# Pavian walkthrough with metagenomics data from brain infections

*July 28, 2016*

## General workflow

1) Import Kraken and Centrifuge report files (a sample set)
2) Look at the overall statistics across the samples
3) Compare the classification across samples
   - in specific domains
   - interactive table or heatmap
4) Zoom into one sample - flow chart or sunburst diagram
5) Zoom into one pathogen in one sample
   - view the alignment on the genome

## Import Kraken and Centrifuge report files

The first step is loading data. Pavian expects files ending with '.report.csv', generated with kraken-report or centrifuge-report, resp. Click 'Choose files' to select files to upload to Pavian. The uploaded files will be added to a new sample set with an auto-generated name. You can change the name of the samples and the sample set using the interface that appears, once the files are loaded.

In this walkthrough we use data from Salzberg et al. (2016). To load this data, click the 'Load example data' button. This study sequenced brain or spinal cord biopsies from 10 patients with suspected central nervous system (CNS) infections. Upon loading the data, the links to 'Results Overview', 'Comparison' and 'Sample' become available in the sidebar, and a table describes the loaded sample set (see Figure 1).

### Brain biopsies data

Salzberg et al. (2016) used sequencing to detect the presence of pathogenic microbes in brain or spinal cord biopsies from 10 patients with neurologic problems indicating possible infection, but for whom conventional clinical and microbiology studies yielded negative or inconclusive results. Direct DNA and RNA sequencing of brain tissue biopsies generated 8.3 million to 29.1 million sequence reads per sample.

Every samples is from a diseased patient. There are no healthy controls, and for all but one patient (PT8) there is only one sample available. The patients had different diseases; thus it was not expected that the same bug was the cause in the patients. The samples are their own control - ubiquitously present microbes were disregarded as sequencing or laboratory contaminants. Most of the reads in each sample are human, or from common contaminants - bacteria that always appear but are unlikely to cause the infection, such as P. acnes.

The FASTQ files for this study are available at http://www.ncbi.nlm.nih.gov/bioproject/PRJNA314149. The reads were classified with Kraken. The Kraken report files from these sampels are available at ADDLINK.

## Look at the overall statistics across the samples

Select 'Results Overview' in the sidebar to load the report files. This will show a vire with the number of reads in the different samples, as well as the overall classification in different categories (see Figure 2).
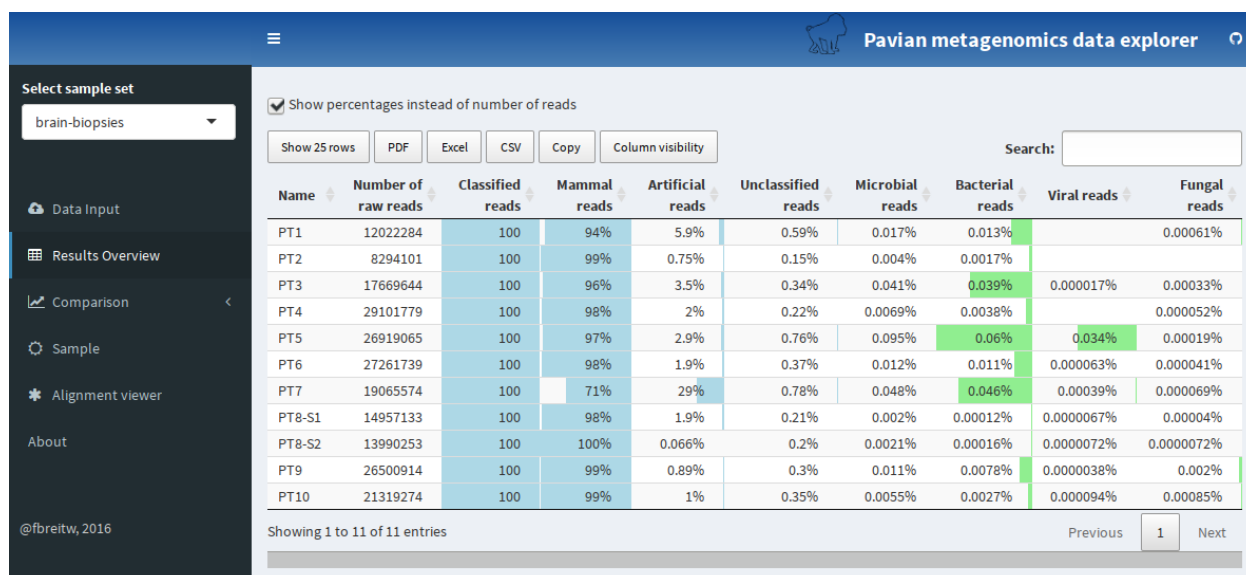
Figure 1: Pavian interface with the sample files loaded.

