

Analyzing microbial ChIP-Seq data with Pique

Russell Y. Neches^{1,3*}, Elizabeth G. Wilbanks^{1,3}, Phillip M. Seitzer,^{2,3} and Marc T. Facciotti^{2,3}

¹Microbiology Graduate Group, University of California, Davis.

²Department of Biomedical Engineering, University of California, Davis.

³Genome Center, University of California, Davis.

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

ABSTRACT

Motivation:

While numerous effective peak finders have been developed for eukaryotic systems, we have found that the approaches used can be surprisingly error prone when run on high-coverage bacterial and archaeal ChIP-Seq datasets.

Results:

We have developed Pique, a conceptually simple, easy to use ChIP-Seq peak finding application for bacterial and archaeal ChIP-Seq experiments. The software is cross-platform and Open Source, and based on only Open Source dependencies. Output is provided in standardized file formats, and may be easily imported by the Gaggles Genome Browser for manual curation and data exploration, or into statistical and graphics software such as R for further analysis.

Availability:

The software is available under the BSD-3 license at <http://github.com/ryneches/pique>.

A tutorial and test data are included with the documentation.

Contact: ryneches@ucdavis.edu

1 INTRODUCTION

Next generation sequencing coupled with chromatin immunoprecipitation (ChIP-Seq) is revolutionizing our ability to genomically map protein-DNA interaction. The growing popularity of ChIP-Seq has spurred the development of over thirty peak picking algorithms (an extensive survey of these packages was conducted by Wilbanks *et al.* [1]). The relative performance of representative peak detection algorithms on eukaryotic data, and methods to assess performance have been recently reviewed by several authors [2, 3, 4, 5, 6].

While ChIP-Seq has been predominantly used to interrogate protein-DNA interactions in eukaryotic systems, it is an especially powerful tool for studying microbes. The small genomes and rapid growth rates, as well as the extensive repertoire of experimental genetic tools available for microbial systems permit ChIP-Seq to provide a particularly clear picture of the state of a cell's transcriptional regulatory machinery.

However, only one ChIP-Seq analysis tool, CSDeconv [7], has been explicitly developed for microbial data. This MATLAB package successfully finds peaks in microbial ChIP-Seq data, but

its application is limited by its dependency on costly proprietary software, slow performance, lack of support for manual curation. Herein we describe Pique, a conceptually simple, Python-based peak finding package that enables easy and rapid peak finding in bacterial and archaeal ChIP-Seq datasets.

2 APPROACH

ChIP-Seq in bacteria and archaea yields coverage several orders of magnitude larger than in eukaryotic systems, resulting in continuous coverage rather than the sparse coverage typically present in eukaryotic ChIP-Seq data. This feature of microbial ChIP-Seq experiments permits simpler, faster algorithms to be used. We have based our algorithm on classic noise reduction techniques from signal processing.

Pique is designed for use in systems that have genomic complexities such as IS elements, gene dosage polymorphisms and accessory genomes that cause coverage variations unrelated to ChIP, or in cases where the organism under study is not identical to the reference genome. The resulting enrichment “pedestals” and “holes” can be problematic for detecting peaks and calculating enrichment levels. If the user provides a map of these features, the software will automatically perform a segmented analysis.

The wide variety of microbial systems, target proteins, protocols, and experimental conditions calls for tailored statistical approaches to ChIP-Seq. Rather than attempting to anticipate each of these (and their combinations) with a very large number of statistical and heuristic parameters, we have chosen to focus on the aspects of the analysis that are common to all ChIP-Seq experiments; finding putative peaks, estimating binding coordinates and binding affinities. The determination of statistical significance is typically straightforward for any particular experiment, but is quite difficult to robustly generalize.

Pique allows users to create high-quality peak lists in two ways. First, for each peak reported we report metrics that can be used to ascertain which peaks are statistically significant (usually, this involves little more than sorting the table and choosing a cutoff). Second, we provide integrated support for curation using the Gaggles Genome Browser. This permits interactive curation of the peak list and analysis of the ChIP-Seq data in the context of other Gaggles-enabled resources. Interactive curation of a microbial ChIP-Seq data set can typically be completed in a few minutes.

