for real data sets 7-12 . . . . .	31
7	Example discrepant differential abundance calls . . . . .	32

## Tables

Setting	Method	Median specificity	Percent of data sets below 95% specificity	Percent of data sets below 50% specificity
Microbial	ALDEx2	0.89	58%	18%
Microbial	DESeq2	0.86	78%	15%
Microbial	scran	0.90	72%	10%
Cell transcriptomic	ALDEx2	0.92	56%	5%
Cell transcriptomic	DESeq2	0.87	76%	11%
Cell transcriptomic	scran	0.92	67%	6%
Bulk transcriptomic	ALDEx2	0.92	56%	2%
Bulk transcriptomic	DESeq2	0.86	77%	1%
Bulk transcriptomic	scran	0.91	66%	4%

Table 1: Low specificity is a common outcome in simulated data. A majority of simulated data sets have false positive rates in excess of 5% and in the setting with lowest feature number and highest volatility (the Microbial setting), a substantial proportion of data sets had false positive rates in excess of 50%.

Feature number	Method	Median specificity	Percent of data sets below 95% specificity	Percent of data sets below 50% specificity
100	ALDEx2	0.83	67%	27%
100	DESeq2	0.84	82%	16%
100	scran	0.89	75%	12%
1000	ALDEx2	0.91	60%	10%
1000	DESeq2	0.86	80%	13%
1000	scran	0.90	71%	8%
5000	ALDEx2	0.92	61%	3%
5000	DESeq2	0.88	72%	9%
5000	scran	0.92	64%	4%

Table 2: Specificity in simulated data. Per-method results are grouped and sorted by increasing feature number.

Method	Sensitivity prediction $R^2$	Specificity prediction $R^2$
ALDEx2	0.817	0.856
DESeq2	0.826	0.941
scran	0.809	0.907

Table 3: Performance of random forest predictive models of sensitivity and specificity on held out data sets.

Source	Description	Number sequence variants	Number samples (per-condition)	Percent zeros	Approx. fold change	Approx. percent differential features
Vieira-Silva et al. (2019)	16S metagenomics from human gut samples of control and Crohn's disease patients	263	14, 54	72%	1.8	7%
Muraro et al. (2016)	single cell RNA-seq of pancreatic islet cells	13,764	100, 100	64%	1.1	13%
Hagai et al. (2018)	bulk RNA sequencing of both unstimulated and mock-viral infected mouse fibroblasts	13,848	10, 29	34%	1.5	19%
Hashimshony et al. (2016)	single cell RNA-sequencing of quiescent and cycling mouse fibroblasts	12,075	31, 38	22%	1.2	23%
Grün et al. (2014)	mouse embryonic stem cells cultured in serum and a two-inhibitor solution	11,572	76, 56	47%	1.3	40%
Kimmerling et al. (2018)	cycling, stimulated CD8+ T cells	9,930	79, 79	35%	2.0	74%
Song et al. (2021)	nCounter array of human primary lung cancer vs. brain metastases	773	13, 15	0%	1.0	39%
Barlow et al. (2020)	16S metagenomics from ketogenic diet and control mice	103	17, 18	41%	3.4	43%
Monaco et al. (2019)	immune cell profiling in human humans via bulk RNA-seq	20,576	4, 8	15%	3.4	45%
Yu et al. (2014)	single cell expression profiling of rat brain and liver tissue	14,085	32, 32	14%	2.0	85%
Klein et al. (2015)	single cell RNA-sequencing of normally developing and leukemia inhibitory factor-treated mouse ESCs	23,316	100, 100	67%	1.4	62%
Owens et al. (2016)	single cell sequencing of zebrafish embryos; early vs. late time course samples drawn	35,794	24, 35	17%	3.7	89%

Table 4: Real 16S metabarcoding, bulk RNA-seq, and single cell RNA-seq data sets corresponding to the abundances shown in Figures 3 and 4. The percent differential features are those significantly differential to a negative binomial GLM in the nominal abundances.



