

Clipper: p-value-free FDR control on high-throughput data from two conditions

Xin Zhou Ge^{1,†}, Yiling Elaine Chen^{1,†}, Dongyuan Song², MeiLu McDermott^{3,4}, Kyla Woyshner³, Antigoni Manousopoulou³, Ning Wang², Wei Li⁵, Leo D. Wang³, and Jingyi Jessica Li^{1,2,6,7,8,*}

¹Department of Statistics, University of California, Los Angeles, CA 90095

²Interdepartmental Program in Bioinformatics, University of California, Los Angeles, CA 90095

³Beckman Research Institute, City of Hope National Medical Center, Duarte, CA 91010

⁴The Quantitative and Computational Biology section, University of Southern California, Los Angeles, CA 90089

⁵Division of Computational Biomedicine, Department of Biological Chemistry, School of Medicine, University of California, Irvine, CA 92697

⁶Department of Human Genetics, University of California, Los Angeles, CA 90095

⁷Department of Computational Medicine, University of California, Los Angeles, CA 90095

⁸Department of Biostatistics, University of California, Los Angeles, CA 90095

[†]These authors contributed equally to this work.

*To whom correspondence should be addressed. Email: jli@stat.ucla.edu

Abstract

High-throughput biological data analysis commonly involves identifying “interesting” features (e.g., genes, genomic regions, and proteins), whose values differ between two conditions, from numerous features measured simultaneously. The most widely-used criterion to ensure the analysis reliability is the false discovery rate (FDR), the expected proportion of uninteresting features among the identified ones. Existing bioinformatics tools primarily control the FDR based on p-values. However, obtaining valid p-values relies on either reasonable assumptions of data distribution or large numbers of replicates under both conditions, two requirements that are often unmet in biological studies. To address this issue, we propose Clipper, a general statistical framework for FDR control without relying on p-values or specific data distributions. Clipper is applicable to identifying both enriched and differential features from high-throughput biological data of diverse types. In comprehensive simulation and real-data benchmarking, Clipper outperforms existing generic FDR control methods and specific bioinformatics tools designed for various tasks, including peak calling from ChIP-seq data, differentially expressed gene identification from bulk or single-cell RNA-seq data, differentially interacting chromatin region identification from Hi-C data, and peptide identification from mass spectrometry data. Notably, our benchmarking results for peptide identification are based on the first mass spectrometry data standard with a realistic dynamic range. Our results demonstrate Clipper’s flexibility and reliability for FDR control, as well as its broad applications in high-throughput data analysis.

Significance Statement

The reproducibility crisis has been increasingly alarming in biomedical research, which often involves high-throughput data analysis to identify targets for downstream experimental validation. False discovery rate (FDR) is the state-of-the-art criterion to guard reproducibility in such biological data analysis. Existing bioinformatics tools control the FDR using p-values, which are usually ill-posed, leading to failed FDR control or poor power. Clipper is a flexible, powerful FDR-control framework that removes the need for high-resolution, well-calibrated p-values. Applicable to various bioinformatics analyses,

Clipper outperforms popular bioinformatics tools, including identifying peaks from ChIP-seq data, differentially expressed genes from bulk or single-cell RNA-seq data, and differentially interacting chromatin regions from Hi-C data. Clipper is a significant computational advance to addressing the reproducibility crisis in biomedical research.

Introduction

High-throughput technologies are widely used to measure system-wide biological features, such as genes, genomic regions, and proteins (“high-throughput” means the number of features is large, at least in thousands). The most common goal of analyzing high-throughput data is to contrast two conditions so as to reliably screen “interesting features,” where “interesting” means “enriched” or “differential.” “Enriched features” are defined to have higher expected measurements (without measurement errors) under the experimental (i.e., treatment) condition than the background (i.e., the negative control) condition. The detection of enriched features is called “enrichment analysis.” For example, typical enrichment analyses include calling protein-binding sites in a genome from chromatin immunoprecipitation sequencing (ChIP-seq) data [1, 2] and identifying peptides from mass spectrometry (MS) data [3]. In contrast, “differential features” are defined to have different expected measurements between two conditions, and their detection is called “differential analysis.” For example, popular differential analyses include the identification of differentially expressed genes (DEGs) from genome-wide gene expression data (e.g., microarray and RNA sequencing (RNA-seq) data [4–10]) and differentially interacting chromatin regions (DIRs) from Hi-C data [11–13] (Fig. 1a). In most scientific research, the interesting features only constitute a small proportion of all features, and the remaining majority is referred to as “uninteresting features.”

The identified features, also called the “discoveries” from enrichment or differential analysis, are subject to further investigation and validation. Hence, to reduce experimental validation that is often laborious or expensive, researchers demand reliable discoveries that contain few false discoveries. Accordingly, the false discovery rate (FDR) [14] has been developed as a statistical criterion for ensuring discoveries’ reliability. The FDR technically is defined as the expected proportion of uninteresting features among the discoveries under the frequentist statistical paradigm. In parallel, under the Bayesian paradigm, other criteria have been developed, including the Bayesian false discovery rate [15], the local false discovery rate (local fdr) [16], and the local false sign rate [17]. Among all these frequentist and Bayesian criteria, the FDR is the dominant criterion for setting thresholds in biological data analysis [1, 10, 18–24] and is thus the focus of this paper.

FDR control refers to the goal of finding discoveries such that the FDR is under a pre-specified threshold (e.g., 0.05). Existing computational methods for FDR control primarily rely on p-values, one per feature. Among the p-value-based methods, the most classic and popular ones are the Benjamini-Hochberg (BH) procedure [14] and the Storey’s q-value [25]; later development introduced methods that incorporate feature weights [26] or covariates (e.g., independent hypothesis weighting (IHW) [27], adaptive p-value thresholding [28], and Boca and Leek’s FDR regression [29]) to boost the detection power. All these methods set a p-value cutoff based on the pre-specified FDR threshold. However, the calculation of p-values requires either distributional assumptions, which are often questionable, or large numbers of replicates, which are often unachievable in biological studies (see Results). Due to these limitations of p-value-based methods in high-throughput biological data analysis, bioinformatics tools often output ill-posed p-values. This issue is evidenced by serious concerns about the widespread miscalculation and misuse of p-values in the scientific community [30]. As a result, bioinformatics tools using questionable p-values either cannot reliably control the FDR to a target level [23] or lack power

to make discoveries [31]; see Results. Therefore, p-value-free control of FDR is desirable, as it would make data analysis more transparent and thus improve the reproducibility of scientific research.

Although p-value-free FDR control has been implemented in the MACS2 method for ChIP-seq peak calling [1] and the SAM method for microarray DEG identification [32], these two methods are restricted to specific applications and lack theoretical guarantee for FDR control¹. More recently, the Barber-Candès (BC) procedure has been proposed to achieve theoretical FDR control without using p-values [35], and it has been shown to perform comparably to the BH procedure with well-calibrated p-values [36]. The BC procedure is advantageous because it does not require well-calibrated p-values, so it holds tremendous potential in various high-throughput data analyses where p-value calibration is challenging [37]. For example, a recent paper has implemented a generalization of the BC procedure to control the FDR in peptide identification from MS data [38].

Inspired by the BC procedure, we propose a general statistical framework Clipper to provide reliable FDR control for high-throughput biological data analysis, without using p-values or relying on specific data distributions. Clipper is a robust and flexible framework that applies to both enrichment and differential analyses and that works for high-throughput data with various characteristics, including data distributions, replicate numbers (from one to multiple), and outlier existence.

Results

Clipper consists of two main steps: construction and thresholding of contrast scores. First, Clipper defines a contrast score for each feature, as a replacement of a p-value, to summarize that feature's measurements between two conditions and to describe the degree of interestingness of that feature. Second, as its name suggests, Clipper establishes a cutoff on features' contrast scores and calls as discoveries the features whose contrast scores exceed the cutoff (see Online Methods and Supplementary). Clipper is a flexible framework that only requires a minimal input: all features' measurements under two conditions and a target FDR threshold (e.g., 5%) (Fig. 1b).

Clipper only relies on two fundamental statistical assumptions of biological data analysis: (1) measurement errors (i.e., differences between measurements and their expectations, with the expectations including both biological signals and batch effects) are independent across all features and replicates; (2) every uninteresting feature has measurement errors identically distributed across all replicates under both conditions. These two assumptions are used in almost all bioinformatics tools and are commonly referred to as the "measurement model" in statistical genomics [39]. With these two assumptions, Clipper has a theoretical guarantee for FDR control under both enrichment and differential analyses with any number of replicates (see Online Methods and Supp. Section S2).

To verify Clipper's performance, we designed comprehensive simulation studies to benchmark Clipper against existing generic FDR control methods (Supp. Section S1). We also benchmarked Clipper against bioinformatics tools in studies including peak calling from ChIP-seq data, peptide identification from mass spectrometry data, DEG identification from bulk or single-cell RNA-seq data, and DIR identification from Hi-C data. Notably, our benchmarking results for peptide identification are based on the first MS data standard with a realistic dynamic range.

Clipper has verified FDR control and power advantage in simulation

Simulation is essential because we can generate numerous datasets from the same distribution with known truths to calculate the FDR, which is not observable from real data. Our simulation covers both

¹ Although later works have studied some theoretical properties of SAM, they are not about the exact control of the FDR [33, 34].

enrichment and differential analyses. In enrichment analysis, we consider four “experimental designs”: 1vs1 design (one replicate per condition), 2vs1 design (two and one replicates under the experimental and background conditions, respectively), 3vs3 design (three replicates per condition), and 10vs10 design (ten replicates per condition). In differential analysis, since Clipper requires that at least one condition has two replicates, we only consider the 2vs1 and 3vs3 designs. For each analysis and design, we simulated data from three “distributional families”—Gaussian, Poisson, and negative binomial—for individual features under two “background scenarios” (i.e., scenarios of the background condition): homogeneous and heterogeneous. Under the homogeneous scenario, all features’ measurements follow the same distribution under the background condition; otherwise, we are under the heterogeneous scenario, which is ubiquitous in applications, e.g., identifying DEGs from RNA-seq data and calling protein-binding sites from ChIP-seq data. By simulation setting, we refer to a combination of an experimental design, a distributional family, and a background scenario. The details of simulation settings are described in Supp. Section S4.

For both enrichment and differential analyses and each simulation setting, we compared Clipper against generic FDR control methods, including p-value-based methods and local-fdr-based methods. The p-value-based methods include BH-pair, BH-pool, qvalue-pair, and qvalue-pool, where “BH” and “qvalue” stand for p-value thresholding procedures, and “pair” and “pool” represent the paired and pooled p-value calculation approaches, respectively. The local-fdr-based methods include locfdr-emp and locfdr-swap, where “emp” and “swap” represent the empirical null and swapping null local-fdr calculation approaches, respectively. See Online Methods for detail.

The comparison results are in Fig. 2 and Supp. Figs. S1–S11. A good FDR control method should have actual FDR no larger than the target FDR threshold and achieve high power. The results show that Clipper controls the FDR and is overall more powerful than other methods, excluding those that fail to control the FDR, under all settings. Clipper is also shown to be more robust to the number of features and the existence of outliers than other methods. In detail, in both enrichment analyses (1vs1, 2vs1, 3vs3, and 10vs10 designs) and differential analyses (2vs1 and 3vs3 designs), Clipper consistently controls the FDR, and it is more powerful than the generic methods in most cases under the realistic, heterogeneous background, where features do not follow the same distribution under the background condition. Under the idealistic, homogeneous background, Clipper is still powerful and only second to BH-pool and qvalue-pool, which, however, cannot control the FDR under the heterogeneous background.

Here we summarize the performance of the generic FDR control methods. First, the two p-value-based methods using the pooled approach, BH-pool and qvalue-pool, are the most powerful under the idealistic, homogeneous background, which is their inherent assumption; however, they cannot control the FDR under the heterogeneous background (Fig. 2b). Besides, they cannot control the FDR when the number of features is small (Fig. 2a and Supp. Fig. S1). These results show that the validity of BH-pool and qvalue-pool requires a large number of features and the homogeneous background assumption, two requirements that rarely hold in biological applications.

Second, the four p-value-based methods using the paired approach with misspecified models or misspecified tests (BH-pair-mis, qvalue-pair-mis, BH-pair-2as1, and qvalue-pair-2as1; see Online Methods) fail to control the FDR by a large margin in most cases, and rarely when they control the FDR, they lack power (Fig. 2c–d and Supp. Figs. S1–S8). These results confirm that the BH-pair and qvalue-pair rely on the correct model specification to control the FDR; however, the correct model specification is hardly achievable with no more than three replicates per condition.

Third, even when models are correctly specified (an idealistic scenario), the p-value-based methods that use the paired approach—BH-pair-correct and qvalue-pair-correct (see Online Methods)—fail to control the FDR in the existence of outliers (Fig. 2e and Supp. Figs. S3 and S7) or for the negative

binomial distribution with unknown dispersion (Fig. 2f and Supp. Fig. S9). It is worth noting that even when they control the FDR, they are less powerful than Clipper in most cases except for the 3vs3 differential analysis with the Poisson distribution (Fig. 2d and Supp. Figs. S4 and S8).

Fourth, the two local-fdr-based methods—locfdr-emp and locfdr-swap—achieve the FDR control under all designs and analyses; however, they are less powerful than Clipper in most cases (Supp. Figs. S1–S4).

Fifth, when the numbers of replicates are large (10vs10 design), non-parametric tests become applicable. We compared Clipper with three BH-pair methods that use different statistical tests: BH-pair-Wilcoxon (the non-parametric Wilcoxon rank-sum test), BH-pair-permutation (the non-parametric permutation test), and BH-pair-parametric (the parametric test based on the correct model specification, equivalent to BH-pair-correct). Although all the three methods control the FDR, they are less powerful than Clipper (Supp. Fig. S10).

Moreover, the above five phenomena are consistently observed across the three distributions (Gaussian, Poisson, and negative binomial) that we have examined, further confirming the robustness of Clipper.

In addition, for the 3vs3 enrichment analysis, we also varied the proportion of interesting features as 10%, 20%, and 40%. The comparison results in Supp. Fig. S3 (columns 1 and 3 for 10%) and Supp. Fig. S12 (for 20% and 40%) show that the performance of Clipper is robust to the proportion of interesting features.

The above results are all based on simulations with independent features. To examine the robustness of Clipper, we introduced feature correlations to our simulated data, on which we compared Clipper with other generic FDR control methods. The comparison results in Supp. Fig. S11 show that even when the feature independence assumption is violated, Clipper still demonstrates strong performance in both FDR control and power.

Clipper has broad applications in omics data analyses

We then demonstrate the use of Clipper in four omics data applications: peak calling from ChIP-seq data, peptide identification from MS data, DEG identification from bulk or single-cell RNA-seq data, and DIR identification from Hi-C data. The first two applications are enrichment analyses, and the last two are differential analyses. In each application, we compared Clipper with mainstream bioinformatics methods to demonstrate Clipper's superiority in FDR control and detection power.

Peak calling from ChIP-seq data (enrichment analysis I)

ChIP-seq is a genome-wide experimental assay for measuring binding intensities of a DNA-associated protein [40], often a transcription factor that activates or represses gene expression [41, 42]. ChIP-seq data are crucial for studying gene expression regulation, and the indispensable analysis is to identify genomic regions with enriched sequence reads in ChIP-seq data. These regions are likely to be bound by the target protein and thus of biological interest. The identification of these regions is termed “peak calling” in ChIP-seq data analysis.

As the identified peaks are subject to experimental validation that is often expensive [43], it is essential to control the FDR of peak identification to reduce unnecessary costs. The two most highly-cited peak-calling methods are MACS2 [1] and [2], both of which claim to control the FDR for their identified peaks. Specifically, both MACS2 and HOMER assume that the read counts for each putative peak (one count per sample/replicate) follow the Poisson distribution, and they use modified paired approaches to assign each putative peak a p-value and a corresponding Storey's q-value. Then given a target FDR

threshold $0 < q < 1$, they call the putative peaks with q-values $\leq q$ as identified peaks. Despite being popular, MACS2 and HOMER have not been verified for their FDR control, to our knowledge.

To verify the FDR control of MACS2 and HOMER (Supp. Section S5.1), we used ENCODE ChIP-seq data of cell line GM12878 [44] and ChiPulate [45], a ChIP-seq data simulator, to generate semi-synthetic data with spiked-in peaks (Supp. Section S6.1). We examined the actual FDR and power of MACS2 and HOMER in a range of target FDR thresholds: $q = 1\%, 2\%, \dots, 10\%$. Fig. 4a shows that MACS2 and HOMER cannot control the FDR as standalone peak-calling methods. However, with Clipper as an add-on (Supp. Section S7.1), both MACS2 and HOMER can guarantee the FDR control. This result demonstrates the flexibility and usability of Clipper for reducing false discoveries in peak calling analysis.

Technically, the failed FDR control by MACS2 and HOMER is attributable to the likely model misspecification and test misformulation in their use of the paired approach. Both MACS2 and HOMER assume the Poisson distribution for read counts in a putative peak; however, it has been widely acknowledged that read counts are over-dispersed and thus better modeled by the negative binomial distribution [46]. Besides, MACS2 uses one-sample tests to compute p-values when two-sample tests should have been performed. As a result, the p-values of MACS2 and HOMER are questionable, so using their p-values for FDR control would not lead to success. (Note that MACS2 does not use p-values to control the FDR but instead swaps experimental and background samples to calculate the empirical FDR; yet, we emphasize that controlling the empirical FDR does not guarantee the FDR control.) As a remedy, Clipper strengthens both methods to control the FDR while maintaining high power.

It is known that uninteresting regions tend to have larger read counts in the control sample than in the experimental (ChIP) sample, making them more likely to have negative contrast scores than positive ones. However, this phenomenon does not violate Clipper's theoretical assumption (Lemma 1(a) in Supp. Section S2), which can be relaxed as we note in Methods.

Peptide identification from MS data (enrichment analysis II)

The state-of-the-art proteomics studies use MS experiments and database search algorithms to identify and quantify proteins in biological samples. In a typical proteomics experiment, a protein mixture sample is first digested into peptides and then measured by tandem MS technology as mass spectra, which encode peptide sequence information. "Peptide identification" is the process that decodes mass spectra and converts mass spectra into peptide sequences in a protein sequence database via search algorithms. The search process matches each mass spectrum to peptide sequences in the database and outputs the best match, called a "peptide-spectrum match" (PSM). The identified PSMs are used to infer and quantify proteins in a high-throughput manner.

False PSMs could occur when mass spectra are matched to wrong peptide sequences due to issues such as low-quality spectra, data-processing errors, and incomplete protein databases, causing problems in the downstream protein identification and quantification [47]. Therefore, a common goal of database search algorithms is to simultaneously control the FDR and maximize the number of identified PSMs, so as to maximize the number of proteins identified in a proteomics study [3, 48]. A widely used FDR control strategy is the target-decoy search, where mass spectra of interest are matched to peptide sequences in both the original (target) database and a decoy database that contains artificial false protein sequences. The resulting PSMs are called the target PSMs and decoy PSMs, respectively. The decoy PSMs, i.e., matched mass spectrum and decoy peptide pairs, are known to be false and thus used by database search algorithms to control the FDR. Mainstream database search algorithms output a q-value for each PSM, target or decoy. Discoveries are the target PSMs whose q-values are no greater than the target FDR threshold q .

We used the first comprehensive benchmark dataset from an archaea species *Pyrococcus furiosus*

to examine the FDR control and power of a popular database search algorithm SEQUEST [3] (Supp. Section S5.2). Using this benchmark dataset (Supp. Section S6.2), we demonstrate that, as an add-on, Clipper improves the power of SEQUEST. Specifically, Clipper treats mass spectra as features. For each mass spectrum, Clipper considers its measurement under the experimental condition as the $-\log_{10}$ -transformed q-value of the target PSM that includes it, and its measurement under the background condition as the $-\log_{10}$ -transformed q-value of the decoy PSM that includes it. Then Clipper decides which mass spectra and their corresponding target PSMs are discoveries (Supp. Section S7.2). Based on the benchmark dataset, we examined the empirical FDR, i.e., the FDP calculated based on the true positives and negatives, and the power of SEQUEST with or without Clipper as an add-on, for a range of target FDR thresholds: $q = 1\%, 2\%, \dots, 10\%$. Fig. 4b shows that although SEQUEST and SEQUEST+Clipper both control the FDR, SEQUEST+Clipper consistently improves the power, thus enhancing the peptide identification efficiency of proteomics experiments.

While preparing this manuscript, we found a recent work [38] that used a similar idea to identify PSMs without using p-values. Clipper differs from this work in two aspects: (1) Clipper is directly applicable as an add-on to any existing database search algorithms that output q-values; (2) Clipper is not restricted to the peptide identification application.

DEG identification from bulk RNA-seq data (differential analysis I)

RNA-seq data measure genome-wide gene expression levels in biological samples. An important use of RNA-seq data is the DEG analysis, which aims to discover genes whose expression levels change between two conditions. The FDR is a widely used criterion in DEG analysis [4–9].

We compared Clipper with two popular DEG identification methods: edgeR [4] and DESeq2 [5] (Supp. Section S5.3). Specifically, when we implemented Clipper, we first performed the trimmed mean of M values (TMM) normalization [49] to correct for batch effects; then we treated genes as features and their normalized expression levels as measurements under two conditions (Supp. Section S7.3). We also implemented two versions of DESeq2 and edgeR: with or without IHW, a popular procedure for boosting the power of p-value-based FDR control methods by incorporating feature covariates [27]. In our implementation of the two versions of DESeq2 and edgeR, we used their standard pipelines, including normalization, model fitting, and gene filtering (edgeR only). To verify the FDR control, we generated four realistic synthetic datasets from two real RNA-seq datasets—one from classical and non-classical human monocytes [50] and the other from yeasts with or without *snf2* knockout [51]—using simulation strategies 1 and 2 (Supp. Section S6.3).

In detail, in simulation strategy 1, we used bulk RNA-seq samples from two conditions to compute a fold change for every gene between the two conditions; then we defined true DEGs as the genes whose fold changes exceeded a threshold; next, we randomly drew three RNA-seq samples and treated them as replicates from each condition ($m = n = 3$ as in Methods); using those subsampled replicates of two conditions, we preserved the true DEGs' read counts and permuted the read counts of the true non-DEGs, i.e., the genes other than true DEGs, between conditions. In summary, simulation strategy 1 guarantees that the measurements of true non-DEGs are i.i.d., an assumption that Clipper relies on for theoretical FDR control.

In simulation strategy 2, borrowed from a benchmark study [52], we first randomly selected at most 30% genes as true DEGs; next, we randomly drew six RNA-seq samples from one condition (classical human monocytes and yeasts without knockout) and split the samples into two “synthetic conditions,” each with three replicates ($m = n = 3$ as in Methods); then for each true DEG, we multiplied its read counts under one of the two synthetic conditions (randomly picked independently for each gene) by a randomly generated fold change (see Supp. Section S6.3); finally, for the true non-DEGs, we preserved

their read counts in the six samples. In summary, simulation strategy 2 preserves batch effects, if existent in real data, for the true non-DEGs (the majority of genes). As a result, the semi-synthetic data generated under strategy 2 may violate the Clipper assumption for theoretical FDR control and thus can help evaluate the robustness of Clipper on real data.

The four semi-synthetic datasets have ground truths (true DEGs and non-DEGs) to evaluate each DEG identification method's FDR and power for a range of target FDR thresholds: $q = 1\%, 2\%, \dots, 10\%$. Our results in Fig. 3a and Supp. Figs. S15a–S17a show that Clipper consistently controls the FDR and achieves high power on all four semi-synthetic datasets. In contrast, DESeq2 and edgeR cannot consistently control the FDR except for the yeast semi-synthetic dataset generated under simulation strategy 2. Given the fact that DESeq2 and edgeR do not consistently perform well on the three other semi-synthetic datasets, we hypothesize that their parametric distributional assumptions, if violated on real data, hinder valid FDR control, in line with our motivation for developing Clipper. By examining whether true non-DEGs' p-values calculated by DESeq2 or edgeR follow the theoretical Uniform[0, 1] distribution, we find that the answer is no for many non-DEGs, as indicated by the small p-values (one per non-DEG) of uniformity tests (Supp. Fig. S30); this issue is more serious for DESeq2, consistent with the worse FDR control of DESeq2 (Fig. 3a and Supp. Figs. S15a–S17a). Furthermore, we observe that adding IHW to edgeR and DESeq2 has negligible effects on the four semi-synthetic datasets.

To further explain why DESeq2 fails to control the FDR, we examined the p-value distributions of 16 non-DEGs that were most frequently identified (from the 100 semi-synthetic datasets generated from the human monocyte dataset using simulation strategy 1) by DESeq2 at the target FDR threshold $q = 0.05$. Our results in Supp. Fig. S18 show that the 16 non-DEGs' p-values are non-uniformly distributed with a mode close to 0. Such unusual enrichment of overly small p-values makes these non-DEGs mistakenly called discoveries by DESeq2.

In addition, we compared the DEG ranking by Clipper, edgeR, and DESeq2 in two ways. First, for true DEGs, we compared their ranking by each method with their true ranking based on true expression fold changes (from large to small, as in semi-synthetic data generation in Supp. Section S6.3). Specifically, we ranked true DEGs using Clipper's contrast scores (from large to small), edgeR's p-values (from small to large), or DESeq2's p-values (from small to large). Our results in Fig. 3b and Supp. Figs. S15b–S17b show that Clipper's contrast scores exhibit the most consistent ranking with the ranking based on true fold changes. Second, to compare the power of Clipper, edgeR, and DESeq2 based on their DEG rankings instead of nominal p-values, we calculated their power under the actual FDRs, which only depend on gene rankings (for the definition of actual FDR, see Supp. Section S6.3). Fig. 3a and Supp. Figs. S15a–S17a show that, when Clipper, edgeR, and DESeq2 have the same actual FDR, Clipper consistently outperforms edgeR and DESeq2 in terms of power, i.e., Clipper has the most true DEGs in its top ranked genes.

We also compared the reproducibility of Clipper, edgeR, and DESeq2 in the presence of sampling randomness. Specifically, we used two semi-synthetic datasets (generated independently from the same procedure in Supp. Section S6.3) as technical replicates and computed Clipper's contrast scores and edgeR's and DESeq2's p-values on each dataset. For each method, we evaluated its reproducibility between the two semi-synthetic datasets by computing three criteria—the irreproducibility discovery rate (IDR) [53], Pearson correlation, and Spearman correlation—using its contrast scores or negative \log_{10} transformed p-values. Fig. 3c and Supp. Figs. S15–S17c show that Clipper's contrast scores have higher reproducibility by all three criteria compared to edgeR's and DESeq2's p-values.

Finally, we compared Clipper with DESeq2 and edgeR on the real RNA-seq data of classical and non-classical human monocytes [50]. In this dataset, gene expression changes are expected to be associated with the immune response process. We input three classical and three non-classical samples into Clipper, DESeq2, and edgeR for DEG identification. Fig. 5a shows that edgeR identifies the fewest

DEGs, while DESeq2 identifies the most DEGs, followed by Clipper. Notably, most DEGs identified by DESeq2 are not identified by Clipper or edgeR. To investigate whether DESeq2 makes too many false discoveries and whether the DEGs found by Clipper but missed by DESeq2 or edgeR are biologically meaningful, we performed functional analysis on the set of DEGs identified by each method. We first performed the gene ontology (GO) analysis on the three sets of identified DEGs using the R package `clusterProfiler` [54]. Fig. 5b (“Total”) shows that more GO terms are enriched (with enrichment q -values ≤ 0.01) in the DEGs identified by Clipper than in the DEGs identified by DESeq2 or edgeR. For the GO terms enriched in all three sets of identified DEGs, Fig. 5c shows that they are all related to the immune response and thus biologically meaningful. Notably, these biologically meaningful GO terms have more significant enrichment in the DEGs identified by Clipper than in those identified by edgeR and DESeq2. We further performed GO analysis on the DEGs uniquely identified by one method in pairwise comparisons of Clipper vs. DESeq2 and Clipper vs. edgeR. Fig. 5b and Supp. Fig. S20 show that multiple immune-related GO terms are enriched in Clipper-specific DEGs, while no GO terms are enriched in edgeR-specific or DESeq2-specific DEGs. In addition, we examined the DEGs that were identified by Clipper only but missed by both edgeR and DESeq2. Fig. 5d and Supplementary Table show that these genes include multiple key immune-related genes, including *CD36*, *DUSP2*, and *TNFAIP3*. We further performed pathway analysis on these genes and the DEGs that were identified by DESeq2 only but missed by both edgeR and Clipper, using the R package `limma` [10]. Supp. Fig. S21a shows that the DEGs that were only identified by Clipper have significant enrichment for immune-related pathways including phagosome, a key function of monocytes and macrophages. On the contrary, Supp. Fig. S21b shows that fewer immune-related pathways are enriched in DEGs that were only identified by DESeq2. Altogether, these results confirm the capacity of Clipper in real-data DEG analysis, and they are consistent with our simulation results that edgeR lacks power, while DESeq2 fails to control the FDR.

DEG identification from single-cell RNA-seq data (differential analysis II)

Single-cell RNA sequencing (scRNA-seq) technologies have revolutionized biomedical sciences by enabling genome-wide profiling of gene expression levels at an unprecedented single-cell resolution. DEG analysis is widely applied to scRNA-seq data for discovering genes whose expression levels change between two conditions or between two cell types. Compared with bulk RNA-seq data, scRNA-seq data have many more “replicates” (i.e., cells, whose number is often in hundreds) under each condition or within each cell type.

We compared Clipper with edgeR [4], MAST [55], Monocle3 [56], the two-sample t test, and the Wilcoxon rank-sum test, five methods that are either popular or reported to have comparatively top performance from a previous benchmark study [57]. To verify the FDR control, we used scDesign2, a flexible probabilistic simulator to generate scRNA-seq count data with known true DEGs [58]. scDesign2 offers three key advantages that enable the generation of realistic synthetic scRNA-seq count data: (1) it captures distinct marginal distributions of different genes; (2) it preserves gene-gene correlations; (3) it adapts to various scRNA-seq protocols. Using scDesign2, we generated two synthetic scRNA-seq datasets from two real scRNA-seq datasets of peripheral blood mononuclear cells (PBMCs) [59]: one using 10x Genomics [60] and the other using Drop-seq [61]. Each synthetic dataset contains two cell types, CD4⁺ T cells and cytotoxic T cells, which we treated as two conditions. Having true DEGs known, the synthetic datasets allow us to evaluate Clipper’s and the other five methods’ FDRs and power for a range of target FDR thresholds: $q = 1\%, 2\%, \dots, 10\%$. Fig. 4d and Supp. Fig. S19 show that on both 10x Genomics and Drop-seq synthetic datasets, Clipper consistently controls the FDR and remains the most powerful among all the methods that achieve FDR control. These results demonstrate Clipper’s robust performance in scRNA-seq DEG analysis.

DIR analysis of Hi-C data (differential analysis III)

Hi-C experiments are widely used to investigate spatial organizations of chromosomes and map chromatin interactions across the genome. A Hi-C dataset is often processed and summarized into an interaction matrix, whose rows and columns represent manually binned chromosomal regions and whose (i, j) -th entry represents the measured contact intensity between the i -th and j -th binned regions. The DIR analysis aims to identify pairs of genomic regions whose contact intensities differ between conditions. Same as DEG analysis, DIR analysis also uses the FDR as a decision criterion [11–13].

We compared Clipper with three popular DIR identification methods: diffHic [13], FIND [12], and multiHiCcompare [11] (Supp. Section S5.5). Specifically, we applied Clipper to DIR identification by treating pairs of genomic regions as features and contact intensities as measurements. To verify the FDR control of Clipper (Supp. Section S7.5), diffHic, FIND, and multiHiCcompare, we generated realistic semi-synthetic data from real interaction matrices of ENCODE cell line GM12878 [44] with true spiked-in DIRs to evaluate the FDR and power (Supp. Section S6.5). We examined the actual FDR and power in a range of target FDR thresholds: $q = 1\%, 2\%, \dots, 10\%$. Fig. 4d shows that Clipper and diffHic are the only two methods that consistently control the FDR, while multiHiCcompare and FIND fail by a large margin. In terms of power, Clipper outperforms diffHic except for $q = 0.01$ and 0.02 , even though Clipper has not been optimized for Hi-C data analysis. This result demonstrates Clipper's general applicability and strong potential for DIR analysis.

Discussion

In this paper, we proposed a new statistical framework, Clipper, for identifying interesting features with FDR control from high-throughput data. Clipper avoids the use of p-values and makes FDR control more reliable and flexible. We used comprehensive simulation studies to verify the FDR control by Clipper under various settings. We demonstrate that Clipper outperforms existing generic FDR control methods by having higher power and greater robustness to model misspecification. We further applied Clipper to four popular bioinformatics analyses: peak calling from ChIP-seq data, peptide identification from MS data, DEG identification from RNA-seq data, and DIR identification from Hi-C data. Our results indicate that Clipper can provide a powerful add-on to existing bioinformatics tools to improve the reliability of FDR control and thus the reproducibility of scientific discoveries.

Clipper's FDR control procedures (BC and GZ procedures in Methods) are motivated by the Barber-Candès (BC)'s knockoff paper [35] and the Gimenez-Zou's multiple knockoff paper [62], but we do not need to construct knockoffs in enrichment analysis when two conditions have the same number of replicates; the reason is that the replicates under the background condition serve as natural negative controls. For differential analysis and enrichment analysis with unequal numbers of replicates, in order to guarantee the theoretical assumptions for FDR control, Clipper uses permutations instead of the complicated knockoff construction because Clipper only examines features marginally and does not concern about features' joint distribution.

We validated the FDR control by Clipper using extensive and concrete simulations, including both model-based and real-data-based data generation with ground truths, which are widely used to validate newly developed computational frameworks [63]. In contrast, in most bioinformatics method papers, the FDR control was merely mentioned but rarely validated. Many of them assumed that using the BH procedure on p-values would lead to valid FDR control; however, the reality is often otherwise because p-values would be invalid when model assumptions were violated or the p-value calculation was problematic. Here we voice the importance of validating the FDR control in bioinformatics method development, and we use this work as a demonstration. We believe that Clipper provides a powerful

booster to this movement. As a p-value-free alternative to the classic p-value-based BH procedure, Clipper relies less on model assumptions and is thus more robust to model misspecifications, making it an appealing choice for FDR control in diverse high-throughput biomedical data analyses.

Clipper is a flexible framework that is easily generalizable to identify a variety of interesting features. The core component of Clipper summarizes each feature's measurements under each condition into an informative statistic (e.g., the sample mean); then Clipper combines each feature's informative statistics under two conditions into a contrast score to enable FDR control. The current implementation of Clipper only uses the sample mean as the informative statistic to identify the interesting features that have distinct expected values under two conditions. However, by modifying the informative statistic, we can generalize Clipper to identify the features that are interesting in other aspects, e.g., having different variances between two conditions. Regarding the contrast score, Clipper currently makes careful choices between two contrast scores, minus and maximum, based on the number of replicates and the analysis task (enrichment or differential).

Notably, Clipper achieves FDR control and high power using those two simple contrast scores, which are calculated for individual features without borrowing information from other features. However, Clipper does leverage the power of multiple testing by setting a contrast score threshold based on all features' contrast scores. This is a likely reason why Clipper achieves good power even with simple contrast scores. An advantage of Clipper is that it allows other definitions of contrast scores, such as the two-sample t statistic that considers within-condition variances. Empirical evidence (Supp. Figs. S13 and S14) shows that the Clipper variant using the two-sample t statistic is underpowered by the default Clipper, which uses the minus summary statistic (difference of two conditions' sample means) as the contrast score in the 3vs3 enrichment analysis or as the degree of interestingness in the 3vs3 differential analysis (see Methods). Here is our current interpretation of this seemingly counter-intuitive result.

- First, both the minus statistic and the t statistic satisfy Clipper's theoretical conditions (Lemmas 1 and 3 in Supp. Section S2), which guarantee the FDR control by the BC and GZ procedures; this is confirmed in Supp. Figs. S13 and S14. Hence, from the FDR control perspective, Clipper does not require the adjustment for within-condition variances by using a t statistic.
- Second, Clipper is different from the two-sample t test or the regression-based t test, where the t statistic was purposely derived as a pivotal statistic so that its null distribution (the t distribution) does not depend on unknown parameters. Since Clipper does not require a null distribution for each feature, the advantage of the t statistic being pivotal no longer matters.
- Third, the minus statistic only requires estimates of two conditions' mean parameters, while the t statistic additionally requires estimates of the two conditions' variances. Hence, when the sample sizes (i.e., the numbers of replicates) are small, the two more parameters that need estimation in the t statistic might contribute to the observed power loss of the Clipper t statistic variant. Indeed, the power difference between the two statistics diminishes as the sample sizes increase from 3vs3 in Supp. Figs. S13–S14 to 10vs10 in Supp. Figs. S10 (where we compared the default Clipper with BH-pair-parametric, which is based on the two-sample t test and is highly similar to the Clipper t statistic variant).
- Fourth, we observe empirically that a contrast score would have better power if its distribution (based on its values of all features) has a larger range and a heavier right tail (in the positive domain). Compared to the minus statistic, the t statistic has a smaller range and a lighter right tail due to its adjustment for features' within-condition variances (Supp. Fig. S28). This observation is consistent with the power difference of the two statistics.

Beyond our current interpretation, however, we admit that future studies are needed to explore alternative contrast scores and their power with respect to data characteristics and analysis tasks. Furthermore, we may generalize Clipper to be robust against sample batch effects by constructing the contrast score as a regression-based test statistic that has batch effects removed.

Our current version of Clipper allows the identification of interesting features between two conditions. However, there is a growing need to generalize our framework to identify features across more than two conditions. For example, temporal analysis based on scRNA-seq data aims to identify genes whose expression levels change along time [31]. To tailor Clipper for such analysis, we could define a new contrast score that differentiates the genes with stationary expression (uninteresting features) from the other genes with varying expression (interesting features). Further studies are needed to explore the possibility of extending Clipper to the regression framework so that Clipper can accommodate data of multiple conditions or even continuous conditions, as well as adjusting for confounding covariates.

We have demonstrated the broad application potential of Clipper in various bioinformatics data analyses. Specifically, when used as an add-on to established, popular bioinformatics methods such as MACS2 for peak calling and SEQUEST for peptide identification, Clipper guaranteed the desired FDR control and in some cases boosted the power. However, many more careful thoughts are needed to escalate Clipper into standalone bioinformatics methods for specific data analyses, for which data processing and characteristics (e.g., peak lengths, GC contents, proportions of zeros, and batch effects) must be appropriately accounted for before Clipper is used for the FDR control [57, 64]. We expect that the Clipper framework will propel future development of bioinformatics methods by providing a flexible p-value-free approach to control the FDR, thus improving the reliability of scientific discoveries.

After finishing this manuscript, we were informed of the work by He et al. [65], which is highly similar to the part of Clipper for differential analysis, as both work use permutation for generating negative controls and the GZ procedure for thresholding (test statistics in He et al. and contrast scores in Clipper). However, the test statistics used in He et al. are the two-sample t statistic and the two-sample Wilcoxon statistic, both of which are different from the minus and maximum contrast scores used in Clipper. While we leave the optimization of contrast scores to future work, we note that the two-sample Wilcoxon statistic, though being a valid contrast score for differential analysis, requires a large sample size to achieve good power. For this reason, we did not consider it as a contrast score in the current Clipper implementation, whose focus is on sample-sample-size high-throughput biological data.

Methods

Notations and assumptions

We first introduce notations and assumptions used in Clipper. While the differential analysis treats the two conditions symmetric, the enrichment analysis requires one condition to be the experimental condition (i.e., the condition of interest) and the other condition to be the background condition (i.e., the negative control). For simplicity, we use the same set of notations for both analyses. For two random vectors $\mathbf{X} = (X_1, \dots, X_m)^\top$ and $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$, we write $\mathbf{X} \perp \mathbf{Y}$ if X_i is independent of Y_j for all $i = 1, \dots, m$ and $j = 1, \dots, n$. To avoid confusion, we use $\text{card}(A)$ to denote the cardinality of a set A and $|c|$ to denote the absolute value of a scalar c . We define $a \vee b := \max(a, b)$.

Clipper only requires two inputs: the target FDR threshold $q \in (0, 1)$ and the input data. Regarding the input data, we use d to denote the number of features with measurements under two conditions, and we use m and n to denote the numbers of replicates under the two conditions. For each feature $j = 1, \dots, d$, we use $\mathbf{X}_j = (X_{j1}, \dots, X_{jm})^\top \in \mathbb{R}_{\geq 0}^m$ and $\mathbf{Y}_j = (Y_{j1}, \dots, Y_{jn})^\top \in \mathbb{R}_{\geq 0}^n$ to denote its

measurements under the two conditions, where $\mathbb{R}_{\geq 0}$ denotes the set of non-negative real numbers. We assume that all measurements are non-negative, as in the case of most high-throughput experiments. (If this assumption does not hold, transformations can be applied to make data satisfy this assumption.)

Clipper has the following assumptions on the joint distribution of $X_1, \dots, X_d, Y_1, \dots, Y_d$. For $j = 1, \dots, d$, Clipper assumes that X_{j1}, \dots, X_{jm} are identically distributed, so are Y_{j1}, \dots, Y_{jn} . Let $\mu_{Xj} = \mathbb{E}[X_{j1}]$ and $\mu_{Yj} = \mathbb{E}[Y_{j1}]$ denote the expected measurement of feature j under the two conditions, respectively. Then conditioning on $\{\mu_{Xj}\}_{j=1}^d$ and $\{\mu_{Yj}\}_{j=1}^d$,

$$\begin{aligned} X_{j1}, \dots, X_{jm}, Y_{j1}, \dots, Y_{jn} \text{ are mutually independent;} \\ \mathbf{X}_j \perp \mathbf{X}_k, \mathbf{Y}_j \perp \mathbf{Y}_k \text{ and } \mathbf{X}_j \perp \mathbf{Y}_k, \forall j, k = 1, \dots, d. \end{aligned} \quad (1)$$

An enrichment analysis aims to identify interesting features with $\mu_{Xj} > \mu_{Yj}$ (with \mathbf{X}_j and \mathbf{Y}_j defined as the measurements under the experimental and background conditions, respectively), while a differential analysis aims to call interesting features with $\mu_{Xj} \neq \mu_{Yj}$. We define $\mathcal{N} := \{j : \mu_{Xj} = \mu_{Yj}\}$ as the set of uninteresting features and denote $N := \text{card}(\mathcal{N})$. In both analyses, Clipper further assumes that an uninteresting feature j satisfies

$$X_{j1}, \dots, X_{jm}, Y_{j1}, \dots, Y_{jn} \text{ are identically distributed, } \forall j \in \mathcal{N}. \quad (2)$$

Clipper consists of two main steps: construction and thresholding of contrast scores. First, Clipper computes contrast scores, one per feature, as summary statistics that reflect the extent to which features are interesting. Second, Clipper establishes a contrast-score cutoff and calls as discoveries the features whose contrast scores exceed the cutoff.

To construct contrast scores, Clipper uses two summary statistics $t(\cdot, \cdot) : \mathbb{R}_{\geq 0}^m \times \mathbb{R}_{\geq 0}^n \rightarrow \mathbb{R}$ to extract data information regarding whether a feature is interesting or not:

$$t^{\text{minus}}(\mathbf{x}, \mathbf{y}) := \bar{x} - \bar{y}; \quad (3)$$

$$t^{\text{max}}(\mathbf{x}, \mathbf{y}) := \max(\bar{x}, \bar{y}) \cdot \text{sign}(\bar{x} - \bar{y}), \quad (4)$$

where $\mathbf{x} = (x_1, \dots, x_m)^\top \in \mathbb{R}_{\geq 0}^m$, $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}_{\geq 0}^n$, $\bar{x} = \sum_{i=1}^m x_i/m$, $\bar{y} = \sum_{i=1}^n y_i/n$, and $\text{sign}(\cdot) : \mathbb{R} \rightarrow \{-1, 0, 1\}$ with $\text{sign}(x) = 1$ if $x > 0$, $\text{sign}(x) = -1$ if $x < 0$, and $\text{sign}(x) = 0$ otherwise.