*to whom correspondence should be addressed

3 METHODS

Pique requires BAM files as input[8]. Therefore, prior to using Pique, reads should first be quality filtered, quality trimmed, and aligned to a reference genome (ideally, all contigs of the reference genome should be used as the mapping target).

By default, Pique treats each contig as a single analysis region, but the user may designate regions within a contig for separate analysis. This may be useful when coverage levels are systematically altered over large regions. Pique supports three features types; analysis regions, masking regions, and normalization regions. Masking regions are simply masked out of their respective analysis regions, and are useful for removing coverage variation due to repetitive DNA. Normalization regions selected within an analysis region are used to compensate for total coverage discrepancies between the background and ChIP alignments.

The user launches the analysis by providing alignment an file for the ChIP data, an alignment file for the control data, and a coverage feature map. The primary analysis proceeds as follows :

- The alignment files are digested into numeric coverage tracks, and the analysis regions are initialized in memory. Masking regions are applied.
- High- k noise is removed using a Wiener-Kolmogorov filter. The filter delay α is chosen to approximate to the expected footprint size of the targeted protein. The choice of filter implies the existence of two inputs; a “true” signal, and a noise source. Both are assumed to be stationary stochastic processes combined additively.
- A Blackman window of a diameter equal to the read length is convolved with the filtered coverage track to remove features smaller than one read. This reduces the effect of fragmentation position bias, and may be especially useful when transposon-based library construction is used.
- The noise threshold in the ChIP coverage track is measured by comparing the coverage distribution in the ChIP track to the control track within user-annotated non-peak regions. Features that exceed the noise threshold are identified.
- Because read orientations are constrained by the fragment size, binding events cause a offset enrichment between the forward and reverse strands. Pique exploits this by requiring that the stop coordinate of the forward strand enrichment envelope must fall between the coordinates of the reverse strand enrichment envelope, and that the start coordinate of the reverse strand enrichment envelope fall between the coordinates of the forward strand enrichment envelope. (We call this the overlap criterion.)

For each putative peak, Pique calculates the enrichment ratio of the ChIP alignment to the control alignment, the binding coordinate, and the enrichment normalization factor for that analysis region.

4 RESULTS

Performance of Pique was benchmarked against CSDeconv two ChIP-Seq experiments. The first experiment targeted TfbD (transcription initiation factor IIB 4) DNA-binding events in late stationary phase of *Halobacterium salinarum* sp. NRC1 described in Wilbanks *et al.* [1]. Peak lists for both Pique and CSDeconv were compared with experimentally verified transcription start sites collected by Koide *et al* [11] and scanned for the presence of putative binding motifs using MAST[10]. An *ab initio* search for putative binding motifs in each peak list was conducted, and the recovered motifs compared using STAMP [12].

The second experiment targeted the DosR (dormancy survival regulator) DNA-binding events after two hours exposure to oxygen after growth to early log phase of *Mycobacterium tuberculosis* described in Lun *et al.* [7]. To the best of our knowledge,

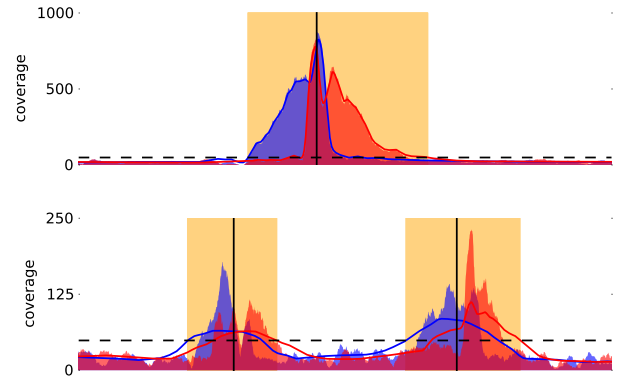


Fig. 1. Peaks found in the included sample dataset, derived from ChIP-seq of tfbD in *Halobacterium salinarum* sp. NRC1. Blue and red shading indicate coverage of reads aligned from the ChIP-derived data to the forward and reverse strands, respectively. Blue and red lines represent the filtered coverage levels for the forward and reverse strands, respectively. The dashed line is the detected noise threshold for the region. Detected peaks are indicated in orange boxes.