Data set	Method	Sensitivity	Specificity
Vieira-Silva et al.	ALDEx2	0.053	0.000
Vieira-Silva et al.	DESeq2	0.263	0.041
Vieira-Silva et al.	scran	0.105	0.074
Muraro et al.	ALDEx2	0.370	0.021
Muraro et al.	DESeq2	0.325	0.009
Muraro et al.	scran	0.696	0.097
Hagai et al.	ALDEx2	0.105	0.024
Hagai et al.	DESeq2	0.475	0.032
Hagai et al.	scran	0.578	0.137
Hashimshony et al.	ALDEx2	0.000	0.000
Hashimshony et al.	DESeq2	0.058	0.007
Hashimshony et al.	scran	0.063	0.004
Grün et al.	ALDEx2	0.182	0.050
Grün et al.	DESeq2	0.420	0.040
Grün et al.	scran	0.519	0.184
Kimmerling et al.	ALDEx2	0.003	0.000
Kimmerling et al.	DESeq2	0.002	0.002
Kimmerling et al.	scran	0.082	0.006
Song et al.	ALDEx2	0.577	0.141
Song et al.	DESeq2	0.626	0.212
Song et al.	scran	0.630	0.049
Barlow et al.	ALDEx2	0.409	0.288
Barlow et al.	DESeq2	0.500	0.153
Barlow et al.	scran	0.545	0.305
Monaco et al.	ALDEx2	0.056	0.011
Monaco et al.	DESeq2	0.379	0.165
Yu et al.	ALDEx2	0.775	0.389
Yu et al.	DESeq2	0.931	0.117
Yu et al.	scran	0.875	0.378
Klein et al.	ALDEx2	0.237	0.377
Klein et al.	DESeq2	0.842	0.311
Klein et al.	scran	0.978	0.763
Owens et al.	ALDEx2	0.708	0.801
Owens et al.	DESeq2	0.814	0.919
Owens et al.	scran	0.801	0.853

Table 5: Observed sensitivities and specificities for all methods in real data sets.