Notably, other summary statistics can also be used to construct contrast scores. For example, an alternative summary statistic is the t statistic from the two-sample t test:

$$t^t(\mathbf{x}, \mathbf{y}) := \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sum_{i=1}^m (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - \bar{y})^2}{m+n-2}}}. \quad (5)$$

Then we introduce how Clipper works in three analysis tasks: the enrichment analysis with equal numbers of replicates under two conditions ($m = n$), the enrichment analysis with different numbers of replicates under two conditions ($m \neq n$), and the differential analysis (when $m + n > 2$).

Enrichment analysis with equal numbers of replicates ($m = n$)

Under the enrichment analysis, we assume that $\mathbf{X}_j \in \mathbb{R}_{\geq 0}^m$ and $\mathbf{Y}_j \in \mathbb{R}_{\geq 0}^n$ are the measurements of feature j , $j = 1, \dots, d$, under the experimental and background conditions with m and n replicates, respectively. We start with the simple case when $m = n$. Clipper defines a contrast score C_j of feature

j in one of two ways:

$$C_j := t^{\text{minus}}(X_j, Y_j) \quad \text{minus contrast score,} \quad (6)$$

or

$$C_j := t^{\text{max}}(X_j, Y_j) \quad \text{maximum contrast score.} \quad (7)$$

Fig. 6a shows a cartoon illustration of contrast scores when $m = n = 1$. Accordingly, a large positive value of C_j bears evidence that $\mu_{X_j} > \mu_{Y_j}$. Motivated by Barber and Candès [35], Clipper uses the following procedure to control the FDR under the target level $q \in (0, 1)$.

Definition 1 (Barber-Candès (BC) procedure for thresholding contrast scores [35]) *Given contrast scores $\{C_j\}_{j=1}^d$, $\mathcal{C} = \{|C_j| : C_j \neq 0; j = 1, \dots, d\}$ is defined as the set of non-zero absolute values of C_j 's. The BC procedure finds a contrast-score cutoff T^{BC} based on the target FDR threshold $q \in (0, 1)$ as*

$$T^{\text{BC}} := \min \left\{ t \in \mathcal{C} : \frac{\text{card}(\{j : C_j \leq -t\}) + 1}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (8)$$

and outputs $\{j : C_j \geq T^{\text{BC}}\}$ as discoveries.

Enrichment analysis with any numbers of replicates m and n

When $m \neq n$, Clipper constructs contrast scores via permutation of replicates across conditions. The idea is that, after permutation, every feature becomes uninteresting and can serve as its own negative control.

Definition 2 (Permutation) *We define σ as permutation, i.e., a bijection from the set $\{1, \dots, m+n\}$ onto itself, and we rewrite the data $X_1, \dots, X_d, Y_1, \dots, Y_d$ into a matrix \mathbf{W} :*

$$\mathbf{W} = \begin{bmatrix} W_{11} & \cdots & W_{1m} & W_{1(m+1)} & \cdots & W_{1(m+n)} \\ & & \vdots & & & \vdots \\ W_{d1} & \cdots & W_{dm} & W_{d(m+1)} & \cdots & W_{d(m+n)} \end{bmatrix} := \begin{bmatrix} X_{11} & \cdots & X_{1m} & Y_{11} & \cdots & Y_{1n} \\ & & \vdots & & & \vdots \\ X_{d1} & \cdots & X_{dm} & Y_{d1} & \cdots & Y_{dn} \end{bmatrix}.$$

We then apply σ to permute the columns of \mathbf{W} and obtain

$$\mathbf{W}_\sigma := \begin{bmatrix} W_{1\sigma(1)} & \cdots & W_{1\sigma(m)} & W_{1\sigma(m+1)} & \cdots & W_{1\sigma(m+n)} \\ & & \vdots & & & \vdots \\ W_{d\sigma(1)} & \cdots & W_{d\sigma(m)} & W_{d\sigma(m+1)} & \cdots & W_{d\sigma(m+n)} \end{bmatrix},$$

from which we obtain the permuted measurements $\{(X_j^\sigma, Y_j^\sigma)\}_{j=1}^d$, where

$$\begin{aligned} X_j^\sigma &:= (W_{j\sigma(1)}, \dots, W_{j\sigma(m)})^\top, \\ Y_j^\sigma &:= (W_{j\sigma(m+1)}, \dots, W_{j\sigma(m+n)})^\top. \end{aligned} \quad (9)$$

In the enrichment analysis, if two permutations σ and σ' satisfy that

$$\{\sigma(1), \dots, \sigma(m)\} = \{\sigma'(1), \dots, \sigma'(m)\},$$

then we define σ and σ' to be in one equivalence class. That is, permutations in the same equivalence class lead to the same division of $m+n$ replicates (from the two conditions) into two groups with sizes m and n . In total, there are $\binom{m+n}{m}$ equivalence classes of permutations.

We define σ_0 as the identity permutation such that $\sigma_0(i) = i$ for all $i \in \{1, \dots, m+n\}$. In addition, Clipper randomly samples h equivalence classes $\sigma_1, \dots, \sigma_h$ with equal probabilities without replacement from the other $h_{\max} := \binom{m+n}{m} - 1$ equivalence classes (after excluding the equivalence class containing σ_0). Note that h_{\max} is the maximum value h can take.

Clipper then obtains $\{(\mathbf{X}_j^{\sigma_0}, \mathbf{Y}_j^{\sigma_0}), (\mathbf{X}_j^{\sigma_1}, \mathbf{Y}_j^{\sigma_1}), \dots, (\mathbf{X}_j^{\sigma_h}, \mathbf{Y}_j^{\sigma_h})\}_{j=1}^d$, where $(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})$ are the permuted measurements based on σ_ℓ , $\ell = 0, 1, \dots, h$. Then Clipper computes $T_j^{\sigma_\ell} := t^{\min}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})$ to indicate the degree of “interestingness” of feature j reflected by $(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})$. Note that Clipper chooses t^{\min} instead of t^{\max} because empirical evidence shows that t^{\min} leads to better power. Sorting $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$ gives

$$T_j^{(0)} \geq T_j^{(1)} \geq \dots \geq T_j^{(h)}.$$

Then Clipper defines the contrast score of feature j , $j = 1, \dots, d$, in one of two ways:

$$C_j := \begin{cases} T_j^{(0)} - T_j^{(1)} & \text{if } T_j^{(0)} = T_j^{\sigma_0} \\ T_j^{(1)} - T_j^{(0)} & \text{otherwise} \end{cases} \quad \text{minus contrast score,} \quad (10)$$

or

$$C_j := \begin{cases} |T_j^{(0)}| & \text{if } T_j^{(0)} = T_j^{\sigma_0} > T_j^{(1)} \\ 0 & \text{if } T_j^{(0)} = T_j^{(1)} \\ -|T_j^{(0)}| & \text{otherwise} \end{cases} \quad \text{maximum contrast score.} \quad (11)$$

The intuition behind the contrast scores is that, if $C_j < 0$, then $T_j^{(0)} \neq T_j^{\sigma_0}$, which means that at least one of $T_j^{\sigma_1}, \dots, T_j^{\sigma_h}$ (after random permutation) is greater than $T_j^{\sigma_0}$ calculated from the original data (identity permutation), suggesting that feature j is likely an uninteresting feature in enrichment analysis. Fig. 6b (right) shows a cartoon illustration of contrast scores when $m = 2$ and $n = 1$. Motivated by Gimenez and Zou [62], we propose the following procedure for Clipper to control the FDR under the target level $q \in (0, 1)$.

Definition 3 (Gimenez-Zou (GZ) procedure for thresholding contrast scores [62]) Given $h \in \{1, \dots, h_{\max}\}$ and contrast scores $\{C_j\}_{j=1}^d$, $\mathcal{C} = \{|C_j| : C_j \neq 0; j = 1, \dots, d\}$ is defined as the set of non-zero absolute values of C_j 's. The GZ procedure finds a contrast-score cutoff T^{GZ} based on the target FDR threshold $q \in (0, 1)$ as:

$$T^{\text{GZ}} := \min \left\{ t \in \mathcal{C} : \frac{\frac{1}{h} + \frac{1}{h} \text{card}(\{j : C_j \leq -t\})}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (12)$$

and outputs $\{j : C_j \geq T^{\text{GZ}}\}$ as discoveries.

Differential analysis with $m + n > 2$

For differential analysis, Clipper also uses permutation to construct contrast scores. When $m \neq n$, the equivalence classes of permutations are defined the same as for the enrichment analysis with $m \neq n$. When $m = n$, there is a slight change in the definition of equivalence classes of permutations: if σ and σ' satisfy that

$$\{\sigma(1), \dots, \sigma(m)\} = \{\sigma'(1), \dots, \sigma'(m)\} \text{ or } \{\sigma'(m+1), \dots, \sigma'(2m)\},$$

then we say that σ and σ' are in one equivalence class. In total, there are $h_{\text{total}} := \binom{m+n}{m}$ (when $m \neq n$) or $\binom{2m}{m}/2$ (when $m = n$) equivalence classes of permutations. Hence, to have more than one

equivalence class, we cannot perform differential analysis with $m = n = 1$; in other words, the total number of replicates $m + n$ must be at least 3.

Then Clipper randomly samples $\sigma_1, \dots, \sigma_h$ with equal probabilities without replacement from the $h_{\max} := h_{\text{total}} - 1$ equivalence classes that exclude the class containing σ_0 , i.e., the identity permutation. Note that h_{\max} is the maximum value h can take. Next, Clipper computes $T_j^{\sigma_\ell} := |t^{\text{minus}}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})|$, where $\mathbf{X}_j^{\sigma_\ell}$ and $\mathbf{Y}_j^{\sigma_\ell}$ are the permuted data defined in (9), and it defines C_j as the contrast score of feature j , $j = 1, \dots, d$, in the same ways as in (10) or (11). Fig. 6b (left) shows a cartoon illustration of contrast scores when $m = 2$ and $n = 1$.

Same as in the enrichment analysis with $m \neq n$, Clipper also uses the GZ procedure [62] to set a cutoff on contrast scores to control the FDR under the target level $q \in (0, 1)$.

Granted, when we use permutations to construct contrast scores in the GZ procedure, we can convert contrast scores into permutation-based p-values (see Supp. S1.1.2). However, when the numbers of replicates are small, the number of possible permutations is small, so permutation-based p-values would have a low resolution (e.g., when $m = 2$ and $n = 1$, the number of non-identity permutations is only 2). Hence, applying the BH procedure to the permutation-based p-values would result in almost no power. Although Yekutieli and Benjamini proposed another thresholding procedure for permutation-based p-values [66], it still requires the number of permutations to be large to obtain a reliable FDR control. Furthermore, if we apply the SeqStep+ procedure by Barber and Candés [35] to permutation-based p-values, it would be equivalent to our application of the GZ procedure to contrast scores (Supp. Section S1.1.2).

For both differential and enrichment analyses, the two contrast scores (minus and maximum) can both control the FDR. Based on the power comparison results in Supp. Section S3 and Supp. Figs. S22–S25, Clipper has the following default choice of contrast score: for the enrichment analysis when two conditions have the same number of replicates (“Enrichment analysis with equal numbers of replicates ($m = n$)” in Methods), Clipper uses the BC procedure with the minus contrast score; for the enrichment analysis when two conditions have different numbers of replicates (“Enrichment analysis with any numbers of replicates m and n ” in Methods) or the differential analysis (“Differential analysis with $m + n > 2$ ” in Methods), Clipper uses the GZ procedure with maximum contrast score.

Generic FDR control methods

In our simulation analysis, we compared Clipper against generic FDR control methods including p-value-based methods and local-fdr-based methods. Briefly, each p-value-based method is a combination of a p-value calculation approach and a p-value thresholding procedure. We use either the “paired” or “pooled” approach (see next paragraph) to calculate p-values of features and then threshold the p-values using the BH procedure (Supp. Definition S1) or Storey’s qvalue procedure (Supp. Definition S2) to make discoveries (Supp. Section S1.1). As a result, we have four p-value-based methods: BH-pair, BH-pool, qvalue-pair, and qvalue-pool (Fig. 1b).

Regarding the existing p-value calculation approaches in bioinformatics tools, we categorize them as “paired” or “pooled.” The paired approach has been widely used to detect DEGs and protein-binding sites [1, 2, 4, 5]. It examines one feature at a time and compares the feature’s measurements between two conditions using a statistical test. In contrast, the pooled approach is popular in proteomics for identifying peptide sequences from MS data [67]. For every feature, it defines a test statistic and estimates a null distribution by pooling all features’ observed test statistic values under the background condition. Finally, it calculates a p-value for every feature under the experimental condition based on the feature’s observed test statistic and the null distribution.

In parallel to p-value-based methods, local-fdr-based methods estimate local fdrs of features and

then threshold the local fdrs using the locfdr procedure (Supp. Definition S5) to make discoveries. The estimation of local fdrs takes one of two approaches: (1) empirical null, which is estimated parametrically from the test statistic values that are likely drawn from the null distribution, and (2) swapping null, which is constructed by swapping measurements between experimental and background conditions. The resulting two local-fdr-based-methods are referred to as locfdr-emp and locfdr-swap (Figs. 1b and 2). Supp. Section S1 provides a detailed explanation of these generic methods and how we implemented them in this work.

Specific to the p-value-based methods, for the paired approach, besides the ideal implementation that uses the correct model to calculate p-values (BH-pair-correct and qvalue-pair-correct), we also consider common mis-implementations. The first mis-implementations is misspecification of the distribution (BH-pair-mis and qvalue-pair-mis). An example is the detection of protein-binding sites from ChIP-seq data. A common assumption is that ChIP-seq read counts in a genomic region (i.e., a feature) follow the Poisson distribution [1, 2], which implies that the counts have the variance equal to the mean. However, if only two replicates are available, it is impossible to check whether this Poisson distribution is reasonably specified. The second mis-implementation is the misspecification of a two-sample test as a one-sample test (BH-pair-2as1 and qvalue-pair-2as1), which ignores the sampling randomness of replicates under one condition. This issue is implicit but widespread in bioinformatics methods [1, 68].

To summarize, we compared Clipper against the following implementations of generic FDR control methods:

- **BH-pool** or **qvalue-pool**: p-values calculated by the pooled approach and thresholded by the BH or qvalue procedure.
- **BH-pair-correct** or **qvalue-pair-correct**: p-values calculated by the paired approach with the correct model specification and thresholded by the BH or qvalue procedure.
- **BH-pair-mis** or **qvalue-pair-mis**: p-values calculated by the paired approach with a misspecified model and thresholded by the BH or qvalue procedure.
- **BH-pair-2as1** or **qvalue-pair-2as1**: p-values calculated by the paired approach that misformulates a two-sample test as a one-sample test (2as1) and thresholded by the BH or qvalue procedure.
- **locfdr-emp**: local fdrs calculated by the empirical null approach and thresholded by the locfdr procedure.
- **locfdr-swap**: local fdrs calculated by the swapping approach and thresholded by the locfdr procedure.

Real datasets

- The H3K4me3 ChIP-seq dataset with one experimental sample (GEO accession number GSM733708) and two control samples (GEO accession number GSM733742) from the cell line GM12878 is available at <ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeBroadHistone/>, with the experimental sample under the filename `wgEncodeBroadHistoneGm12878H3k4me3StdA1nRep1.bam` and the two control samples under the filenames `wgEncodeBroadHistoneGm12878ControlStdA1nRep1.bam` and `wgEncodeBroadHistoneGm12878ControlStdA1nRep2.bam`. The processed dataset is available at <https://zenodo.org/record/4404882>.
- The MS benchmark dataset will be published in a future manuscript. Interested readers should contact Dr. Leo Wang at lewang@coh.org.

- The human monocyte RNA-seq dataset is available at <https://www.ncbi.nlm.nih.gov/Traces/study/?acc=srp082682>. The dataset includes 17 samples of classical monocytes and 17 samples of non-classical monocytes, and it is converted to a sample-by-gene count matrix by R package GenomicFeatures (v 1.40.1). The processed count matrix is available at <https://zenodo.org/record/4404882>.
- The Hi-C dataset from the cell line GM12878 is available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63525>. The count matrix is under the filename GSE63525_GM12878_primary_intrachromosomal_contact_matrices.tar.gz, and the matrix corresponding to Chromosome 1 and bin width 1MB is used. The processed dataset is available at <https://zenodo.org/record/4404882>.

Software packages used in this study

- **p.adjust** R function (in R package stats v 4.0.2 with default arguments) [14]: used for BH-pool, BH-pair-correct, BH-pair-mis, and BH-pair-2as1.
- **qvalue** R package (v 2.20.0 with default arguments) [69]: used for qvalue-pool, qvalue-pair-correct, qvalue-pair-mis, and qvalue-pair-2as1.
- **locfdr** R package (v 1.1-8 with default arguments) [70]: used for locfdr-emp.
- **MACS2** software package (v 2.2.6 with default settings) [1]: available at <https://github.com/macs3-project/MACS/releases/tag/v2.2.6>.
- **ChIPulate** software package [45]: available at <https://github.com/vishakad/chipulate>.
- **HOMER** software package (findPeaks v 3.1.9.2 with default settings) [2]: available at <http://www.bcgsc.ca/platform/bioinfo/software/findpeaks/releases/3.1.9.2/findpeaks3-1-9-2-tar.gz>.
- **SEQUEST** in Proteome Discoverer (v 2.3.0.523 with default settings) [3]: commercial software by ThermoScientific.
- **edgeR** R package (v 3.30.0 with default arguments) [4]: available at <https://www.bioconductor.org/packages/release/bioc/html/edgeR.html>.
- **DESeq2** R package (v 1.28.1 with default arguments) [5]: available at <http://bioconductor.org/packages/release/bioc/html/DESeq2.html>.
- **limma** R package (v 3.44.3 with default arguments) [10]: available at <https://www.bioconductor.org/packages/release/bioc/html/limma.html>.
- **MAST** R package (v 1.14.0 with default arguments) [55]: available at <https://www.bioconductor.org/packages/release/bioc/html/MAST.html>.
- **monocle3** R package (v 0.2.3.0 with default arguments) [56]: available at <https://www.bioconductor.org/packages/release/bioc/html/monocle3.html>.
- **MultiHiCcompare** R package (v 1.6.0 with default arguments) [11]: available at <https://bioconductor.org/packages/release/bioc/html/multiHiCcompare.html>.
- **diffHic** R package (v 1.20.0 with default arguments) [13]: available at <https://www.bioconductor.org/packages/release/bioc/html/diffHic.html>.

- **FIND** R package (v 0.99 with default arguments) [12]: available at <https://bitbucket.org/nadhir/find/src/master/>.

Software, code, and video tutorial

- The **Clipper** R package is available at <https://github.com/JSB-UCLA/Clipper/>.
- The code and processed data for reproducing the figures are available at <https://zenodo.org/record/4404882>.
- A video introduction of Clipper is available at <https://youtu.be/-GXyHiJMpLo>.

Funding

This work was supported by the following grants: NIH-NCI T32LM012424 (to Y.E.C.); NCI K08 CA201591, the Gabrielle’s Angel Foundation and Alex’s Lemonade Stand Foundation (to L.D.W.); NIH R01HG007538, R01CA193466, and R01CA228140 (to W.L.); National Science Foundation DBI-1846216, NIH/NIGMS R01GM120507, Johnson & Johnson WiSTEM2D Award, Sloan Research Fellowship, and UCLA David Geffen School of Medicine W.M. Keck Foundation Junior Faculty Award (to J.J.L.).

Conflicts of interests

L.D.W. holds equity in Magenta Therapeutics.

Acknowledgements

The authors would like to thank Dr. Yu-Cheng Yang for his suggestions on the figures and R package. The authors would also like to thank Mr. Nikos Ignatiadis, Dr. Lihua Lei, and Dr. Rina Barber for their insightful comments after we presented this work at the International Seminar on Selective Inference (<https://www.selectiveinferenceseminar.com/past-talks>). The authors also appreciate the comments and feedback from Mr. Tianyi Sun, Ms. Kexin Li, and other members of the Junction of Statistics and Biology at UCLA (<http://jsb.ucla.edu>).

References

- [1] Yong Zhang et al. “Model-based analysis of ChIP-Seq (MACS)”. In: *Genome biology* 9.9 (2008), pp. 1–9.
- [2] Sven Heinz et al. “Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities”. In: *Molecular cell* 38.4 (2010), pp. 576–589.
- [3] Michael P Washburn, Dirk Wolters, and John R Yates. “Large-scale analysis of the yeast proteome by multidimensional protein identification technology”. In: *Nature biotechnology* 19.3 (2001), pp. 242–247.

- [4] Mark D Robinson, Davis J McCarthy, and Gordon K Smyth. “edgeR: a Bioconductor package for differential expression analysis of digital gene expression data”. In: *Bioinformatics* 26.1 (2010), pp. 139–140.
- [5] Michael I Love, Wolfgang Huber, and Simon Anders. “Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2”. In: *Genome biology* 15.12 (2014), p. 550.
- [6] Cole Trapnell et al. “Differential analysis of gene regulation at transcript resolution with RNA-seq”. In: *Nature biotechnology* 31.1 (2013), pp. 46–53.
- [7] Jun Li et al. “Normalization, testing, and false discovery rate estimation for RNA-sequencing data”. In: *Biostatistics* 13.3 (2012), pp. 523–538.
- [8] Thomas J Hardcastle and Krystyna A Kelly. “baySeq: empirical Bayesian methods for identifying differential expression in sequence count data”. In: *BMC bioinformatics* 11.1 (2010), pp. 1–14.
- [9] Gordon K Smyth. “Linear models and empirical bayes methods for assessing differential expression in microarray experiments”. In: *Statistical applications in genetics and molecular biology* 3.1 (2004).
- [10] Matthew E Ritchie et al. “limma powers differential expression analyses for RNA-sequencing and microarray studies”. In: *Nucleic acids research* 43.7 (2015), e47–e47.
- [11] John C Stansfield, Kellen G Cresswell, and Mikhail G Dozmorov. “multiHiCcompare: joint normalization and comparative analysis of complex Hi-C experiments”. In: *Bioinformatics* 35.17 (2019), pp. 2916–2923.
- [12] Mohamed Nadhir Djekidel, Yang Chen, and Michael Q Zhang. “FIND: diffERential chromatin INteractions Detection using a spatial Poisson process”. In: *Genome research* 28.3 (2018), pp. 412–422.
- [13] Aaron TL Lun and Gordon K Smyth. “diffHic: a Bioconductor package to detect differential genomic interactions in Hi-C data”. In: *BMC bioinformatics* 16.1 (2015), pp. 1–11.
- [14] Yoav Benjamini and Yosef Hochberg. “Controlling the false discovery rate: a practical and powerful approach to multiple testing”. In: *Journal of the Royal statistical society: series B (Methodological)* 57.1 (1995), pp. 289–300.
- [15] Bradley Efron and Robert Tibshirani. “Empirical Bayes methods and false discovery rates for microarrays”. In: *Genetic epidemiology* 23.1 (2002), pp. 70–86.
- [16] Bradley Efron et al. “Empirical Bayes analysis of a microarray experiment”. In: *Journal of the American statistical association* 96.456 (2001), pp. 1151–1160.
- [17] Matthew Stephens. “False discovery rates: a new deal”. In: *Biostatistics* 18.2 (2017), pp. 275–294.
- [18] John D Storey and Robert Tibshirani. “Statistical significance for genomewide studies”. In: *Proceedings of the National Academy of Sciences* 100.16 (2003), pp. 9440–9445.
- [19] Anat Reiner, Daniel Yekutieli, and Yoav Benjamini. “Identifying differentially expressed genes using false discovery rate controlling procedures”. In: *Bioinformatics* 19.3 (2003), pp. 368–375.
- [20] Bing Yang et al. “Identification of cross-linked peptides from complex samples”. In: *Nature methods* 9.9 (2012), pp. 904–906.
- [21] James Robert White, Niranjana Nagarajan, and Mihai Pop. “Statistical methods for detecting differentially abundant features in clinical metagenomic samples”. In: *PLoS Comput Biol* 5.4 (2009), e1000352.
- [22] Andrey A Shabalin. “Matrix eQTL: ultra fast eQTL analysis via large matrix operations”. In: *Bioinformatics* 28.10 (2012), pp. 1353–1358.

- [23] Stijn Hawinkel et al. “A broken promise: microbiome differential abundance methods do not control the false discovery rate”. In: *Briefings in bioinformatics* 20.1 (2019), pp. 210–221.
- [24] Ye Zheng and Sündüz Keleş. “FreeHi-C simulates high-fidelity Hi-C data for benchmarking and data augmentation”. In: *Nature Methods* 17.1 (2020), pp. 37–40.
- [25] John D Storey. “A direct approach to false discovery rates”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64.3 (2002), pp. 479–498.
- [26] Yoav Benjamini and Yosef Hochberg. “Multiple hypotheses testing with weights”. In: *Scandinavian Journal of Statistics* 24.3 (1997), pp. 407–418.
- [27] Nikolaos Ignatiadis et al. “Data-driven hypothesis weighting increases detection power in genome-scale multiple testing”. In: *Nature methods* 13.7 (2016), pp. 577–580.
- [28] Lihua Lei and William Fithian. “Adapt: an interactive procedure for multiple testing with side information”. In: *arXiv preprint arXiv:1609.06035* (2016).
- [29] Simina M Boca and Jeffrey T Leek. “A direct approach to estimating false discovery rates conditional on covariates”. In: *PeerJ* 6 (2018), e6035.
- [30] Joses Ho et al. “Moving beyond P values: data analysis with estimation graphics”. In: *Nature methods* 16.7 (2019), pp. 565–566.
- [31] Dongyuan Song and Jingyi Jessica Li. “PseudotimeDE: inference of differential gene expression along cell pseudotime with well-calibrated p-values from single-cell RNA sequencing data”. In: *bioRxiv* (2020).
- [32] Virginia Goss Tusher, Robert Tibshirani, and Gilbert Chu. “Significance analysis of microarrays applied to the ionizing radiation response”. In: *Proceedings of the National Academy of Sciences* 98.9 (2001), pp. 5116–5121.
- [33] Jesse Hemerik and Jelle J Goeman. “False discovery proportion estimation by permutations: confidence for significance analysis of microarrays”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80.1 (2018), pp. 137–155.
- [34] Jesse Hemerik, Aldo Solari, and Jelle J Goeman. “Permutation-based simultaneous confidence bounds for the false discovery proportion”. In: *Biometrika* 106.3 (2019), pp. 635–649.
- [35] Rina Foygel Barber and Emmanuel J Candès. “Controlling the false discovery rate via knockoffs”. In: *The Annals of Statistics* 43.5 (2015), pp. 2055–2085.
- [36] Ery Arias-Castro and Shiyun Chen. “Distribution-free multiple testing”. In: *Electronic Journal of Statistics* 11.1 (2017), pp. 1983–2001.
- [37] Yoav Benjamini. “Selective inference: The silent killer of replicability”. In: *Issue 2.4* 2.4 (2020).
- [38] Kristen Emery et al. “Multiple Competition-Based FDR Control and Its Application to Peptide Detection”. In: *International Conference on Research in Computational Molecular Biology*. Springer. 2020, pp. 54–71.
- [39] Abhishek K Sarkar and Matthew Stephens. “Separating measurement and expression models clarifies confusion in single cell RNA-seq analysis”. In: *BioRxiv* (2020).
- [40] Peter J Park. “ChIP-seq: advantages and challenges of a maturing technology”. In: *Nature reviews genetics* 10.10 (2009), pp. 669–680.
- [41] Pamela J Mitchell and Robert Tjian. “Transcriptional regulation in mammalian cells by sequence-specific DNA binding proteins”. In: *Science* 245.4916 (1989), pp. 371–378.

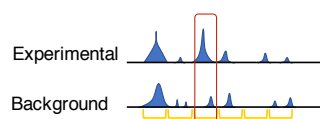
- [42] Mark Ptashne and Alexander Gann. “Transcriptional activation by recruitment”. In: *Nature* 386.6625 (1997), pp. 569–577.
- [43] Timothy Bailey et al. “Practical guidelines for the comprehensive analysis of ChIP-seq data”. In: *PLoS Comput Biol* 9.11 (2013), e1003326.
- [44] ENCODE Project Consortium et al. “An integrated encyclopedia of DNA elements in the human genome”. In: *Nature* 489.7414 (2012), pp. 57–74.
- [45] Vishaka Datta, Sridhar Hannenhalli, and Rahul Siddharthan. “ChIPulate: A comprehensive ChIP-seq simulation pipeline”. In: *PLoS computational biology* 15.3 (2019), e1006921.
- [46] Aaron Diaz et al. “Normalization, bias correction, and peak calling for ChIP-seq”. In: *Statistical applications in genetics and molecular biology* 11.3 (2012).
- [47] Boris Bogdanow, Henrik Zauber, and Matthias Selbach. “Systematic errors in peptide and protein identification and quantification by modified peptides”. In: *Molecular & Cellular Proteomics* 15.8 (2016), pp. 2791–2801.
- [48] Marshall Bern, Yong J Kil, and Christopher Becker. “Byonic: advanced peptide and protein identification software”. In: *Current protocols in bioinformatics* 40.1 (2012), pp. 13–20.
- [49] Mark D Robinson and Alicia Oshlack. “A scaling normalization method for differential expression analysis of RNA-seq data”. In: *Genome biology* 11.3 (2010), pp. 1–9.
- [50] Claire R Williams et al. “Empirical assessment of analysis workflows for differential expression analysis of human samples using RNA-Seq”. In: *BMC bioinformatics* 18.1 (2017), p. 38.
- [51] Marek Gierliński et al. “Statistical models for RNA-seq data derived from a two-condition 48-replicate experiment”. In: *Bioinformatics* 31.22 (2015), pp. 3625–3630.
- [52] Keegan Korthauer et al. “A practical guide to methods controlling false discoveries in computational biology”. In: *Genome biology* 20.1 (2019), pp. 1–21.
- [53] Qunhua Li et al. “Measuring reproducibility of high-throughput experiments”. In: *The annals of applied statistics* 5.3 (2011), pp. 1752–1779.
- [54] Guangchuang Yu et al. “clusterProfiler: an R package for comparing biological themes among gene clusters”. In: *Omics: a journal of integrative biology* 16.5 (2012), pp. 284–287.
- [55] Greg Finak et al. “MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data”. In: *Genome biology* 16.1 (2015), pp. 1–13.
- [56] Xiaojie Qiu et al. “Single-cell mRNA quantification and differential analysis with Censur”. In: *Nature methods* 14.3 (2017), pp. 309–315.
- [57] Charlotte Sonesson and Mark D Robinson. “Bias, robustness and scalability in single-cell differential expression analysis”. In: *Nature methods* 15.4 (2018), p. 255.
- [58] Tianyi Sun et al. “scDesign2: an interpretable simulator that generates high-fidelity single-cell gene expression count data with gene correlations captured”. In: *bioRxiv* (2020).
- [59] Jiarui Ding et al. “Systematic comparison of single-cell and single-nucleus RNA-sequencing methods”. In: *Nature biotechnology* 38.6 (2020), pp. 737–746.
- [60] Grace XY Zheng et al. “Massively parallel digital transcriptional profiling of single cells”. In: *Nature communications* 8.1 (2017), pp. 1–12.
- [61] Evan Z Macosko et al. “Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets”. In: *Cell* 161.5 (2015), pp. 1202–1214.

- [62] Jaime Roquero Gimenez and James Zou. “Improving the Stability of the Knockoff Procedure: Multiple Simultaneous Knockoffs and Entropy Maximization”. In: *arXiv preprint arXiv:1810.11378* (2018).
- [63] Ning Wang et al. “Identifying the combinatorial control of signal-dependent transcription factors”. In: *PLOS Computational Biology* 17.6 (2021), e1009095.
- [64] Jonathan Thorsen et al. “Large-scale benchmarking reveals false discoveries and count transformation sensitivity in 16S rRNA gene amplicon data analysis methods used in microbiome studies”. In: *Microbiome* 4.1 (2016), p. 62.
- [65] Kun He et al. “Null-free False Discovery Rate Control Using Decoy Permutations for Multiple Testing”. In: *arXiv preprint arXiv:1804.08222* (2018).
- [66] Daniel Yekutieli and Yoav Benjamini. “Resampling-based false discovery rate controlling multiple test procedures for correlated test statistics”. In: *Journal of Statistical Planning and Inference* 82.1-2 (1999), pp. 171–196.
- [67] Alexey I Nesvizhskii. “A survey of computational methods and error rate estimation procedures for peptide and protein identification in shotgun proteomics”. In: *Journal of proteomics* 73.11 (2010), pp. 2092–2123.
- [68] Wei Li et al. “MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens”. In: *Genome biology* 15.12 (2014), p. 554.
- [69] John D. Storey et al. *qvalue: Q-value estimation for false discovery rate control*. R package version 2.20.0. 2020. URL: <http://github.com/jdstorey/qvalue>.
- [70] Bradley Efron. “Large-scale simultaneous hypothesis testing: the choice of a null hypothesis”. In: *Journal of the American Statistical Association* 99.465 (2004), pp. 96–104.

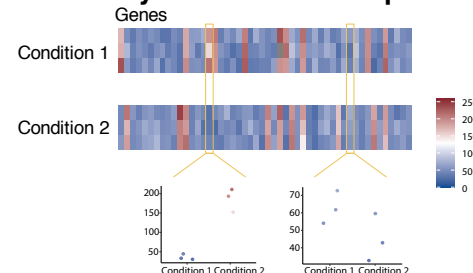
Figures

a High-throughput omics data analyses

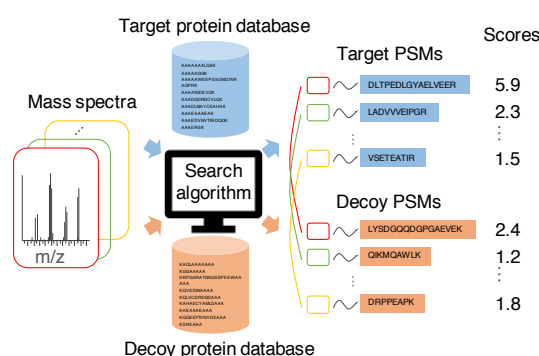
Peak calling from ChIP-seq data



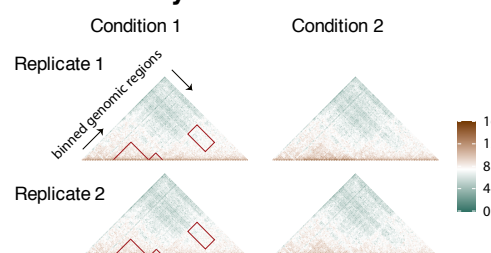
DEG analysis from RNA-seq data



Peptide identification from MS data



DIR analysis from Hi-C data



b Generic FDR control methods

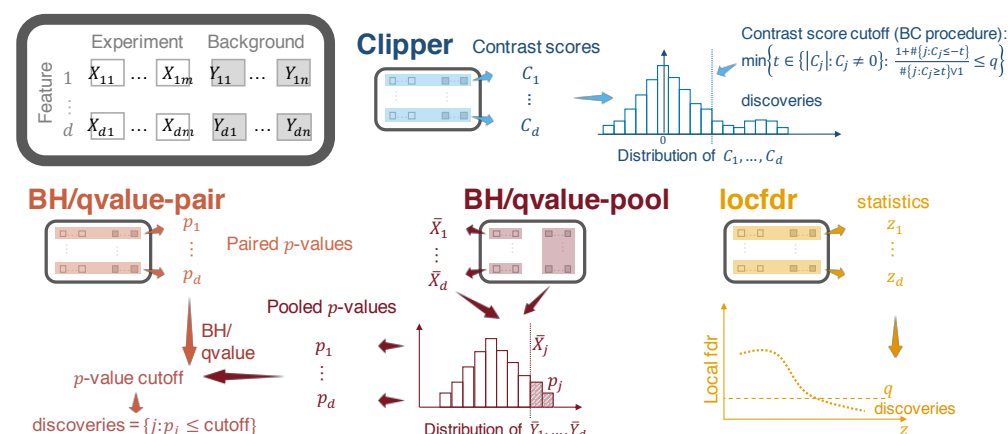


Figure 1: High-throughput omics data analyses and generic FDR control methods. (a) Illustration of four common high-throughput omics data analyses: peak calling from ChIP-seq data, peptide identification from MS data, DEG analysis from RNA-seq data, and DIR analysis from Hi-C data. In these four analyses, the corresponding features are genomic regions (yellow intervals), peptide-spectrum matches (PSMs; a pair of a mass spectrum and a peptide sequence), genes (columns in the heatmaps), and chromatin interacting regions (entries in the heatmaps). (b) Illustration of Clipper and five generic FDR control methods: BH-pair (and qvalue-pair), BH-pool (and qvalue-pool), and locfdr. The input data are d features with m and n repeated measurements under the experimental and background conditions, respectively. Clipper computes a contrast score for each feature based on the feature's m and n measurements, decides a contrast-score cutoff, and calls the features with contrast scores above the cutoff as discoveries. (This illustration is Clipper for enrichment analysis with $m = n$.) BH-pair or qvalue-pair computes a p-value for each feature based on the feature's m and n measurements, sets a p-value cutoff, and calls the features with p-values below the cutoff as discoveries. BH-pool or qvalue-pool constructs a null distribution from the d features' average (across the n replicates) measurements under the background condition, calculates a p-value for each feature based on the null distribution and the feature's average (across the m replicates) measurements under the experimental condition, sets a p-value cutoff, and calls the features with p-values below the cutoff as discoveries. The locfdr method computes a summary statistic for each feature based on the feature's m and n measurements, estimates the empirical null distribution and the empirical distribution of the statistic across features, computes a local fdr for each feature, sets a local fdr cutoff, and calls the features with local fdr below the cutoff as discoveries.

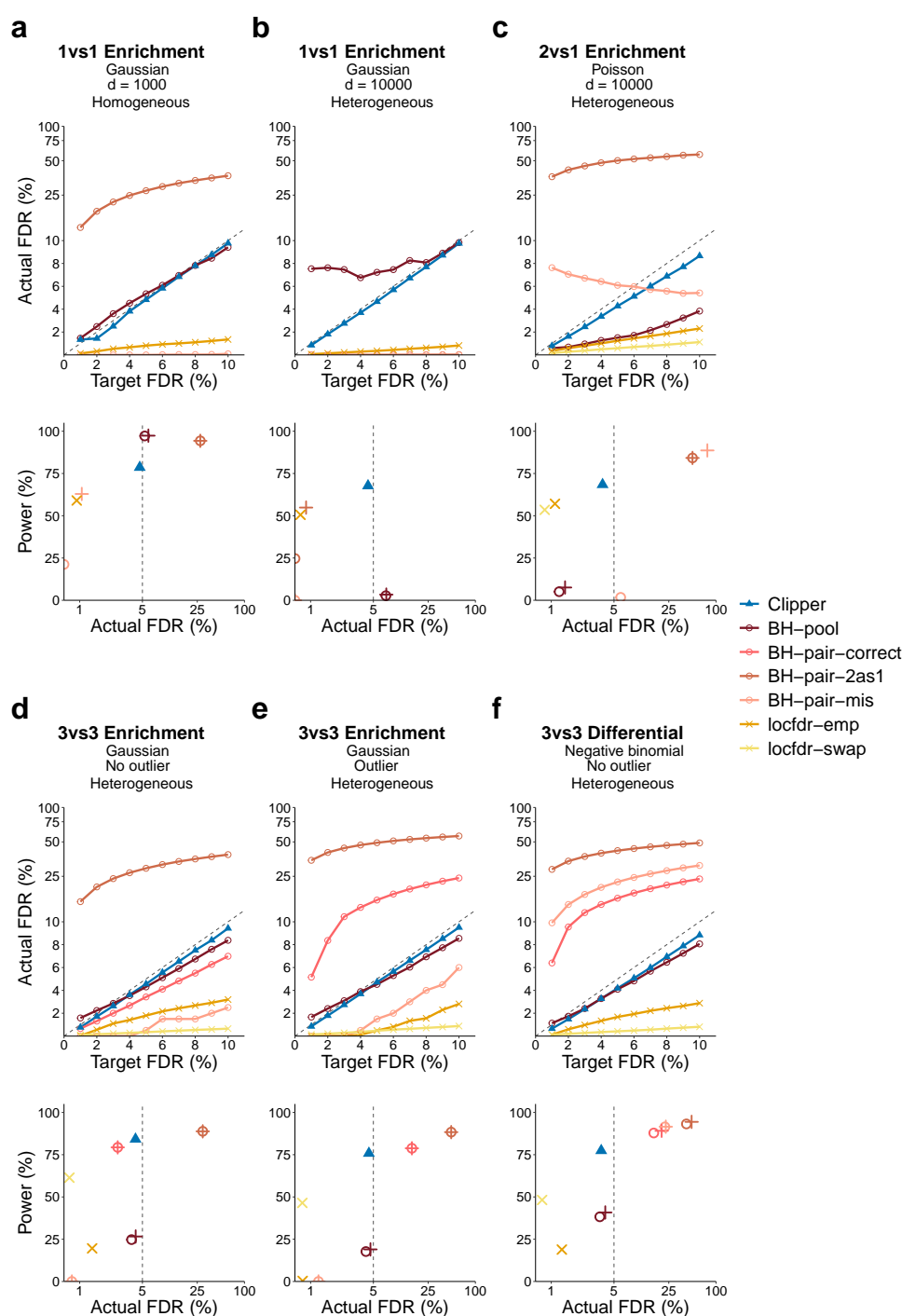


Figure 2: Comparison of Clipper with generic FDR control methods in terms of their FDR control and power in six example simulation studies. (a) 1vs1 enrichment analysis with 1000 features generated from the Gaussian distribution with a homogeneous background; (b) 1vs1 enrichment analysis with 10,000 features generated from the Gaussian distribution with a heterogeneous background; (c) 2vs1 enrichment analysis with 10,000 features generated from the Poisson distribution with a heterogeneous background; (d) 3vs3 enrichment analysis with 10,000 features generated from the Gaussian distribution without outliers and with a heterogeneous background; (e) 3vs3 enrichment analysis with 10,000 features generated from the Gaussian distribution without outliers and with a heterogeneous background; (f) 3vs3 differential analysis with 10,000 features generated from the negative binomial distribution with a heterogeneous background. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are approximated by the averages of false discovery proportions (see Eq. (S14) in the Supplementary) and power evaluated on 200 simulated datasets. In each panel, the top row shows each method's actual FDRs at target FDR thresholds: whenever the actual FDR is larger than the target FDR (the solid line is higher than the dashed line), FDR control is failed; the bottom row shows each method's actual FDRs and power at the target FDR threshold $q = 5\%$: whenever the actual FDR is greater than q (on the right of the vertical dashed line), FDR control is failed. Under the FDR control, the larger the power, the better. Note that BH-pair-correct is not included in (a)–(c) because it is impossible to correctly specify the model with only one replicate per condition; locfdr-swap is not included in (a)–(b) because it is inapplicable to the 1vs1 design.