Found by	Peaks	TSS	Motif	TSS & Motif
Pique or CSDeconv	610	332	197	69
Pique	417	225	128	50
CSDeconv	449	213	138	50
Pique & CSDeconv	257	106	69	31
Pique only	160	119	59	19
CSDeconv only	192	107	69	19

Table 1. Putative binding motifs and transcript start sites detected in peaks recovered by Pique and CSDeconv.

transcription start sites have not been systematically verified experimentally in this organism, and so putative transcription start sites were identified using the putative binding motif found constructed by Gerasimova *et al.* [?] using MAST. An *ab initio* search for putative binding motifs was conducted, and the recovered motifs compared using STAMP.

4.1 TfbD in *Halobacterium salinarum* sp. NRC1

The putative TfbD binding motif for this organism is analogous to [??] a similar archaeon’s TFB protein in photocrosslinking study [9], and is thought to consist of a TFB recognition element (BRE), a TATA box, a proximal promoter element (PPE) and a transcription start site (TSS). This putative promoter motif was scanned against each peak region using MAST [10] with an e -value cutoff of 100 and the default motif p -value cutoff of 10^{-4} . Peaks found by each software package were also for presence of experimentally determined transcript start sites (TSS) [11] (Table 1).

Enrichment ratios (the ratio of the integrated coverage to the background) for each peak were calculated, and peaks from each software package for each category were then placed in rank order (Fig. 2). It was found that Pique recovered more peaks than CSDeconv, and more peaks with experimentally verified transcript start sites. It was found that Pique and CSDeconv recovered peaks

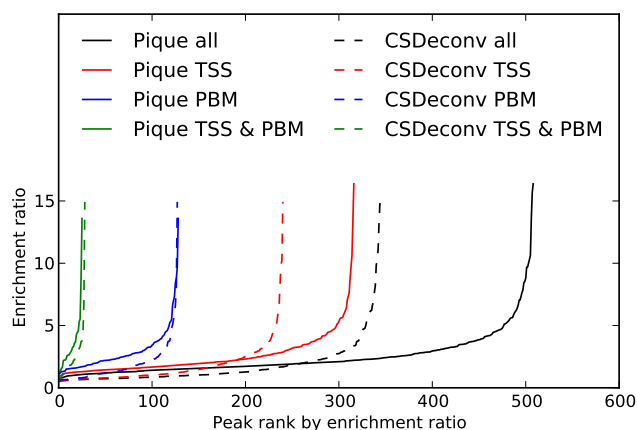


Fig. 2. Performance of Pique and CSDeconv on a ChIP-seq experiment (*Halobacterium salinarum* sp. NRC1 transcription initiation factor TFB tfbD in late stationary phase) was studied by comparing lists of peaks. Peak lists are shown here ranked by enrichment ratio. Pique recovers more peaks than CSDeconv, and more of these peak regions contain an experimentally determined transcription start site (TSS). Both software packages recovered the same number of peaks containing the predicted binding motif (PBM). Both software packages recovered the same number of peaks in which a transcript start site was found less than 50bp downstream from a binding motif. TSS coordinates were experimentally verified (Koide *et al.*) [11] and motifs were predicted using Motif Catcher (Seitzer *et al.*, publication pending).

containing the full putative binding motif (BRE-TATA-PPE-TSS) at the same rate.

In order to ascertain the quality of the predicted peak regions uniquely predicted by each software package, an *ab initio* search for the putative binding motif was conducted for each list of peaks.