## Figures

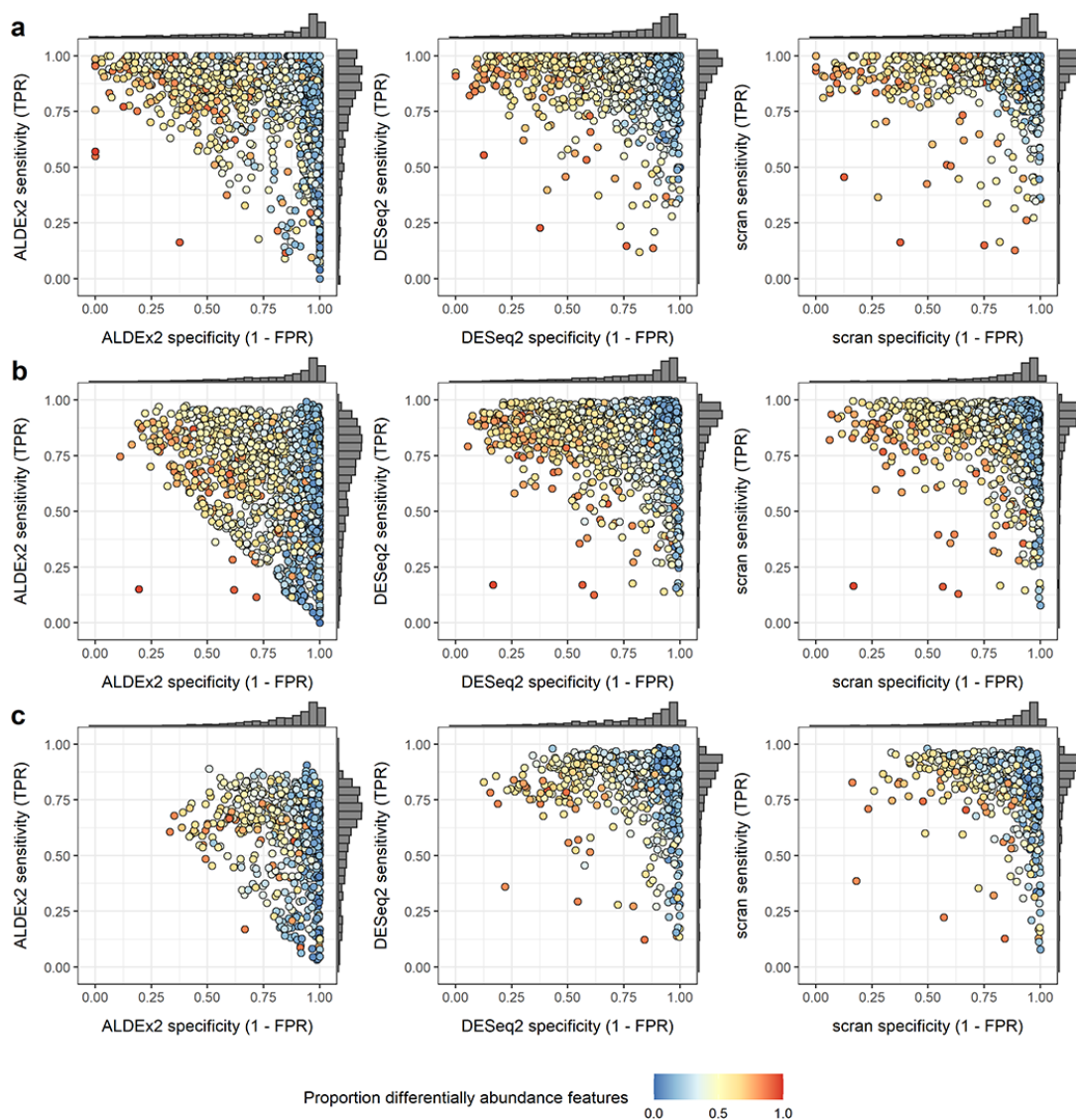


Figure 1: Sensitivity and specificity for three differential abundance calling methods in three experimental settings. Data sets are labeled by proportion of features with at least a 50% increase or decrease in abundance between conditions. Results shown are for all methods in the a) Microbial, b) Bulk transcriptomic, and c) Cell transcriptomic settings.