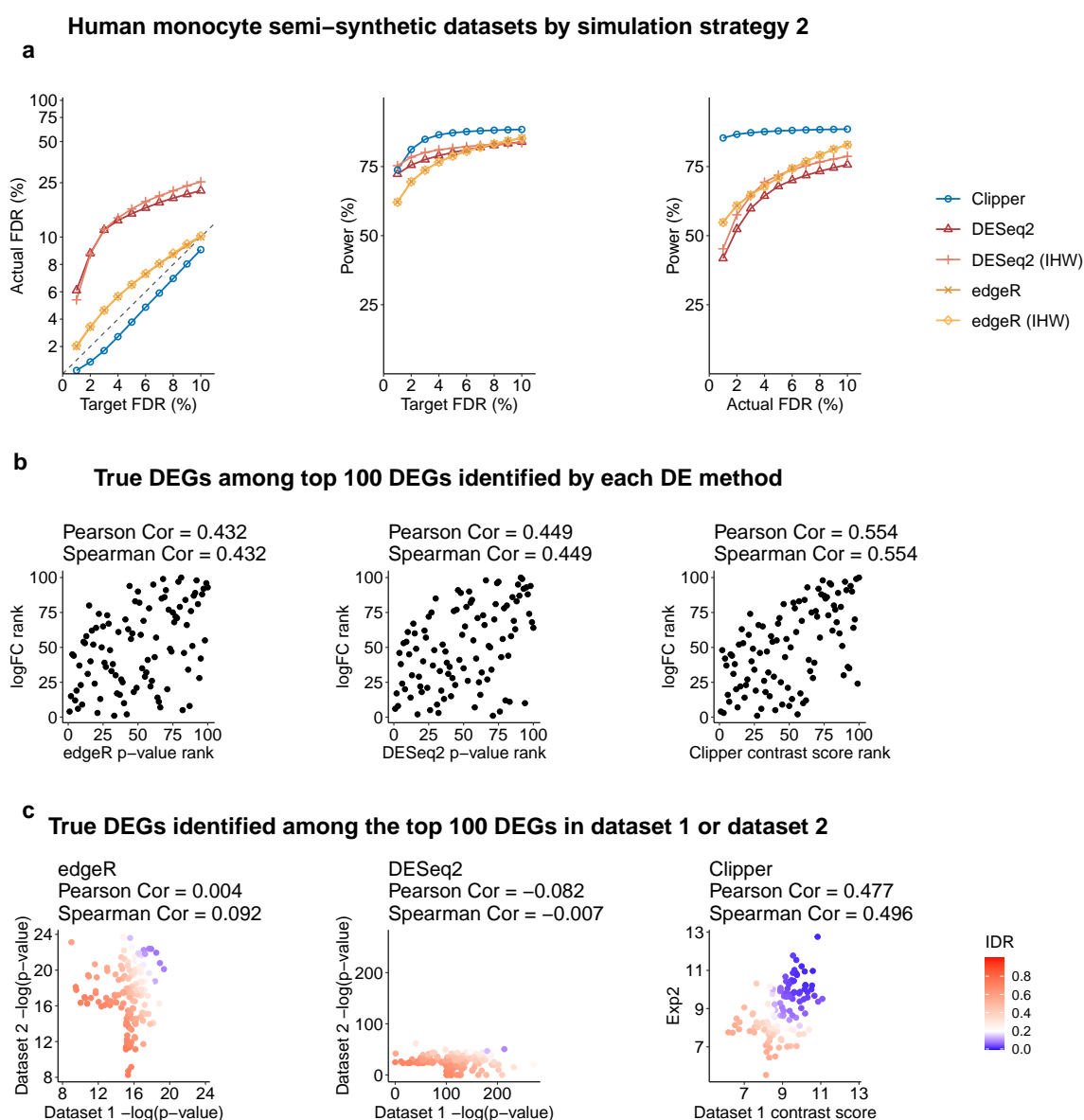
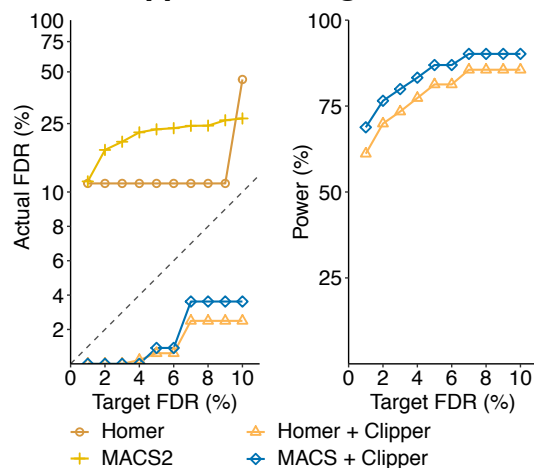
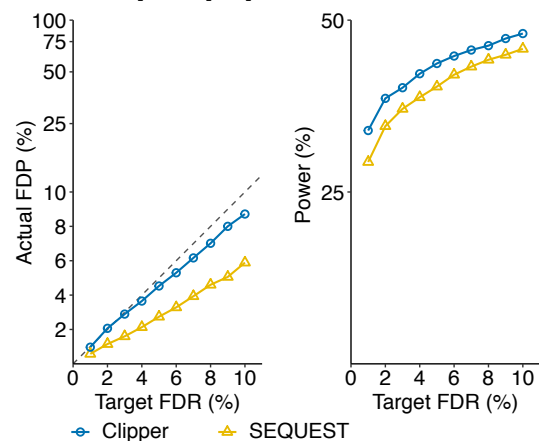


Figure 3: Comparison of Clipper and two popular DEG identification methods—edgeR and DESeq2—in DEG analysis on semi-synthetic bulk RNA-seq data (generated from human monocyte real data using simulation strategy 2 in Supp. Section S6.3). (a) FDR control, power given the same target FDR, and power given the same actual FDR. (b) Ranking consistency of the true DEGs among the top 100 DEGs identified by each method. The consistency is defined between the genes' ranking based on edgeR/DESeq2's p-values or Clipper's contrast scores and their ranking based on true expression fold changes. (c) Reproducibility between two semi-synthetic datasets as technical replicates. Three reproducibility criteria are used: the IDR, Pearson correlation, and Spearman correlation. Each criterion is calculated for edgeR/DESeq2's p-values or Clipper's contrast scores on the two semi-synthetic datasets. Among the three methods, only Clipper controls the FDR, and Clipper achieves the highest power, the best gene ranking consistency, and the best reproducibility.

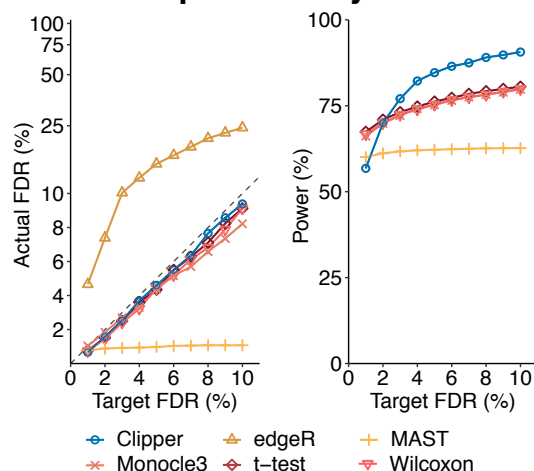
a ChIP-seq peak calling



b Mass-spec peptide identification



c scRNA-seq DEG analysis



d Hi-C DIR analysis

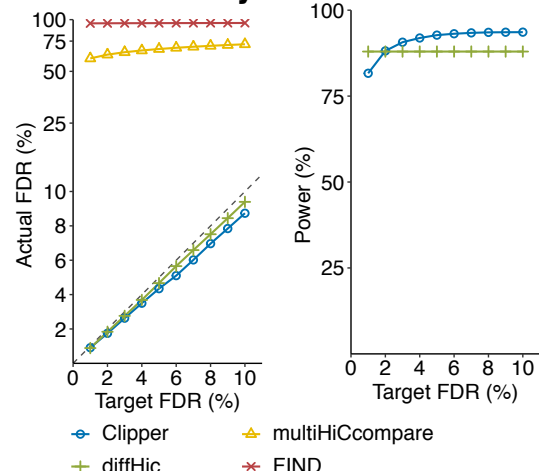


Figure 4: Comparison of Clipper and popular bioinformatics methods in terms of FDR control and power. (a) peaking calling analysis on semi-synthetic ChIP-seq data; (b) peptide identification on real proteomics data; (c) DEG analysis on synthetic 10x Genomics scRNA-seq data; (d) DIR analysis on semi-synthetic Hi-C data. In all four panels, the target FDR threshold q ranges from 1% to 10%. In the “Actual FDR vs. Target FDR” plot of each panel, points above the dashed diagonal line indicate failed FDR control; when this happens, the power of the corresponding methods is not shown, including HOMER in (a), MACS2 for target FDR less than 5% in (a), edgeR in (c), and multiHiCcompare, and FIND in (d). In all four applications, Clipper controls the FDR while maintaining high power, demonstrating Clipper’s broad applicability in high-throughput data analyses.

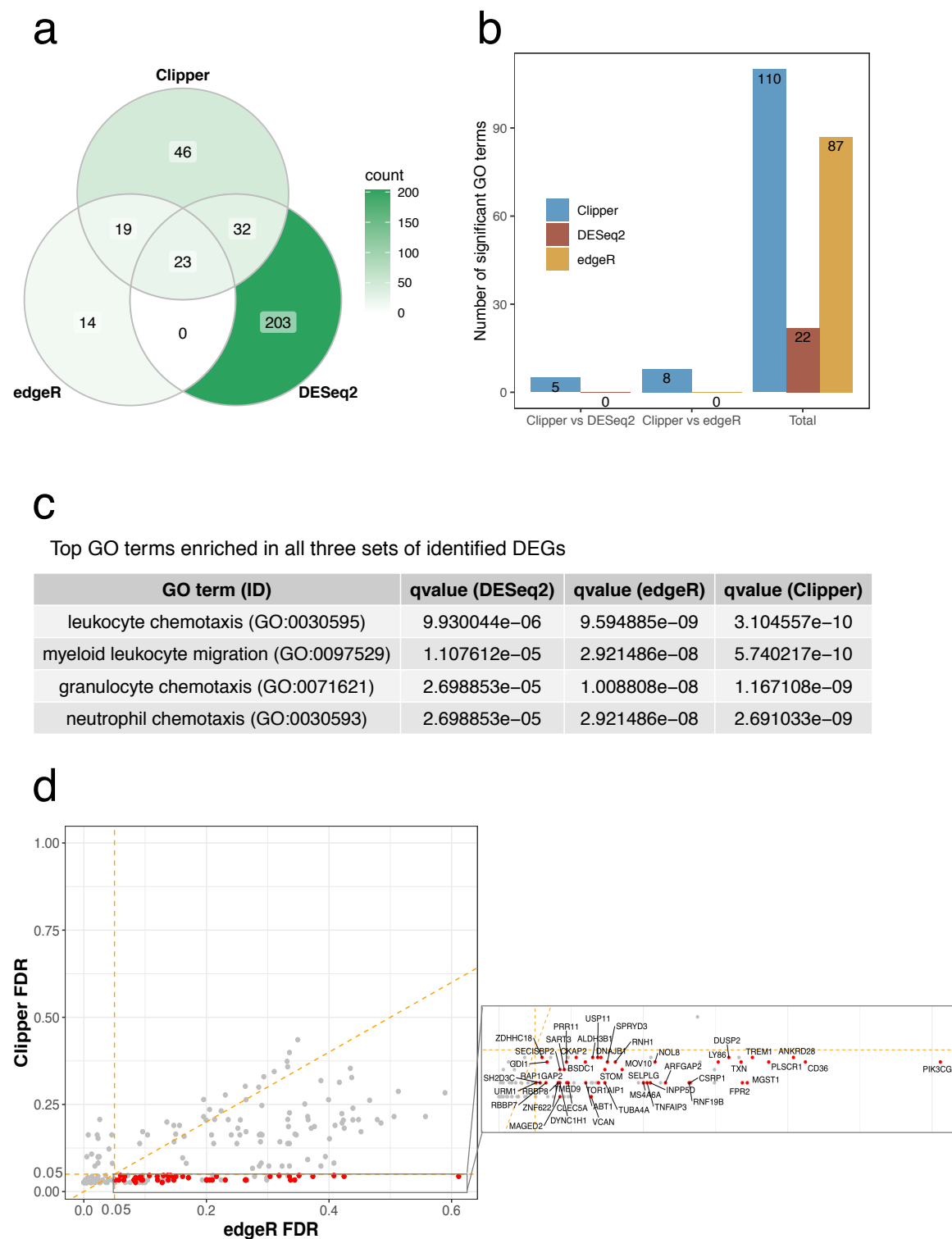
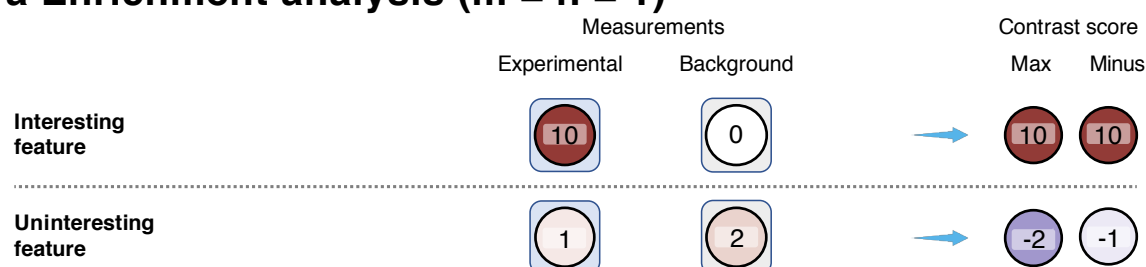


Figure 5: Application of Clipper, DESeq2, and edgeR to identifying DEGs from the classical and non-classical human monocyte dataset. (a) A Venn diagram showing the overlaps of the identified DEGs (at the FDR threshold $q = 5\%$) by the three DE methods. (b) Numbers of GO terms enriched (with enrichment q -values < 0.01) in the DEGs found by Clipper, DESeq2 and edgeR (column 3), or in the DEGs specifically identified by Clipper or DESeq2/edgeR in the pairwise comparison between Clipper and DESeq2 (column 1) or between Clipper and edgeR (column 2). More GO terms are enriched in the DEGs identified by Clipper than in those identified by edgeR or DESeq2. (c) Enrichment q -values of four GO terms that are found enriched (with enrichment q -values < 0.01) in all three sets of identified DEGs, one set per method. All the four terms are most enriched in the DEGs identified by Clipper. (d) A scatterplot of the claimed FDR of Clipper against that of edgeR for all the DEGs identified by Clipper, edgeR or DESeq2. The 46 DEGs only identified by Clipper are highlighted with red.

a Enrichment analysis ($m = n = 1$)



b Differential analysis and enrichment analysis ($m = 2, n = 1$)

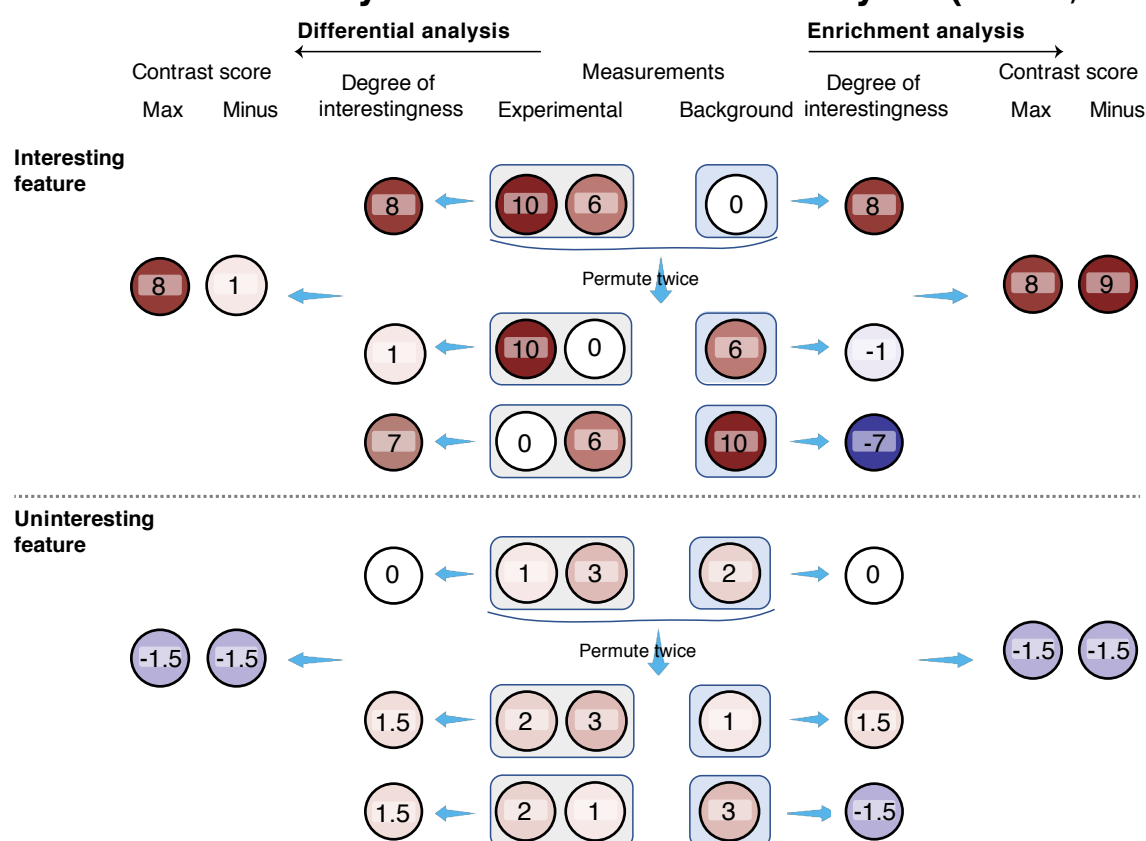


Figure 6: Illustration of the construction of contrast scores. (a) 1vs1 enrichment analysis; (b) 2vs1 differential analysis (left) or enrichment analysis (right). In each panel, an interesting feature (top) and an uninteresting feature (bottom) are plotted for contrast; both features have measurements under the experimental and background conditions. In (a), each feature's measurements are summarized into a maximum (max) contrast score or a minus contrast score. In (b), each feature's measurements are permuted across the two conditions, resulting in two sets of permuted measurements. Then for each feature, we calculate its degrees of interestingness (as the difference that equals the average of experimental measurements minus the average of background measurements (in enrichment analysis; right), or the absolute value of the difference (in differential analysis; left)) from its original measurements and permuted measurements, respectively. Finally, we summarize each feature's degrees of interestingness into a maximum (max) contrast score or a minus contrast score.

Supplementary Material

S1 Review of generic FDR control methods

To facilitate our discussion, we introduce the notations for data. For feature $j = 1, \dots, d$, we use $\mathbf{X}_j = (X_{j1}, \dots, X_{jm})^\top \in \mathbb{R}_{\geq 0}^m$ and $\mathbf{Y}_j = (Y_{j1}, \dots, Y_{jn})^\top \in \mathbb{R}_{\geq 0}^n$ to denote its measurements under the experimental and background conditions, respectively. We assume that X_{j1}, \dots, X_{jm} are identically distributed, so are Y_{j1}, \dots, Y_{jn} . Let $\mu_{Xj} = \mathbb{E}[X_{j1}]$ and $\mu_{Yj} = \mathbb{E}[Y_{j1}]$ denote the expected measurement of feature j under the two conditions, respectively. Then we denote by \bar{X}_j the sample average of X_{j1}, \dots, X_{jm} and by \bar{Y}_j the sample average of Y_{j1}, \dots, Y_{jn} .

S1.1 P-value-based methods

Here we describe the details of p-value-based FDR control methods, including BH-pair, BH-pool, qvalue-pair, and qvalue-pool. Each of these four methods first computes p-values using either the pooled approach or the paired approach, and it then relies on the BH procedure [1] or Storey's qvalue procedure [2] for FDR control. In short, every p-value-based method is a combination of a p-value calculation approach and a p-value thresholding procedure. Below we introduce two p-value calculation approaches (paired and pooled) and two p-value thresholding procedures (BH and Storey's qvalue).

S1.1.1 P-value calculation approaches

The paired approach. The paired approach examines one feature at a time and compares its measurements between two conditions. Besides the ideal implementation, i.e., the *correct paired approach* that uses the correct model to calculate p-values, we also include commonly-used flawed implementations that either misspecify the distribution, i.e., the *misspecified paired approach*, or misformulate the two-sample test as a one-sample test, i.e., the *2as1 paired approach*.

Here we use the negative binomial distribution as an example to demonstrate the ideas of the correct, misspecified, and 2as1 paired approaches. Suppose that for each feature j , its measurements under each condition follow a negative binomial distribution, and the two distributions under the two conditions have the same dispersion; that is, $X_{j1}, \dots, X_{jm} \stackrel{\text{i.i.d.}}{\sim} \text{NB}(\mu_{Xj}, \theta_j)$; $Y_{j1}, \dots, Y_{jn} \stackrel{\text{i.i.d.}}{\sim} \text{NB}(\mu_{Yj}, \theta_j)$, where θ_j is the dispersion parameter such that the variance $\text{Var}(X_{ji}) = \mu_{Xj} + \theta_j \mu_{Xj}^2$.

- The correct paired approach assumes that the two negative binomial distributions have the same dispersion parameter θ_j , and it uses the two-sample test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ (enrichment analysis) or $H_1 : \mu_{Xj} \neq \mu_{Yj}$ (differential analysis).
- The misspecified paired approach misspecifies the negative binomial distribution as Poisson, and it uses the two-sample test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ (enrichment analysis) or $H_1 : \mu_{Xj} \neq \mu_{Yj}$ (differential analysis).
- The 2as1 paired approach bluntly assumes $\mu_{Yj} = \bar{Y}_j$, and it performs the one-sample test based on X_{j1}, \dots, X_{jm} for the null hypotheses $H_0 : \mu_{Xj} = \bar{Y}_j$ against the alternative hypothesis $H_1 : \mu_{Xj} > \bar{Y}_j$ (enrichment analysis) or $H_1 : \mu_{Xj} \neq \bar{Y}_j$ (differential analysis).

The pooled approach. The pooled approach pools all features' average measurements under the background condition $\{\bar{Y}_j\}_{j=1}^d$ to form a null distribution, and it calculates a p-value for each feature j

by comparing \bar{X}_j to the null distribution. Specifically, in enrichment analysis, the p-value of feature j is computed as:

$$p_j = \frac{\text{card}(\{k : \bar{Y}_k \geq \bar{X}_j\})}{d}.$$

In differential analysis, the p-value of feature j is computed as:

$$p_j = 2 \cdot \min\left(\frac{\text{card}(\{k : \bar{Y}_k \geq \bar{X}_j\})}{d}, \frac{\text{card}(\{k : \bar{Y}_k \leq \bar{X}_j\})}{d}\right).$$

S1.1.2 P-value thresholding procedures for FDR control

Definition S1 (BH procedure for thresholding p-values [1]) *The features' p-values p_1, \dots, p_d are sorted in an ascending order $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(d)}$. Given the target FDR threshold q , the Benjamini–Hochberg (BH) procedure finds a p-value cutoff T^{BH} as*

$$T^{\text{BH}} := p_{(k)}, \text{ where } k = \max\left\{j = 1, \dots, d : p_{(j)} \leq \frac{j}{d}q\right\}. \quad (\text{S1})$$

Then BH outputs $\{j : p_j \leq T^{\text{BH}}\}$ as discoveries.

Definition S2 (Storey's qvalue procedure for thresholding p-values [2]) *The features' p-values p_1, \dots, p_d are sorted in an ascending order $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(d)}$. Let $\hat{\pi}_0$ denote an estimate of the probability $P(\text{the } i\text{-th feature is uninteresting})$ (see Storey [2] for details). Storey's qvalue procedure defines the q-value for $p_{(d)}$ as*

$$\hat{q}(p_{(d)}) := \frac{\hat{\pi}_0 \cdot d \cdot p_{(d)}}{\text{card}(\{k : p_k \leq p_{(d)}\})} = \hat{\pi}_0 \cdot p_{(d)}.$$

Then for $j = d-1, d-2, \dots, 1$, the q-value for $p_{(j)}$ is defined as:

$$\hat{q}(p_{(j)}) := \min\left(\hat{q}(p_{(j+1)}), \frac{\hat{\pi}_0 \cdot d \cdot p_{(j)}}{\text{card}(\{k : p_k \leq p_{(j)}\})}\right).$$

Then Storey's qvalue procedure outputs $\{j : \hat{q}(p_j) \leq q\}$ as discoveries.

We use function `qvalue` from R package `qvalue` (v 2.20.0; with default estimate $\hat{\pi}_0$) to calculate q-values.

Definition S3 (SeqStep+ procedure for thresholding p-values [3]) *Define H_0^j as the null hypothesis for feature j and p_j as the p-value for H_0^j , $j = 1, \dots, d$. Order the null hypotheses H_0^1, \dots, H_0^d from the most to the least promising (here more promising means more likely to be interesting) and denote the resulting null hypotheses and p-values as $H_0^{(1)}, \dots, H_0^{(d)}$ and $p_{(1)}, \dots, p_{(d)}$. Given any target FDR threshold q , a pre-specified constant $s \in (0, 1)$, and subset $\mathcal{K} \subseteq \{1, \dots, d\}$, the SeqStep+ procedure finds a cutoff \hat{j} as*

$$\hat{j} := \max\left\{j \in \mathcal{K} : \frac{1 + \text{card}(\{k \in \mathcal{K}, k \leq j : p_{(k)} > s\})}{\text{card}(\{k \in \mathcal{K}, k \leq j : p_{(k)} \leq s\}) \vee 1} \leq \frac{1-s}{s}q\right\} \quad (\text{S2})$$

Then SeqStep+ rejects $\{H_0^{(j)} : p_{(j)} \leq s, j \leq \hat{j}, j \in \mathcal{K}\}$. If the orders of the null hypotheses are independent of the p-values, the SeqStep+ procedure ensures FDR control.

The GZ procedure (Definition 3) used in Clipper is a special case of the SeqStep+ procedure with $s = 1/(h+1)$. Recall that given the number of non-identical permutations $h \in \{1, \dots, h_{\max}\}$ and contrast

scores $\{C_j\}_{j=1}^d$, the GZ procedure sorts $\{|C_j|\}_{j=1}^d$ in a decreasing order:

$$|C_{(1)}| \geq |C_{(2)}| \geq \cdots \geq |C_{(d)}|. \quad (\text{S3})$$

To see the connection between the GZ procedure and SeqStep+, we consider the null hypothesis for the j -th ordered feature, $j = 1, \dots, d$, as $H_0^{(j)} : \mu_{X(j)} = \mu_{Y(j)}$ and define the corresponding p-value $p_{(j)} := \frac{r(T_{(j)}^{\sigma_0})}{h+1}$, where $r(T_{(j)}^{\sigma_0})$ is the rank of $T_{(j)}^{\sigma_0}$ in $\{T_{(j)}^{\sigma_0}, \dots, T_{(j)}^{\sigma_h}\}$ in a descending order. We also define $\mathcal{K} := \{j = 1, \dots, d : C_j \neq 0\}$ as the subset of features with non-zero C_j 's. Finally, we input the p-values, null hypothesis orders in (S3), $s = 1/(h+1)$, q and \mathcal{K} into the SeqStep+ procedure, and we obtain the GZ procedure.

The BC procedure (Definition 1) is a further special case with $h = 1$, $p_{(j)} := (\mathbb{1}(C_{(j)} > 0) + 1)/2$, and $\mathcal{K} := \{j = 1, \dots, d : C_j \neq 0\}$.

S1.2 Local-fdr-based methods

The FDR is statistical criterion that ensures the reliability of discoveries as a whole. In contrast, the local fdr focuses on the reliability of each discovery. The definition of the local fdr relies on some pre-computed summary statistics z_j for feature j , $j = 1, \dots, d$. In the calculation of local fdr, $\{z_1, \dots, z_d\}$ are assumed to be realizations of an abstract random variable Z that represents any feature. Let p_0 or p_1 denote the prior probability that any feature is uninteresting or interesting, with $p_0 + p_1 = 1$. Let $f_0(z) := \mathbb{P}(Z = z | \text{uninteresting})$ or $f_1(z) := \mathbb{P}(Z = z | \text{interesting})$ denote the conditional probability density of Z at z given that Z represents an uninteresting or interesting feature. Thus by Bayes' theorem, the posterior probability of any feature being uninteresting given its summary statistic $Z = z$ is

$$\mathbb{P}(\text{uninteresting} | Z = z) = p_0 f_0(z) / f(z), \quad (\text{S4})$$

where $f(z) := p_0 f_0(z) + p_1 f_1(z)$ is the marginal probability density of Z . Accordingly, the local fdr of feature j is defined as follows.

Definition S4 (Local fdr [4]) *Given notations defined above, the local fdr of feature j is defined as*

$$\text{local-fdr}_j := f_0(z_j) / f(z_j).$$

Because $p_0 \leq 1$, local-fdr_j is an upper bound of the posterior probability of feature j being uninteresting given its summary statistic z_j , defined in (S4).

Note that another definition of the local fdr is the posterior probability $\mathbb{P}(\text{uninteresting} | z)$ in (S4) [5]. Although this other definition is more reasonable, we do not use it but choose Definition S4 because the estimation of p_0 is usually difficult. Another reason is that uninteresting features are the dominant majority in high-throughput biological data, so p_0 is often close to 1.

We define local-fdr-based methods as a type of FDR control methods by thresholding local fdrs of features under the target FDR threshold q . Although the local fdr is different from FDR, it has been shown that thresholding the local fdrs at q will approximately control the FDR under q [4]. This makes local-fdr-based methods competitors against Clipper and p-value-based methods.

Every local-fdr-based method is a combination of a local fdr calculation approach and a local fdr thresholding procedure. Below we introduce two local fdr calculation approaches (empirical null and swapping) and one local fdr thresholding procedure. After the combination, we have two local-fdr-based methods: locfdr-emp and locfdr-swap.

S1.2.1 Local fdr calculation approaches

With z_1, \dots, z_d , the calculation of local fdr defined in Definition S4 requires the estimation of f_0 and f , two probability densities. f is estimated by nonparametric density estimation, and f_0 is estimated by either the empirical null approach [4] or the swapping approach, which shuffles replicates between conditions [5]. With the estimated \hat{f} and \hat{f}_0 , the estimated local fdr of feature j is

$$\widehat{\text{local-fdr}}_j := \hat{f}_0(z_j) / \hat{f}(z_j). \quad (\text{S5})$$

The empirical null approach. This approach assumes a parametric distribution, typically the Gaussian distribution, to estimate f_0 . Then with the density estimate \hat{f} , the local fdr is estimated for each feature j . The implementation of this approach depends on the numbers of replicates.

- In 1vs1 enrichment and differential analyses, we define z_j as

$$z_j := \frac{D_j}{\sqrt{\frac{1}{d} \sum_{j=1}^d (D_j - \bar{D})^2}},$$

where $D_j = X_{j1} - Y_{j1}$ and $\bar{D} = \sum_{j=1}^d D_j / d$.

- In 2vs1 enrichment and differential analyses, we define z_j as

$$z_j := \frac{\bar{X}_j - Y_{j1}}{\sqrt{\frac{s_{X_j}^2}{2}}},$$

where $s_{X_j}^2 = \sum_{i=1}^2 (X_{ji} - \bar{X}_j)^2$.

- In m vs n enrichment and differential analyses with $m, n \geq 2$, we define z_j as the two-sample t-statistic with unequal variances:

$$z_j := \frac{\bar{X}_j - \bar{Y}_j}{\sqrt{\frac{s_{X_j}^2}{m} + \frac{s_{Y_j}^2}{n}}},$$

where $s_{X_j}^2 = \frac{1}{m-1} \sum_{i=1}^m (X_{ji} - \bar{X}_j)^2$ and $s_{Y_j}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_{ji} - \bar{Y}_j)^2$ are the sample variances of feature j under the experimental and background conditions.

Then $\{\widehat{\text{locfdr}}_j\}_{j=1}^d$ are estimated from $\{z_j\}_{j=1}^d$ by function `locfdr` in R package `locfdr` (v 1.1-8; with default arguments).

The swapping approach. This approach swaps $\lceil m/2 \rceil$ replicates under the experimental condition with $\lceil n/2 \rceil$ replicates under the background condition. Then it calculates the summary statistic for each feature on the swapped data, obtaining z'_1, \dots, z'_d . Finally, it estimates f_0 and f by applying kernel density estimation to z'_1, \dots, z'_d and z_1, \dots, z_d , respectively (by function `kde` in R package `ks`). With \hat{f}_0 and \hat{f} , $\{\widehat{\text{locfdr}}_j\}_{j=1}^d$ are calculated by Definition S4.

The implementation of this approach depends on the numbers of replicates. Below are three special cases included in this work.

- In 1vs1 enrichment and differential analyses, the swapping approach is inapplicable because interesting features would not become uninteresting after the swapping.

- In 2vs1 enrichment and differential analyses, we define z_j and z'_j as

$$z_j = \frac{X_{j1} + X_{j2}}{2} - Y_{j1},$$

$$z'_j = \frac{X_{j1} + Y_{j1}}{2} - X_{j2}.$$

- In 3vs3 enrichment and differential analyses with, we define z_j and z'_j as

$$z_j = \frac{X_{j1} + X_{j2}}{2} - \frac{Y_{j1} + Y_{j2}}{2},$$

$$z'_j = \frac{X_{j1} + Y_{j1}}{2} - \frac{X_{j2} + Y_{j2}}{2}.$$

Then we apply kernel density estimation to $\{z_j\}_{j=1}^d$ and $\{z'_j\}_{j=1}^d$ to obtain \hat{f} and \hat{f}_0 , respectively. By (S5), we calculate $\{\widehat{\text{locfdr}}_j\}_{j=1}^d$.

S1.2.2 The local fdr thresholding procedure

Definition S5 (locfdr procedure) Given the local fdr estimates $\{\widehat{\text{local-fdr}}_j\}_{j=1}^d$ and the target FDR threshold q , the locfdr procedure outputs $\{j = 1, \dots, d : \widehat{\text{local-fdr}}_j \leq q\}$ as discoveries.

S2 The Clipper methodology

Clipper is a flexible framework that reliably controls the FDR without using p-values in high-throughput data analysis with two conditions. Clipper has two functionalities: (I) enrichment analysis, which identifies the “interesting” features that have higher expected measurements (i.e., true signals) under the experimental condition than the background, a.k.a. negative control condition (if the goal is to identify the interesting features with smaller expected measurements under the experimental condition, enrichment analysis can be applied after the values are negated); (II) differential analysis, which identifies the interesting features that have different expected measurements between the two conditions. For both functionalities, uninteresting features are defined as those that have equal expected measurements under the two conditions.

Clipper only relies on two fundamental statistical assumptions of biological data analysis: (1) measurement errors (i.e., differences between measurements and their expectations, with the expectations including biological signals and batch effects) are independent across all features and experiments; (2) every uninteresting feature has measurement errors identically distributed across all experiments. These two assumptions are used in almost all bioinformatics tools and commonly referred to as the “measurement model” in statistical genomics [6].

In the following subsections, we will first introduce notations and assumptions used in Clipper. Then we will detail how Clipper works and discuss its theoretical guarantee in three analysis tasks: the enrichment analysis with equal numbers of replicates under two conditions ($m = n$), the enrichment analysis with different numbers of replicates under two conditions ($m \neq n$), and the differential analysis (when $m + n > 2$).

S2.1 Notations and assumptions

To facilitate our discussion, we first introduce the following mathematical notations. For two random vectors $\mathbf{X} = (X_1, \dots, X_m)^\top$ and $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$, or two sets of random variables $\mathcal{X} = \{X_1, \dots, X_m\}$

and $\mathcal{Y} = \{Y_1, \dots, Y_n\}$, we write $\mathbf{X} \perp \mathbf{Y}$ or $\mathcal{X} \perp \mathcal{Y}$ if X_i is independent of Y_j for all $i = 1, \dots, m$ and $j = 1, \dots, n$. To avoid confusion, we use $\text{card}(A)$ to denote the cardinality of a set A and $|c|$ to denote the absolute value of a scalar c . We define $a \vee b := \max(a, b)$.

Clipper only requires two inputs: the target FDR threshold $q \in (0, 1)$ and the input data. Regarding the input data, we use d to denote the number of features with measurements under two conditions, and we use m and n to denote the numbers of replicates under the two conditions. For each feature $j = 1, \dots, d$, we use $\mathbf{X}_j = (X_{j1}, \dots, X_{jm})^\top \in \mathbb{R}_{\geq 0}^m$ and $\mathbf{Y}_j = (Y_{j1}, \dots, Y_{jn})^\top \in \mathbb{R}_{\geq 0}^n$ to denote its measurements under the two conditions, where $\mathbb{R}_{\geq 0}$ denotes the set of non-negative real numbers. We assume that all measurements are non-negative, as in the case of most high-throughput experiments. (If this assumption does not hold, transformations can be applied to make data satisfy this assumption.)

Clipper has the following assumptions on the joint distribution of $\mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{Y}_1, \dots, \mathbf{Y}_d$. For $j = 1, \dots, d$, Clipper assumes that X_{j1}, \dots, X_{jm} are identically distributed, so are Y_{j1}, \dots, Y_{jn} . Let $\mu_{Xj} = \mathbb{E}[X_{j1}]$ and $\mu_{Yj} = \mathbb{E}[Y_{j1}]$ denote the expected measurement of feature j under the two conditions, respectively. Then conditioning on $\{\mu_{Xj}\}_{j=1}^d$ and $\{\mu_{Yj}\}_{j=1}^d$,

$$\begin{aligned} X_{j1}, \dots, X_{jm}, Y_{j1}, \dots, Y_{jn} \text{ are mutually independent;} \\ \mathbf{X}_j \perp \mathbf{X}_k, \mathbf{Y}_j \perp \mathbf{Y}_k \text{ and } \mathbf{X}_j \perp \mathbf{Y}_k, \forall j, k = 1, \dots, d. \end{aligned} \quad (\text{S6})$$

An enrichment analysis aims to identify interesting features with $\mu_{Xj} > \mu_{Yj}$ (with \mathbf{X}_j and \mathbf{Y}_j defined as the measurements under the experimental and background conditions, respectively), while a differential analysis aims to call interesting features with $\mu_{Xj} \neq \mu_{Yj}$. We define $\mathcal{N} := \{j : \mu_{Xj} = \mu_{Yj}\}$ as the set of uninteresting features and denote $N := \text{card}(\mathcal{N})$. In both analyses, Clipper further assumes that an uninteresting feature j satisfies

$$X_{j1}, \dots, X_{jm}, Y_{j1}, \dots, Y_{jn} \text{ are identically distributed, } \forall j \in \mathcal{N}. \quad (\text{S7})$$

Clipper consists of two main steps: construction and thresholding of contrast scores. First, Clipper computes contrast scores, one per feature, as summary statistics that reflect the extent to which features are interesting. Second, Clipper establishes a contrast-score cutoff and calls as discoveries the features whose contrast scores exceed the cutoff.

To construct contrast scores, Clipper uses two summary statistics $t(\cdot, \cdot) : \mathbb{R}_{\geq 0}^m \times \mathbb{R}_{\geq 0}^n \rightarrow \mathbb{R}$ to extract data information regarding whether a feature is interesting or not:

$$t^{\text{minus}}(\mathbf{x}, \mathbf{y}) := \bar{x} - \bar{y}; \quad (\text{S8})$$

$$t^{\text{max}}(\mathbf{x}, \mathbf{y}) := \max(\bar{x}, \bar{y}) \cdot \text{sign}(\bar{x} - \bar{y}), \quad (\text{S9})$$

where $\mathbf{x} = (x_1, \dots, x_m)^\top \in \mathbb{R}_{\geq 0}^m$, $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}_{\geq 0}^n$, $\bar{x} = \sum_{i=1}^m x_i/m$, $\bar{y} = \sum_{i=1}^n y_i/n$, and $\text{sign}(\cdot) : \mathbb{R} \rightarrow \{-1, 0, 1\}$ with $\text{sign}(x) = 1$ if $x > 0$, $\text{sign}(x) = -1$ if $x < 0$, and $\text{sign}(x) = 0$ otherwise.

Notably, other summary statistics can also be used to construct contrast scores. For example, an alternative summary statistic is the t statistic from the two-sample t test:

$$t^t(\mathbf{x}, \mathbf{y}) := \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sum_{i=1}^m (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - \bar{y})^2}{m+n-2}}}.$$

S2.2 Enrichment analysis with equal numbers of replicates ($m = n$)

Under the enrichment analysis, we assume that $\mathbf{X}_j \in \mathbb{R}_{\geq 0}^m$ and $\mathbf{Y}_j \in \mathbb{R}_{\geq 0}^n$ are the measurements of feature j , $j = 1, \dots, d$, under the experimental and background conditions with m and n replicates, respectively. We start with the simple case when $m = n$. Clipper defines a contrast score C_j of feature j in one of two ways:

$$C_j := t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j) \quad \text{minus contrast score,} \quad (\text{S10})$$

or

$$C_j := t^{\text{max}}(\mathbf{X}_j, \mathbf{Y}_j) \quad \text{maximum contrast score.} \quad (\text{S11})$$

Accordingly, a large positive value of C_j bears evidence that $\mu_{Xj} > \mu_{Yj}$. Motivated by Barber and Candès [3] and Arias-Castro and Chen [7], Clipper proposes the following BC procedure to control the FDR under the target level $q \in (0, 1)$.

Definition S6 (Barber-Candès (BC) procedure for thresholding contrast scores [3]) *Given contrast scores $\{C_j\}_{j=1}^d$, $\mathcal{C} = \{|C_j| : C_j \neq 0; j = 1, \dots, d\}$ is defined as the set of non-zero absolute values of C_j 's. The BC procedure finds a contrast-score cutoff T^{BC} based on the target FDR threshold $q \in (0, 1)$ as*

$$T^{\text{BC}} := \min \left\{ t \in \mathcal{C} : \frac{\text{card}(\{j : C_j \leq -t\}) + 1}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (\text{S12})$$

and outputs $\{j : C_j \geq T^{\text{BC}}\}$ as discoveries.

Theorem 1 *Suppose that the input data satisfy the Clipper assumptions (S6)–(S7) and $m = n$. Then for any $q \in (0, 1)$ and either definition of contrast scores in (S10) or (S11), the contrast-score cutoff T^{BC} found by the BC procedure guarantees that the discoveries have the FDR under q :*

$$\text{FDR} = \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \right] \leq q,$$

where $\mathcal{N} = \{j : \mu_{Xj} = \mu_{Yj}\}$ denotes the set of uninteresting features.

The proof of Theorem 1 (Supp. Section S8) requires two key ingredients: Lemma 1, which states important properties of contrast scores, and Lemma 2 from [8], which states a property of a Bernoulli process with independent but not necessarily identically distributed random variables. The cutoff T^{BC} can be viewed as a stopping time of a Bernoulli process.

Lemma 1 *Suppose that the input data that satisfy the Clipper assumptions (S6)–(S7) and $m = n$, and that Clipper constructs contrast scores $\{C_j\}_{j=1}^d$ based on (S10) or (S11). Denote $S_j = \text{sign}(C_j) \in \{-1, 0, 1\}$. Then $\{S_j\}_{j=1}^d$ satisfy the following properties:*

- (a) S_1, \dots, S_d are mutually independent ;
- (b) $\mathbb{P}(S_j = 1) = \mathbb{P}(S_j = -1)$ for all $j \in \mathcal{N}$;
- (c) $\{S_j\}_{j \in \mathcal{N}} \perp \mathcal{C}$.

Notably, Lemma 1(a) can be relaxed as $\mathbb{P}(S_j = 1) \leq \mathbb{P}(S_j = -1)$ for all $j \in \mathcal{N}$. Then Lemma 2 still holds, and so does Theorem 1, making Clipper still have theoretical FDR control.

Lemma 2 Suppose that Z_1, \dots, Z_d are independent with $Z_j \sim \text{Bernoulli}(\rho_j)$, and $\min_j \rho_j \geq \rho > 0$. Let J be a stopping time in reverse time with respect to the filtration $\{\mathcal{F}_j\}$, where

$$\mathcal{F}_j = \sigma((Z_1 + \dots + Z_j), Z_{j+1}, \dots, Z_d), \quad (\text{S13})$$

with $\sigma(\cdot)$ denoting a σ -algebra. Then

$$\mathbb{E} \left[\frac{1 + J}{1 + Z_1 + \dots + Z_J} \right] \leq \rho^{-1}.$$

Here we give a brief intuition about how Lemma 2 bridges Lemma 1 and Theorem 1 for FDR control. First we note that the false discovery proportion (FDP), whose expectation is the FDR, satisfies

$$\text{FDP} := \frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \quad (\text{S14})$$

$$= \frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot \frac{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \quad (\text{S15})$$

$$\leq \frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot \frac{\text{card}(\{j : C_j \leq -T^{\text{BC}}\}) + 1}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \quad (\text{S16})$$

$$\leq \frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot q, \quad (\text{S17})$$

where the last inequality follows from the definition of T^{BC} (S12).