Sequence entries were created by extracting 100-bp stretches of sequence centered at each respective called peak center. Each sequence data set was subjected to an iterative MC MAST/MEME MotifCatcher search [3] with 100 seeds. This search was conducted for five groups of sequences corresponding to all peak regions found by Pique, regions found exclusively by Pique, all peak regions found by CSDeconv, regions found by CSDeconv exclusively, and peak regions found by both software packages. Motifs could be discovered on the forward and reverse strand, and could be anywhere from 20 to 40 nucleotides in length (Supplementary Figure 2). A putative binding motif was discovered in all five sets of peaks.

The pairwise distances between motifs was computed by the average log likelihood ratio (ALLR), and clustered using UPGMA (program defaults). All motifs that were significantly similar to the putative canonical TFB motif were clustered, and for cases where more than two significant motif matches were discovered in a given data set, motifs were clustered and familial profiles were computed. Membership clustering thresholds were tried from 0 to 1.0 in increments of 0.10, and the familial profile motif with the lowest *e*-value was taken from each dataset as its “representative” motif profile. For cases where two or fewer significant motif matches were determined, the motif match with the lower *e*-value was taken as the representative motif profile of its respective data set. Representative

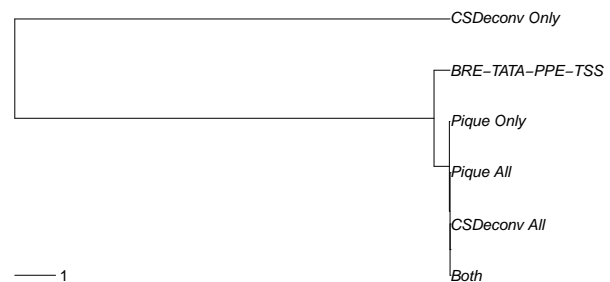


Fig. 3. Motifs recovered by Pique and CSDeconv are qualitatively different.

Peak list	Motif <i>e</i> -value
CSDeconv all	3.3×10^{-144}
CSDeconv only	5.2×10^{-35}
Pique & CSDeconv	6.9×10^{-77}
Pique all	4.6×10^{-130}
Pique only	4.9×10^{-82}

Table 2. Statistical support for motifs computed from lists of peaks found by Pique and CSDeconv.

motifs were compared using STAMP [12] with the ALLR motif comparison measure and UPGMA clustering (Fig. 3).

It was found that the motifs computed from the list of peak regions recovered by Pique and CSDeconv clustered closely with the canonical binding motif, as did motifs computed from the list of peak regions in common to both software packages. However, motifs computed from the list of peak regions recovered *exclusively* by Pique clustered with the canonical motif, but the motifs computed from the list of peak regions recovered *exclusively* by CSDeconv did not (Fig. 3). Furthermore, the *e*-value of the best motif computed from the CSDeconv-unique peak regions is the lowest among the five, suggesting that a higher proportion of false positives may be among this list.

4.2 DosR in *Mycobacterium tuberculosis*

5 DISCUSSION

Pique does not attempt to filter peaks that are statistically insignificant. We have found that this part of the analysis is usually specific to the data and to the experiment, and can be highly idiosyncratic. Pique is designed to achieve a low false-negative rate. This allows Pique to work without modification on many different kinds of experiments, but at the cost of some post-filtering. Pique provides the user with output that can be used to support a variety of such statistical tests.