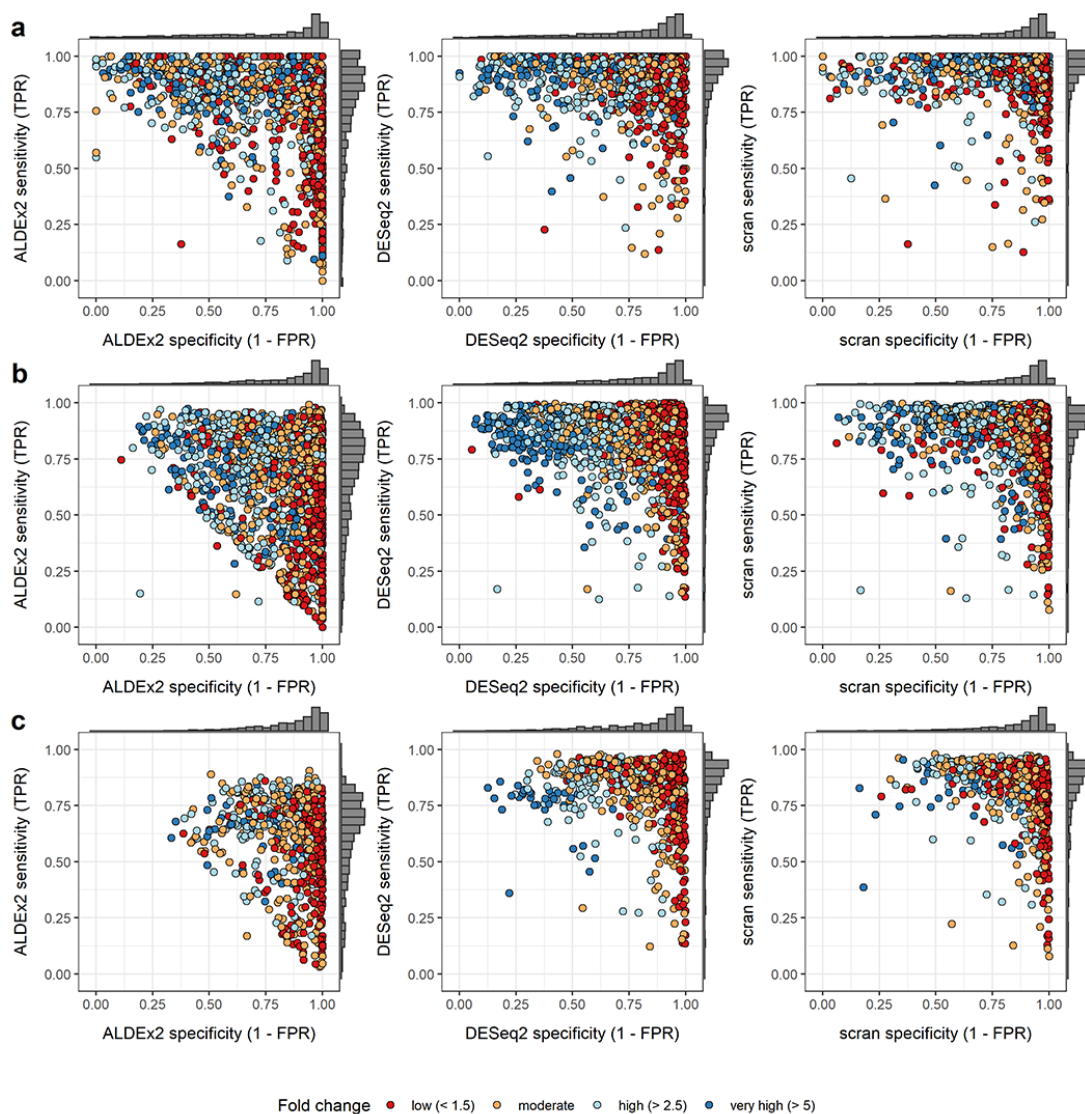


Figure 2: Sensitivity and specificity for three differential abundance calling methods in three compositional complexity settings. Data sets are labeled by absolute mean fold change across conditions. Results shown are for all methods in the **a**) Microbial, **b**) Bulk transcriptomic, and **c**) Cell transcriptomic settings.