By its definition, if T^{BC} exists, it is positive. This implies that Clipper would never call the features with $C_j = 0$ as discoveries. Here we sketch the idea of proving Theorem 1 by considering a simplified case where \mathcal{C} is fixed instead of being random; that is, we assume the features with non-zero contrast scores to be known. Then, without loss of generality, we assume $\mathcal{C} = \{1, \dots, d\}$. Then we order the absolute values of uninteresting features' contrast scores, i.e., elements in $\{|C_j| : j \in \mathcal{N}\}$, from the largest to the smallest, denoted by $|C_{(1)}| \geq |C_{(2)}| \geq \dots \geq |C_{(N)}|$. Let $J = \sum_{j \in \mathcal{N}} \mathbb{1}(|C_j| \geq T^{\text{BC}})$, the number of uninteresting features whose contrast scores have absolute values no less than T^{BC} . When $J > 0$, $|C_{(1)}| \geq \dots \geq |C_{(J)}| \geq T^{\text{BC}}$. Define $Z_k = \mathbb{1}(C_{(k)} < 0)$, $k = 1, \dots, N$. Then for each order k , the following holds

$$\begin{aligned} C_{(k)} \geq T^{\text{BC}} &\iff |C_{(k)}| \geq T^{\text{BC}} \text{ and } C_{(k)} > 0 \iff k \leq J \text{ and } Z_k = 0; \\ C_{(k)} \leq -T^{\text{BC}} &\iff |C_{(k)}| \geq T^{\text{BC}} \text{ and } C_{(k)} < 0 \iff k \leq J \text{ and } Z_k = 1. \end{aligned}$$

Then the upper bound of FDP becomes

$$\begin{aligned} \frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot q &= \frac{\sum_{k=1}^N \mathbb{1}(C_{(k)} \geq T^{\text{BC}})}{1 + \sum_{k=1}^N \mathbb{1}(C_{(k)} \leq -T^{\text{BC}})} \cdot q \\ &= \frac{\sum_{k=1}^J \mathbb{1}(C_{(k)} \geq T^{\text{BC}})}{1 + \sum_{k=1}^J \mathbb{1}(C_{(k)} \leq -T^{\text{BC}})} \cdot q \\ &= \frac{(1 - Z_1) + \dots + (1 - Z_J)}{1 + Z_1 + \dots + Z_J} \cdot q \\ &= \left(\frac{1 + J}{1 + Z_1 + \dots + Z_J} - 1 \right) \cdot q. \end{aligned}$$

By Lemma 1(a)–(b), $Z_k \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(0.5)$, which together with Lemma 1(c) satisfy the condition of

Lemma 2 and make $\rho = 0.5$. Then by Lemma 2, we have

$$\text{FDR} = \mathbb{E}[\text{FDP}] \leq \mathbb{E} \left[\frac{1 + J}{1 + Z_1 + \dots + Z_J} - 1 \right] \cdot q \leq (\rho^{-1} - 1) \cdot q = q,$$

which is the statement of Theorem 1. The complete proof of Theorem 1 is in Supp. Section S8.

S2.2.1 An optional, heuristic fix if the BC procedure makes no discoveries

Although the BC procedure has theoretical guarantee of FDR control, it lacks power when the number of replicates $m = n$, the target FDR threshold q , and the number of features d are all small (e.g., $m = n = 1$, $q = 0.01$ and $d = 1000$ in Fig. S22). As a result, the BC procedure may lead to no discoveries. In that case, Clipper implements a heuristic fix—an approximate p-value Benjamini-Hochberg (aBH) procedure—to increase the power. The aBH procedure constructs an empirical null distribution of contrast scores by additionally assuming that uninteresting features' contrast scores follow a symmetric distribution around zero; it then computes approximate p-values of features based on the empirical null distribution, and finally it uses the BH procedure [1] to threshold the approximate p-values.

Definition S7 (The aBH procedure) *Given contrast scores $\{C_j\}_{j=1}^d$, an empirical null distribution is defined on $\mathcal{E} := \{C_j : C_j < 0; j = 1, \dots, d\} \cup \{-C_j : C_j < 0; j = 1, \dots, d\}$. The aBH procedure defines the approximate p-value of feature j as*

$$p_j := \frac{\sum_{c \in \mathcal{E}} \mathbb{1}(c \geq C_j)}{\text{card}(\mathcal{E}) \vee 1}.$$

Then it applies the BH procedure with the target FDR threshold q to $\{p_j\}_{j=1}^d$ to call discoveries.

S2.3 Enrichment analysis with any numbers of replicates m and n

When $m \neq n$, the BC procedure cannot guarantee FDR control because Lemma 1 no longer holds. To control the FDR in a more general setting ($m = n$ or $m \neq n$), Clipper constructs contrast scores via permutation of replicates across conditions. The idea is that, after permutation, every feature becomes uninteresting and can serve as its own negative control.

Definition S8 (Permutation) *We define σ as permutation, i.e., a bijection from the set $\{1, \dots, m + n\}$ onto itself, and we rewrite the data $\mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{Y}_1, \dots, \mathbf{Y}_d$ into a matrix \mathbf{W} :*

$$\mathbf{W} = \begin{bmatrix} W_{11} & \dots & W_{1m} & W_{1(m+1)} & \dots & W_{1(m+n)} \\ \vdots & & & & & \\ W_{d1} & \dots & W_{dm} & W_{d(m+1)} & \dots & W_{d(m+n)} \end{bmatrix} := \begin{bmatrix} X_{11} & \dots & X_{1m} & Y_{11} & \dots & Y_{1n} \\ \vdots & & & & & \\ X_{d1} & \dots & X_{dm} & Y_{d1} & \dots & Y_{dn} \end{bmatrix}.$$

We then apply σ to permute the columns of \mathbf{W} and obtain

$$\mathbf{W}_\sigma := \begin{bmatrix} W_{1\sigma(1)} & \dots & W_{1\sigma(m)} & W_{1\sigma(m+1)} & \dots & W_{1\sigma(m+n)} \\ \vdots & & & & & \\ W_{d\sigma(1)} & \dots & W_{d\sigma(m)} & W_{d\sigma(m+1)} & \dots & W_{d\sigma(m+n)} \end{bmatrix},$$

from which we obtain the permuted measurements $\{(X_j^\sigma, Y_j^\sigma)\}_{j=1}^d$, where

$$\begin{aligned} X_j^\sigma &:= (W_{j\sigma(1)}, \dots, W_{j\sigma(m)})^\top, \\ Y_j^\sigma &:= (W_{j\sigma(m+1)}, \dots, W_{j\sigma(m+n)})^\top. \end{aligned} \quad (\text{S18})$$

In the enrichment analysis, if two permutations σ and σ' satisfy that

$$\{\sigma(1), \dots, \sigma(m)\} = \{\sigma'(1), \dots, \sigma'(m)\},$$

then we define σ and σ' to be in one equivalence class. That is, permutations in the same equivalence class lead to the same division of $m + n$ replicates (from the two conditions) into two groups with sizes m and n . In total, there are $\binom{m+n}{m}$ equivalence classes of permutations.

We define σ_0 as the identity permutation such that $\sigma_0(i) = i$ for all $i \in \{1, \dots, m+n\}$. In addition, Clipper randomly samples h equivalence classes $\sigma_1, \dots, \sigma_h$ with equal probabilities without replacement from the other $h_{\max} := \binom{m+n}{m} - 1$ equivalence classes (after excluding the equivalence class containing σ_0). Note that h_{\max} is the maximum value h can take.

Clipper then obtains $\{(X_j^{\sigma_0}, Y_j^{\sigma_0}), (X_j^{\sigma_1}, Y_j^{\sigma_1}), \dots, (X_j^{\sigma_h}, Y_j^{\sigma_h})\}_{j=1}^d$, where $(X_j^{\sigma_\ell}, Y_j^{\sigma_\ell})$ are the permuted measurements based on σ_ℓ , $\ell = 0, \dots, h$. Then Clipper computes $T_j^{\sigma_\ell} := t^{\min}_{\text{minus}}(X_j^{\sigma_\ell}, Y_j^{\sigma_\ell})$ to indicate the degree of “interestingness” of feature j reflected by $(X_j^{\sigma_\ell}, Y_j^{\sigma_\ell})$. Note that Clipper chooses t^{\min}_{minus} instead of t^{\max} because empirical evidence shows that t^{\min}_{minus} leads to better power. Sorting $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$ gives

$$T_j^{(0)} \geq T_j^{(1)} \geq \dots \geq T_j^{(h)}.$$

Then Clipper defines the contrast score of feature j , $j = 1, \dots, d$, in one of two ways:

$$C_j := \begin{cases} T_j^{(0)} - T_j^{(1)} & \text{if } T_j^{(0)} = T_j^{\sigma_0} \\ T_j^{(1)} - T_j^{(0)} & \text{otherwise} \end{cases} \quad \text{minus contrast score}, \quad (\text{S19})$$

or

$$C_j := \begin{cases} |T_j^{(0)}| & \text{if } T_j^{(0)} = T_j^{\sigma_0} > T_j^{(1)} \\ 0 & \text{if } T_j^{(0)} = T_j^{(1)} \\ -|T_j^{(0)}| & \text{otherwise} \end{cases} \quad \text{maximum contrast score}. \quad (\text{S20})$$

The intuition behind the contrast scores is that, if $C_j < 0$, then $\mathbb{1}(T_j^{(0)} = T_j^{\sigma_0}) = 0$, which means that at least one of $T_j^{\sigma_1}, \dots, T_j^{\sigma_h}$ (after random permutation) is greater than $T_j^{\sigma_0}$ calculated from the original data (identity permutation), suggesting that feature j is likely an uninteresting feature in enrichment analysis. Motivated by Gimenez and Zou [9], we propose the following procedure for Clipper to control the FDR under the target level $q \in (0, 1)$.

Definition S9 (Gimenez-Zou (GZ) procedure for thresholding contrast scores [9]) Given $h \in \{1, \dots, h_{\max}\}$ and contrast scores $\{C_j\}_{j=1}^d$, $\mathcal{C} = \{|C_j| : C_j \neq 0; j = 1, \dots, d\}$ is defined as the set of non-zero absolute values of C_j 's. The GZ procedure finds a contrast-score cutoff T^{GZ} based on the target FDR threshold $q \in (0, 1)$ as:

$$T^{\text{GZ}} := \min \left\{ t \in \mathcal{C} : \frac{\frac{1}{h} + \frac{1}{h} \text{card}(\{j : C_j \leq -t\})}{\text{card}(\{j : C_j \geq t\}) \vee 1} \leq q \right\} \quad (\text{S21})$$

and outputs $\{j : C_j \geq T^{\text{GZ}}\}$ as discoveries.

Theorem 2 Suppose that the input data that satisfy the Clipper assumptions (S6)–(S7). Then for any $q \in (0, 1)$ and either definition of contrast scores in (S19) or (S20), the contrast-score cutoff T^{GZ} found by the GZ procedure (S21) guarantees that the discoveries have the FDR under q :

$$\text{FDR} = \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j : C_j \geq T^{\text{GZ}}\}) \vee 1} \right] \leq q,$$

where \mathcal{N} denotes the set of uninteresting features.

The proof of Theorem 2 (Supp. Section S8) is similar to that of Theorem 1 and requires two key ingredients: Lemma 2, which is also used in the proof of Theorem 1, and Lemma 3, which is similar to Lemma 1 and is about the properties of signs of $\{C_j\}_{j=1}^d$. The cutoff T^{GZ} can also be viewed as a stopping time of a Bernoulli process.

Lemma 3 For input data that satisfy the Clipper assumptions (S6) and (S7), Clipper constructs contrast scores $\{C_j\}_{j=1}^d$ based on (S20) or (S19). Denote $S_j = \text{sign}(C_j) \in \{-1, 0, 1\}$. Then $\{S_j\}_{j=1}^d$ and $\{C_j\}_{j=1}^d$ satisfy the following properties:

- (a) S_1, \dots, S_d are mutually independent ;
- (b) $\mathbb{P}(S_j = 1) \leq \frac{1}{h+1}$ for all $j \in \mathcal{N}$;
- (c) $\{S_j\}_{j \in \mathcal{N}} \perp \mathcal{C}$.

We note that the GZ procedure is also applicable to the enrichment analysis with equal numbers of replicates, i.e., $m = n$ (Section S2.2). We will compare the GZ procedure against the BC procedure in our results.

S2.4 Differential analysis with $m + n > 2$

For differential analysis, Clipper also uses permutation to construct contrast scores. When $m \neq n$, the equivalence classes of permutations are defined the same as for the enrichment analysis with $m \neq n$. When $m = n$, there is a slight change in the definition of equivalence classes of permutations: if σ and σ' satisfy that

$$\{\sigma(1), \dots, \sigma(m)\} = \{\sigma'(1), \dots, \sigma'(m)\} \text{ or } \{\sigma'(m+1), \dots, \sigma'(2m)\},$$

then we say that σ and σ' are in one equivalence class. In total, there are $h_{\text{total}} := \binom{m+n}{m}$ (when $m \neq n$) or $\binom{2m}{m}/2$ (when $m = n$) equivalence classes of permutations. Hence, to have more than one equivalence class, we cannot perform differential analysis with $m = n = 1$; in other words, the total number of replicates $m + n$ must be at least 3.

Then Clipper randomly samples $\sigma_1, \dots, \sigma_h$ with equal probabilities without replacement from the $h_{\text{max}} := h_{\text{total}} - 1$ equivalence classes that exclude the class containing σ_0 , i.e., the identity permutation. Note that h_{max} is the maximum value h can take. Next, Clipper computes $T_j^{\sigma_\ell} := |t^{\text{minus}}(\mathbf{X}_j^{\sigma_\ell}, \mathbf{Y}_j^{\sigma_\ell})|$, where $\mathbf{X}_j^{\sigma_\ell}$ and $\mathbf{Y}_j^{\sigma_\ell}$ are the permuted data defined in (S18), and it defines C_j as the contrast score of feature j , $j = 1, \dots, d$, in the same ways as in (S19) or (S20).

Same as in the enrichment analysis with $m \neq n$, Clipper also uses the GZ procedure [9] to set a cutoff on contrast scores to control the FDR under the target level $q \in (0, 1)$, following Theorem 2.

S2.5 Clipper variant algorithms

For nomenclature, we assign the following names to Clipper variant algorithms, each of which combines a contrast score definition with a thresholding procedure.

- **Clipper-minus-BC**: minus contrast score $C_j = t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j)$ (S10) and BC procedure (Definition S6);
- **Clipper-minus-aBH**: minus contrast score $C_j = t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j)$ and aBH procedure (Definition S7);
- **Clipper-minus-GZ**: minus contrast score $\tau_j = T_j^{(0)} - T_j^{(1)}$ (S19) and GZ procedure (Definition S9);
- **Clipper-max-BC**: maximum contrast score $C_j = t^{\text{max}}(\mathbf{X}_j, \mathbf{Y}_j)$ (S11) and BC procedure;
- **Clipper-max-aBH**: maximum contrast score $C_j = t^{\text{max}}(\mathbf{X}_j, \mathbf{Y}_j)$ and aBH procedure;
- **Clipper-max-GZ**: maximum contrast score $\tau_j = T_j^{(0)}$ (S20) and GZ procedure.

S2.6 R package “Clipper”

In the R package `Clipper`, the default implementation is as follows. Based on the power comparison results in Section S3 and Figs. S22, S23, S24, and S25, Clipper uses Clipper-minus-BC as the default algorithm for the enrichment analysis with equal numbers of replicates; when there are no discoveries, Clipper suggests users to increase the target FDR threshold q or to use the Clipper-minus-aBH algorithm with the current q . For the enrichment analysis with different numbers of replicates under two conditions or the differential analysis, Clipper uses the Clipper-max-GZ algorithm by default.

S3 Comparison of Clipper variant algorithms

We compared Clipper variant algorithms applicable to each experimental design. Based on the comparison results, we selected a variant algorithm as the default Clipper implementation for each design.

- **1vs1 enrichment analysis.** Under each of the 12 settings, we compared Clipper-minus-BC, Clipper-minus-aBH, Clipper-max-BC, and Clipper-max-aBH (Section S2.5), the only four Clipper variant algorithms applicable to 1vs1 enrichment analysis. The results in Fig. S22 show that, regardless of the contrast scores being minus or maximum (max), the BC procedure always guarantees the FDR control under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$. Notably, in terms of power, the two contrast scores consistently have different advantages under the two background scenarios: Clipper-max-BC has higher power under the homogeneous background, while Clipper-minus-BC is more powerful under the heterogeneous background. Considering that the heterogeneous scenario is prevalent in high-throughput biological data, the minus contrast score is preferred. As the power of Clipper-minus-BC drops when q is too small ($q \leq 3\%$) and d is not too large ($d = 1000$), we consider the aBH procedure as an alternative to control the FDR. The results show that Clipper-minus-aBH is indeed more powerful when Clipper-minus-BC lacks power; however, Clipper-minus-aBH cannot guarantee the exact FDR control as Clipper-minus-BC does. Therefore, Clipper uses **Clipper-minus-BC** by default in 1vs1 enrichment analysis, and it recommends users to increase q when too few discoveries are made; if users reject this option, then Clipper would use Clipper-minus-aBH to increase power for the current q .

- **2vs1 enrichment analysis.** Under each of the 6 settings, we compared Clipper-minus-GZ and Clipper-max-GZ (Section S2.5), the only two Clipper variant algorithms applicable to 2vs1 enrichment analysis. For either algorithm, we further compared two numbers of permutation equivalence classes: $h = 1$ or 2 , where the latter is $h_{\max} = \binom{3}{1} - 1$ —the maximum number of equivalence classes that do not include the identity permutation. Note that h is a required input parameter for the GZ procedure. The results in Fig. S23 show that, regardless of h and the contrast score definition—maximum (max) or minus, the GZ procedure always guarantees the FDR control under all target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$. In terms of power, Clipper-max-GZ($h = 1$) is consistently more powerful than the other three Clipper variants under all settings. Therefore, Clipper uses **Clipper-max-GZ**($h = 1$) by default in enrichment analysis with unequal numbers of replicates under two conditions.

- **3vs3 enrichment analysis.** Under each of the 12 settings, we compared five Clipper variant algorithms: Clipper-minus-BC, Clipper-minus-aBH, Clipper-max-BC, Clipper-max-aBH, and Clipper-max-GZ (Section S2.5). Fig. S24 shows the comparison of the first four variants: regardless of the contrast scores being minus or maximum (max), the BC procedure simultaneously guarantees the FDR control and achieves good power under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$. Similar to the results in the 1vs1 enrichment analysis, Clipper-max-BC has higher power under the homogeneous background, while Clipper-minus-BC is more powerful under the heterogeneous background. By the same reasoning—the prevalent heterogeneous scenarios in high-throughput biological data—we prefer the minus contrast score. Unlike the 1vs1 enrichment analysis, here Clipper-minus-BC is consistently as powerful as Clipper-minus-aBH, even when q is small, but Clipper-minus-aBH cannot guarantee the exact FDR control. Therefore, Clipper-minus-BC achieves the overall best performance among the first four Clipper variants. Given that the GZ procedure is also applicable to this setting, we further compared Clipper-minus-BC with Clipper-max-GZ($h = 1$), the most powerful Clipper variant with the GZ procedure and the default Clipper implementation in the 2vs1 enrichment and differential analyses and the 3vs3 differential analysis. The results in Fig. S26 show that while both **Clipper-minus-BC** and Clipper-max-GZ($h = 1$) control the FDR, the former is more powerful. Hence, we will use Clipper-minus-BC as the default when both conditions have more than one and the same number of replicates.

Under the simulation settings from Gaussian distributions, we also compared Clipper-minus-BC with another Clipper variant using the BC procedure and the t statistic as the contrast score (Clipper-t), where the t statistic is from the two-sample t test. Fig. S13 shows that, although Clipper-t always guarantees the FDR control under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, it has lower power compared to Clipper-minus-BC, our default Clipper for enrichment analysis with equal numbers of replicates. Based on this result, we did not consider the t statistic as an alternative contrast score for Clipper.

- **2vs1 differential analysis.** Similar to 2vs1 enrichment analysis, under each of the 6 settings, we compared Clipper-minus-GZ and Clipper-max-GZ (Section S2.5) with $h = 1$ or 2 . The results in Fig. S23 show that, regardless of h and the contrast score definition—maximum (max) or minus, the GZ procedure always guarantees the FDR control under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$. Notably, in terms of power, Clipper-minus-GZ($h = 2$) is the most powerful when q is very small ($q \leq 2\%$) under Poisson and negative binomial settings, while Clipper-max-GZ($h = 1$) is the most powerful otherwise. Considering that Clipper-max-GZ($h = 1$) outperforms the other three Clipper variants in most cases, Clipper uses **Clipper-max-GZ**($h = 1$) by default in 2vs1 differential analysis, and it recommends users to use Clipper-minus-GZ($h = 2$) when too few discoveries are made.

- **3vs3 differential analysis.** Under each of the 12 settings, we compared Clipper-minus-GZ, and Clipper-max-GZ (Section S2.5) with $h = 1, 3$ or 9 , where $h = 9$ is $h_{\max} = \binom{6}{3}/2 - 1$ —the maximum number of equivalence classes that do not include the identity permutation. The results in Fig. S25 show that, regardless of h and the contrast score definition—maximum (max) or minus, the GZ procedure always guarantees the FDR control under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$. In terms of power, Clipper-max-GZ($h = 1$) is consistently more powerful than the other Clipper variant algorithms under all settings. Therefore, Clipper uses **Clipper-max-GZ($h = 1$)** by default in 3vs3 differential analysis.

Under the simulation settings from Gaussian distributions, we also compared Clipper-max-GZ with another Clipper variant using the GZ procedure and the t statistic to calculate the degree of interestingness (Clipper-t), where the t statistic is from the two-sample t test. Fig. S14 shows that, although Clipper-t always guarantees the FDR control under a range of target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, it has lower power compared to Clipper-max-GZ, our default Clipper for differential analysis. Based on this result, we did not consider the t statistic as an alternative contrast scores for Clipper.

In summary, whenever Clipper-minus-BC is applicable (enrichment analysis with equal number of replicates under two conditions), it is chosen as the default Clipper implementation; otherwise, Clipper-max-GZ($h = 1$) is the default.

S4 Data generation and detailed implementation of the paired approach (a p-value calculation approach) in simulation studies

We describe how we simulated data and how we implemented the paired approach in different simulation settings: 1vs1 enrichment analysis, 2vs1 enrichment analysis, 3vs3 enrichment analysis, 2vs1 differential analysis, and 3vs3 differential analysis, combined with three distribution families (Gaussian, Poisson, and negative binomial) and two background scenarios (homogeneous and heterogeneous). Under some settings, we considered different numbers of features and the existence of outliers.

In each simulation setting, we generated 200 simulated datasets, computed an FDP and an empirical power on each dataset, and averaged the 200 FDPs and 200 empirical powers to approximate the FDR and power, respectively. For notation simplicity, we use $N(\mu, \sigma^2)$ to denote the Gaussian distribution with mean μ and variance σ^2 , $\text{Pois}(\lambda)$ to denote the Poisson distribution with mean λ , and $\text{NB}(\mu, \theta)$ to denote the negative binomial distribution with mean μ and dispersion θ (such that its variance equals $\mu + \theta\mu^2$).

For each design and analysis, we compared the default Clipper implementation with other generic FDR control methods. Specifically, seven generic methods (BH-pool, qvalue-pool, BH-pair-mis, qvalue-pair-mis, BH-pair-2as1, qvalue-pair-2as1, and locfdr-emp) are included in all designs and analyses. The two methods relying on correct model specification, BH-pair-correct and qvalue-pair-correct, are only included in the 3vs3 enrichment and differential analyses, because it is almost impossible to correctly specify a model with fewer than three replicates per condition. The permutation-based method, locfdr-swap, is excluded from the 1vs1 enrichment analysis because it requires at least one condition to have more than one replicate.

In addition to the above designs and analyses, we also compared the default Clipper implementation with BH-pair methods that use parametric or non-parametric tests to calculate p-values when the numbers of replicates are 10 under both conditions for enrichment analysis, i.e., 10vs10 enrichment analysis.

S4.1 1vs1 enrichment analysis

We simulated data with $d = 1000$ and 10,000 features under two background scenarios and three distributional families—a total of 12 settings. In each setting, 10% of the features are interesting ($\mu_{Xj} > \mu_{Yj}$), and the rest are uninteresting (with $\mu_{Xj} = \mu_{Yj}$). Recall that \mathcal{N} denotes the set of uninteresting features.

Gaussian distribution

We simulated data from Gaussian using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 0$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 0$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- We independently generated X_{j1} from $N(\mu_{Xj}, 1)$ and Y_{j1} from $N(\mu_{Yj}, 1)$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that the null distribution of $X_{j1} - Y_{j1}$, $j = 1, \dots, d$ is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{1}{d-1} \sum_{j=1}^d \left(X_{j1} - \frac{1}{d} \sum_{j=1}^d X_{j1} \right)^2 + \frac{1}{d-1} \sum_{j=1}^d \left(Y_{j1} - \frac{1}{d} \sum_{j=1}^d Y_{j1} \right)^2.$$

This is a misspecified model that assumes that μ_{Xj} 's are all equal and so are μ_{Yj} 's. Then we computed the p-value of feature $j = 1, \dots, d$ as the right tail probability of $X_{j1} - Y_{j1}$ in $N(0, \hat{\sigma}^2)$, i.e., $1 - \Phi\left(\frac{X_{j1} - Y_{j1}}{\hat{\sigma}}\right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $N(Y_{j1}, 1)$ conditioning on the observed Y_{j1} as the null distribution of X_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of X_{j1} in $N(Y_{j1}, 1)$, i.e., $1 - \Phi(X_{j1} - Y_{j1})$.

Poisson distribution

We simulated data from Poisson using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{Pois}(20)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- We independently generated X_{j1} from $\text{Pois}(\mu_{Xj})$ and Y_{j1} from $\text{Pois}(\mu_{Yj})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we first defined a log-transformation $f(x) = \log(x + 0.01)$, which we applied to X_{j1} and Y_{j1} , $j = 1, \dots, d$. We assumed that the null distribution of $f(X_{j1}) - f(Y_{j1})$, $j = 1, \dots, d$ is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{1}{d-1} \sum_{j=1}^d \left(f(X_{j1}) - \frac{1}{d} \sum_{j=1}^d f(X_{j1}) \right)^2 + \frac{1}{d-1} \sum_{j=1}^d \left(f(Y_{j1}) - \frac{1}{d} \sum_{j=1}^d f(Y_{j1}) \right)^2.$$

This model misspecifies the Poisson distribution as the log-normal distribution.

Then we computed the p-value of feature $j = 1, \dots, d$ as the right tail probability of $f(X_{j1}) - f(Y_{j1})$ in $N(0, \hat{\sigma}^2)$, i.e., $1 - \Phi\left(\frac{f(X_{j1}) - f(Y_{j1})}{\hat{\sigma}}\right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $\text{Pois}(Y_{j1})$ conditioning on the observed Y_{j1} as the null distribution of X_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of X_{j1} in $\text{Pois}(Y_{j1})$, i.e., $\mathbb{P}(Z \geq X_{j1})$ where $Z \sim \text{Pois}(Y_{j1})$.

Negative binomial distribution

We simulated data from negative binomial using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{NB}(20, 20^{-1})$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- We independently generated X_{j1} from $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ and Y_{j1} from $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that for each uninteresting feature j , Y_{j1} and X_{j1} follow the same Poisson distribution. We calculated the p-value of feature j from a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ using function `poisson.test` in R package `stats`.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $\text{NB}(Y_{j1}, Y_{j1}^{-1})$ conditioning on the observed Y_{j1} as the null distribution of X_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of X_{j1} in $\text{NB}(Y_{j1}, Y_{j1}^{-1})$.

S4.2 2vs1 enrichment analysis

We simulated data with $d = 10,000$ features under two background scenarios and three distributional families—a total of 6 settings. In each setting, 10% of the features are interesting ($\mu_{Xj} > \mu_{Yj}$) and the rest are uninteresting (with $\mu_{Xj} = \mu_{Yj}$). Recall that \mathcal{N} denotes the set of uninteresting features.

Gaussian distribution

We simulated data from Gaussian using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 0$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 0$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j \in \mathcal{N}}$ i.i.d. from $N(0, 2^2)$ and set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. We next generated $\{\mu_{Yj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(0, 2^2)$ and $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- We independently generated X_{j1} and X_{j2} from $N(\mu_{Xj}, 1)$ and Y_{j1} from $N(\mu_{Yj}, 1)$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that the null distribution of $\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}$, $j = 1, \dots, d$, is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{1}{2(2d-1)} \sum_{j=1}^d \sum_{i=1}^2 \left(X_{ji} - \frac{1}{2d} \sum_{j=1}^d \sum_{i=1}^2 X_{ji} \right)^2 + \frac{1}{d-1} \sum_{j=1}^d \left(Y_{j1} - \frac{1}{d} \sum_{j=1}^d Y_{j1} \right)^2.$$

This is a misspecified model that assumes μ_{Xj} 's are all equal and so are μ_{Yj} 's. Then we computed the p-value of feature $j = 1, \dots, d$ as the right tail probability of $\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}$ in $N(0, \hat{\sigma}^2)$, i.e., $1 - \Phi\left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{\hat{\sigma}}\right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $N(Y_{j1}, 1/2)$ conditioning on the observed Y_{j1} as the null distribution of $\frac{1}{2}(X_{j1} + X_{j2})$. Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of $\frac{1}{2}(X_{j1} + X_{j2})$ in $N(Y_{j1}, 1/2)$, i.e., $1 - \Phi\left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{1/\sqrt{2}}\right)$.

Poisson distribution

We simulated data from Poisson using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{Pois}(20)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- We independently generated X_{j1} and X_{j2} from $\text{Pois}(\mu_{Xj})$ and Y_{j1} from $\text{Pois}(\mu_{Yj})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we first defined a log-transformation $f(x) = \log(x + 0.01)$, which we applied to X_{j1} and Y_{j1} , $j = 1, \dots, d$. We assumed that the null distribution of $f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})$, $j = 1, \dots, d$ is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{6}{d-1} \sum_{j=1}^d \left(f(Y_{j1}) - \frac{1}{d} \sum_{j=1}^d f(Y_{j1}) \right)^2.$$

This model misspecifies the Poisson distribution as the log-normal distribution.

Then we computed the p-value of feature $j = 1, \dots, d$ as the right tail probability of $f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})$ in $N(0, \hat{\sigma}^2)$, i.e., $1 - \Phi\left(\frac{f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})}{\hat{\sigma}}\right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we assumed that for each uninteresting feature j , X_{j1} and X_{j2} independently follow $\text{Pois}(Y_{j1})$ conditioning on the observed Y_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ by performing a one-sample Poisson test using the R function `poisson.test` for the null hypothesis $H_0 : \mu_{Xj} = Y_{j1}$ against the alternative hypothesis $H_1 : \mu_{Xj} > Y_{j1}$.

Negative binomial distribution

We simulated data from negative binomial using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{NB}(20, 20^{-1})$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- We independently generated X_{j1} and X_{j2} from $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ and Y_{j1} from $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that for each uninteresting feature j , X_{ji} , $i = 1, 2$ and Y_{j1} follow the same Poisson distribution. We calculated the p-value of feature j from a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ using the function `poisson.test` in R package `stats`.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $\text{NB}(2Y_{j1}, (2Y_{j1})^{-1})$ conditioning on the observed Y_{j1} as the null distribution of $X_{j1} + X_{j2}$. Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of $X_{j1} + X_{j2}$ in $\text{NB}(2Y_{j1}, (2Y_{j1})^{-1})$.

S4.3 3vs3 enrichment analysis

We simulated data with and without outliers under two background scenarios and three distributional families—a total of 12 settings. In each setting, we generated $d = 10,000$ features, among which 10% are interesting (with $\mu_{Xj} > \mu_{Yj}$) and the rest are uninteresting (with $\mu_{Xj} = \mu_{Yj}$). For the results in Fig. S12, we simulated data without outliers under two background scenarios and three distributional families using two more proportions of interesting features: 20% and 40%. The data generation under the Gaussian, Poisson, and negative binomial distributions is the same as the settings with 10% interesting features.

Under the settings with outliers, we generated $\{O_{ji}^X : j = 1, \dots, d; i = 1, \dots, 3\}$ and $\{O_{ji}^Y : j = 1, \dots, d; i = 1, \dots, 3\}$ i.i.d. from $\text{Bernoulli}(0.1)$, where $O_{ji}^X = 1$ or $O_{ji}^Y = 1$ indicates X_{ji} or Y_{ji} is an outlier, respectively. Under settings without outliers, $O_{ji}^X = O_{ji}^Y = 0$ for all $j = 1, \dots, d; i = 1, \dots, 3$.

Gaussian distribution

- Under the homogeneous background scenario, we set $\mu_{Yj} = 0$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 0$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $N(5, 1)$.
- We independently generated X_{ji} from $N(\mu_{Xj}, 1)$ if $O_{ji}^X = 0$ or from the top 1% percentile of $N(\mu_{Xj}, 1)$ if $O_{ji}^X = 1$, $j = 1, \dots, d; i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $N(\mu_{Yj}, 1)$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $N(\mu_{Yj}, 1)$ if $O_{ji}^Y = 1$, $j = 1, \dots, d; i = 1, \dots, 3$.
- For the results in Supp. Fig. S13, under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$, and $\{s_j\}_{j=1}^d$ i.i.d. from a uniform distribution $U(0.5, 2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$

i.i.d. from $N(5, 1)$. We then independently generated X_{ji} from $N(\mu_{Xj}, s_j^2)$ if $O_{ji}^X = 0$ or from the top 1% percentile of $N(\mu_{Xj}, s_j^2)$ if $O_{ji}^X = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $N(\mu_{Yj}, s_j^2)$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $N(\mu_{Yj}, s_j^2)$ if $O_{ji}^Y = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$.

- For the results in Supp. Fig. S11, we generated correlated features. We first selected 10 groups of features (2 groups of interesting features and 8 groups of uninteresting features), with each group containing 200 features. For each group k , we used k_1, \dots, k_{200} to denote the indices of the 200 features within that group and generated $\{X_{ki}\}_{i=1}^{200}$ from a multivariate Gaussian distribution $N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, where $\boldsymbol{\mu}_k = (\mu_{Xk_1}, \dots, \mu_{Xk_{200}})$ and $\boldsymbol{\Sigma}_k$ is a matrix with diagonal entries as 1 and other entries as a fixed correlation. In our simulation, the fixed correlation took two values: 0.2 and 0.4.

To implement the correct paired approach (as in BH-pair-correct and qvalue-pair-correct), we calculated the p-value of feature j from a two-sample t-test with equal variance for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we calculated the p-value of feature j from a two-sample t-test with unequal variance for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we assumed that for each uninteresting feature j , X_{ji} , $i = 1, \dots, 3$ are i.i.d. Gaussian with mean \bar{Y}_j conditioning on the observed \bar{Y}_j and unknown variance. We calculated the p-value of feature j using a one-sample t-test for the null hypothesis $H_0 : \mu_{Xj} = \bar{Y}_j$ against the alternative hypothesis $H_1 : \mu_{Xj} > \bar{Y}_j$.

Poisson distribution

We simulated data from Poisson using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{Pois}(20)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{Pois}(40)$.
- We independently generated X_{ji} from $\text{Pois}(\mu_{Xj})$ if $O_{ji}^X = 0$ or from the top 1% percentile of $\text{Pois}(\mu_{Xj})$ if $O_{ji}^X = 1$, $j = 1, \dots, d$, $i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $\text{Pois}(\mu_{Yj})$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $\text{Pois}(\mu_{Yj})$ if $O_{ji}^Y = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$.

To implement the correct paired approach (as in BH-pair-correct and qvalue-pair-correct), we calculated the p-value of feature j by performing a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ using the function `poisson.test` in R package `stats`.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we first defined a log-transformation $f(x) = \log(x + 0.01)$, which we applied to X_{ji} and Y_{ji} , $j = 1, \dots, d$; $i = 1, \dots, 3$. We assumed that for each uninteresting feature j , $\{f(X_{ji})\}_{i=1}^3$ and $\{f(Y_{ji})\}_{i=1}^3$ follow Gaussian distributions with mean $\mu_{f(Xj)}$ and $\mu_{f(Yj)}$, respectively. Then we computed the p-value of feature j using a two-sample equal variance t-test for the null hypothesis $H_0 : \mu_{f(Xj)} = \mu_{f(Yj)}$ against the alternative hypothesis $H_1 : \mu_{f(Xj)} > \mu_{f(Yj)}$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we assumed that for each uninteresting feature j , $\{X_{ji}\}_{i=1}^3$ follow $\text{Pois}(\bar{Y}_j)$ conditioning on the observed \bar{Y}_j . We calculated the p-value of feature j by performing a one-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \bar{Y}_j$ against the alternative hypothesis $H_1 : \mu_{Xj} > \bar{Y}_j$ using R function `poisson.test` from package `stats`.

Negative binomial distribution

We simulated data from negative binomial using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{NB}(20, 20^{-1})$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For interesting features, we generated $\{\mu_{Xj}\}_{j \notin \mathcal{N}}$ i.i.d. from $\text{NB}(45, 45^{-1})$.
- We independently generated X_{ji} from $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ if $O_{ji}^X = 0$ or from the top 1% percentile of $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ if $O_{ji}^X = 1$, $j = 1, \dots, d$, $i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$ if $O_{ji}^Y = 1$, $j = 1, \dots, d$, $i = 1, \dots, 3$.

To implement the correct paired approach (as in BH-pair-correct and qvalue-pair-correct), we performed a two-sample negative binomial test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative $H_1 : \mu_{Xj} > \mu_{Yj}$ using $T_j := \sum_{i=1}^3 X_{ji} - \sum_{i=1}^3 Y_{ji}$ as the test statistic. We computed the p-value of feature j as the right tail probability

$$\mathbb{P}(T_j \geq t_j) = \sum_{k_1=0}^{\infty} \sum_{k_2=k_1+t_j}^{\infty} \mathbb{P}\left(\sum_{i=1}^3 X_{ji} \geq k_2\right) \mathbb{P}\left(\sum_{i=1}^3 Y_{ji} = k_1\right),$$

where t_j is the realization of T_j , $\mathbb{P}(\sum_{i=1}^3 X_{ji} \geq k_2)$ and $\mathbb{P}(\sum_{i=1}^3 Y_{ji} = k_1)$ can be estimated from the null distribution of X_{ji} and Y_{ji} , $j = 1, \dots, d$; $i = 1, \dots, 3$. As $\sum_{i=1}^3 X_{ji}$ and $\sum_{i=1}^3 Y_{ji}$ follow the same distribution under null, we estimated μ_{Xj} and μ_{Yj} as $\hat{\mu}_{Xj} = \hat{\mu}_{Yj} := (\sum_{i=1}^3 X_{ji} + \sum_{i=1}^3 Y_{ji})/6$. Then, we calculated $\mathbb{P}(\sum_{i=1}^3 X_{ji} \geq k_2)$ and $\mathbb{P}(\sum_{i=1}^3 Y_{ji} = k_1)$ using the estimated distribution of X_{ji} and Y_{ji} as $\text{NB}(\hat{\mu}_{Xj}, (\hat{\mu}_{Xj})^{-1})$ and $\text{NB}(\hat{\mu}_{Yj}, (\hat{\mu}_{Yj})^{-1})$, respectively, $j = 1, \dots, d$; $i = 1, \dots, 3$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that for each uninteresting feature j , $\{X_{ji}\}_{j=1}^3$ and $\{Y_{ji}\}_{j=1}^3$ follow the same Poisson distribution. We calculated the p-value of feature j from a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} > \mu_{Yj}$ using function `poisson.test` in R package `stats`.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $\text{NB}(\sum_{i=1}^3 Y_{ji}, (\sum_{i=1}^3 Y_{ji})^{-1})$ conditioning on the observed $\sum_{i=1}^3 Y_{ji}$ as the null distribution of $\sum_{i=1}^3 X_{ji}$. Then we calculated the p-value of feature $j = 1, \dots, d$ as the right tail probability of $\sum_{i=1}^3 X_{ji}$ in $\text{NB}(\sum_{i=1}^3 Y_{ji}, (\sum_{i=1}^3 Y_{ji})^{-1})$.

S4.4 10vs10 enrichment analysis

We simulated data without outliers under heterogeneous background scenario and three distributional families—a total of 3 settings. In each setting, we generated $d = 10,000$ features, among which 10% are

interesting (with $\mu_{Xj} > \mu_{Yj}$) and the rest are uninteresting (with $\mu_{Xj} = \mu_{Yj}$).

The data generation under the Gaussian, Poisson, and negative binomial distributions is the same as in the 3vs3 enrichment analysis (Section S4.3) except that we set the number of replicates to 10 under each condition, and we did not generate outliers.

The correct paired approaches in BH-pair-parametric are the same as the corresponding BH-pair-correct in the 3vs3 enrichment analysis (Section S4.3) except that, under the negative binomial distribution, the test statistic T_j and its null distribution should have the number of replicates changed from 3 to 10. The misspecified and 2as1 paired approaches (BH-pair-mis and BH-pair-2as1) are also the same as the corresponding approaches in the 3vs3 enrichment analysis (Section S4.3).

To implement the non-parametric paired approaches, we calculated the p-value of feature j from the one-sided two-sample Wilcoxon rank-sum test (using R function `wilcox.test` in package `stats`) in BH-pair-Wilcoxon and from the one-sided two-sample permutation test (using R function `oneway.test` in package `coin`) in BH-pair-permutation.

S4.5 2vs1 differential analysis

We simulated data with $d = 10,000$ features under two background scenarios and three distributional families—a total of 6 settings. In each setting, we set 10% features as “up-regulated” with $\mu_{Xj} > \mu_{Yj}$ and another 10% features as “down-regulated” with $\mu_{Xj} < \mu_{Yj}$.

Gaussian distribution

We simulated data from Gaussian using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 0$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 0$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $N(5, 1)$. For down-regulated features, generated μ_{Xj} i.i.d. from $N(-5, 1)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $N(5, 1)$. For down-regulated features, generated μ_{Xj} i.i.d. from $N(-5, 1)$.
- We independently generated X_{j1} and X_{j2} from $N(\mu_{Xj}, 1)$ and Y_{j1} from $N(\mu_{Yj}, 1)$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that the null distribution of $\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}$, $j = 1, \dots, d$, is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{1}{2(2d-1)} \sum_{j=1}^d \sum_{i=1}^2 \left(X_{ji} - \frac{1}{2d} \sum_{j=1}^d \sum_{i=1}^2 X_{ji} \right)^2 + \frac{1}{d-1} \sum_{j=1}^d \left(Y_{j1} - \frac{1}{d} \sum_{j=1}^d Y_{j1} \right)^2.$$

This is a misspecified model assuming that μ_{Xj} 's are all equal and so are μ_{Yj} 's. Then we computed the p-value of feature $j = 1, \dots, d$ as the two-sided tail probability of $\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}$ in $N(0, \hat{\sigma}^2)$, i.e., $2 \cdot \min \left(1 - \Phi \left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{\hat{\sigma}} \right), \Phi \left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{\hat{\sigma}} \right) \right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $N(Y_{j1}, 1)$ conditioning on the observed Y_{j1} as the null distribution of X_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ as the two-sided tail probability of $\frac{1}{2}(X_{j1} + X_{j2})$ in $N(Y_{j1}, 1/2)$, i.e., $2 \cdot \min \left(1 - \Phi \left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{1/\sqrt{2}} \right), \Phi \left(\frac{\frac{1}{2}(X_{j1} + X_{j2}) - Y_{j1}}{1/\sqrt{2}} \right) \right)$.