Some recommended filtering might include eliminating peaks that are significantly narrower than the size range of the sequencing library, peaks with a normalized enrichment ratio below unity, or peaks that have predicted binding sites that are very skewed from the center of the enriched region. Depending on how many peaks are recovered, the user may wish to try one or all of these, perhaps with clustering. However, if a “perfect” peak list is required, we have found that heuristic filtering is inadequate regardless of the software used. To facilitate manual curation, Pique outputs a track

file of the coverage, a quantitative positional data of the estimated binding sites, and a bookmark file annotating the peaks. These files are simple to process by a variety of tools, and can be loaded directly into the Gaggle Genome Browser. False positives are easy to recognize visually, and can be easily deleted.¹

6 CONCLUSION

We conclude that Pique provides a rapid, open source platform for the sensitive detection of transcription factor binding sites in bacterial and archaeal ChIP-seq experiments. We leverage standard signal processing algorithms to rapidly identify binding sites. Downstream analysis is supported via integration with statistical and graphics software such as R, and curation via integration with the user-friendly Gaggle Genome Browser and the suite of Gaggle tools.

We note that Pique should also work well with eukaryotic datasets provided they are gathered with greater coverage than has been previously reported.

ACKNOWLEDGEMENT

Funding: This project was funded by UC Davis startup funds to M.T.F., NSF Graduate Research Fellowship awarded to E.G.W. and DARPA award number HR0011-05-1-0057 to R.Y.N.

REFERENCES

- [1]Elizabeth G. Wilbanks and Marc T. Facciotti. Evaluation of Algorithm Performance in ChIP-Seq Peak Detection. *PLoS ONE*, 5(7):e11471+, July 2010.
- [2]Shirley Pepke, Barbara Wold, and Ali Mortazavi. Computation for ChIP-seq and RNA-seq studies. *Nature Methods*, 6(11s):S22–S32, October 2009.
- [3]Teemu Laajala, Sunil Raghav, Soile Tuomela, Riitta Lahesmaa, Tero Aittokallio, and Laura Elo. A practical comparison of methods for detecting transcription factor binding sites in ChIP-seq experiments. *BMC Genomics*, 10(1):618+, December 2009.
- [4]Adam M. Szalkowski and Christoph D. Schmid. Rapid innovation in ChIP-seq peak-calling algorithms is outdistancing benchmarking efforts. *Briefings in Bioinformatics*, November 2010.
- [5]Xin Feng, Robert Grossman, and Lincoln Stein. PeakRanger: A cloud-enabled peak caller for ChIP-seq data. *BMC Bioinformatics*, 12(1):139+, 2011.
- [6]Morten B. Rye, Pål Sætrom, and Finn Drabløs. A manually curated ChIP-seq benchmark demonstrates room for improvement in current peak-finder programs. *Nucleic Acids Research*, 39(4):e25, March 2011.
- [7]Desmond Lun, Ashley Sherrid, Brian Weiner, David Sherman, and James Galagan. A blind deconvolution approach to high-resolution mapping of transcription factor binding sites from ChIP-seq data. *Genome Biology*, 10(12):R142+, December 2009.
- [8]Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–2079, August 2009.
- [9]Matthew B. Renfrow, Nikolai Naryshkin, L. Michelle Lewis, Hung-Ta T. Chen, Richard H. Ebright, and Robert A. Scott. Transcription factor B contacts promoter DNA near the transcription start site of the archaeal transcription initiation complex. *The Journal of biological chemistry*, 279(4):2825–2831, January 2004.
- [10]T. L. Bailey and M. Gribskov. Combining evidence using p-values: application to sequence homology searches. *Bioinformatics*, 14(1):48–54, January 1998.
- [11]Tie Koide, David J. Reiss, J. Christopher Bare, Wyming L. Pang, Marc T. Facciotti, Amy K. Schmid, Min Pan, Bruz Marzolf, Phu T. Van, Fang-Yin Lo, Abhishek Pratap, Eric W. Deutsch, Amelia Peterson, Dan Martin, and Nitin S.

¹ See supplementary figure.

Baliga. Prevalence of transcription promoters within archaeal operons and coding sequences. *Molecular Systems Biology*, 5(1), June 2009.

[12]Shaun Mahony, Philip E. Auron, and Panayiotis V. Benos. DNA familial binding profiles made easy: comparison of various motif alignment and clustering strategies. *PLoS computational biology*, 3(3):e61+, March 2007.