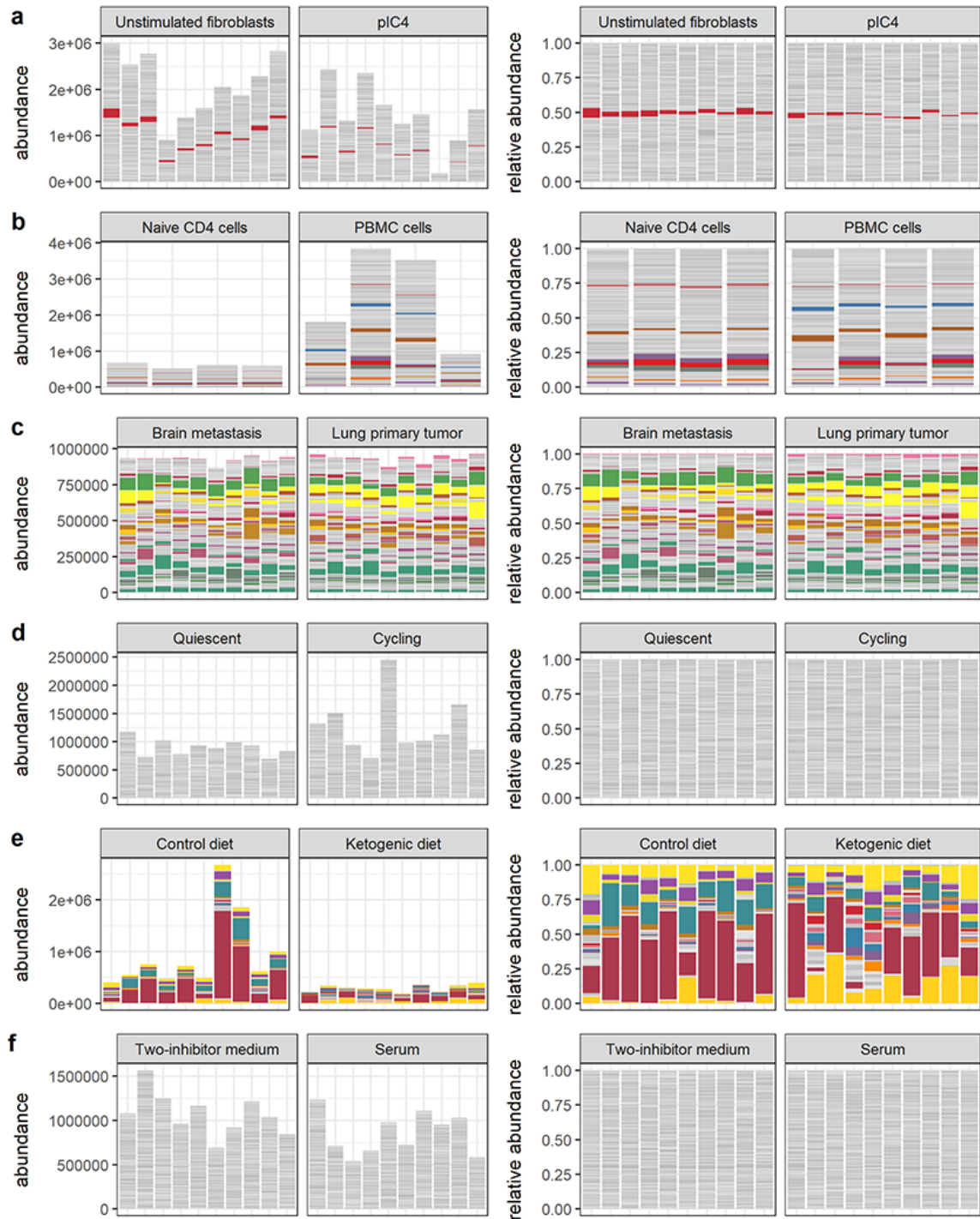


Figure 3: Visual summaries of nominal abundances (left panels) and relative abundances (right panels) for a) Vieira-Silva et al. b) Muraro et al. c) Hagai et al. d) Hashimshony et al. e) Grün et al. f) Kimmerling et al. Features (genes or bacterial sequence variants) with at least 1% relative abundance across all samples are colored; all other features are gray.