Poisson distribution

We simulated data from Poisson using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(60)$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(5)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{Pois}(20)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(60)$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(5)$.
- We independently generated X_{j1} and X_{j2} from $\text{Pois}(\mu_{Xj})$ and Y_{j1} from $\text{Pois}(\mu_{Yj})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we first defined a log-transformation $f(x) = \log(x + 0.01)$, which we applied to X_{j1} and Y_{j1} , $j = 1, \dots, d$. We assumed that the null distribution of $f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})$, $j = 1, \dots, d$ is $N(0, \hat{\sigma}^2)$, where

$$\hat{\sigma}^2 = \frac{6}{d-1} \sum_{j=1}^d \left(f(Y_{j1}) - \frac{1}{d} \sum_{j=1}^d f(Y_{j1}) \right)^2.$$

This model misspecifies the Poisson distribution as the log-normal distribution. Then we computed the p-value of feature $j = 1, \dots, d$ as the two-sided tail probability of $f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})$ in $N(0, \hat{\sigma}^2)$, i.e., $2 \cdot \min \left(1 - \Phi \left(\frac{f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})}{\hat{\sigma}} \right), \Phi \left(\frac{f(X_{j1}) + f(X_{j2}) - 2f(Y_{j1})}{\hat{\sigma}} \right) \right)$, where Φ is the cumulative distribution function of $N(0, 1)$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we assumed that for each uninteresting feature j , X_{j1} and X_{j2} independently follow $\text{Pois}(Y_{j1})$ conditioning on the observed Y_{j1} . Then we calculated the p-value of feature $j = 1, \dots, d$ by performing a one-sample Poisson test using the R function `poisson.test` for the null hypothesis $H_0 : \mu_{Xj} = Y_{j1}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq Y_{j1}$.

Negative binomial distribution

We simulated data from negative binomial using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 30$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 30$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(70, 70^{-1})$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(7, 7^{-1})$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{NB}(30, 30^{-1})$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(70, 70^{-1})$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(7, 7^{-1})$.
- We independently generated X_{j1} and X_{j2} from $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ and Y_{j1} from $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$, $j = 1, \dots, d$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that for each uninteresting feature j , X_{j1} , X_{j2} , and Y_{j1} follow the same Poisson distribution. We calculated the p-value of feature j from a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$ using the function `poisson.test` in R package `stats`.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $NB(2Y_{j1}, (2Y_{j1})^{-1})$ conditioning on the observed Y_{j1} as the null distribution of $X_{j1} + X_{j2}$. Then we calculated the p-value of feature $j = 1, \dots, d$ as the two-sided tail probability of $X_{j1} + X_{j2}$ in $NB(2Y_{j1}, (2Y_{j1})^{-1})$, i.e., twice the smaller of the left-tail and right-tail probabilities.

S4.6 3vs3 differential analysis

We simulated data with or without outliers under two background scenarios and three distributional families—a total of 12 settings. In each setting, we generated $d = 10,000$ features, among which 10% features were “up-regulated features” with $\mu_{Xj} > \mu_{Yj}$ and another 10% were “down-regulated features” with $\mu_{Xj} < \mu_{Yj}$.

Under the settings with outliers, we generated $\{O_{ji}^X : j = 1, \dots, d; i = 1, \dots, 3\}$ and $\{O_{ji}^Y : j = 1, \dots, d; i = 1, \dots, 3\}$ i.i.d. from Bernoulli(0.1), where $O_{ji}^X = 1$ or $O_{ji}^Y = 1$ indicates X_{ji} or Y_{ji} is an outlier, respectively. Under settings without outliers, $O_{ji}^X = O_{ji}^Y = 0$ for all $j = 1, \dots, d; i = 1, \dots, 3$.

Gaussian distribution

- Under the homogeneous background scenario, we set $\mu_{Yj} = 0$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 0$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $N(5, 1)$. For down-regulated features, generated μ_{Xj} i.i.d. from $N(-5, 1)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $N(5, 1)$. For down-regulated features, generated μ_{Xj} i.i.d. from $N(-5, 1)$.
- We independently generated X_{ji} from $N(\mu_{Xj}, 1)$ if $O_{ji}^X = 0$ or from the top 1% percentile of $N(\mu_{Xj}, 1)$ if $O_{ji}^X = 1$, $j = 1, \dots, d; i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $N(\mu_{Yj}, 1)$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $N(\mu_{Yj}, 1)$ if $O_{ji}^Y = 1$, $j = 1, \dots, d; i = 1, \dots, 3$.
- For the results in Supp. Fig. S14, under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $N(0, 2^2)$, and $\{s_j\}_{j=1}^d$ i.i.d. from a uniform distribution $U(0.5, 2)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $N(5, 1)$. For down-regulated features, generated μ_{Xj} i.i.d. from $N(-5, 1)$. We then independently generated X_{ji} from $N(\mu_{Xj}, s_j^2)$ if $O_{ji}^X = 0$ or from the top 1% percentile of $N(\mu_{Xj}, s_j^2)$ if $O_{ji}^X = 1$, $j = 1, \dots, d; i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $N(\mu_{Yj}, s_j^2)$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $N(\mu_{Yj}, s_j^2)$ if $O_{ji}^Y = 1$, $j = 1, \dots, d; i = 1, \dots, 3$.

To implement the correct paired approach (as in BH-pair-correct and qvalue-pair-correct), we calculated the p-value of feature j from a two-sample t-test with equal variance for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we calculated the p-value of feature j from a two-sample t-test with unequal variance for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we treated $N(\bar{Y}_j, 1)$ conditioning on observed \bar{Y}_j as the null distribution of X_{ji} , $i = 1, \dots, 3$. We calculated the p-value of feature j using a one-sample t-test for the null hypothesis $H_0 : \mu_{Xj} = \bar{Y}_j$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \bar{Y}_j$.

Poisson distribution

We simulated data from Poisson using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 20$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 20$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(40)$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(5)$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{Pois}(20)$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(40)$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{Pois}(5)$.
- We independently generated X_{ji} from $\text{Pois}(\mu_{Xj})$ if $O_{ji}^X = 0$ or from the top 1% percentile of $\text{Pois}(\mu_{Xj})$ if $O_{ji}^X = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $\text{Pois}(\mu_{Yj})$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $\text{Pois}(\mu_{Yj})$ if $O_{ji}^Y = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$.

To implement the correct paired approach (as in BH-pair-correct and qvalue-pair-correct), we calculated the p-value of feature j by performing a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$ using function `poisson.test` in R package `stats`.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we first defined a log-transformation $f(x) = \log(x + 0.01)$, which we applied to X_{ji} and Y_{ji} , $j = 1, \dots, d$; $i = 1, \dots, 3$. We assumed that for each uninteresting feature j , $\{f(X_{ji})\}_{i=1}^3$ and $\{f(Y_{ji})\}_{i=1}^3$ follow Gaussian distributions with mean $\mu_{f(Xj)}$ and $\mu_{f(Yj)}$, respectively. Then we computed the p-value of feature j using a two-sample equal variance t-test for the null hypothesis $H_0 : \mu_{f(Xj)} = \mu_{f(Yj)}$ against the alternative hypothesis $H_1 : \mu_{f(Xj)} \neq \mu_{f(Yj)}$.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we assumed that for each uninteresting feature j , $\{X_{ji}\}_{i=1}^3$ follow $\text{Pois}(\bar{Y}_j)$ conditioning on the observed \bar{Y}_j . We calculated the p-value of feature j by performing a one-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \bar{Y}_j$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \bar{Y}_j$ using the function `poisson.test` in R package `stats`.

Negative binomial distribution

We simulated data from negative binomial using the following procedure:

- Under the homogeneous background scenario, we set $\mu_{Yj} = 30$ for all d features. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj} = 30$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(70, 70^{-1})$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(7, 7^{-1})$.
- Under the heterogeneous background scenario, we generated $\{\mu_{Yj}\}_{j=1}^d$ i.i.d. from $\text{NB}(30, 30^{-1})$. For uninteresting features, we set $\mu_{Xj} = \mu_{Yj}$ for $j \in \mathcal{N}$. For up-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(70, 70^{-1})$. For down-regulated features, we generated μ_{Xj} i.i.d. from $\text{NB}(7, 7^{-1})$.
- We independently generated X_{ji} from $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ if $O_{ji}^X = 0$ or from the top 1% percentile of $\text{NB}(\mu_{Xj}, \mu_{Xj}^{-1})$ if $O_{ji}^X = 1$, $j = 1, \dots, d$, $i = 1, \dots, 3$. Similarly, we independently generated Y_{ji} from $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$ if $O_{ji}^Y = 0$ or from the top 1% percentile of $\text{NB}(\mu_{Yj}, \mu_{Yj}^{-1})$ if $O_{ji}^Y = 1$, $j = 1, \dots, d$; $i = 1, \dots, 3$.

To implement the correct paired approach with unknown dispersion (as in BH-pair-correct and qvalue-pair-correct), we performed a two-sample negative binomial test for the null hypothesis $H_0 :$

$\mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$ using the coefficient from the negative binomial regression as the test statistic. Specifically, for each feature j we performed a negative binomial regression by treating the condition labels as a categorical covariate and feature j 's measurements as the response. We implemented this regression analysis using function `glm.nb` in R package `MASS` and extracted the p-value of the coefficient as the p-value of feature j . The dispersion parameter was not pre-specified but estimated by `glm.nb`.

To implement the correct paired approach with known dispersion, we performed a similar negative binomial regression but with the pre-specified dispersion parameter 30^{-1} for each feature j . Then we computed the feature's p-value as the p-value of the coefficient of the condition covariate. We implemented this regression analysis using function `glm` in R package `stats`.

To implement the misspecified paired approach (as in BH-pair-mis and qvalue-pair-mis), we assumed that for each uninteresting feature j , $\{X_{ji}\}_{j=1}^3$ and $\{Y_{ji}\}_{j=1}^3$ follow the same Poisson distribution. We calculated the p-value of feature j from a two-sample Poisson test for the null hypothesis $H_0 : \mu_{Xj} = \mu_{Yj}$ against the alternative hypothesis $H_1 : \mu_{Xj} \neq \mu_{Yj}$ using function `poisson.test` in R package `stats`.

To implement the 2as1 paired approach (as in BH-pair-2as1 and qvalue-pair-2as1), we first used function `glm.nb` in R package `MASS` to estimate $\hat{\mu}_{Yj}$ and $\hat{\theta}_{Yj}$ from $\{Y_{ji}\}_{j=1}^3$. Then we computed the p-value of feature j by treating $NB(3\hat{\mu}_{Yj}, (3\hat{\theta}_{Yj})^{-1})$ as the null distribution of $\sum_{i=1}^3 X_{ji}$ and calculated its two-sided tail probability, i.e., twice the smaller of the left-tail and right-tail probabilities.

S5 Bioinformatic methods with FDR control functionality

S5.1 Peak calling methods for ChIP-seq data

MACS2 MACS2 [10] uses sliding windows with a fixed length across the genome and identifies peaks by using a Poisson distribution to model the read counts within each window, which has one read count per replicate. Specifically, for each region (which is combined from sliding windows), MACS2 performs a one-sample Poisson test to calculate a p-value, where the null distribution is set to be Poisson with its parameter estimated from the background. By thresholding p-values, MACS2 identifies a set of candidate peaks. It also estimates for each candidate peak a q-value by swapping the experimental sample with the background (negative control) sample, and the q-values are used for FDR control. We used MACS2 software (version 2.2.6) with its default settings.

HOMER We used `findPeaks`, a program in HOMER [11], to perform peak calling on ChIP-seq data. The p-value calculation in `findPeaks` is similar to that in MACS2; that is, `findPeaks` also uses the Poisson distribution as the null distribution of read counts in a genomic region, and it also estimates the Poisson mean from the background sample. Then `findPeaks` identifies peaks by setting thresholds on p-values and fold-changes (the fold change of a region is defined as the observed read count under the experimental sample divided by the estimated Poisson mean from the background sample). We used `findPeaks` version 3.1.9.2.

S5.2 SEQUEST for peptide identification from MS data

SEQUEST SEQUEST uses probability-based scoring to identify PSMs from mass-spectrometry data. We ran SEQUEST in Proteome Discoverer 2.3.0.523 (ThermoScientific) with the following settings: 10 ppm precursor tolerance; 0.6 Da fragment tolerance; static modifications: methylthio (C); dynamic modifications: deamination (NQ), oxidation (M). We ran Percolator [12] in conjunction with SEQUEST with

the target/decoy selection mode set to “separate.” For SEQUEST, for a range of target FDR thresholds ($q \in \{1\%, 2\%, \dots, 10\%\}$), we identified the target PSMs with SEQUEST q-values no greater than q as discoveries. To prepare the input for Clipper, we set peptide and protein FDRs to 100% to obtain the entire lists of target PSMs and decoy PSMs with their SEQUEST q-values.

S5.3 Differentially expressed gene (DEG) methods for bulk RNA-seq data

edgeR edgeR models each gene's read counts by using a negative binomial regression, where the condition is incorporated as an indicator covariate, and the condition's coefficient represents the gene-wise differential expression effect [13]. We used R package edgeR version 3.30.0.

DESeq2 DESeq2 uses a similar negative binomial regression as edgeR to model each gene's read counts under two conditions. DESeq2 differs from edgeR mainly in their estimation of the dispersion parameter in the negative binomial distribution [14]. We used R package DESeq2 version 1.28.1.

S5.4 Differentially expressed gene (DEG) methods for scRNA-seq data

MAST MAST models each gene's log read counts (TPM) by using a two-part generalized regression model. Each gene's expression rate was modeled using logistic regression and, conditioning on a cell expressing the gene, the gene's expression level was modeled as Gaussian [15]. We used R package MAST version 1.14.0.

Monocle3 Monocle3 uses a generalized linear model to model each gene's normalized expression value, with other information included as covariates (time, treatment, and so on) [16]. We used R package monocle3 version 0.2.3.0.

S5.5 Differentially interacting chromatin regions (DIR) methods for Hi-C data

MultiHiCcompare MultiHiCcompare relies on a non-parametric method to jointly normalize multiple Hi-C interaction matrices [17]. It uses a generalized linear model to detect DIRs. MultiHiCcompare is an extension of the HiCcompare package [18]. We used R package multiHiCcompare version 1.6.0.

diffHic diffHic uses the statistical framework of the edgeR package to model biological variability and to test for significant differences between conditions [19]. We used R package diffHic version 1.20.0.

FIND FIND uses a spatial Poisson process to detect chromosomal regions that display a significant difference between two regions' contact intensity and their neighbouring contact intensities [20]. We used R package FIND version 0.99.

S6 Benchmark data generation in omics data applications

S6.1 ChIP-seq data with synthetic spike-in peaks

We used two control samples (which we refer to as Control 1 and Control 2) from H3K4me3 ChIP-seq data in Chromosome 1 of the cell line GM12878 [21].

- (i) We created two semi-synthetic experimental samples by adding synthetic true peaks to Control 1. To mimic real H3K4me3 ChIP-seq data, where peaks are located predominantly in promoter regions, we added synthetic true peaks to promoter regions annotated from Ensembl BioMart (Ensemble hg 19, regulation 104) [22]. Specifically, we randomly sampled 585 genes' promoter regions from Chromosome 1. We then used ChIPulate to simulate reads from these promoter regions (for each simulation, extraction efficiency parameter and PCR efficiency parameter were randomly sampled from a uniform distribution between 0 to 1; binding energy parameters were randomly sampled from a uniform distribution between 0 and 2; sequencing depth parameter was set to 50). Then we added the simulated reads to Control 1. We repeated this procedure for twice to obtain two semi-synthetic experimental samples (i.e., two replicates under the experimental condition).
- (ii) We repeated Step (i) for 20 times to generate 20 sets of semi-synthetic experimental samples. For each set of experimental samples, we paired them with Control 2, which was treated as the background sample (i.e., one replicate under the background condition). Hence, we obtained 20 semi-synthetic ChIP-seq datasets, each containing 585 synthetic true peaks.
- (iii) After applying a peak calling method to these 20 semi-synthetic datasets, we evaluated the method's 20 FDPs and 20 empirical power, which were then averaged as the method's approximate FDR and power. In the evaluation, a called peak was a true positive if it overlapped with a synthetic true peak; otherwise, it was a false positive.

S6.2 Real MS benchmark data

We purchased the complex proteomics standard (CPS) (part number 400510) from Agilent (Agilent, Santa Clara, CA, USA). The CPS contains soluble proteins extracted from the archaeon *Pyrococcus furiosus* (*Pfu*), which has a complete protein database; that is, all proteins from *Pfu* were catalogued into its protein database with known protein sequences. We subjected the CPS to a shotgun proteomics analysis that generated mass spectra of *Pfu*.

To generate a benchmark dataset, we first generated a reference database by concatenating the Uniprot *Pyrococcus furiosus* (*Pfu*) database, the Uniprot Human database, and two contaminant databases: the CRAPome [23] and the contaminant database from MaxQuant [24]. During the process, we purified the reference database by first performing *in silico* digestion of *Pfu* proteins and then removing human proteins that contained *Pfu* peptides from the reference database. We then input the *Pfu* mass spectra (from the CPS) and the purified reference database into SEQUEST. We considered a target PSM as true if SEQUEST reported its protein as from *Pfu* or the two contaminants; otherwise (if from Human), we considered the target PSM as false. The *in silico* digestion was performed in Python using the `pyteomics.parser` function from `pyteomics` with the following settings: Trypsin digestion, two allowed missed cleavages, minimum peptide length of six [25, 26].

S6.3 Bulk RNA-seq data with synthetic spike-in DEGs

We generated four sets of realistic semi-synthetic data from two real RNA-seq datasets. The first one is a human monocyte RNA-seq dataset including 17 samples of classical monocytes and 17 samples of non-classical monocytes [27]. Each sample contains expression levels of $d = 52,376$ transcripts.

The second one is a yeast RNA-seq dataset including 48 samples of a *snf2* knockout mutant cell line and 48 samples of negative control (without the knockout) [28]. Each sample contains expression levels of $d = 7126$ genes. We preprocessed this dataset by removing low-quality replicates (replicates 6, 13, 25, 35 from the knockout; replicates 21, 22, 25, 28, 34, 36 from the control) identified by the original paper Gierliński et al. [28], leaving us with 44 replicates under the knockout condition and 42 replicates under the negative control.

Here we describe our **simulation strategy 1**. Given either the human monocyte dataset or the yeast dataset, we performed the following steps.

- (i) We first performed normalization on all samples across two conditions using the edgeR normalization method trimmed mean of M-values (TMM) [29]. We denote the resulting normalized read count matrix of classical human monocytes or yeasts without the knockout by \mathbf{X}^1 and the normalized read count matrix of non-classical human monocytes or yeast with the knockout by \mathbf{X}^2 , respectively. Following the convention in bioinformatics, the columns and rows of \mathbf{X}^1 and \mathbf{X}^2 represent biological samples and genes, respectively.
- (ii) To define true DEGs, we first computed the fold change of gene j by $FC_j = [(\bar{\mathbf{X}}_j^2 + 1)/(\bar{\mathbf{X}}_j^1 + 1)]$ for $j = 1, \dots, d$, where \mathbf{X}_j^1 and \mathbf{X}_j^2 denote the j -th row vector of \mathbf{X}^1 and \mathbf{X}^2 respectively and $\bar{\cdot}$ denotes the average of elements in a vector. We added the pseudo-count of 1 to avoid division by 0. We defined true DEGs as those with $|\log_2 FC_j| \geq 4$ for the human monocyte dataset and with $|\log_2 FC_j| \geq 1.5$ for the yeast dataset, resulting 191 true human DEGs (transcripts) and 152 true yeast DEGs.
- (iii) We generated semi-synthetic data with 3 samples under both the experimental and background conditions, a typical design in bulk RNA-seq experiments. Specifically, if gene j is a true DEG, we randomly sampled without replacement 3 values from \mathbf{X}_j^1 as counts under the experimental condition, and another 3 values from \mathbf{X}_j^2 as counts under the background condition. If gene j is not a true DEG, we randomly sampled 6 values without replacement from $(\mathbf{X}_j^1, \mathbf{X}_j^2)$ and randomly split them into 3 and 3 counts under two conditions. Doing so guaranteed that a non-DEG's read counts are i.i.d. regardless of condition.
- (iv) We repeated Step (iii) for 100 times to generate 100 semi-synthetic datasets.

Next, we describe our **simulation strategy 2**. Let us now re-use notations \mathbf{X}^1 to denote the original read count matrix of classical human monocytes or yeast without the knockout, and \mathbf{X}^2 to denote the original read count matrix of non-classical human monocytes or yeast with the knockout. Both \mathbf{X}^1 and \mathbf{X}^2 have rows as genes or transcripts and columns as biological samples. Given either the human monocyte dataset or the yeast dataset, we performed the following steps.

- (i) We first identified genes whose read counts are positive in all samples under both conditions and denote the number of such genes by d_p . Then from these identified genes, we randomly sampled without replacement $\min(d_p, 0.3d)$ genes as true DEGs. The remaining $d - \min(d_p, 0.3d)$ genes were considered true non-DEGs.
- (ii) To generate fold changes of true DEGs, we first computed the fold change of gene j by $FC_j = [(\bar{\mathbf{X}}_j^2 + 1)/(\bar{\mathbf{X}}_j^1 + 1)]$ for $j = 1, \dots, d$, where \mathbf{X}_j^1 and \mathbf{X}_j^2 denote the j -th row vector of \mathbf{X}^1 and \mathbf{X}^2

respectively and $\bar{\cdot}$ denotes the average of elements in a vector. Let \mathcal{W} denote $\{\text{FC}_j : \text{FC}_j \geq 16, j = 1, \dots, d\}$ for the human monocyte dataset and $\{\text{FC}_j : \text{FC}_j \geq 1.5, j = 1, \dots, d\}$ for the yeast dataset. We then sorted unique elements in \mathcal{W} and denoted them by $w_{(1)} < \dots < w_{(n_u)}$, where n_u is the number of unique elements in \mathcal{W} . To generate a fold change of a true DEG, say gene j , we randomly generated an integer v with equal probability from $\{1, \dots, n_u - 1\}$ and a value p from $\text{Uniform}(0, 1)$. Then we calculated the fold change as $R_j = w_{(v)} + p(w_{(v+1)} - w_{(v)})$. Using this approach, generated the fold changes independently for all true DEGs.

- (iii) Next, we randomly sampled 6 replicates without replacement from \mathbf{X}^2 and split them into two groups of 3 replicates. We denote the resulting matrices as $\tilde{\mathbf{X}}^1$ and $\tilde{\mathbf{X}}^2$, whose j -th rows are denoted respectively by $\tilde{\mathbf{X}}_j^1$ and $\tilde{\mathbf{X}}_j^2$. If gene j is a true DEG, we generated U_j from $\text{Bernoulli}(1/2)$. Then we set gene j 's expression levels under the two conditions to $R_j \tilde{\mathbf{X}}_j^1$ and $\tilde{\mathbf{X}}_j^2$ if $U_j = 0$ or $\tilde{\mathbf{X}}_j^1$ and $R_j \tilde{\mathbf{X}}_j^2$ if $U_j = 1$. If gene j is not a true DEG, its expression levels under the two conditions would remain unchanged, i.e., $\tilde{\mathbf{X}}_j^1$ and $\tilde{\mathbf{X}}_j^2$. Such data generation strategy has no guarantee of i.i.d. read counts for non-DEGs if the samples in \mathbf{X}^2 have batch effects.

- (iv) We repeated Step (iii) for 100 times to generate 100 semi-synthetic datasets.

The human monocyte RNA-seq dataset is available in the NCBI Sequence Read Archive (SRA) under accession number SRP082682 (<https://www.ncbi.nlm.nih.gov/Traces/study/?acc=srp082682>). The yeast RNA-seq data is available in the European Nucleotide Archive (ENA) archive with project ID PRJEB5348 (<https://www.ebi.ac.uk/ena/browser/view/PRJEB5348>).

S6.4 Single-cell RNA-seq data with synthetic spike-in DEGs

We used scDesign2, a flexible probabilistic simulator to generate realistic scRNA-seq count data with gene correlations captured [30]. Using scDesign2, we generated two sets of semi-synthetic data from two peripheral blood mononuclear cell (PBMC) real datasets [31]: one generated using the 10x Genomics protocol [32] and the other using Drop-seq [33]. Each synthetic dataset contains two types of cells: CD4+ T cells, and cytotoxic T cells, which we treated as two conditions. Starting with the real data generated using either 10x Genomics or Drop-seq, we used the following steps to generate synthetic scRNA-seq data.

- (i) First, we fit the real data count matrices using R function `fit_model_scDesign2` for each cell type by specifying the underlying distribution of each gene as negative binomial. Denote the resulting marginal distributions of gene j as $NB(\hat{\mu}_{j1}, \hat{\theta}_{j1})$ for CD4+ T cells and $NB(\hat{\mu}_{j2}, \hat{\theta}_{j2})$ for cytotoxic T cells, $j = 1, \dots, d$. The gene-gene correlations with each cell type were fitted using a copula model.
- (ii) Let \mathbf{X}^{cd4} and \mathbf{X}^{cyto} denote the read count matrices of CD4+ T cells and cytotoxic T cells. To define true DEGs, we first computed the log fold change of gene j by $\log\text{FC}_j = \log_2 [(\bar{\mathbf{X}}_j^{\text{cd4}} + 1)/(\bar{\mathbf{X}}_j^{\text{cyto}} + 1)]$ for $j = 1, \dots, d$, where $\mathbf{X}_j^{\text{cd4}}$ and $\mathbf{X}_j^{\text{cyto}}$ denote the j -th row vector of \mathbf{X}^{cd4} and \mathbf{X}^{cyto} respectively and $\bar{\cdot}$ denotes the average of elements in a vector. We then selected 1000 genes with the largest absolute fold changes as true DEGs and kept the remaining ones as true non-DEGs.
- (iii) We simulated the semi-synthetic datasets using R function `simulate_count_scDesign2`. Specifically, we set the number of synthetic cells generated by scDesign2 equal to the number of real cells for each cell type. If a gene j is a true DEG, we specify its marginal distributions under the two conditions as $NB(\hat{\mu}_{j1}, \hat{\theta}_{j1})$ and $NB(\hat{\mu}_{j2}, \hat{\theta}_{j2})$ respectively. If a gene j is a true non-DEG, we specify its marginal distribution under both conditions as $NB((\hat{\mu}_{j1} + \hat{\mu}_{j2})/2, (\hat{\theta}_{j1} + \hat{\theta}_{j2})/2)$. We

used the fitted copula models from the two cell types to generate genes' (correlated) expression read counts.

(iv) We repeated Step (iii) for 200 times to generate 200 semi-synthetic datasets.

Both `fit_model_scDesign2` and `simulate_count_scDesign2` come from R package `scDesign2` [30]. The 10x Genomic PBMC dataset and the Drop-seq PBMC dataset are available from the Gene Expression Omnibus (GEO) with accession number GSE132044 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE132044>) and the Single Cell Portal with accession numbers SCP424 (https://singlecell.broadinstitute.org/single_cell/study/SCP424/single-cell-comparison-pbmc-data).

S6.5 Hi-C data with synthetic spike-in DIRs

The real Hi-C interaction matrix contains the pairwise contact intensities of 250 binned genomic regions in Chromosome 1. It is from the cell line GM12878 and available in the NCBI Gene Expression Omnibus (GEO) under accession number GSE63525. We denote the real interaction matrix as \mathbf{X}^{real} . Because \mathbf{X}^{real} is symmetric, we only focus on its upper triangular part.

- (i) Among the $(250 \times 250 - 250)/2 = 31,125$ upper triangular entries (i.e., region pairs), we selected 404 entries as true up-regulated DIRs, and 550 entries as true down-regulated DIRs (Fig. S29).
- (ii) Next, for the (i, j) -th entry, we generated a log fold change, denoted by f_{ij} , between the two conditions as follows. We simulated f_{ij} from truncated Normal($100/|i - j|, 0.5^2$) with support $[0.05, \infty)$ if the (i, j) -th entry is up-regulated, or from truncated Normal($-100/|i - j|, 0.5^2$) with support $(-\infty, -0.05]$ if the (i, j) -th entry is down-regulated; if the (i, j) -th entry is not differential, we set $f_{ij} = 0$.
- (iii) Then we specify the mean measurement of the (i, j) -th entry under the two conditions as $\mu_{X_{ij}} = [\mathbf{X}^{\text{real}}]_{ij}$ and $\mu_{Y_{ij}} = [\mathbf{X}^{\text{real}}]_{ij} \cdot e^{f_{ij}}$, respectively.
- (iv) We generated synthetic read counts of the (i, j) -th entry from $\text{NB}(\mu_{X_{ij}}, 1000^{-1})$ and $\text{NB}(\mu_{Y_{ij}}, 1000^{-1})$ respectively under the two conditions.
- (v) We repeated Step (iv) for 200 times to generate 200 semi-synthetic datasets.

S7 Implementation of Clipper in omics data applications

Below we briefly introduce the implementation of Clipper in the four omics data applications. All the results were obtained by running using R package `Clipper` (see package vignette for details: <https://github.com/JSB-UCLA/Clipper/blob/master/vignettes/Clipper.pdf>).

S7.1 Peak calling from ChIP-seq data

- (i) We consider each genomic location, i.e., a base pair, as a feature and each ChIP-seq sample as a replicate under the experimental or background condition. Then we consider the read count of each location in each sample as the corresponding feature's measurement. Doing so, we summarized ChIP-seq data into a $d \times (m + n)$ matrix, where d is the number of locations, and m and n are the numbers of experimental and control samples, respectively. We then applied Clipper to perform an enrichment analysis to obtain the contrast score C_j for each location j . In our study, $m = n = 1$, so the default Clipper implementation is Clipper-minus-BC.

- (ii) For any target FDR threshold q , Clipper gives a cutoff T_q on contrast scores.
- (iii) We then used existing peak calling methods, e.g., MACS2 and HOMER, to call candidate peaks with the least stringent q-value cutoff. For example, when we used MACS2, we set the q-value cutoff as 1.
- (iv) We computed the contrast score of each candidate peak as the median of the contrast scores of all the locations within.
- (v) The candidate peaks with contrast scores greater than or equal to T_q are called discoveries.

S7.2 Peptide identification from mass spectrometry data

- (i) We consider each mass spectrum as a feature and its target/decoy PSM as a replicate under the experimental/background condition respectively. Then we consider $-\log_{10}(\text{q-value} + 0.01)$ as the measurement of each PSM, where the q-value is output by SEQUEST. Doing so, we summarized the SEQUEST output into a $d \times (m + n)$ matrix, where d is the number of mass spectra, and m and n are the numbers of experimental and control samples, respectively. We then applied Clipper to perform an enrichment analysis to obtain a contrast score C_j for each mass spectrum j . If the mass spectrum has no decoy or background measurement, we set $C_j = 0$. In our study, $m = n = 1$, so the default Clipper implementation is Clipper-minus-BC.
- (ii) For any target FDR threshold q , Clipper gives a cutoff T_q on contrast scores.
- (iii) The target PSMs whose mass spectra have contrast scores greater than or equal to T_q are called discoveries.

S7.3 DEG identification from bulk RNA-seq data

- (i) We consider each gene as a feature and the class label—classical and non-classical human monocytes—as the two conditions. We first performed the TMM normalization method [29]. Then we consider \log_2 -transformed read counts with a pseudocount 1 as measurements. Doing so, we summarized the gene expression matrix into a $d \times (m + n)$ matrix, where d is the number of genes, and m and n are the numbers of samples under the two conditions, respectively. We then applied Clipper to perform a differential analysis to obtain a contrast score C_j for each gene. In our study, $m = n = 3$, so the default Clipper implementation is Clipper-max-GZ with $h = 9$, the maximum number of permutations when we have three replicates under both conditions.
- (ii) For any target FDR threshold q , Clipper gives a cutoff T_q on contrast scores.
- (iii) The genes with contrast scores greater than or equal to T_q are called discoveries.

S7.4 DEG identification from scRNA-seq data

- (i) We consider each gene as a feature and the cell type—CD4+ T cells and cytotoxic T cells—as the two conditions. We first performed the TMM normalization [29]. Then we consider \log_2 -transformed read counts with a pseudocount 1 as measurements. Doing so, we summarized the gene expression matrix into a $d \times (m + n)$ matrix, where d is the number of genes, and m and n are the numbers of samples under the two conditions, respectively. We then applied Clipper to perform differential analysis to obtain a contrast score C_j for each gene j . In our study, $m = 1172$,

$n = 789$ for Drop-seq dataset and $m = 963$, $n = 694$ for 10x Genomics dataset. The default Clipper implementation is Clipper-max-GZ with $h = 1$, the default number of permutations.

- (ii) For any target FDR threshold q , Clipper gives a cutoff T_q on contrast scores.
- (iii) The genes with contrast scores greater than or equal to T_q are called discoveries.

S7.5 DIR identification from Hi-C data

- (i) We consider each pair of genomic regions as a feature and manually created two conditions. Then we consider log-transformed read counts as measurements. Doing so, we summarized the gene expression matrix into a $d \times (m + n)$ matrix, where d is the total pairs of genomic regions, and m and n are the numbers of samples under the two conditions, respectively. We then applied Clipper to perform a differential analysis to obtain a contrast score C_j for each pair of genomic regions. In our study, $m = n = 2$, so the default Clipper implementation is Clipper-max-GZ with $h = 1$.
- (ii) For any target FDR threshold q , Clipper gives a cutoff T_q on contrast scores.
- (iii) The pairs of genomic regions with contrast scores greater than or equal to T_q are called discoveries.

S8 Proofs

S8.1 Proof of Theorem 1

We first prove Theorem 1, which relies on Lemmas 1 and 2. Here we only include the proof of Lemma 1 and defer the proof of Lemma 2 to Section S8.3.

Proof 1 (Proof of Lemma 1) Here we prove that Lemma 1 holds when C_j is constructed using (S10); the proof is similar when C_j is constructed using (S11).

When input data satisfy (S6) and (S7) and $m = n$, properties (a) and (b) can be derived directly. To prove property (c), it suffices to prove that for any $j \in \mathcal{N}$ with $C_j \neq 0$, S_j is independent of $|C_j|$.

Note that \bar{X}_j and \bar{Y}_j are i.i.d for $j \in \mathcal{N}$ when $m = n$. Hence for any measurable set $\mathcal{A} \subset [0, +\infty)$,

$$\mathbb{P}(S_j = 1, |C_j| \in \mathcal{A}) = \mathbb{P}(t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j) \in \mathcal{A}) = \mathbb{P}(t^{\text{minus}}(\mathbf{Y}_j, \mathbf{X}_j) \in \mathcal{A}) = \mathbb{P}(S_j = -1, |C_j| \in \mathcal{A}).$$

The first equality holds because $t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j) = C_j = |C_j|$ when $S_j = 1$. The second equality holds because $t^{\text{minus}}(\mathbf{X}_j, \mathbf{Y}_j)$ and $t^{\text{minus}}(\mathbf{Y}_j, \mathbf{X}_j)$ are identically distributed when $j \in \mathcal{N}$. The third equality holds because $t^{\text{minus}}(\mathbf{Y}_j, \mathbf{X}_j) = -C_j$; if $-C_j \in \mathcal{A}$, then $S_j = -1$.

Because $\mathbb{P}(S_j = 1, |C_j| \in \mathcal{A}) + \mathbb{P}(S_j = -1, |C_j| \in \mathcal{A}) = \mathbb{P}(|C_j| \in \mathcal{A})$, it follows that

$$\mathbb{P}(S_j = 1, |C_j| \in \mathcal{A}) = \frac{1}{2} \mathbb{P}(|C_j| \in \mathcal{A}) = \mathbb{P}(S_j = 1) \mathbb{P}(|C_j| \in \mathcal{A}),$$

where the last equality holds because $\mathbb{P}(S_j = 1) = 1/2$ by property (b).

Hence, S_j and $|C_j|$ are independent $\forall j \in \mathcal{N}$.

Proof 2 (Proof of Theorem 1) Define a random subset of \mathcal{N} as $\mathcal{M} := \mathcal{N} \setminus \{j \in \mathcal{N} : C_j = 0\} = \{j \in \mathcal{N} : S_j \neq 0\}$.

First note that by Lemma 1(b), $\mathbb{P}(S_j = -1) = \mathbb{P}(C_j < 0) = 1/2$ for all $j \in \mathcal{M} \subset \mathcal{N}$. Assume without loss of generality that $\mathcal{M} = \{1, \dots, d'\}$. We order $\{|C_j| : j \in \mathcal{M}\}$, from the largest to the

smallest, denoted by $|C_{(1)}| \geq |C_{(2)}| \geq \dots \geq |C_{(d')}|$. Let $J = \sum_{j \in \mathcal{N}} \mathbb{1}(|C_j| \geq T^{\text{BC}})$, the number of uninteresting features whose contrast scores have absolute values no less than T^{BC} . When $J > 0$, $|C_{(1)}| \geq \dots \geq |C_{(J)}| \geq T^{\text{BC}}$. Define $Z_k = \mathbb{1}(C_{(k)} < 0)$, $k = 1, \dots, d'$. Then for each order k , the following holds

$$\begin{aligned} C_{(k)} \geq T^{\text{BC}} &\iff |C_{(k)}| \geq T^{\text{BC}} \text{ and } C_{(k)} > 0 \iff k \leq J \text{ and } Z_k = 0; \\ C_{(k)} \leq -T^{\text{BC}} &\iff |C_{(k)}| \geq T^{\text{BC}} \text{ and } C_{(k)} < 0 \iff k \leq J \text{ and } Z_k = 1. \end{aligned}$$

Then

$$\begin{aligned} \frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{BC}}\}) + 1} &= \frac{\sum_{k=1}^{d'} \mathbb{1}(C_{(k)} \geq T^{\text{BC}})}{1 + \sum_{k=1}^{d'} \mathbb{1}(C_{(k)} \leq -T^{\text{BC}})} \\ &= \frac{\sum_{k=1}^J \mathbb{1}(C_{(k)} \geq T^{\text{BC}})}{1 + \sum_{k=1}^J \mathbb{1}(C_{(k)} \leq -T^{\text{BC}})} \\ &= \frac{(1 - Z_1) + \dots + (1 - Z_J)}{1 + Z_1 + \dots + Z_J} \\ &= \frac{1 + J}{1 + Z_1 + \dots + Z_J} - 1. \end{aligned}$$

Because $\{S_j\}_{j \in \mathcal{N}}$ is independent of \mathcal{C} (Lemma 1(c)), Lemma 1(a)-(b) still holds after $C_1, \dots, C_{d'}$ are reordered as $C_{(1)}, \dots, C_{(d')}$. Thus $Z_1, \dots, Z_{d'}$ are i.i.d. from Bernoulli(1/2). To summarize, it holds that

$$\{Z_j\}_{j \in \mathcal{M}} \mid \mathcal{M} \stackrel{i.i.d.}{\sim} \text{Bernoulli}(1/2).$$

Then by applying Lemma 2 and making $\rho = 0.5$, we have:

$$\mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{BC}}\}) + 1} \mid \mathcal{M} \right] \leq 1 \quad (\text{S22})$$

Then

$$\begin{aligned} \text{FDR} &= \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \right] \\ &= \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot \frac{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \right] \\ &\leq \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \cdot \frac{\text{card}(\{j : C_j \leq -T^{\text{BC}}\}) + 1}{\text{card}(\{j : C_j \geq T^{\text{BC}}\}) \vee 1} \right] \\ &\leq q \cdot \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{BC}}\}) + 1} \right] \\ &\leq q \cdot \mathbb{E} \left[\mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{BC}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{BC}}\}) + 1} \mid \mathcal{M} \right] \right] \\ &\leq q, \end{aligned}$$

where \mathcal{M} is random subset of \mathcal{N} such that for each $j \in \mathcal{M}$, $|C_j| > 0$. The last inequality follows from (S22).

S8.2 Proof of Theorem 2

We then prove Theorem 2, which relies on Lemmas 2 and 3. Here we introduce the proof of Lemma 3 and defer the proof of Lemma 2 to Section S8.3.

Proof 3 (Proof of Lemma 3) With input data satisfying (S6) and (S7), C_j constructed from (S19) or (S20), property (a) can be derived directly.

To show property (b), note that for each uninteresting feature $j \in \mathcal{N}$, \mathbf{X}_j and \mathbf{Y}_j are from the same distribution; thus $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$ are identically distributed. Define an event $\mathcal{E}_j := \left\{ \sum_{\ell=0}^h \mathbb{1}(T_j^{\sigma_\ell} = T_j^{(0)}) = 1 \right\}$, which indicates that $T_j^{(0)}$, the maximizer of $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$, is unique. Then conditional on \mathcal{E}_j , the maximizer is equally likely to be any of $\{0, \dots, h\}$, and it follows that $\mathbb{P}(S_j = 1 \mid \mathcal{E}_j) = \mathbb{P}(T_j^{\sigma_0} = T_j^{(0)} \mid \mathcal{E}_j) = 1/(h+1)$. Conditioning on that \mathcal{E}_j does not happen, $\mathbb{P}(S_j = 1 \mid \mathcal{E}_j^c) = 0$. Thus $\mathbb{P}(S_j = 1) = \mathbb{P}(S_j = 1 \mid \mathcal{E}_j)\mathbb{P}(\mathcal{E}_j) + \mathbb{P}(S_j = 1 \mid \mathcal{E}_j^c)\mathbb{P}(\mathcal{E}_j^c) \leq 1/(h+1)$.

The proof of property (c) is similar to the Proof of Lemma 1(c). It suffices to show that for any $j \in \mathcal{N}$ with $C_j \neq 0$ (that is, \mathcal{E}_j occurs), S_j is independent of $|C_j|$. As \mathbf{X}_j and \mathbf{Y}_j are from the same distribution, $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$ are identically distributed. Hence for any measurable set $\mathcal{A} \subset [0, +\infty)$,

$$\begin{aligned} \mathbb{P}(S_j = 1, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) &= \mathbb{P}(T_j^{\sigma_0} = T_j^{(0)}, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) \\ &= \frac{1}{h} \mathbb{P}(T_j^{\sigma_0} \neq T_j^{(0)}, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) \\ &= \frac{1}{h} \mathbb{P}(S_j = -1, |C_j| \in \mathcal{A} \mid \mathcal{E}_j). \end{aligned}$$

The first equality holds because $T_j^{\sigma_0} = T_j^{(0)}$ when $S_j = 1$. The second equality holds because $\{T_j^{\sigma_\ell}\}_{\ell=0}^h$ are identically distributed when $j \in \mathcal{N}$. The third equality holds because $T_j^{\sigma_0} \neq T_j^{(0)}$ when $S_j = -1$.

Because $\mathbb{P}(S_j = 1, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) + \mathbb{P}(S_j = -1, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) = \mathbb{P}(|C_j| \in \mathcal{A} \mid \mathcal{E}_j)$, it follows that

$$\mathbb{P}(S_j = 1, |C_j| \in \mathcal{A} \mid \mathcal{E}_j) = \frac{1}{h+1} \mathbb{P}(|C_j| \in \mathcal{A} \mid \mathcal{E}_j) = \mathbb{P}(S_j = 1 \mid \mathcal{E}_j) \mathbb{P}(|C_j| \in \mathcal{A} \mid \mathcal{E}_j),$$

where the last equality holds because $\mathbb{P}(S_j = 1 \mid \mathcal{E}_j) = 1/(h+1)$.

Hence, S_j and $|C_j|$ are independent $\forall j \in \mathcal{N}$ with $C_j \neq 0$.