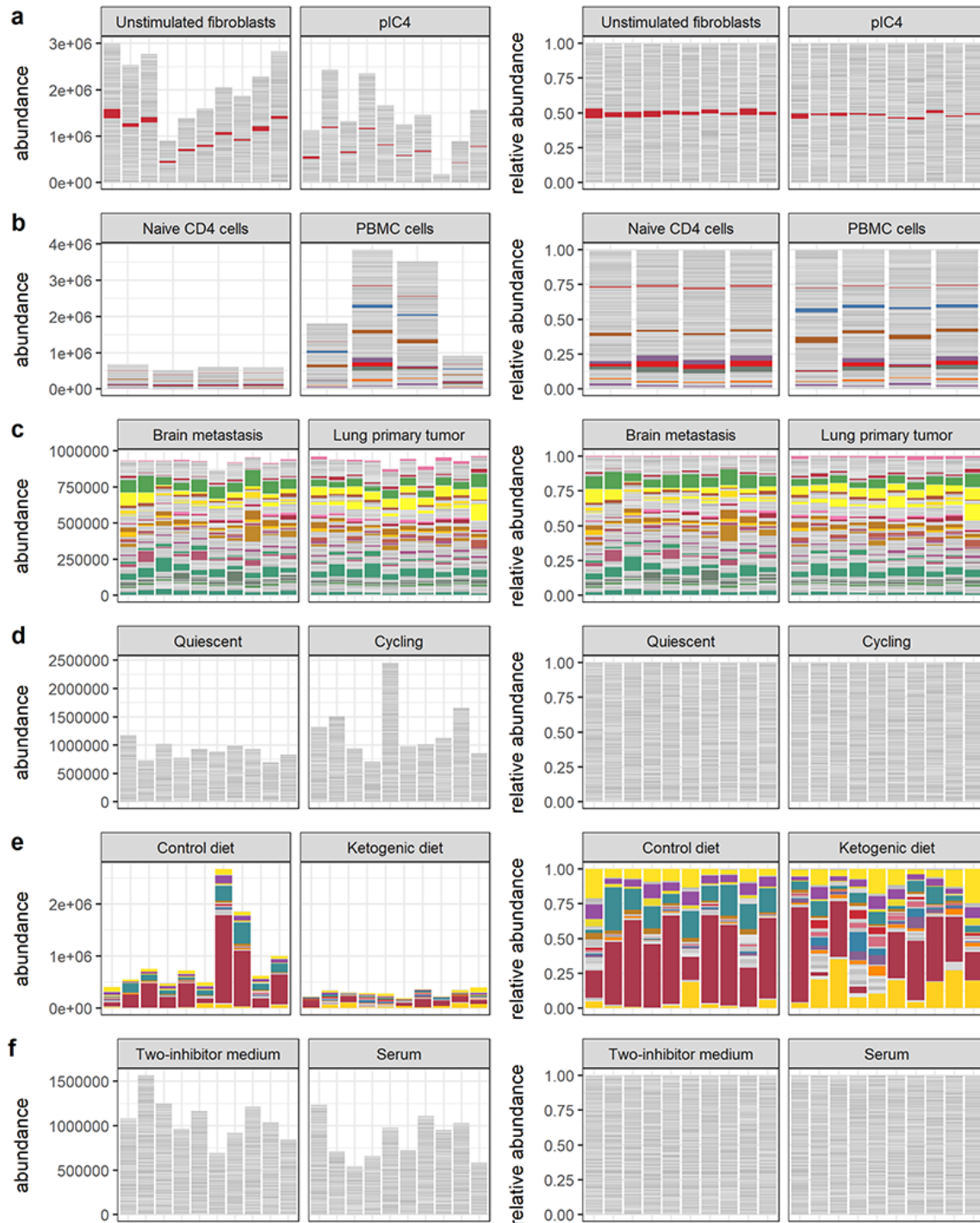


Figure 4: Visual summaries of nominal abundances (left panels) and relative abundances (right panels) for **g)** Song et al. **h)** Barlow et al. **i)** Monaco et al. **j)** Yu et al. **k)** Klein et al. **m)** Owens et al. Features (genes or bacterial sequence variants) with at least 1% relative abundance across all samples are colored; all other features are gray.