Proof 4 (Proof of Theorem 2) Define a random subset of \mathcal{N} as $\mathcal{M} := \mathcal{N} \setminus \{j \in \mathcal{N} : C_j = 0\} = \{j \in \mathcal{N} : S_j \neq 0\}$. Assume without loss of generality that $\mathcal{M} = \{1, \dots, d'\}$. We order $\{|C_j| : j \in \mathcal{M}\}$, from the largest to the smallest, denoted by $|C_{(1)}| \geq |C_{(2)}| \geq \dots \geq |C_{(d')}|$. Let $J = \sum_{j \in \mathcal{N}} \mathbb{1}(|C_j| \geq T^{\text{GZ}})$, the number of uninteresting features whose contrast scores have absolute values no less than T^{GZ} . When $J > 0$, $|C_{(1)}| \geq \dots \geq |C_{(J)}| \geq T^{\text{GZ}}$. Define $Z_k = \mathbb{1}(C_{(k)} < 0)$, $k = 1, \dots, d'$. Then for each order k , the following holds:

$$\begin{aligned} C_{(k)} \geq T^{\text{GZ}} &\iff |C_{(k)}| \geq T^{\text{GZ}} \text{ and } C_{(k)} > 0 \iff k \leq J \text{ and } Z_k = 0; \\ C_{(k)} \leq -T^{\text{GZ}} &\iff |C_{(k)}| \geq T^{\text{GZ}} \text{ and } C_{(k)} < 0 \iff k \leq J \text{ and } Z_k = 1. \end{aligned}$$

Then it follows that

$$\begin{aligned} \frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{GZ}}\}) + 1} &= \frac{\sum_{k=1}^{d'} \mathbb{1}(C_{(k)} \geq T^{\text{GZ}})}{1 + \sum_{k=1}^{d'} \mathbb{1}(C_{(k)} \leq -T^{\text{GZ}})} \\ &= \frac{\sum_{k=1}^J \mathbb{1}(C_{(k)} \geq T^{\text{GZ}})}{1 + \sum_{k=1}^J \mathbb{1}(C_{(k)} \leq -T^{\text{GZ}})} \\ &= \frac{(1 - Z_1) + \dots + (1 - Z_J)}{1 + Z_1 + \dots + Z_J} \\ &= \frac{1 + J}{1 + Z_1 + \dots + Z_J} - 1. \end{aligned}$$

Because $\{S_j\}_{j \in \mathcal{N}}$ is independent of \mathcal{C} (Lemma 1(c)), Lemma 1(a)-(b) still holds after $C_1, \dots, C_{d'}$ are reordered as $C_{(1)}, \dots, C_{(d')}$. Thus $Z_1, \dots, Z_{d'}$ are i.i.d. from $\text{Bernoulli}(\rho_k)$. To summarize, it holds that

$$\{Z_j\}_{j \in \mathcal{M}} \Big| \mathcal{M} \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\rho_k).$$

Then by applying Lemma 2 and making $\rho = h/(h+1)$, we have:

$$\mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{GZ}}\}) + 1} \right] \leq 1/h. \quad (\text{S23})$$

Then

$$\begin{aligned} \text{FDR} &= \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j : C_j \geq T^{\text{GZ}}\}) \vee 1} \right] \\ &= \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{GZ}}\}) + 1} \cdot \frac{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{GZ}}\}) + 1}{\text{card}(\{j : C_j \geq T^{\text{GZ}}\}) \vee 1} \right] \\ &\leq h \cdot \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{GZ}}\}) + 1} \cdot \frac{\frac{1}{h} \text{card}(\{j : C_j \leq -T^{\text{GZ}}\}) + \frac{1}{h}}{\text{card}(\{j : C_j \geq T^{\text{GZ}}\}) \vee 1} \right] \\ &\leq hq \cdot \mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{N} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{N} : C_j \leq -T^{\text{GZ}}\}) + 1} \right] \\ &\leq hq \cdot \mathbb{E} \left[\mathbb{E} \left[\frac{\text{card}(\{j \in \mathcal{M} : C_j \geq T^{\text{GZ}}\})}{\text{card}(\{j \in \mathcal{M} : C_j \leq -T^{\text{GZ}}\}) + 1} \Big| \mathcal{M} \right] \right] \\ &\leq q, \end{aligned}$$

where the second inequality follows from the definition of T_{GZ} (S21) and the last inequality follows from (S23).

S8.3 Proof of Lemma 2

Finally, we derive Lemma 2 by following the same proof same as in [8], which relies on Lemma 4 and Corollary 1.

Lemma 4 Suppose that $Z_1, \dots, Z_d \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\rho)$. Let J be a stopping time in reverse time with respect to the filtration $\{\mathcal{F}_j\}$, where $\mathcal{F}_j = \sigma(\{(Z_1 + \dots + Z_j), Z_{j+1}, \dots, Z_d\})$ with $\sigma(\cdot)$ denoting a σ -algebra, and the variables Z_1, \dots, Z_j are exchangeable with respect to $\{\mathcal{F}_j\}$. Then

$$\mathbb{E} \left[\frac{1 + J}{1 + Z_1 + \dots + Z_J} \right] \leq \rho^{-1}.$$

Proof 5 (Proof of Lemma 4) Define

$$Y_j = Z_1 + \dots + Z_j \in \mathcal{F}_j$$

and define the process

$$M_j = \frac{1+j}{1+Z_1+\dots+Z_j} = \frac{1+j}{1+Y_j} \in \mathcal{F}_j.$$

In [3], it is shown that $\mathbb{E}[M_d] \leq \rho^{-1}$. Therefore, by the optional stopping time theorem it suffices to show that $\{M_j\}$ is a supermartingale with respect to $\{\mathcal{F}_j\}$. As $\{Z_1, \dots, Z_{j+1}\}$ are exchangeable with respect to \mathcal{F}_{j+1} , we have

$$\mathbb{P}(Z_{j+1} = 1 | \mathcal{F}_{j+1}) = \frac{Y_j + 1}{1+j}.$$

Therefore, if $Y_{j+1} > 0$,

$$\begin{aligned} \mathbb{E}[M_j | \mathcal{F}_{j+1}] &= \frac{1+j}{1+Y_{j+1}} \cdot \mathbb{P}(Z_{j+1} = 0 | \mathcal{F}_{j+1}) + \frac{1+j}{1+Y_{j+1}-1} \cdot \mathbb{P}(Z_{j+1} = 1 | \mathcal{F}_{j+1}) \\ &= \frac{1+j}{1+Y_{j+1}} \cdot \frac{1+j-Y_{j+1}}{1+j} + \frac{1+j}{1+Y_{j+1}-1} \cdot \frac{Y_{j+1}}{1+j} \\ &= \frac{1+j-Y_{j+1}}{1+Y_{j+1}} + 1 \\ &= \frac{1+(j+1)}{1+Y_{j+1}} \\ &= M_{j+1}. \end{aligned}$$

If instead $Y_{j+1} = 0$, then trivially $Y_j = 0$, and $M_j = 1+j < 2+j = M_{j+1}$. This proves that $\{M_j\}$ is a supermartingale with respect to $\{\mathcal{F}_j\}$ as desired.

Corollary 1 Suppose that $\mathcal{A} \subseteq \{1, \dots, d\}$ is fixed, while $Z_1, \dots, Z_d \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\rho)$. Let J be a stopping time in reverse time with respect to the filtration $\{\mathcal{F}_j\}$, where $\mathcal{F}_j = \sigma\left(\{\sum_{k \leq j, k \in \mathcal{A}} Z_k\} \cup \{Z_k : j < k < d, k \in \mathcal{A}\}\right)$ with $\sigma(\cdot)$ denoting a σ -algebra, and the variables $\{Z_k : k \leq j, k \in \mathcal{A}\}$ are exchangeable with respect to \mathcal{F}_j . Then

$$\mathbb{E}\left[\frac{1 + \text{card}(\{k : k \leq J, k \in \mathcal{A}\})}{1 + \sum_{k \leq J, k \in \mathcal{A}} Z_k}\right] \leq \rho^{-1}.$$

Proof 6 (Proof of Corollary 1) Let $\mathcal{A} = \{j_1, \dots, j_m\}$ where $1 \leq j_1 < \dots < j_m \leq d$. Then by considering the i.i.d. sequence

$$Z_{j_1}, \dots, Z_{j_m}$$

in place of Z_1, \dots, Z_d , we see that this result is equivalent to Lemma 4.

Proof 7 (Proof of Lemma 2) [From [3]] We may assume $\rho < 1$ to avoid the trivial case. We first introduce a different definition for $\{Z_j\}_{j=1}^d$ by defining a random set $\mathcal{A} \subseteq \{1, \dots, d\}$ where for each j , independently,

$$\mathbb{P}(j \in \mathcal{A}) = \frac{1 - \rho_j}{1 - \rho}.$$

We then define random variables $Q_1, \dots, Q_d \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\rho)$, which are generated independently of the random set \mathcal{A} . Finally, we define

$$Z_j = Q_j \cdot \mathbb{1}(j \in \mathcal{A}) + \mathbb{1}(j \notin \mathcal{A}). \quad (\text{S24})$$

Then $\{Z_j\}_{k=1}^d$ are mutually independent and $\mathbb{P}(Z_j = 1) = 1 - \mathbb{P}(j \in \mathcal{A}) \cdot \mathbb{P}(Q_j = 0) = \rho_j$, that is,

$Z_j \sim \text{Bernoulli}(\rho_j)$. This new definition of $\{Z_j\}_{j=1}^d$ meet all the conditions required by Lemma 2, so that we can apply this new definition in the following proof.

As $Z_j = Q_j \cdot \mathbb{1}(j \in \mathcal{A}) + \mathbb{1}(j \notin \mathcal{A})$ for all j , we have

$$\frac{1 + J}{1 + Z_1 + \dots + Z_J} = \frac{1 + \text{card}(\{j \leq J : j \in \mathcal{A}\}) + \text{card}(\{j \leq J : j \notin \mathcal{A}\})}{1 + \sum_{j \leq J, j \in \mathcal{A}} Q_j + \text{card}(\{j \leq J : j \notin \mathcal{A}\})} \leq \frac{1 + \text{card}(\{j \leq J : j \in \mathcal{A}\})}{1 + \sum_{j \leq J, j \in \mathcal{A}} Q_j}, \quad (\text{S25})$$

where the last step uses the identify $\frac{a+c}{b+c} \leq \frac{a}{b}$ whenever $0 < b \leq a$ and $c \geq 0$. Therefore, it will be sufficient to prove that

$$\mathbb{E} \left[\frac{1 + \text{card}(\{j \leq J : j \in \mathcal{A}\})}{1 + \sum_{j \leq J, j \in \mathcal{A}} Q_j} \middle| \mathcal{A} \right] \leq \rho^{-1}, \quad (\text{S26})$$

To prove (S26), first let $\tilde{Q}_j = Q_j \cdot \mathbb{1}(j \in \mathcal{A})$, and define a filtration $\{\mathcal{F}'_j\}$ where \mathcal{F}'_j is the σ -algebra generated as

$$\mathcal{F}'_j = \sigma \left(\left\{ \tilde{Q}_1 + \dots + \tilde{Q}_j, \tilde{Q}_{j+1}, \dots, \tilde{Q}_d, \mathcal{A} \right\} \right).$$

Next for any j , by (S24) we see that

$$Z_1 + \dots + Z_j, Z_{j+1}, \dots, Z_d \in \mathcal{F}'_j \Rightarrow \mathcal{F}_j \subseteq \mathcal{F}'_j,$$

so J is a stopping time (in reverse time) with respect to \mathcal{F}'_j . Finally, since the Q_j 's are independent of \mathcal{A} , (S26) follows from Corollary 1 after conditioning on \mathcal{A} .

S9 Supplementary figures

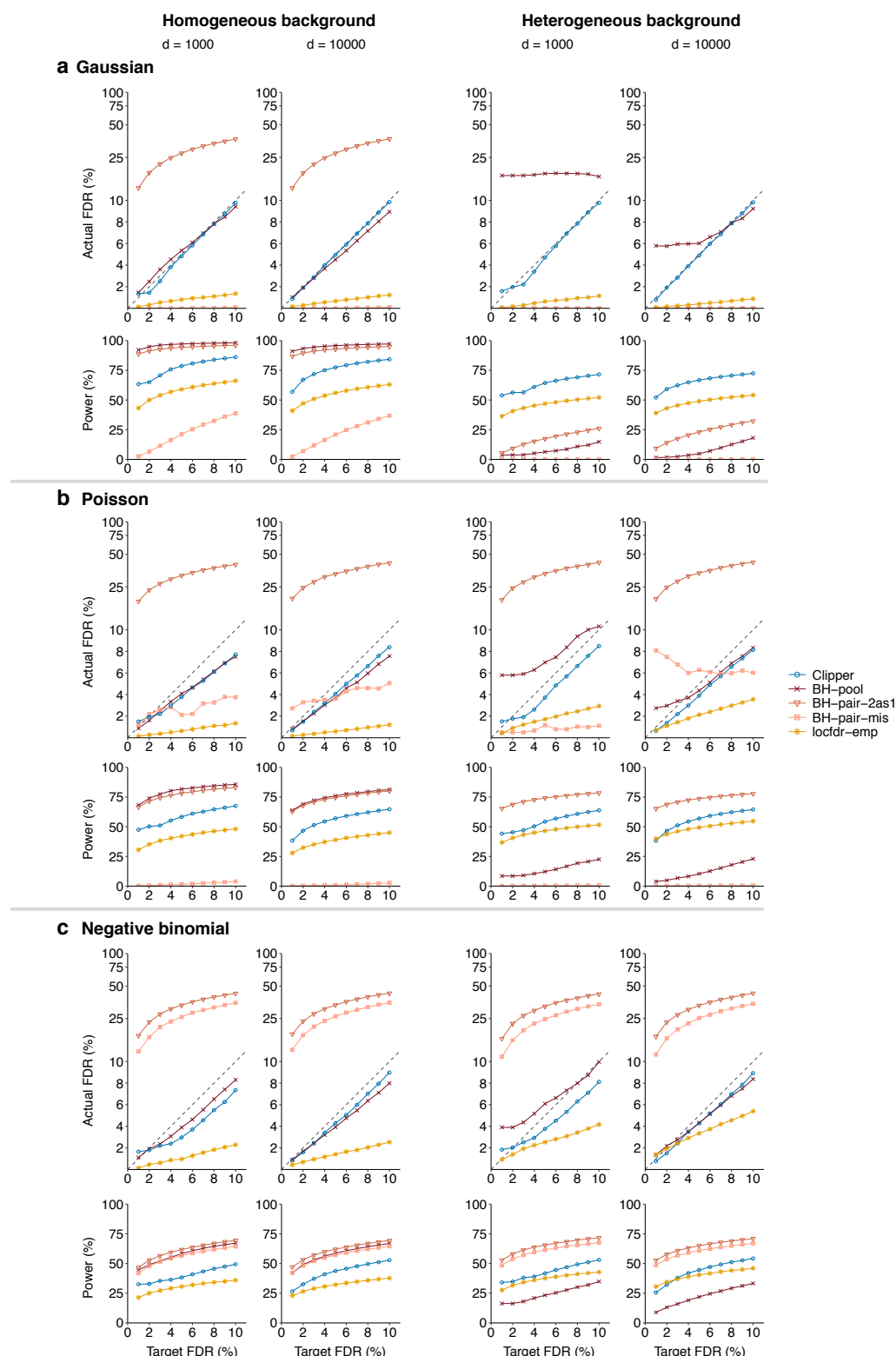


Figure S1: In the 1vs1 enrichment analysis, comparison of Clipper and four other generic FDR control methods (BH-pool, BH-pair-2as1, BH-pair-mis, and locfdr-emp) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 1000$ or $10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

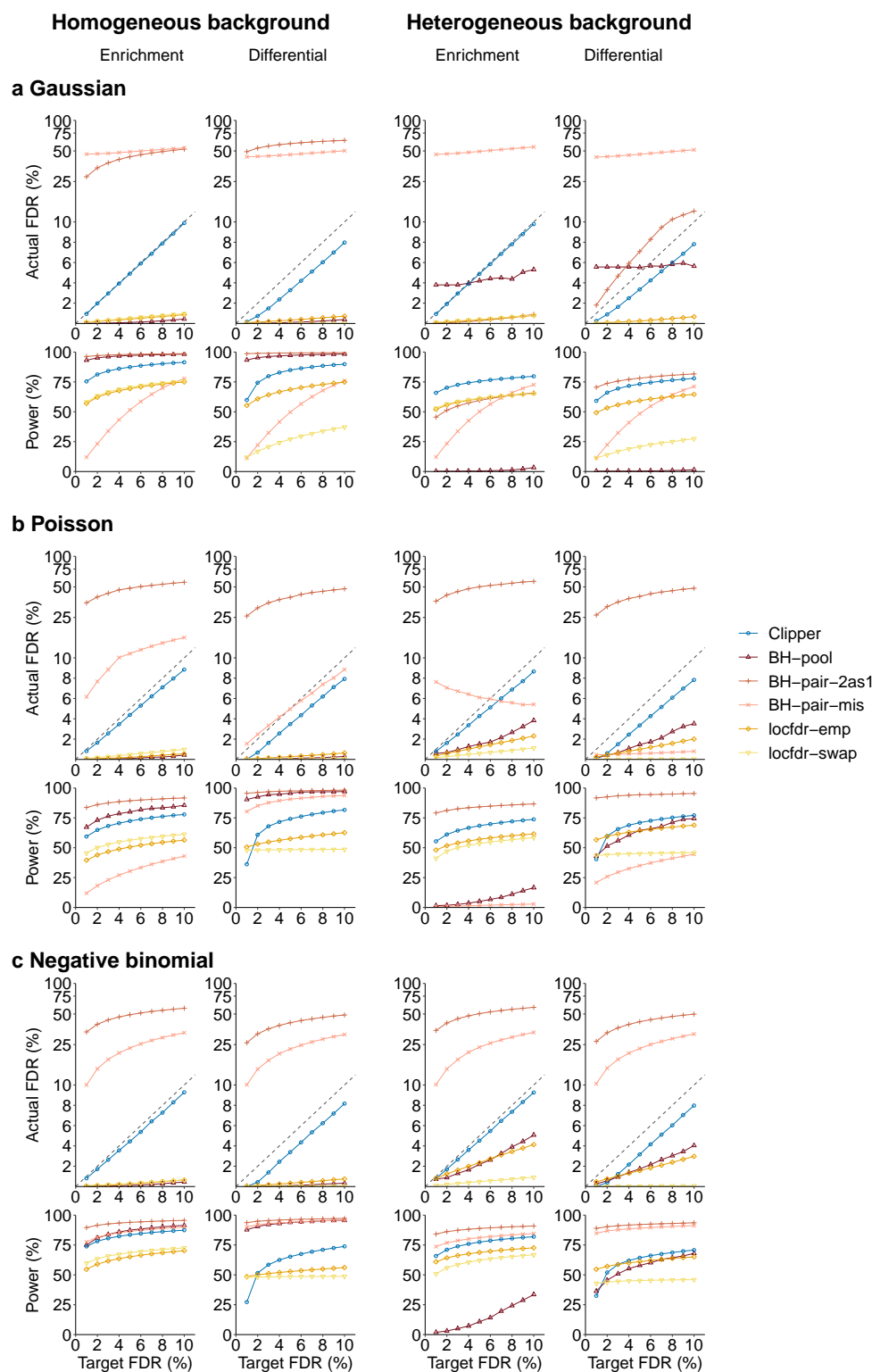


Figure S2: In the 2vs1 enrichment analysis (columns 1 and 3) and differential analysis (columns 2 and 4), comparison of Clipper and five other generic FDR control methods (BH-pooled, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background (except for differential analysis with $q \leq 2\%$).

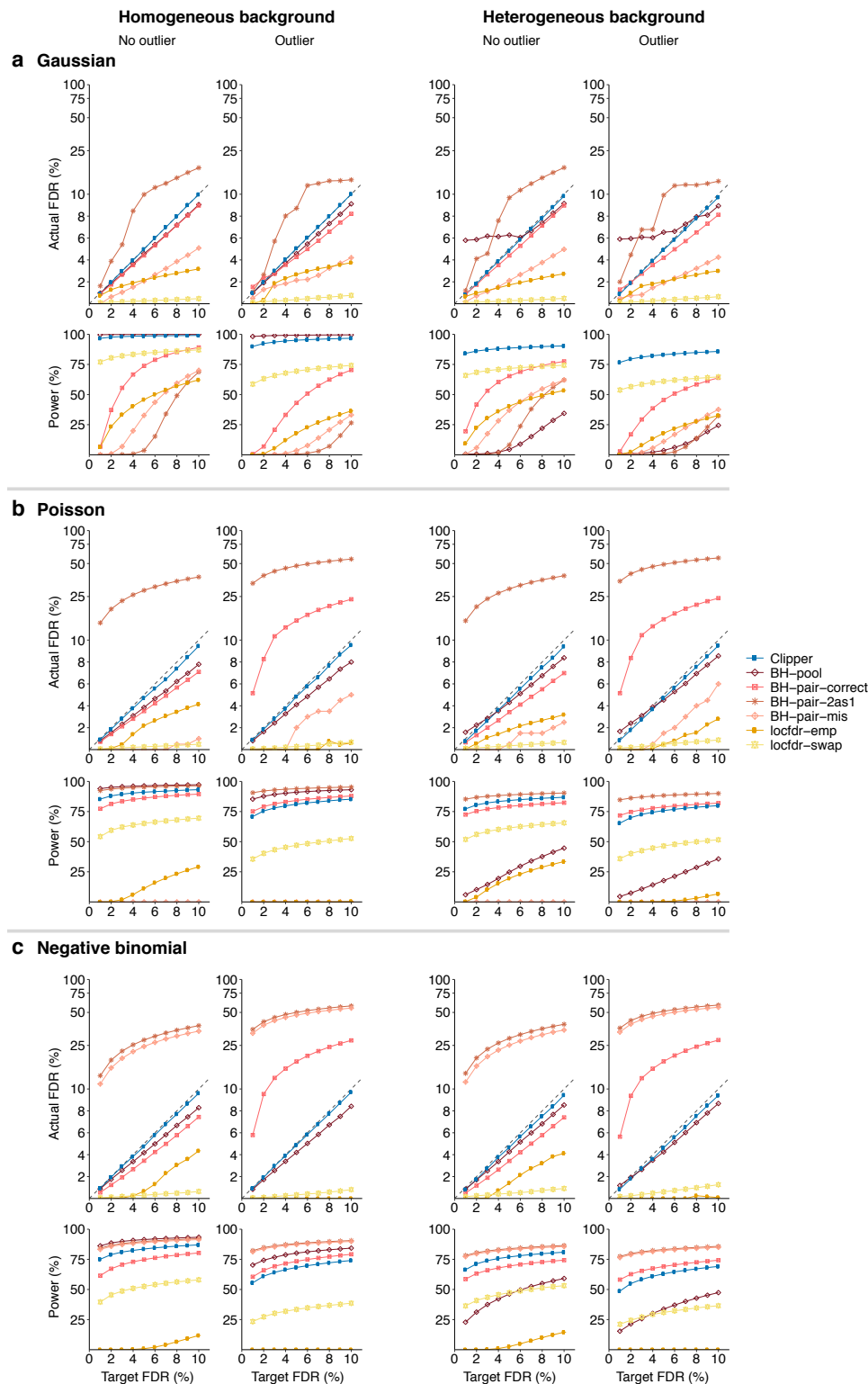


Figure S3: In the 3vs3 enrichment analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of Clipper and six other generic FDR control methods (BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power in 3vs3 enrichment analysis with possible outliers. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. 10% of the features are interesting features. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

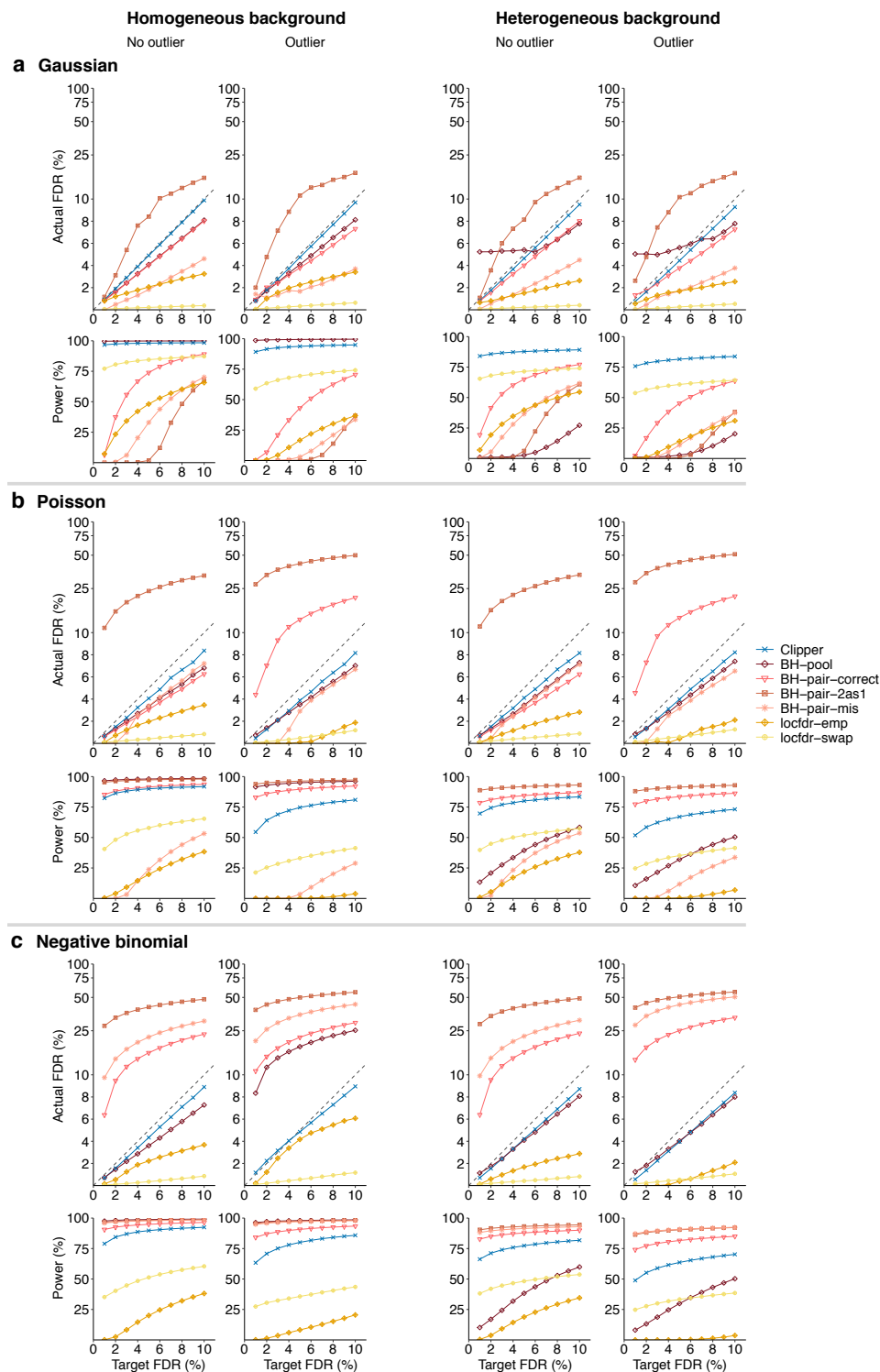


Figure S4: In the 3vs3 differential analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of Clipper and six other generic FDR control methods (BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background (except for Poisson distribution where Clipper is second to BH-pair-correct, an idealistic method).

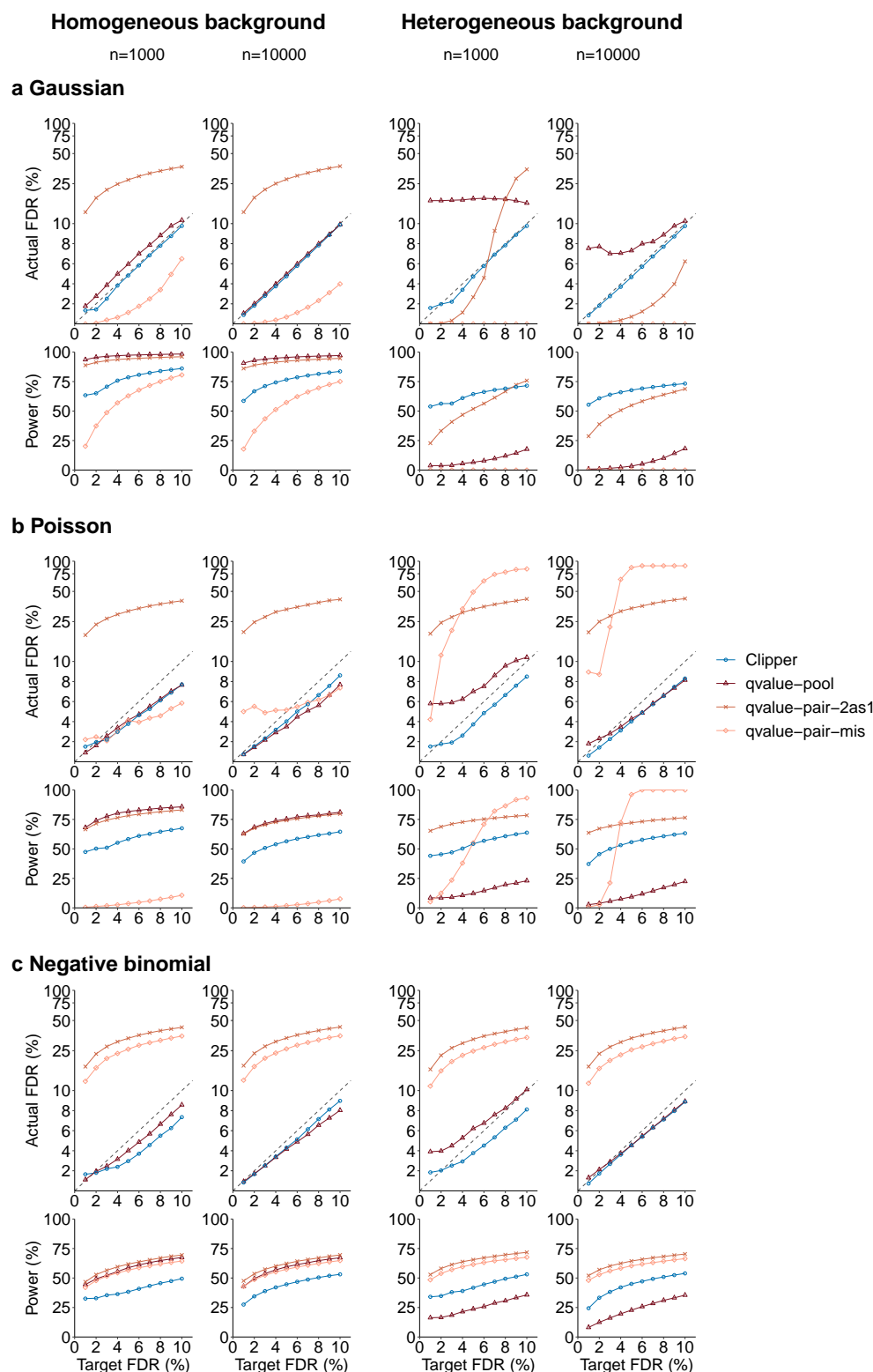


Figure S5: In the 1vs1 enrichment analysis, comparison of Clipper and three other generic FDR control methods using Storey's q-value (qvalue-pool, qvalue-pair-2as1, and qvalue-pair-mis) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 1000$ or $10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

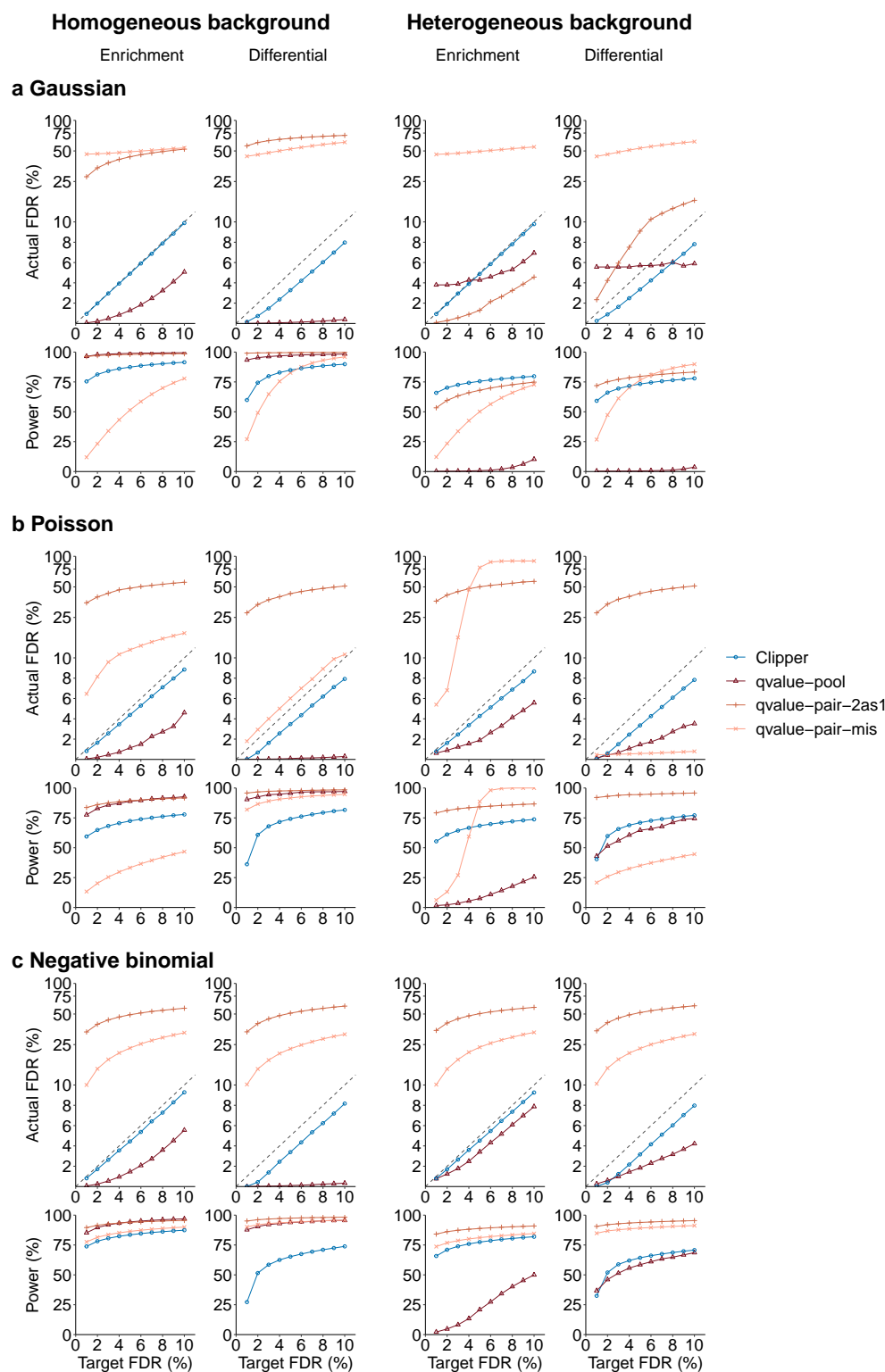


Figure S6: In the 2vs1 enrichment analysis (columns 1 and 3) and differential analysis (columns 2 and 4), comparison of Clipper and three other generic FDR control methods using Storey's q-value (qvalue-pool, qvalue-pair-2as1, and qvalue-pair-mis) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

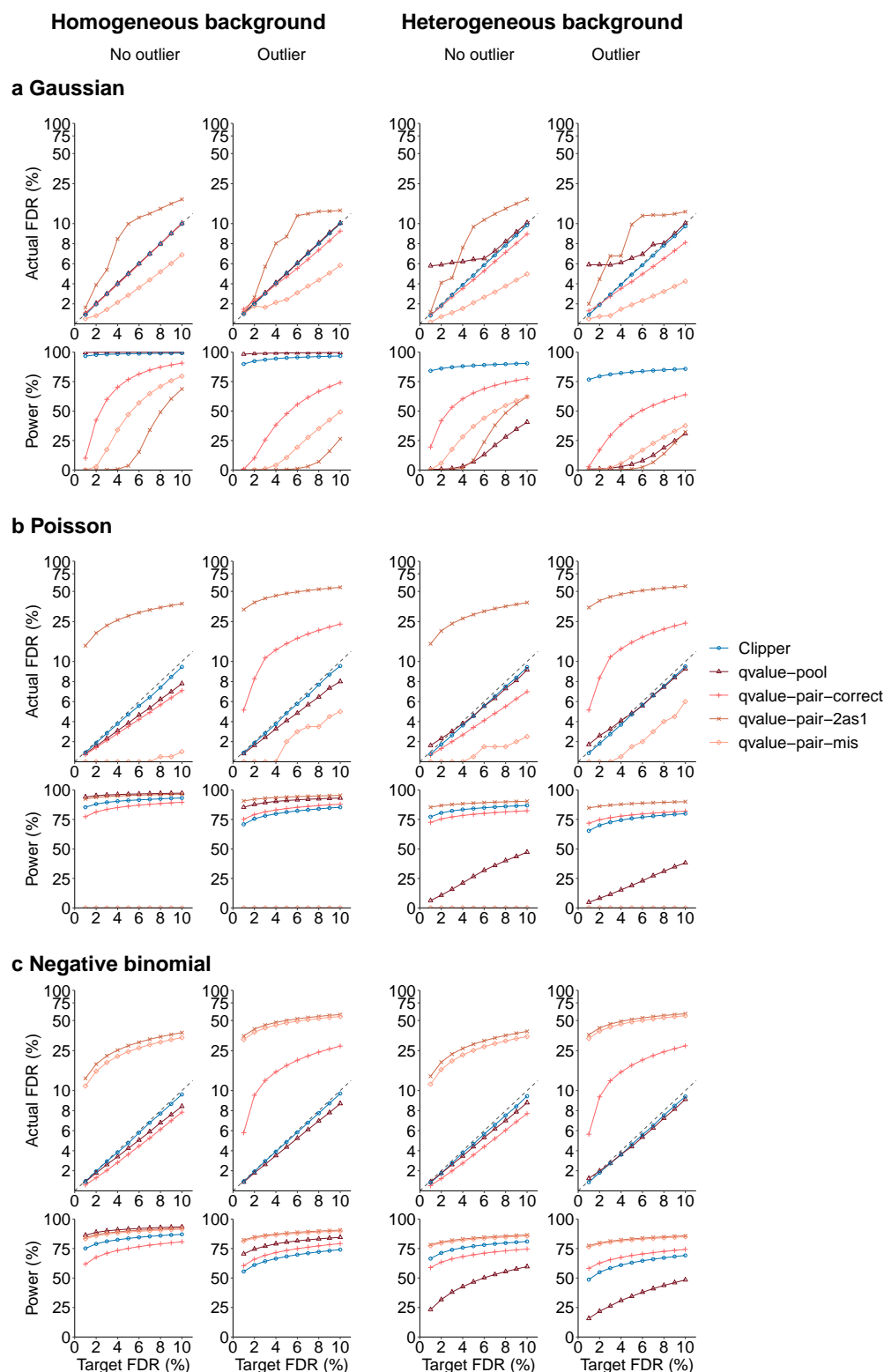


Figure S7: In the 3vs3 enrichment analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of Clipper and four other generic FDR control methods using Storey's q-value (qvalue-pooled, qvalue-pair-correct, qvalue-pair-2as1, and qvalue-pair-mis) in terms of their FDR control and power in 3vs3 enrichment analysis with possible outliers. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

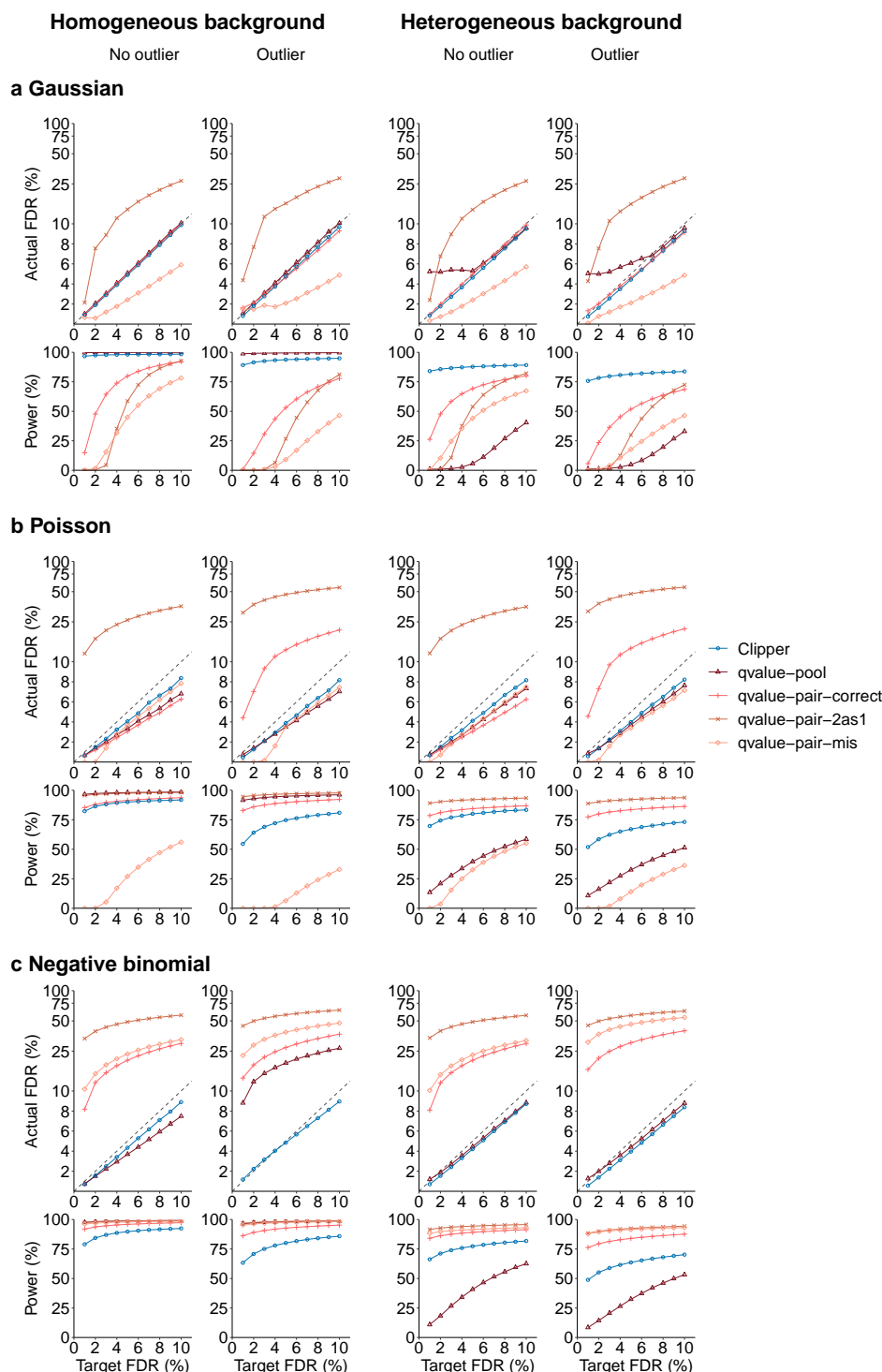


Figure S8: In the 3vs3 differential analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of Clipper and four other generic FDR control methods using Storey's q-value (qvalue-pooled, qvalue-pair-correct, qvalue-pair-2as1, and qvalue-pair-mis) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background (except for Poisson distribution where Clipper is second to qvalue-pair-correct, an idealistic method).

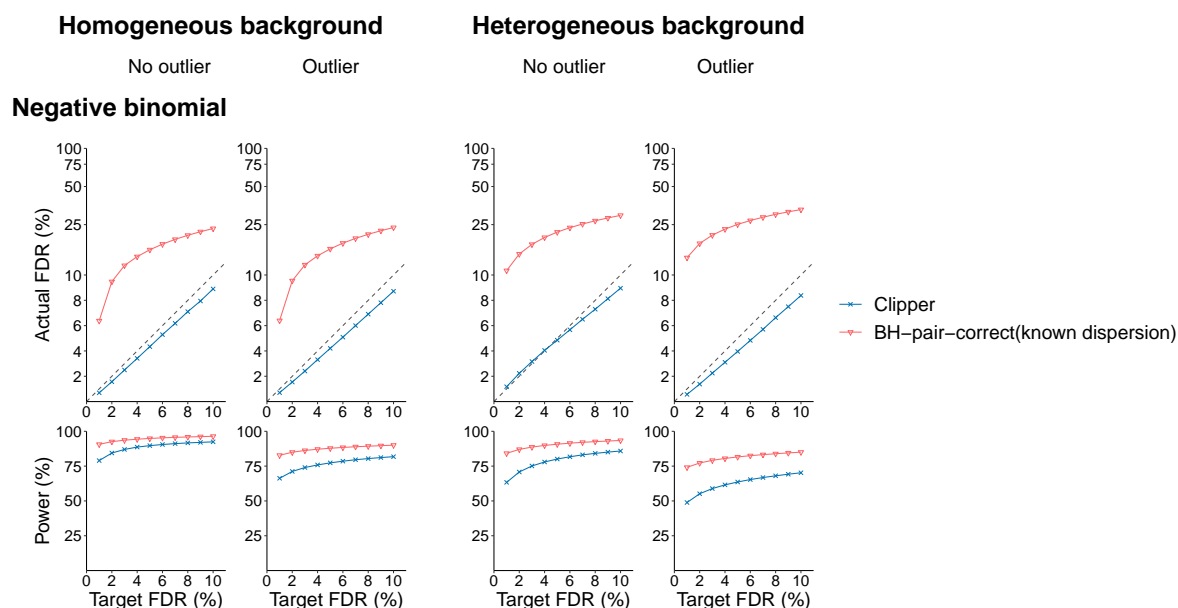


Figure S9: In the 3vs3 differential analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of Clipper, BH-pair-correct (known dispersion), and BH-pair-correct (unknown dispersion) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. BH-pair-correct (unknown dispersion) cannot control the FDR in all settings. In contrast, Clipper is consistently the most powerful for homogeneous and heterogeneous background.

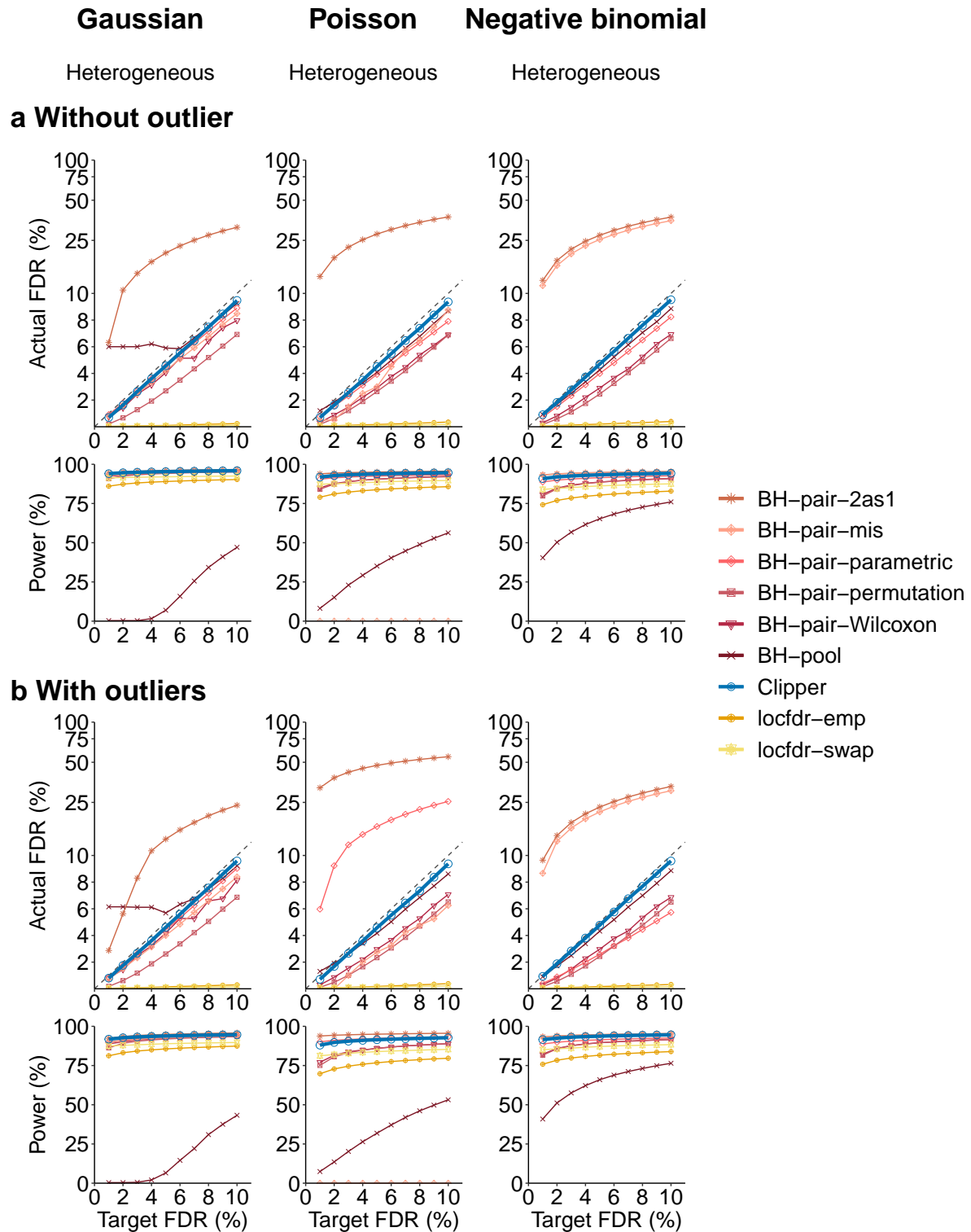


Figure S10: In the 10vs10 enrichment analysis with and without outliers, comparison of Clipper and eight generic FDR control methods (BH-pooled, BH-pair-Wilcoxon, BH-pair-parametric, and BH-pair-permutation, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from the Gaussian distribution (left), the Poisson distribution (middle), or the negative binomial distribution (right) under heterogeneous background scenarios. Clipper achieves the highest power for all three distributions.

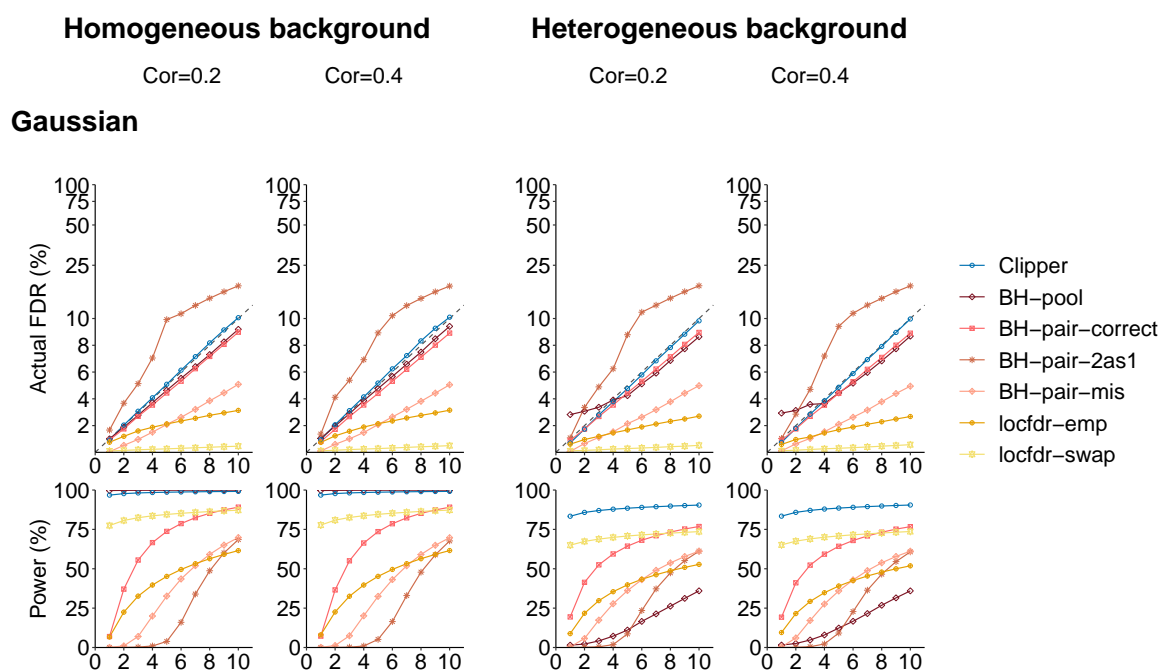


Figure S11: In the 3vs3 enrichment analysis with correlated features, comparison of Clipper and six other generic FDR control methods (BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power in 3vs3 enrichment analysis. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from a multivariate Gaussian distribution with a correlation 0.2 (columns 1 and 3) or 0.4 (columns 2 and 4) between features. Among the methods that control the FDR, Clipper is the second most powerful for homogeneous background and the most powerful for heterogeneous background.

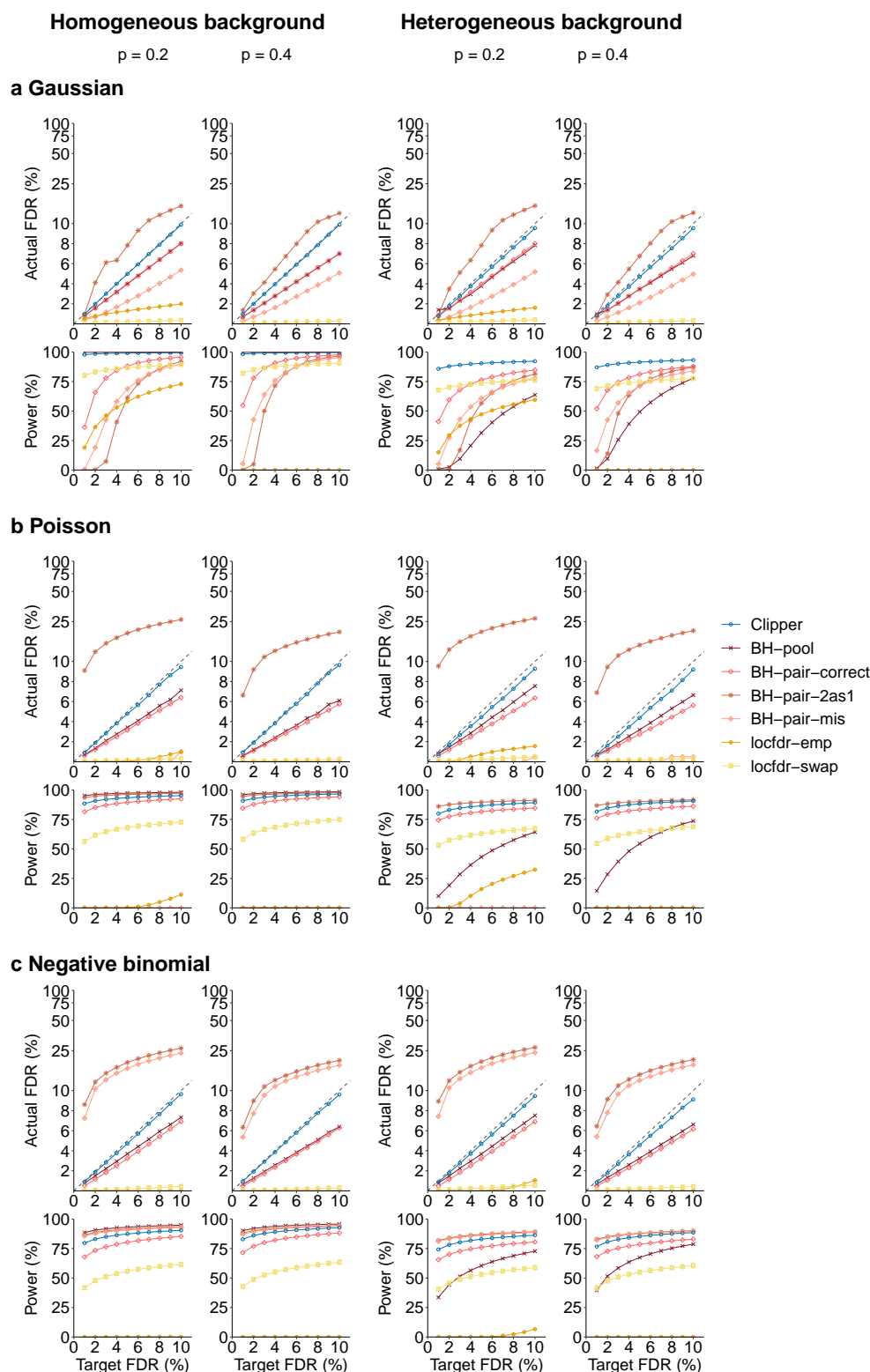


Figure S12: In 3vs3 enrichment analysis with different proportions of interesting features without outliers, comparison of Clipper and six generic FDR control methods (BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from the Gaussian distribution, the Poisson distribution, or the negative binomial distribution, with the proportion of interesting features being 0.2 (columns 1 and 3) or 0.4 (columns 2 and 4) under homogeneous (columns 1 and 2) and heterogeneous (columns 3 and 4) background scenarios. Clipper achieves the highest power for all distributions.

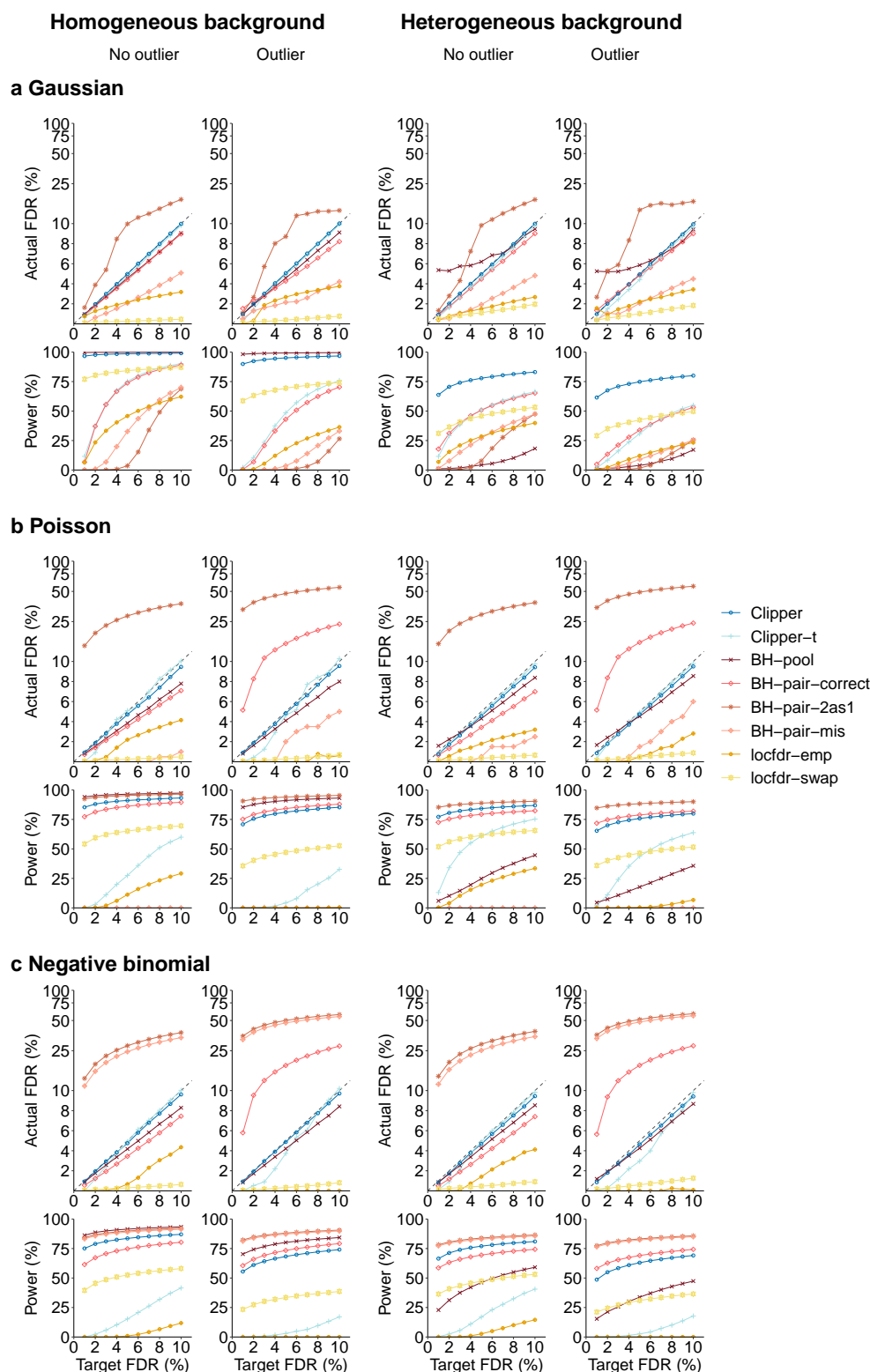


Figure S13: In the 3vs3 enrichment analysis with and without outliers, comparison of the default Clipper, the Clipper variant using the t statistic as the contrast score (Clipper-t), and six generic FDR control methods (Clipper BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from the Gaussian distribution, the Poisson distribution, or the negative binomial distribution under homogeneous (columns 1 and 2) and heterogeneous (columns 3 and 4) background scenarios. Clipper achieves higher power than Clipper-t does.

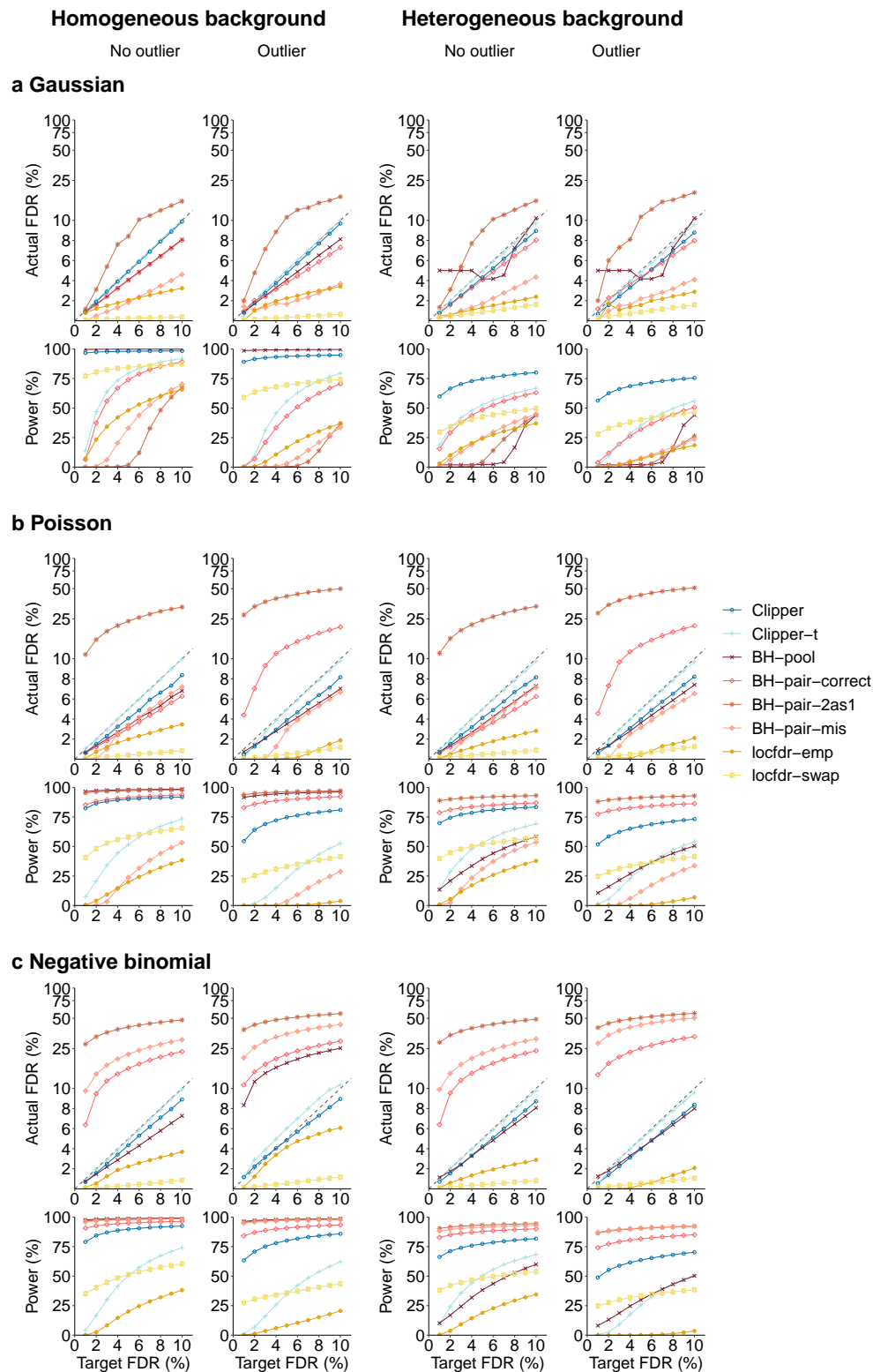


Figure S14: In the 3vs3 differential analysis with and without outliers, comparison of the default Clipper, the Clipper variant using the t statistic to calculate the degree of interestingness (Clipper-t), and six generic FDR control methods (Clipper BH-pooled, BH-pair-correct, BH-pair-2as1, BH-pair-mis, locfdr-emp, and locfdr-swap) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from the Gaussian distribution, the Poisson distribution, or the negative binomial distribution under homogeneous (columns 1 and 2) and heterogeneous (columns 3 and 4) background scenarios. Clipper achieves higher power than Clipper-t does.

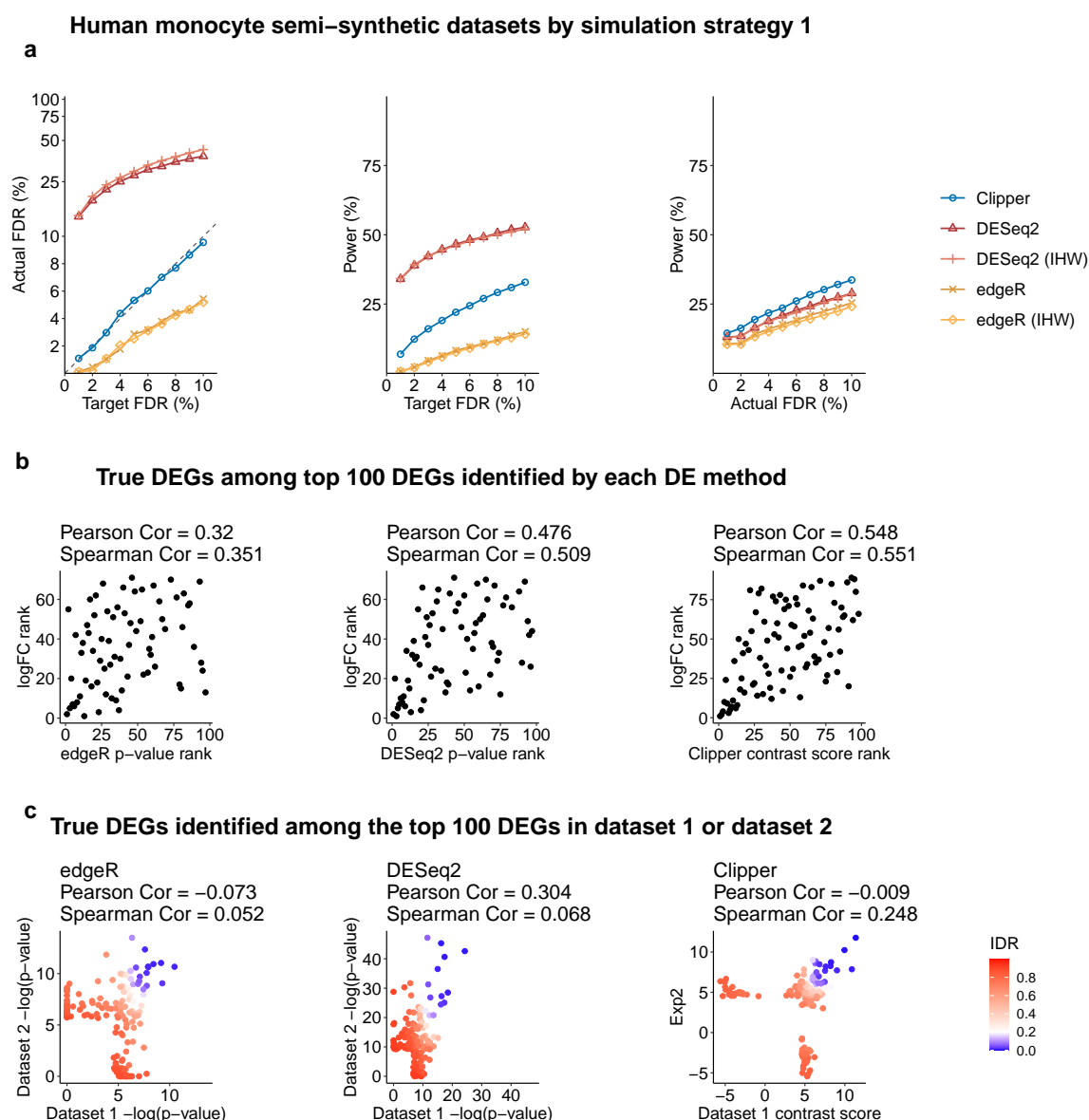


Figure S15: Comparison of Clipper and two popular DEG identification methods—edgeR and DESeq2—in DEG analysis on semi-synthetic bulk RNA-seq data (generated from human monocyte real data using simulation strategy 1 in Supp. Section S6.3). **(a)** FDR control, power given the same target FDR, and power given the same actual FDR. **(b)** Ranking consistency of the true DEGs among the top 100 DEGs identified by each method. The consistency is defined between the genes' ranking based on edgeR/DESeq2's p-values or Clipper's contrast scores and their ranking based on true expression fold changes. **(c)** Reproducibility between two semi-synthetic datasets as technical replicates. Three reproducibility criteria are used: the IDR, Pearson correlation, and Spearman correlation. Each criterion is calculated for edgeR/DESeq2's p-values or Clipper's contrast scores on the two semi-synthetic datasets. Among the three methods, only Clipper controls the FDR, and Clipper achieves the highest power, the best gene ranking consistency, and the best reproducibility.

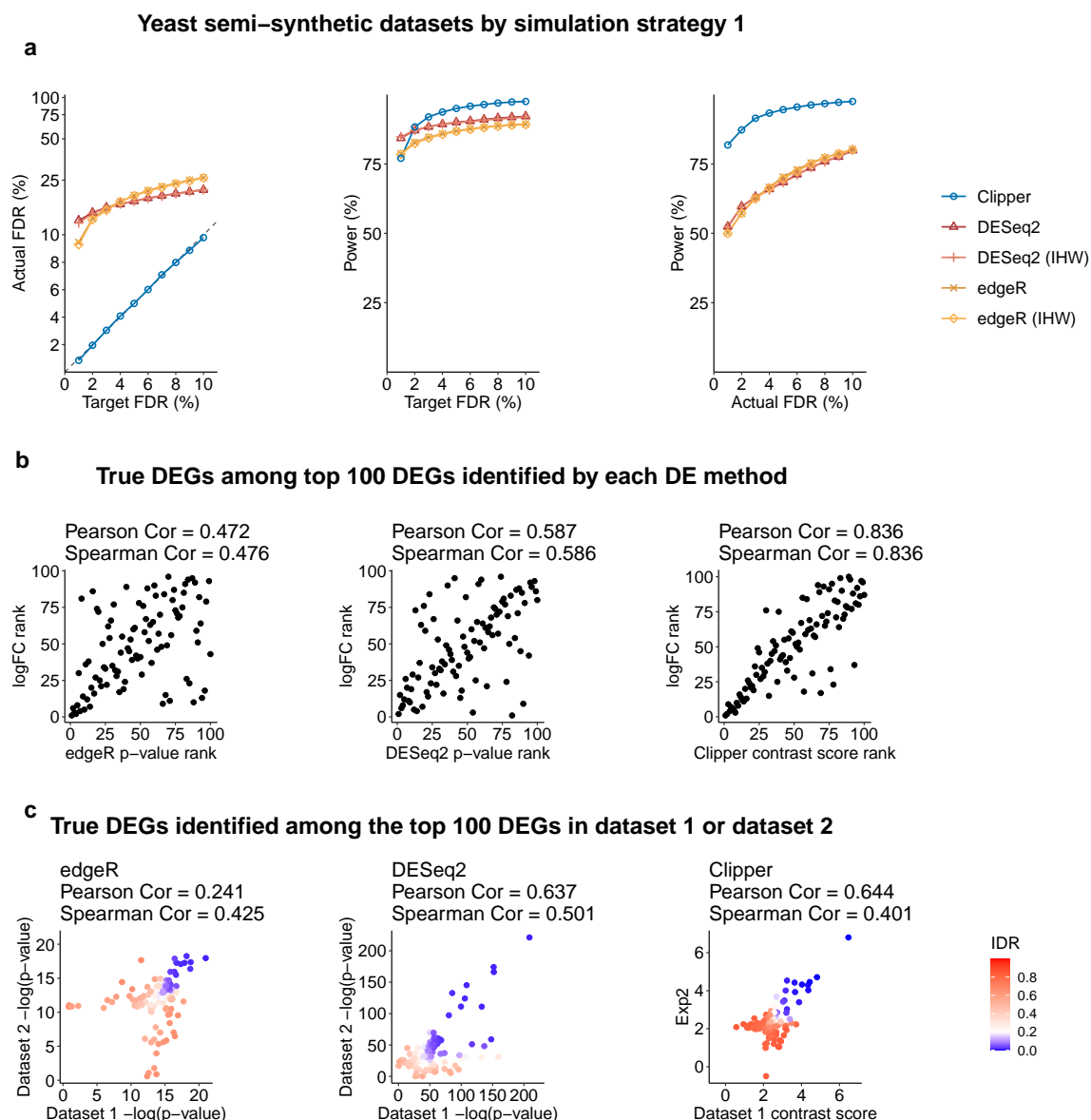


Figure S16: Comparison of Clipper and two popular DEG identification methods—edgeR and DESeq2—in DEG analysis on semi-synthetic bulk RNA-seq data (generated from yeast real data using simulation strategy 1 in Supp. Section S6.3). **(a)** FDR control, power given the same target FDR, and power given the same actual FDR. **(b)** Ranking consistency of the true DEGs among the top 100 DEGs identified by each method. The consistency is defined between the genes' ranking based on edgeR/DESeq2's p-values or Clipper's contrast scores and their ranking based on true expression fold changes. **(c)** Reproducibility between two semi-synthetic datasets as technical replicates. Three reproducibility criteria are used: the IDR, Pearson correlation, and Spearman correlation. Each criterion is calculated for edgeR/DESeq2's p-values or Clipper's contrast scores on the two semi-synthetic datasets. Among the three methods, only Clipper controls the FDR, and Clipper achieves the highest power, the best gene ranking consistency, and the best reproducibility.

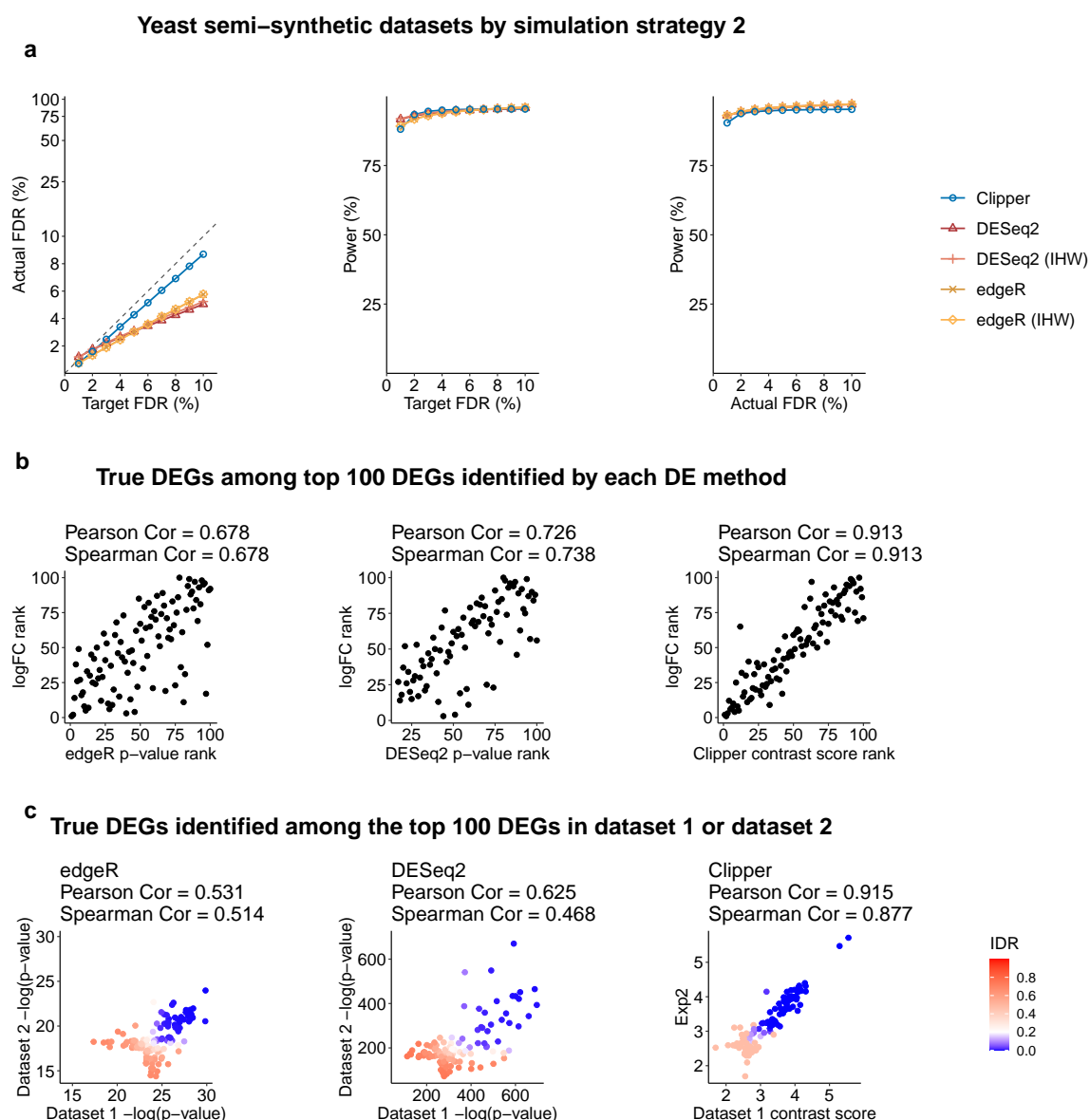


Figure S17: Comparison of Clipper and two popular DEG identification methods—edgeR and DESeq2—in DEG analysis on semi-synthetic bulk RNA-seq data (generated from yeast real data using simulation strategy 2 in Supp. Section S6.3). **(a)** FDR control, power given the same target FDR, and power given the same actual FDR. **(b)** Ranking consistency of the true DEGs among the top 100 DEGs identified by each method. The consistency is defined between the genes' ranking based on edgeR/DESeq2's p-values or Clipper's contrast scores and their ranking based on true expression fold changes. **(c)** Reproducibility between two semi-synthetic datasets as technical replicates. Three reproducibility criteria are used: the IDR, Pearson correlation, and Spearman correlation. Each criterion is calculated for edgeR/DESeq2's p-values or Clipper's contrast scores on the two semi-synthetic datasets. Among the three methods, only Clipper controls the FDR, and Clipper achieves the highest power, the best gene ranking consistency, and the best reproducibility.

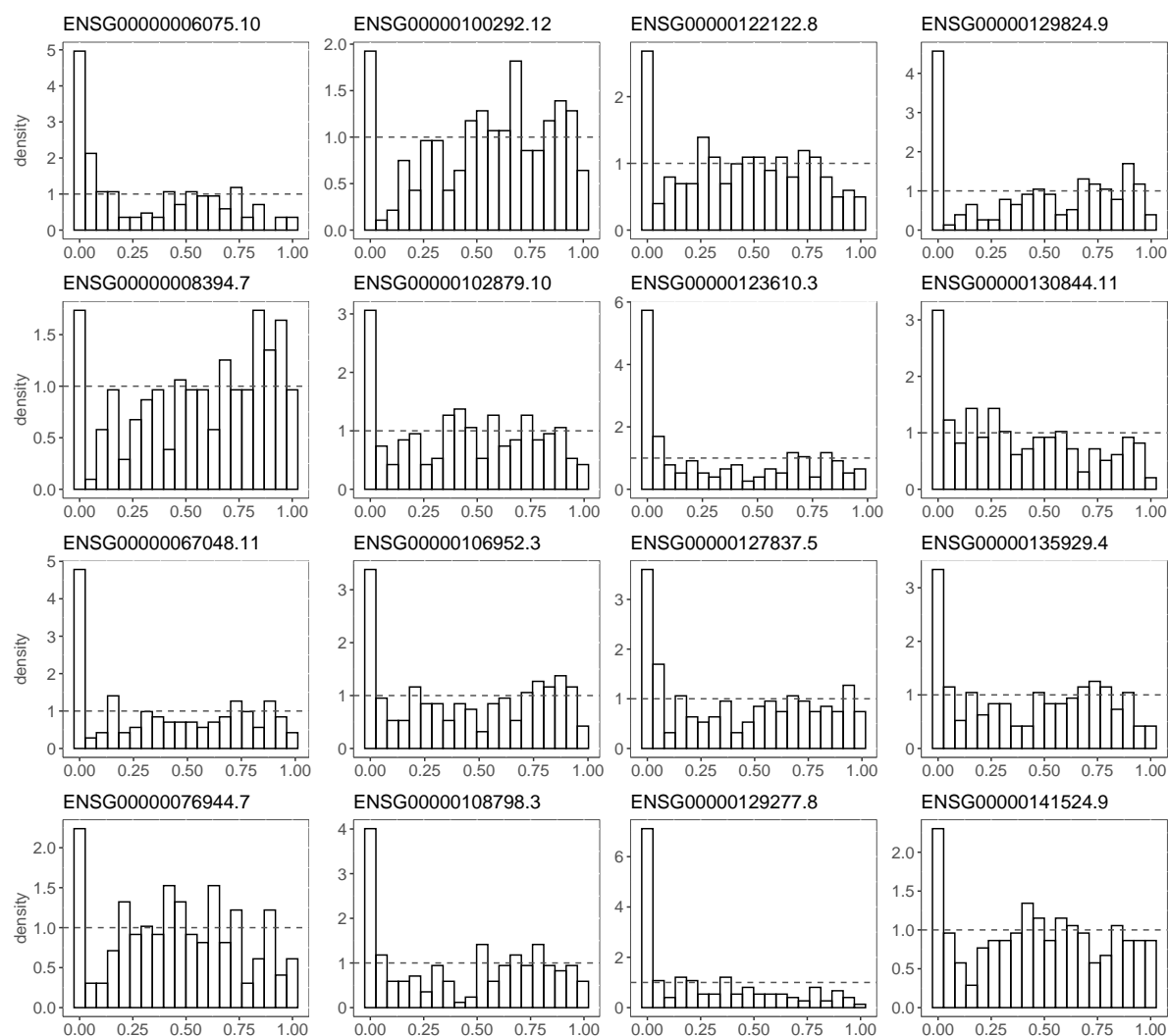


Figure S18: The p-value distributions of 16 non-DEGs that are most frequently identified by DESeq2 at $q = 5\%$ from 200 semi-synthetic datasets. The p-values of these 16 genes tend to be overly small, and their distributions are non-uniform with a mode close to 0.

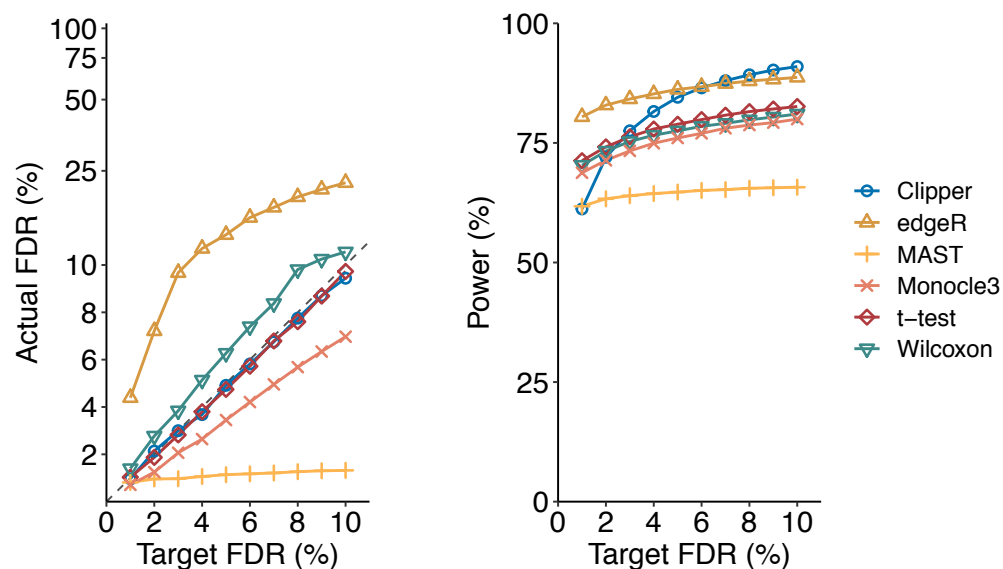


Figure S19: Comparison of Clipper and five scRNA-seq DEG identification methods on synthetic Drop-seq data generated by scDesign2 (based on a real Drop-seq dataset of PBMCs). The target FDR threshold q ranges from 1% to 10%. In the “Actual FDR vs. Target FDR” plot (left), points above the dashed diagonal line indicate failed FDR control. Clipper controls the FDR while maintaining high power, demonstrating Clipper’s good performance in single-cell DE analyses.

a

GO terms enriched in Clipper-specific DEGs in Clipper vs. DESeq2 comparison

GO term (ID)	qvalue (Clipper)
neutrophil activation (GO:0042119)	3.104557e-10
granulocyte activation (GO:0036230)	3.104557e-10
neutrophil degranulation (GO:0043312)	8.587750e-10
neutrophil activation involved in immune response (GO:0002283)	8.591455e-10
neutrophil mediated immunity (GO:0002446)	3.104557e-10

b

GO terms enriched in Clipper-specific DEGs in Clipper vs. edgeR comparison

GO term (ID)	qvalue (Clipper)
neutrophil degranulation (GO:0043312)	8.587750e-10
neutrophil activation involved in immune response (GO:0002283)	8.591455e-10
neutrophil activation (GO:0042119)	3.104557e-10
neutrophil mediated immunity (GO:0002446)	3.104557e-10
granulocyte activation (GO:0036230)	3.104557e-10
cellular response to chemical stress (GO:0062197)	2.157116e-03
response to oxidative stress (GO:0006979)	3.141033e-03
cellular response to oxidative stress (GO:0034599)	2.902893e-03

Figure S20: Enrichment q-values of GO terms that are found enriched in the DEGs that are uniquely identified by Clipper in pairwise comparison of **(a)** Clipper vs. edgeR and **(b)** Clipper vs. DESeq2. These GO terms are all related to immune response and thus biologically meaningful.