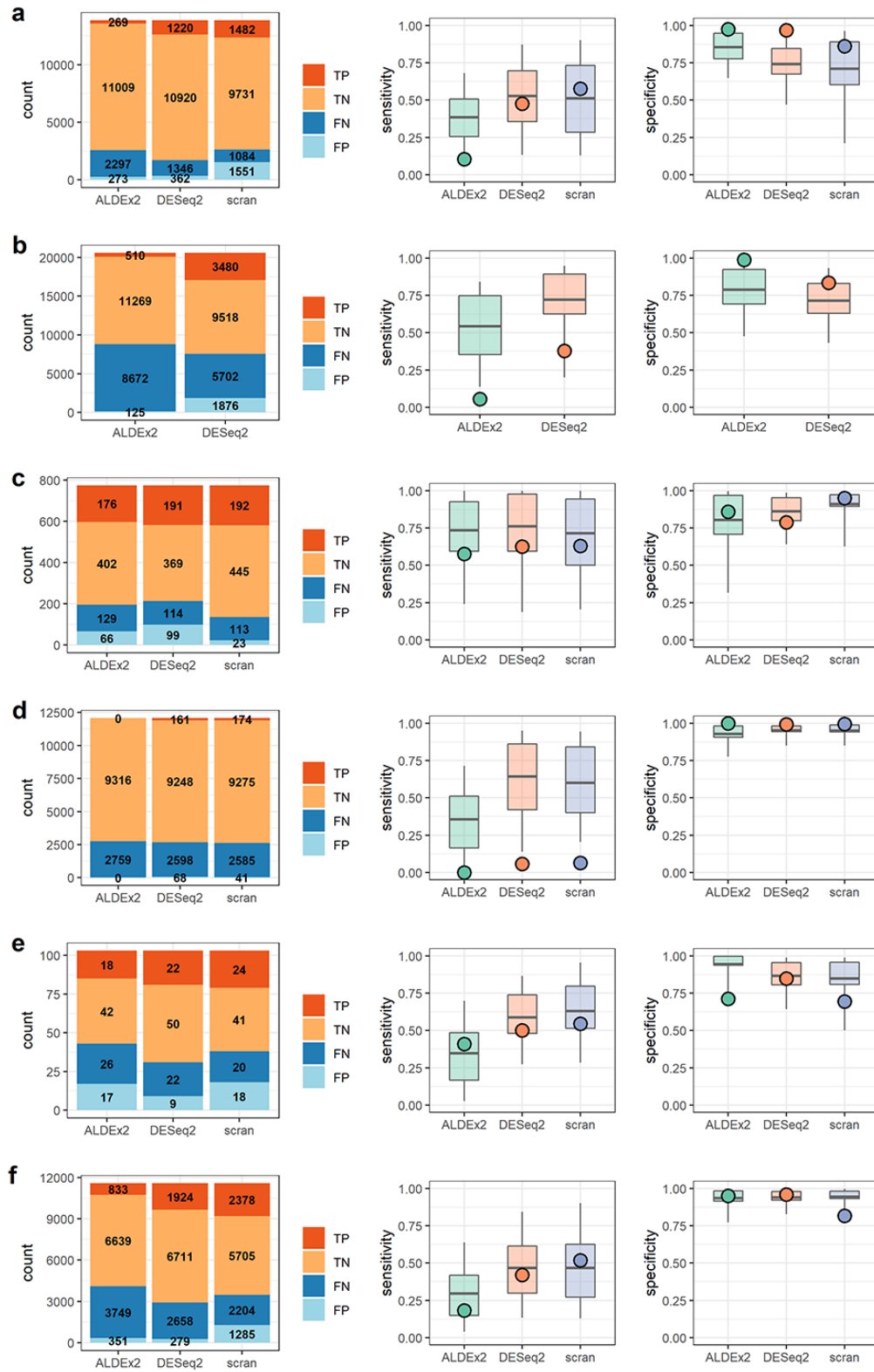


Figure 5: Results summaries for a) Vieira-Silva et al. b) Muraro et al. c) Hagai et al. d) Hashimshony et al. e) Grün et al. f) Kimmerling et al.

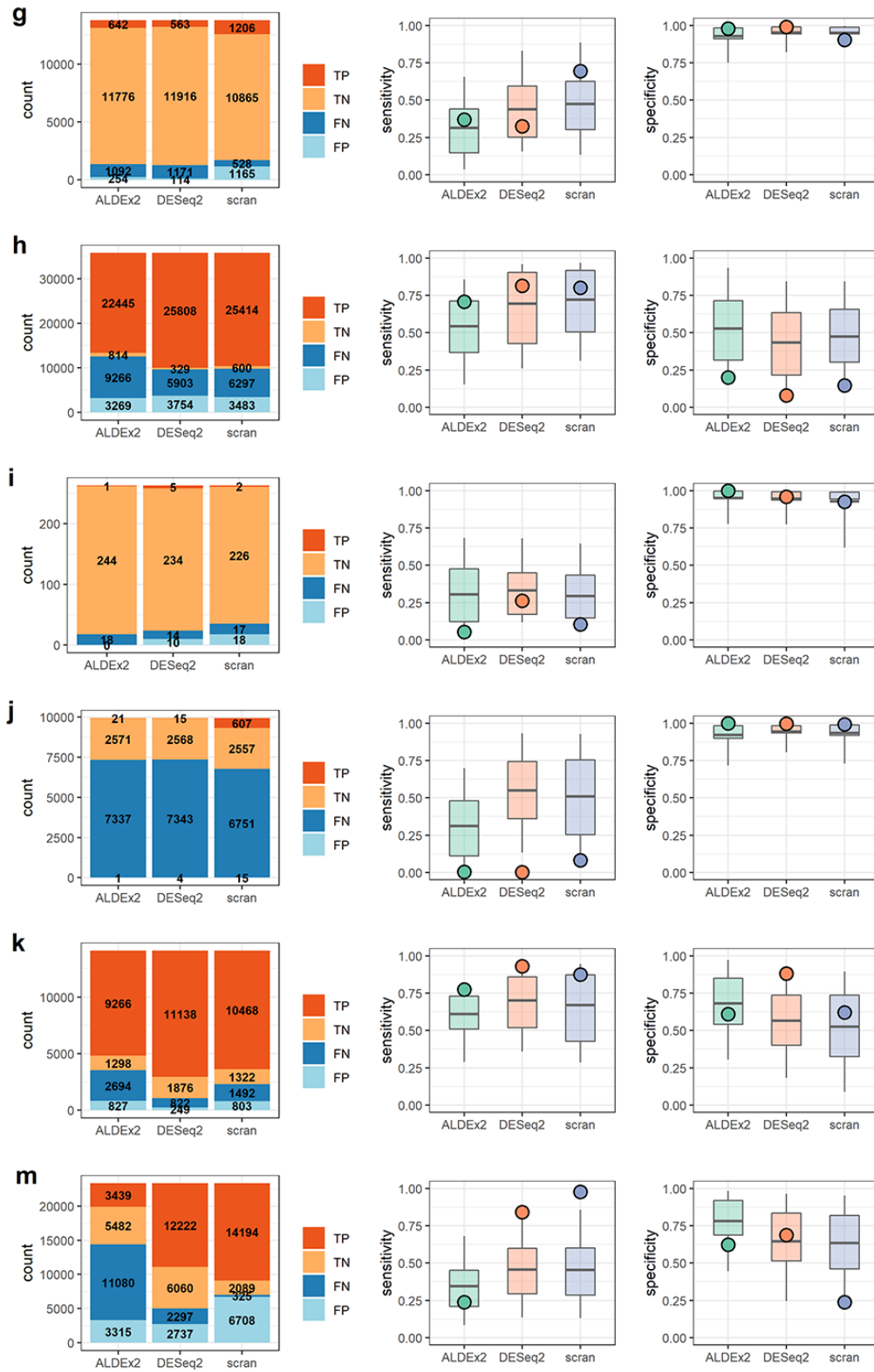


Figure 6: Results summaries for g) Song et al. h) Barlow et al. i) Monaco et al. j) Yu et al. k) Klein et al. m) Owens et al.

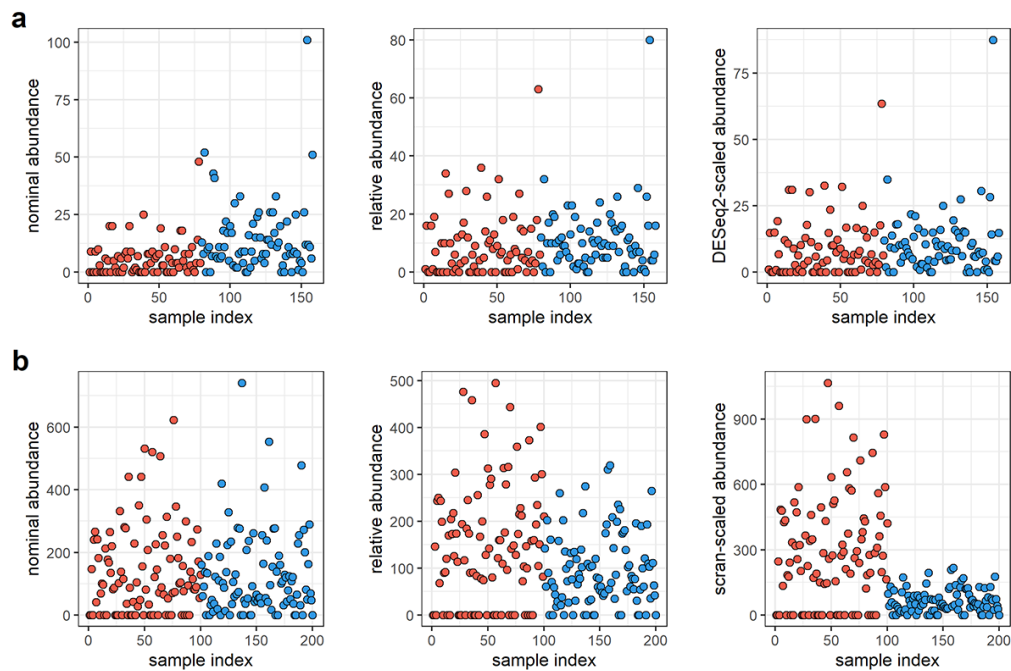


Figure 7: Discrepant calls in nominal and relative abundances. **a)** A typical false negative result in data derived from the experiment of Kimmerling et al. This feature is significantly differentially abundant in the nominal abundances but not in relative abundances or relative abundances rescaled by DESeq2's size factor. **b)** A typical false positive result in data from Klein et al.