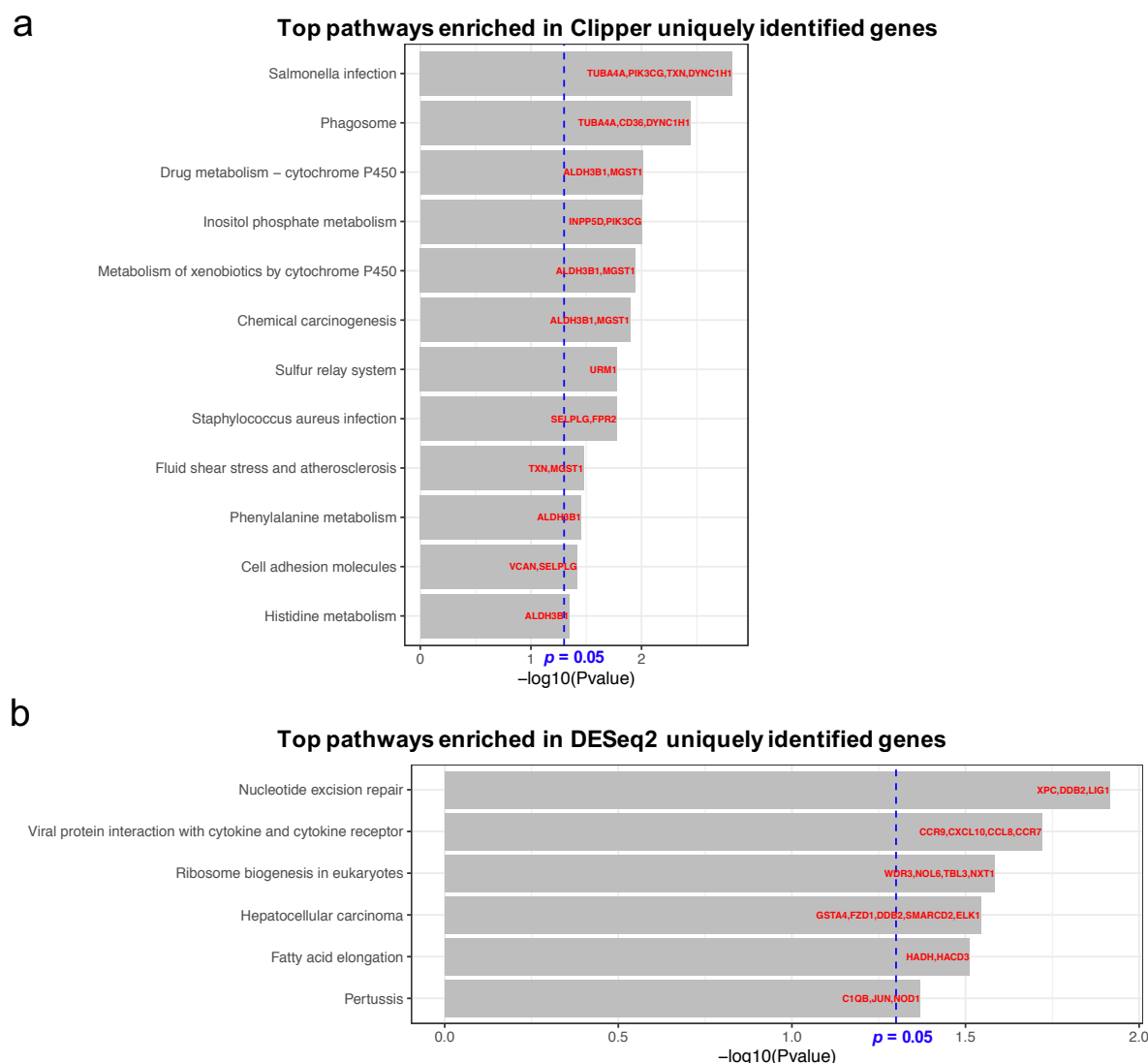


Figure S21: The p-values of the top enriched pathways in the DEGs that are uniquely identified by (a) Clipper and (b) DESeq2; i.e., the DEGs that are only identified by one method by missed by the other two methods. There are more immune-related pathways enriched in (a) than (b).

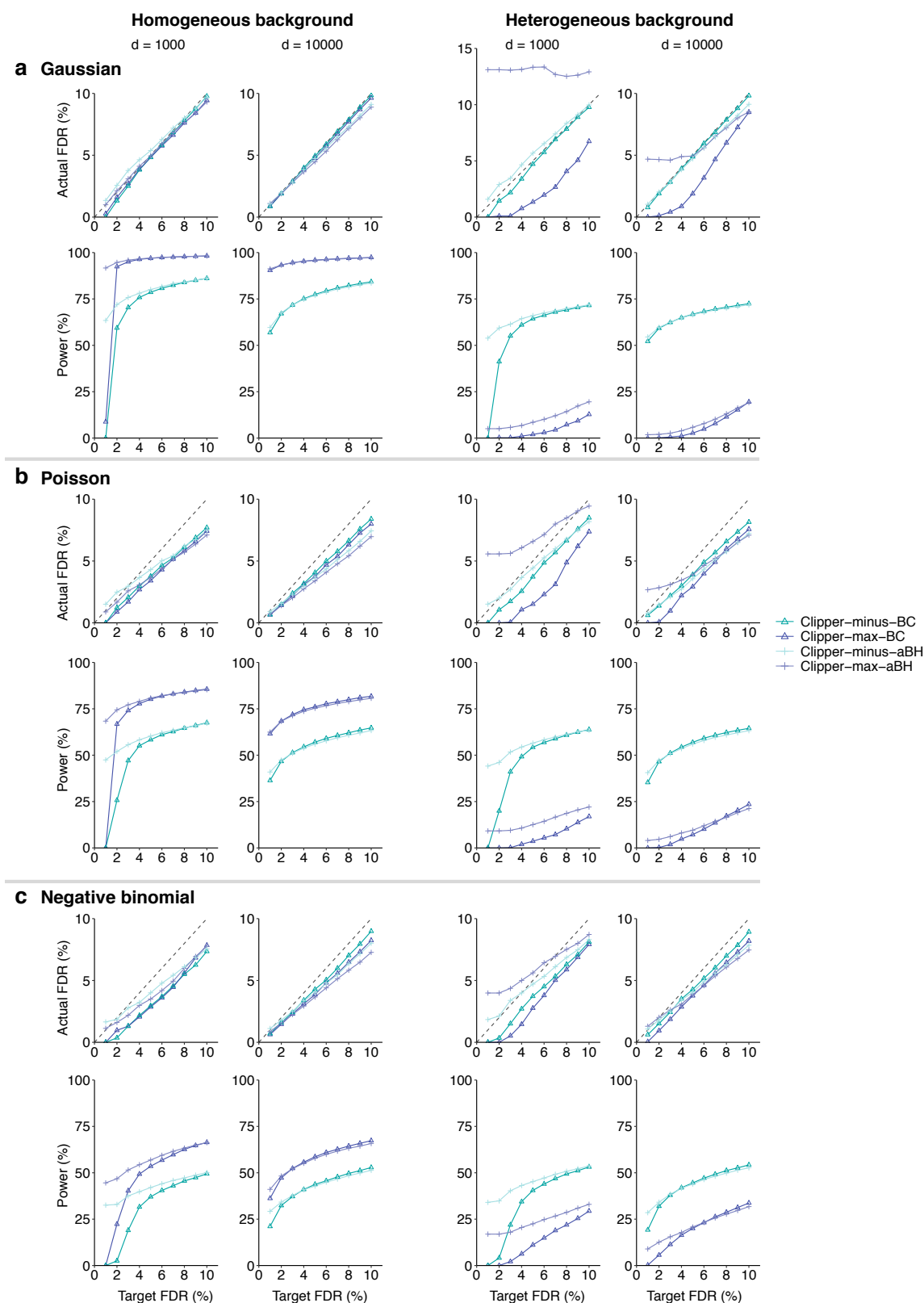


Figure S22: In 1vs1 enrichment analysis, comparison of four Clipper variant algorithms (Clipper-minus-BC, Clipper-minus-aBH, Clipper-max-BC, and Clipper-max-aBH) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 1000$ or $10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Clipper-minus-BC is chosen as the default implementation under this scenario.

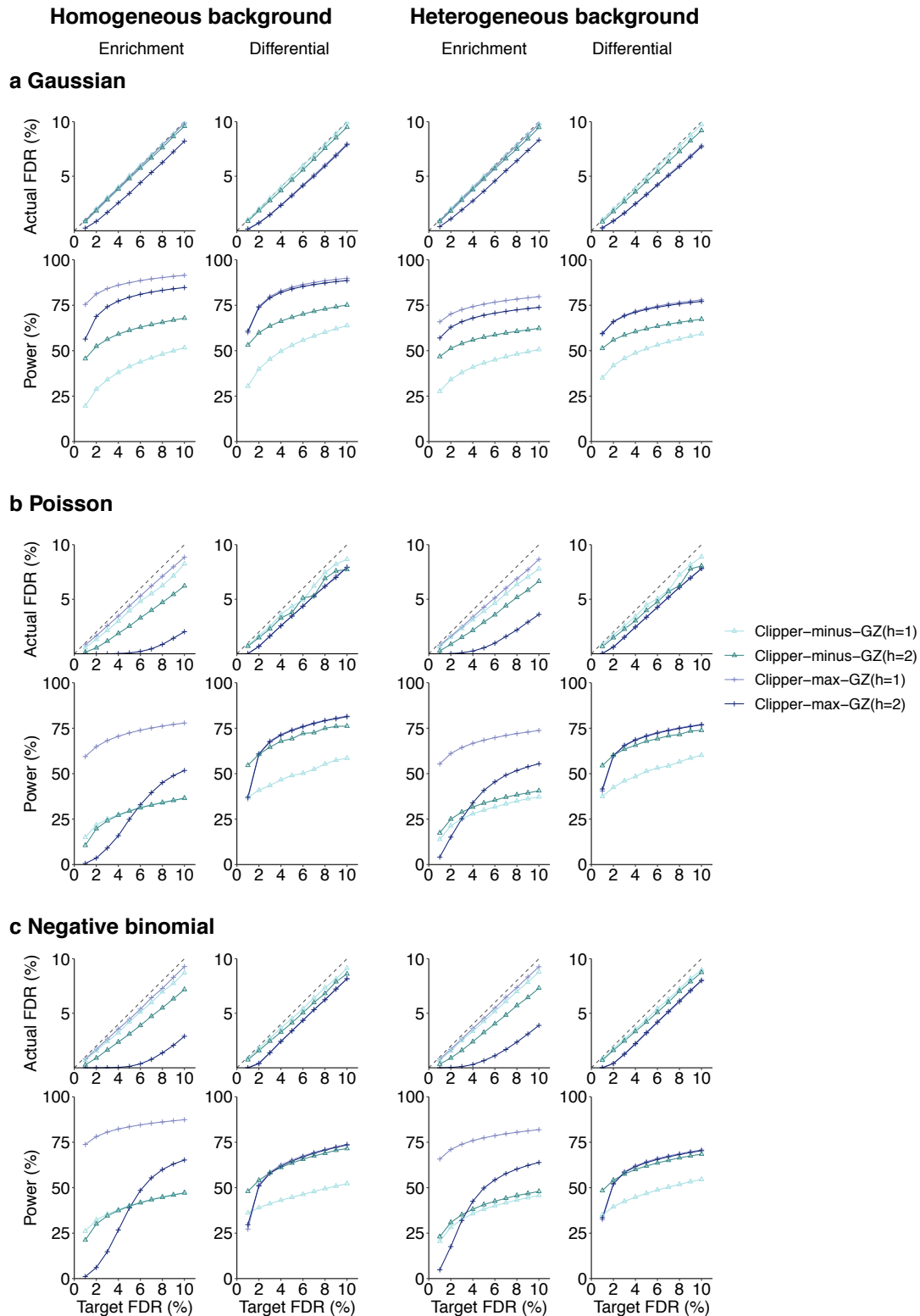


Figure S23: In the 2vs1 enrichment analysis (columns 1 and 3) and differential analysis (columns 2 and 4), comparison of four Clipper variant algorithms (Clipper-minus-GZ(h=1), Clipper-minus-GZ(h=2), Clipper-max-GZ(h=1), and Clipper-max-GZ(h=2)) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Clipper-max-GZ(h=1) is chosen as the default implementation under this scenario.

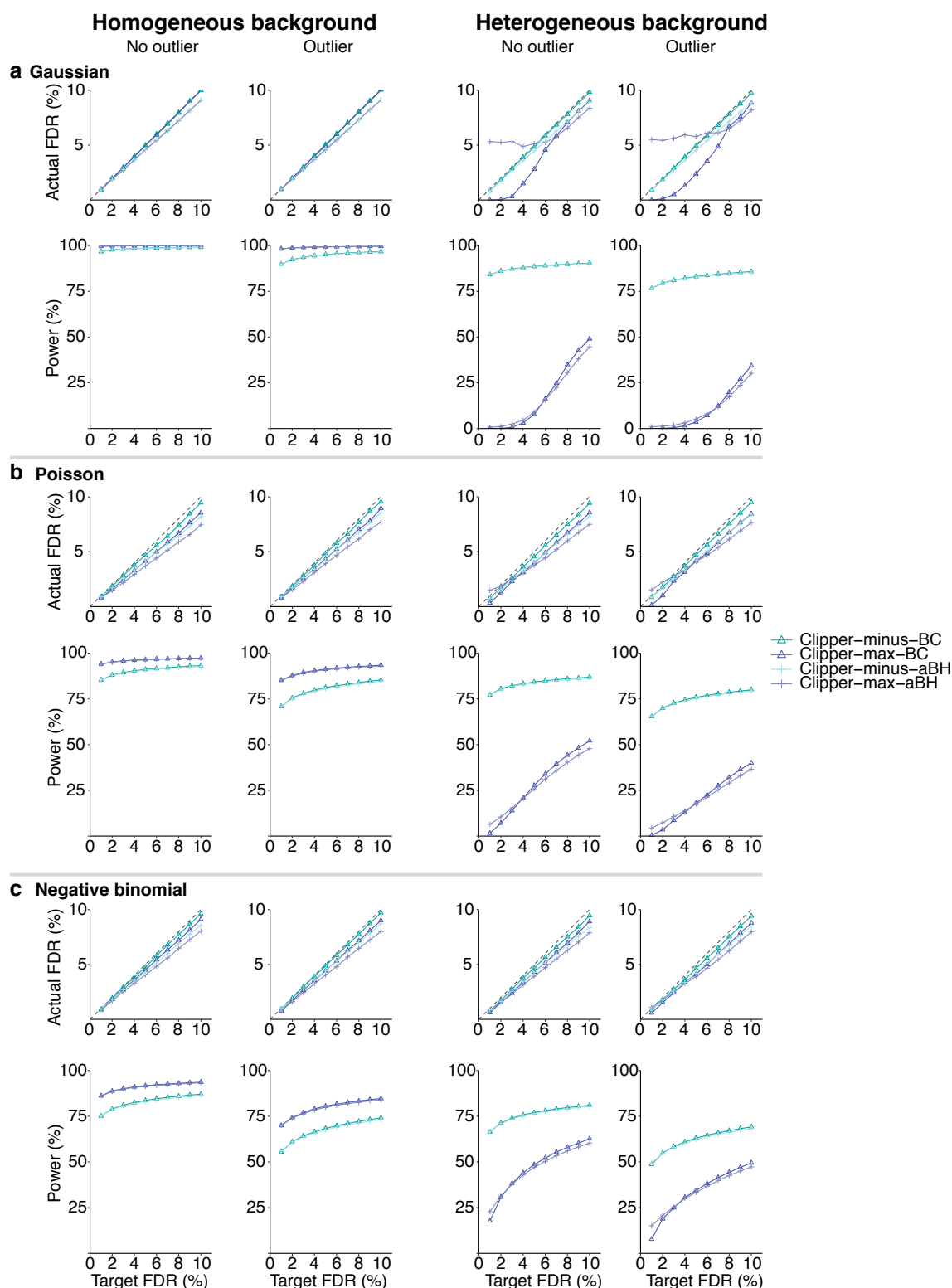


Figure S24: In 3vs3 enrichment analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of four Clipper variant algorithms (Clipper-minus-BC, Clipper-minus-aBH, Clipper-max-BC, and Clipper-max-aBH) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Clipper-minus-BC is chosen as the default implementation under this scenario.

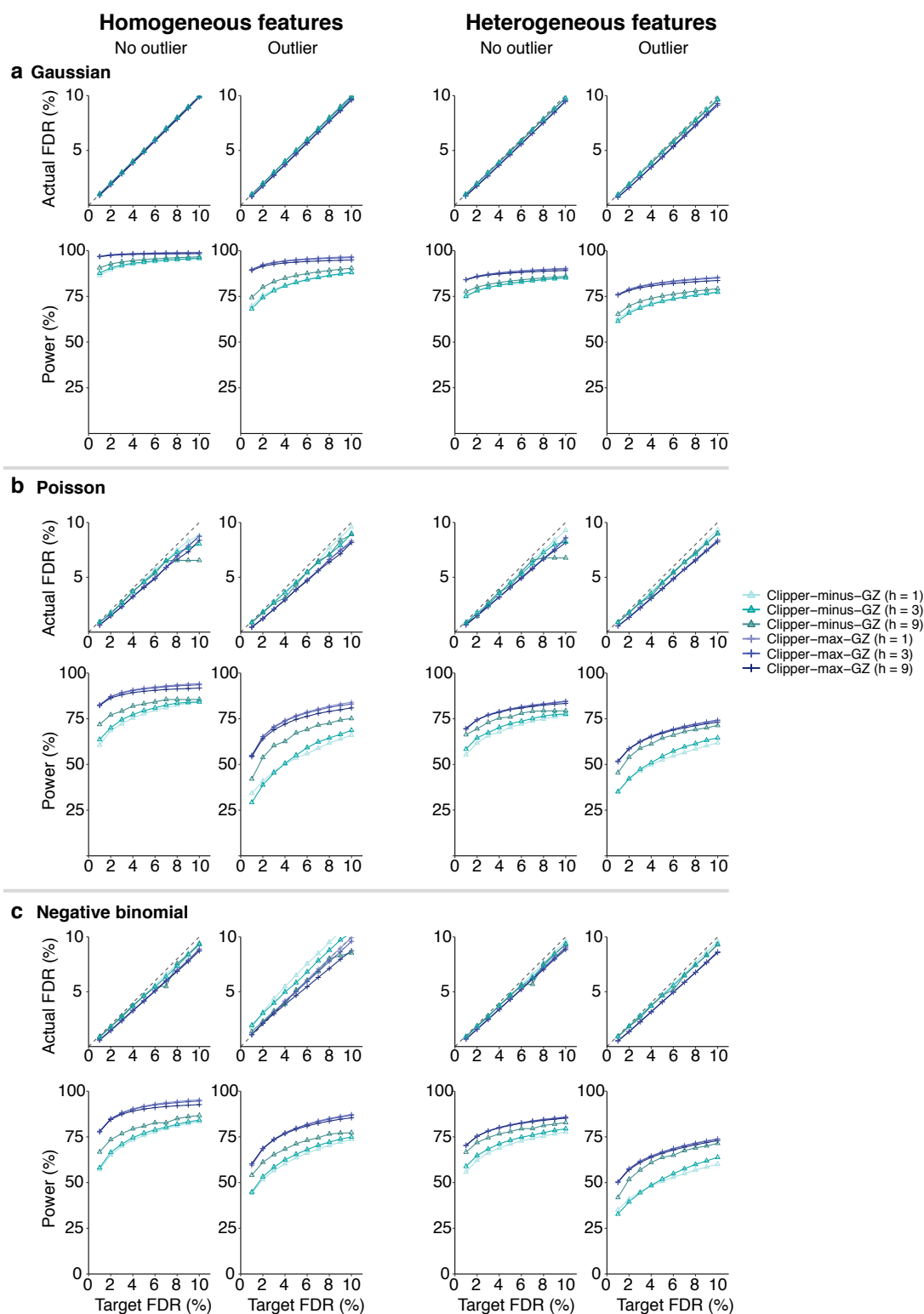


Figure S25: In 3vs3 differential analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of six Clipper variant algorithms (Clipper-minus-GZ($h=1$), Clipper-minus-GZ($h=3$), Clipper-minus-GZ($h=9$), Clipper-max-GZ($h=1$), Clipper-max-GZ($h=3$), and Clipper-max-GZ($h=9$)) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Clipper-max-GZ($h=1$) is chosen as the default implementation under this scenario.

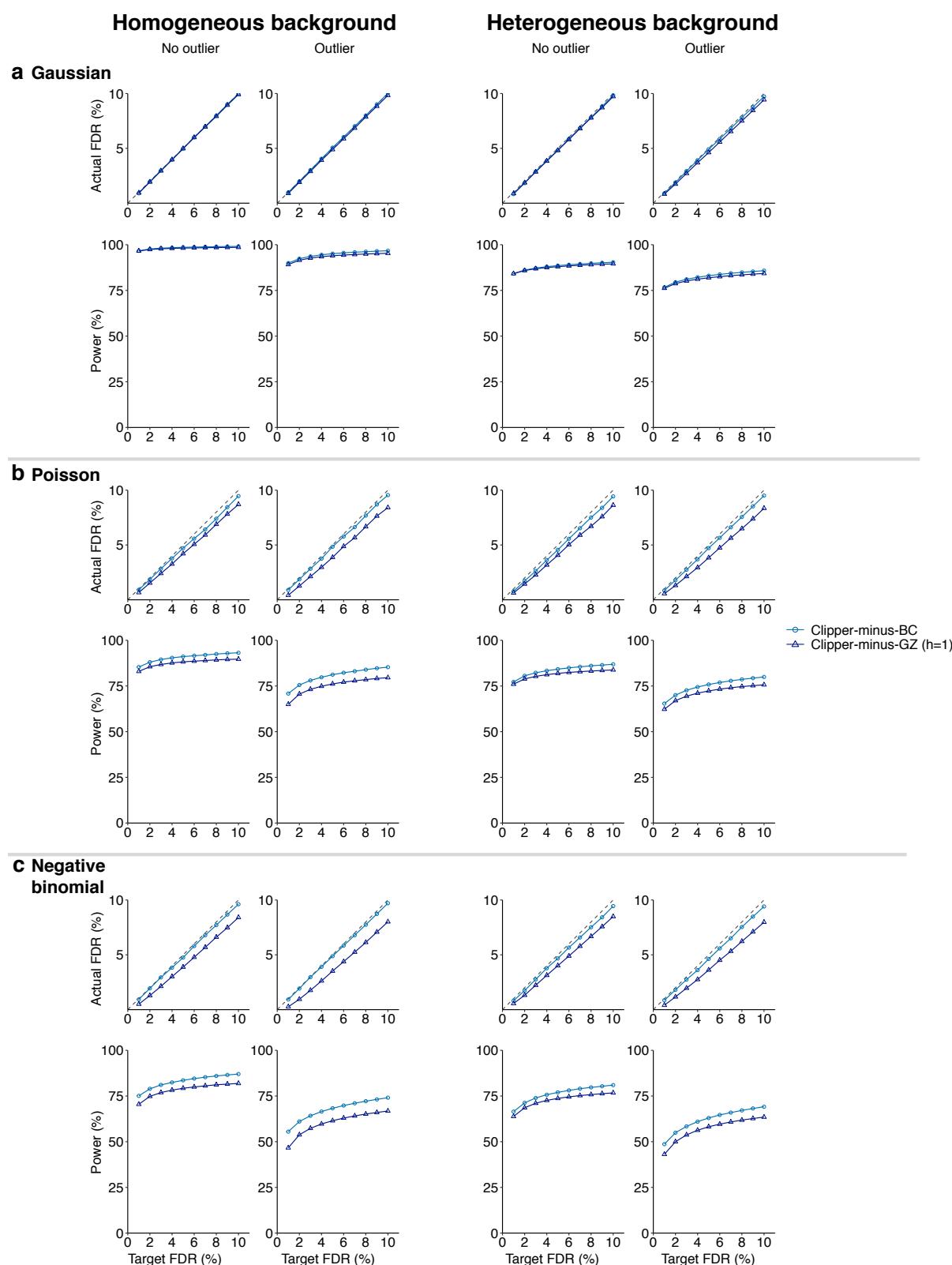


Figure S26: In the 3vs3 enrichment analysis without (columns 1 and 3) or with outliers (columns 2 and 4), comparison of two Clipper variant algorithms (Clipper-minus-BC, Clipper-max-GZ(h=1)) in terms of their FDR control and power. At target FDR thresholds $q \in \{1\%, 2\%, \dots, 10\%\}$, each method's actual FDRs and power are evaluated on 200 simulated datasets with $d = 10,000$ features generated from (a) the Gaussian distribution, (b) the Poisson distribution, or (c) the negative binomial distribution under homogeneous (two left columns) and heterogeneous (two right columns) background scenarios. Clipper-minus-BC is chosen as the default implementation under this scenario.

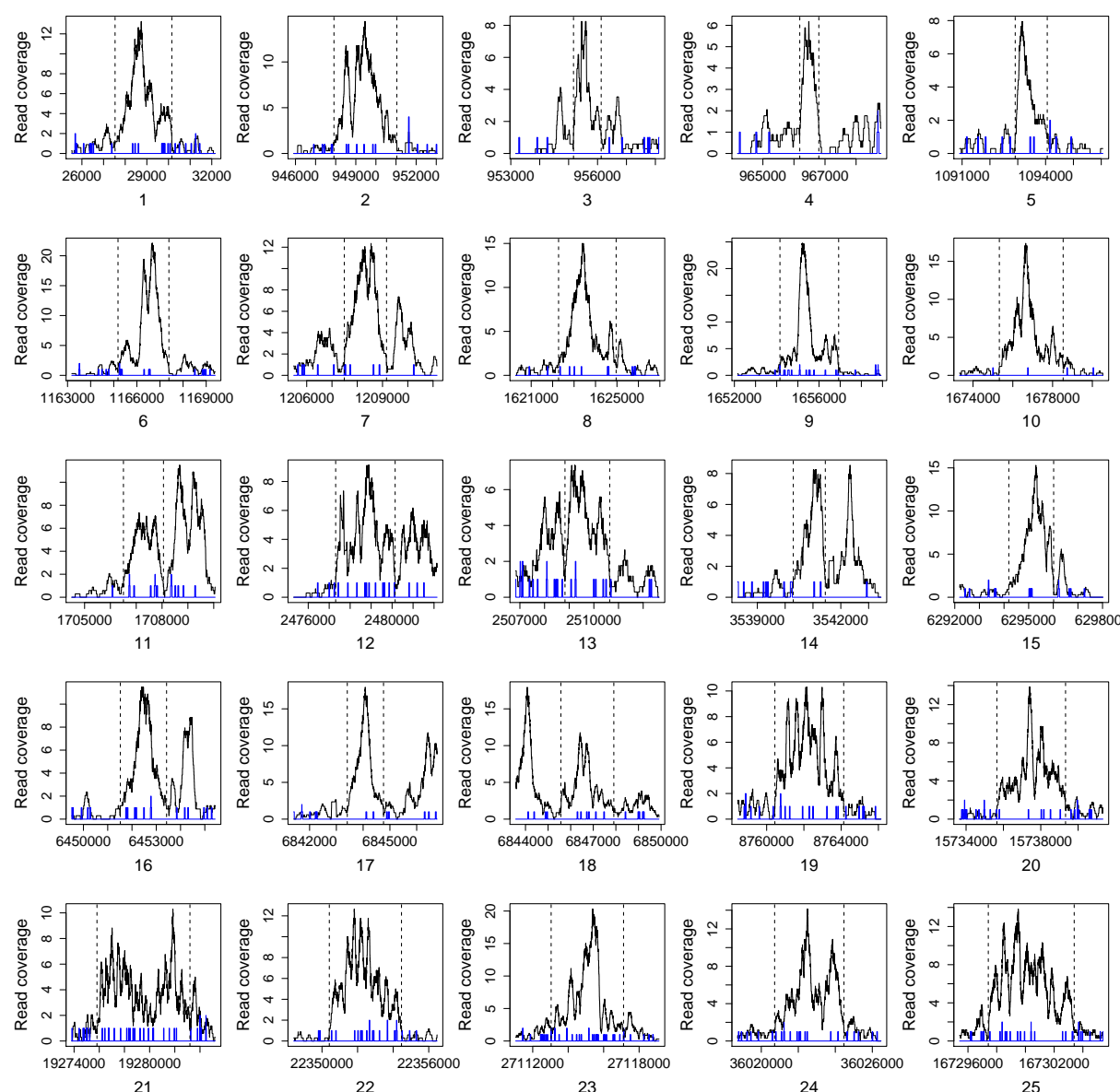


Figure S27: 25 "true peaks" from H3K4me3 ChIP-seq data of cell line GM12878. Black and blue curves indicate the read coverages in the experimental and control samples, respectively. Vertical dashed lines indicate the peak boundaries.

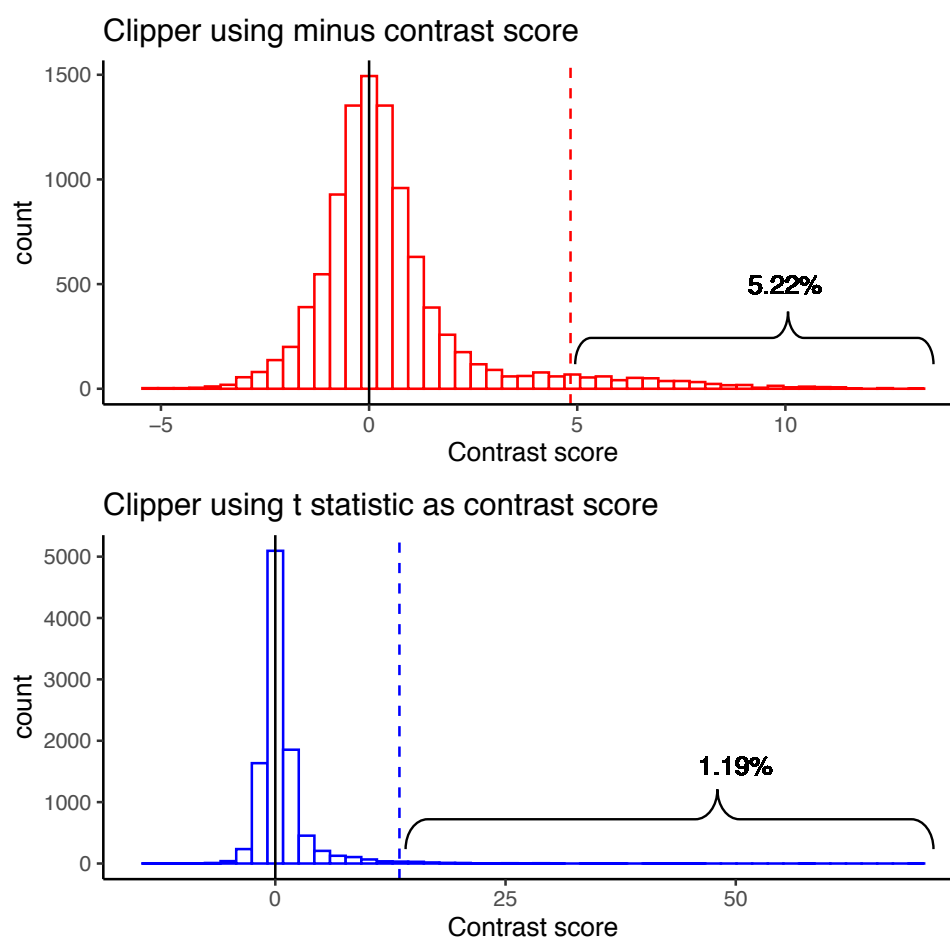


Figure S28: In the 3vs3 enrichment analysis, distributions of contrast scores used by two Clipper variants: the default Clipper using the minus contrast score (top) and the Clipper variant using the two-sample t statistic (bottom). Features are generated from the Gaussian distribution under the heterogeneous background scenario (see Supp. Section S4). The vertical dashed lines indicate the contrast score cutoffs found by the BC procedure at the target FDR threshold $q = 1\%$. The distribution of the minus contrast scores has a heavier right tail (5.22%) than that of the distribution of the t statistic contrast scores (1.19%).

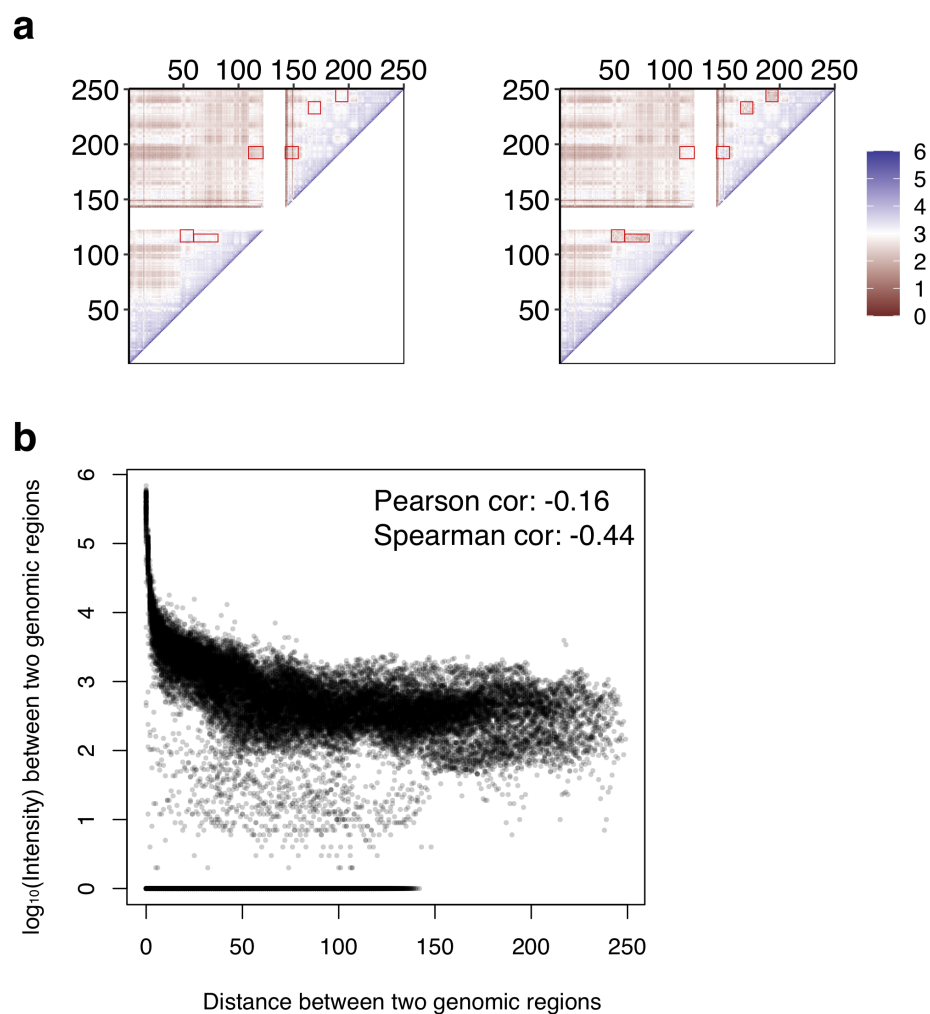


Figure S29: (a) \log_{10} -transformed mean Hi-C interaction matrices (μ_X and μ_Y in Section S6.5) under the two conditions. DIR regions are highlighted in red squares. **(b)** In one randomly picked Hi-C semi-synthetic dataset, closer genomic regions tend to have higher contact intensities.

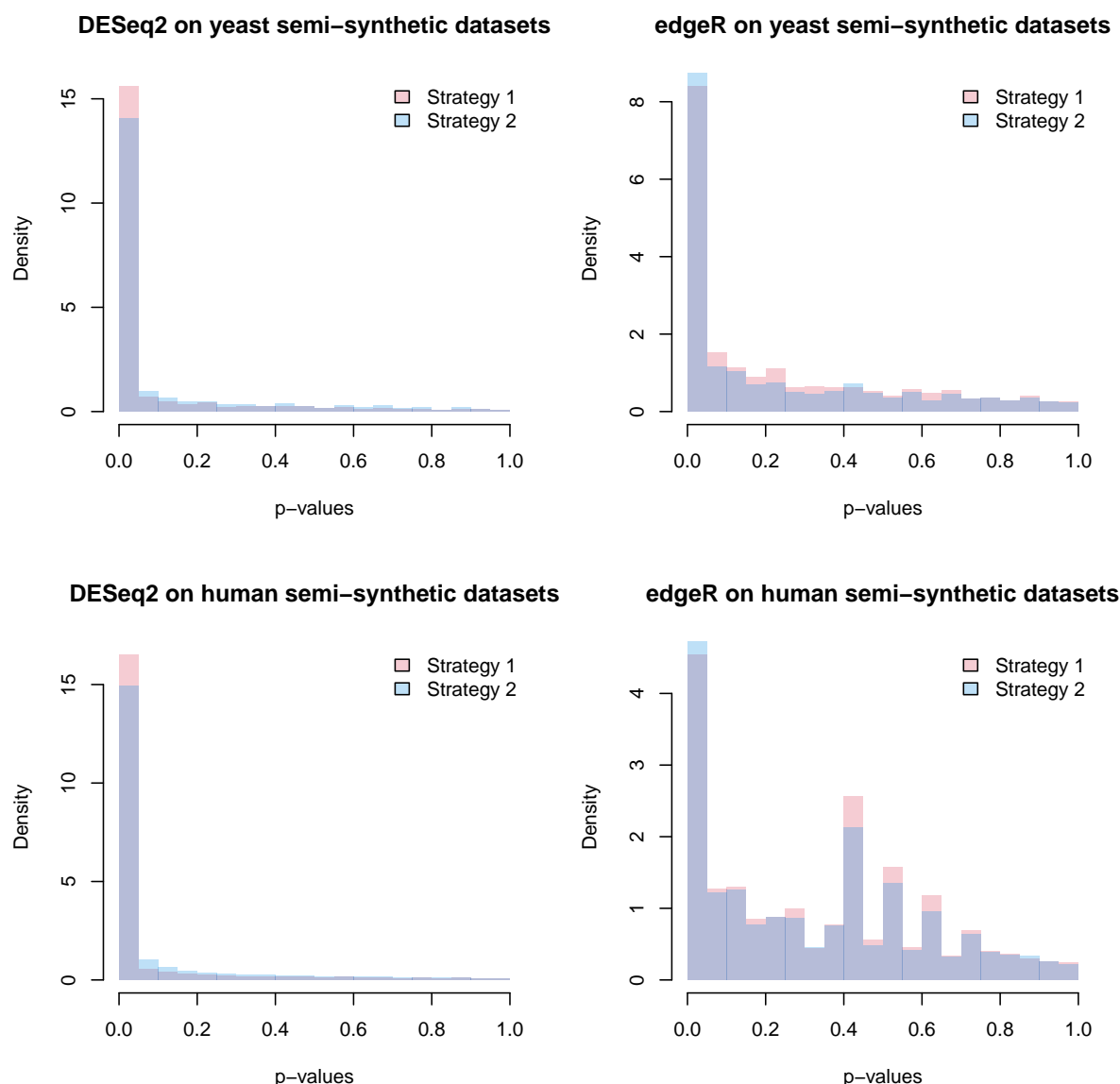


Figure S30: Histograms of p-values (one p-value per non-DEG) that are obtained by testing whether each non-DEG's p-values (output by DESeq2 or edgeR) follow a uniform distribution. For each real dataset (yeasts on the top and human monocytes on the bottom) and each simulation strategy (red for strategy 1 and blue for strategy 2), a histogram is plotted for DESeq2 (left) or edgeR (right); each p-value is calculated across 100 semi-synthetic datasets (excluding NA p-values). In each panel, more right skewed histograms are considered better.

References

- [1] Yoav Benjamini and Yosef Hochberg. “Controlling the false discovery rate: a practical and powerful approach to multiple testing”. In: *Journal of the Royal statistical society: series B (Methodological)* 57.1 (1995), pp. 289–300.
- [2] John D Storey. “A direct approach to false discovery rates”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64.3 (2002), pp. 479–498.
- [3] Rina Foygel Barber and Emmanuel J Candès. “Controlling the false discovery rate via knockoffs”. In: *The Annals of Statistics* 43.5 (2015), pp. 2055–2085.
- [4] Bradley Efron. “Large-scale simultaneous hypothesis testing: the choice of a null hypothesis”. In: *Journal of the American Statistical Association* 99.465 (2004), pp. 96–104.
- [5] Bradley Efron et al. “Empirical Bayes analysis of a microarray experiment”. In: *Journal of the American statistical association* 96.456 (2001), pp. 1151–1160.
- [6] Abhishek K Sarkar and Matthew Stephens. “Separating measurement and expression models clarifies confusion in single cell RNA-seq analysis”. In: *BioRxiv* (2020).
- [7] Ery Arias-Castro and Shiyun Chen. “Distribution-free multiple testing”. In: *Electronic Journal of Statistics* 11.1 (2017), pp. 1983–2001.
- [8] Rina Foygel Barber, Emmanuel J Candès, et al. “A knockoff filter for high-dimensional selective inference”. In: *The Annals of Statistics* 47.5 (2019), pp. 2504–2537.
- [9] Jaime Roquero Gimenez and James Zou. “Improving the Stability of the Knockoff Procedure: Multiple Simultaneous Knockoffs and Entropy Maximization”. In: *arXiv preprint arXiv:1810.11378* (2018).
- [10] Yong Zhang et al. “Model-based analysis of ChIP-Seq (MACS)”. In: *Genome biology* 9.9 (2008), pp. 1–9.
- [11] Sven Heinz et al. “Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities”. In: *Molecular cell* 38.4 (2010), pp. 576–589.
- [12] Marina Spivak et al. “Improvements to the percolator algorithm for peptide identification from shotgun proteomics data sets”. In: *Journal of proteome research* 8.7 (2009), pp. 3737–3745.
- [13] Mark D Robinson, Davis J McCarthy, and Gordon K Smyth. “edgeR: a Bioconductor package for differential expression analysis of digital gene expression data”. In: *Bioinformatics* 26.1 (2010), pp. 139–140.
- [14] Michael I Love, Wolfgang Huber, and Simon Anders. “Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2”. In: *Genome biology* 15.12 (2014), p. 550.
- [15] Greg Finak et al. “MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data”. In: *Genome biology* 16.1 (2015), pp. 1–13.
- [16] Xiaojie Qiu et al. “Single-cell mRNA quantification and differential analysis with Census”. In: *Nature methods* 14.3 (2017), pp. 309–315.
- [17] John C Stansfield, Kellen G Cresswell, and Mikhail G Dozmorov. “multiHiCcompare: joint normalization and comparative analysis of complex Hi-C experiments”. In: *Bioinformatics* 35.17 (2019), pp. 2916–2923.

- [18] John C Stansfield et al. “HiCompare: an R-package for joint normalization and comparison of Hi-C datasets”. In: *BMC bioinformatics* 19.1 (2018), p. 279.
- [19] Aaron TL Lun and Gordon K Smyth. “diffHic: a Bioconductor package to detect differential genomic interactions in Hi-C data”. In: *BMC bioinformatics* 16.1 (2015), pp. 1–11.
- [20] Mohamed Nadhir Djekidel, Yang Chen, and Michael Q Zhang. “FIND: diffERential chromatin INteractions Detection using a spatial Poisson process”. In: *Genome research* 28.3 (2018), pp. 412–422.
- [21] ENCODE Project Consortium et al. “An integrated encyclopedia of DNA elements in the human genome”. In: *Nature* 489.7414 (2012), pp. 57–74.
- [22] Kevin L Howe et al. “Ensembl 2021”. In: *Nucleic acids research* 49.D1 (2021), pp. D884–D891.
- [23] Dattatreya Mellacheruvu et al. “The CRAPome: a contaminant repository for affinity purification–mass spectrometry data”. In: *Nature methods* 10.8 (2013), pp. 730–736.
- [24] Jürgen Cox and Matthias Mann. “MaxQuant enables high peptide identification rates, individualized ppb-range mass accuracies and proteome-wide protein quantification”. In: *Nature biotechnology* 26.12 (2008), pp. 1367–1372.
- [25] Anton A Goloborodko et al. “Pyteomics—a Python framework for exploratory data analysis and rapid software prototyping in proteomics”. In: *Journal of The American Society for Mass Spectrometry* 24.2 (2013), pp. 301–304.
- [26] Lev I Levitsky et al. “Pyteomics 4.0: five years of development of a Python proteomics framework”. In: *Journal of proteome research* 18.2 (2018), pp. 709–714.
- [27] Claire R Williams et al. “Empirical assessment of analysis workflows for differential expression analysis of human samples using RNA-Seq”. In: *BMC bioinformatics* 18.1 (2017), p. 38.
- [28] Marek Gierliński et al. “Statistical models for RNA-seq data derived from a two-condition 48-replicate experiment”. In: *Bioinformatics* 31.22 (2015), pp. 3625–3630.
- [29] Mark D Robinson and Alicia Oshlack. “A scaling normalization method for differential expression analysis of RNA-seq data”. In: *Genome biology* 11.3 (2010), pp. 1–9.
- [30] Tianyi Sun et al. “scDesign2: an interpretable simulator that generates high-fidelity single-cell gene expression count data with gene correlations captured”. In: *bioRxiv* (2020).
- [31] Jiarui Ding et al. “Systematic comparison of single-cell and single-nucleus RNA-sequencing methods”. In: *Nature biotechnology* 38.6 (2020), pp. 737–746.
- [32] Grace XY Zheng et al. “Massively parallel digital transcriptional profiling of single cells”. In: *Nature communications* 8.1 (2017), pp. 1–12.
- [33] Evan Z Macosko et al. “Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets”. In: *Cell* 161.5 (2015), pp. 1202–1214.