Figure 2: Results overview.

The samples of the brain biopsies dataset have between 8.3 and 29 million reads, most of them classified as 'mammal' (a 'mammal' is usually the host in these studies). There is a varying number of reads classifie as artificial - usually below 5%, but two outliers in sample PT1 and PT7 with 5.9 and 29%, respectively. The number of reads classified as microbial is below .1% in all samples.

## Compare the classification across samples

Click 'Comparison' and 'All data' to delve into the data (see Figure 3). The sample comparison view juxtaposes the identification results from the samples in a queryable table with taxa as rows and samples as columns. The table provides a visual guide to the values with inline bars, and the column taxid links to NCBI. You can decide to show percentages or z-scores, only show results at a certain taxon level, and filter uninteresting taxa. By default, the taxa 'Homo sapiens', 'synthetic construct' and 'unclassified' are filtered.

From this view, you can also select to see the data in the table in a heatmap representation. Note that the data displayed always correspond to the table - to show more rows in the heatmap, first select to show more rows in the table. Furthermore, you can look at how the samples are correlated with each other in 'Samples Clustering'.

For the brain biospies dataset, the table is rather wide with ten samples. The fourth column therefore provides an overview of the values in the columns. By default, the rows are sorted by to the mean over all samples. In this sample set, we immediately see an outlier in sample PT5, which is the only one with a substantial number of JC polyomavirus reads.
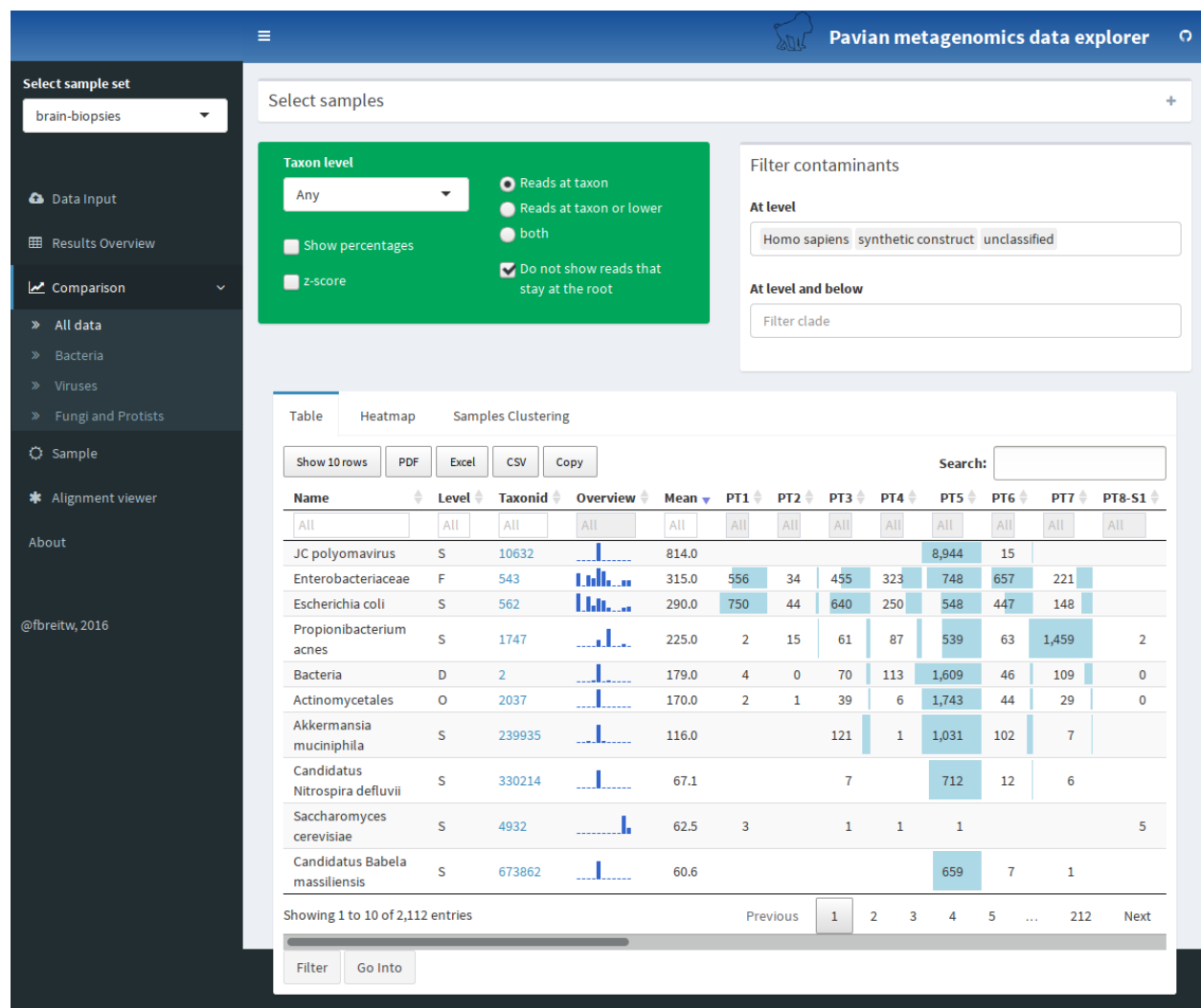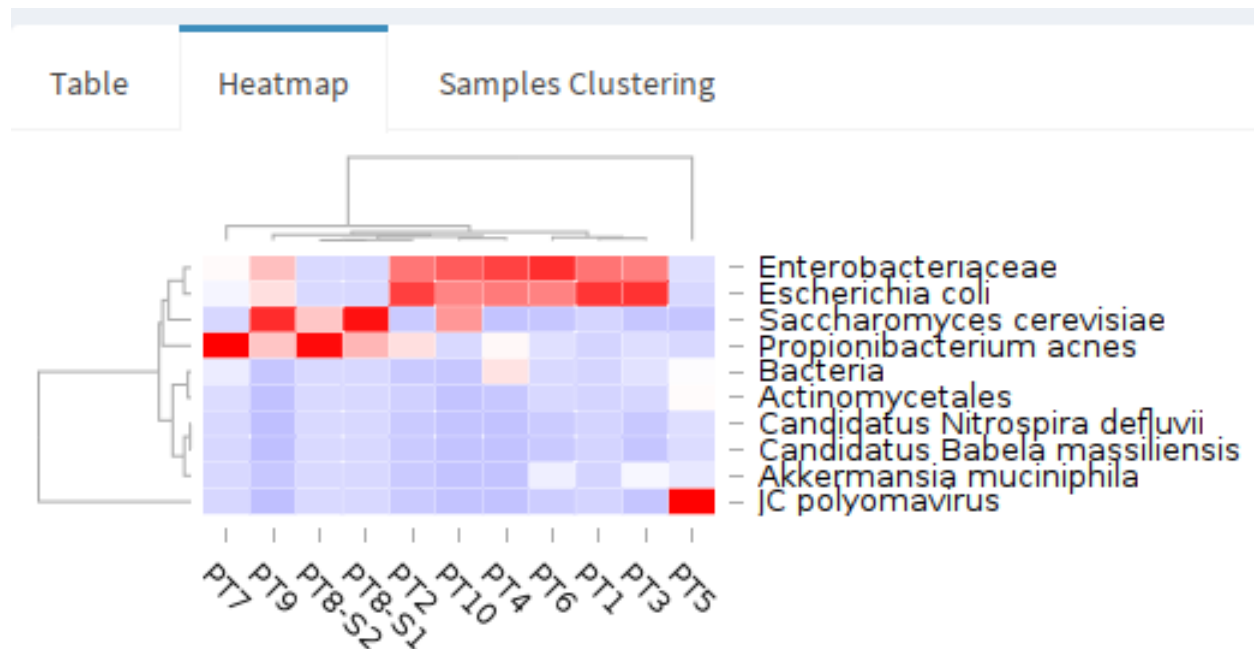
Figure 3: Sample comparision.

Figure 4: Sample comparision.

# Zoom into one sample - flow chart or sunburst diagram

# Zoom into one pathogen in one sample

`- view the alignment on the genome`

Salzberg, Steven L., Florian P. Breitwieser, Anupama Kumar, Haiping Hao, Peter Burger, Fausto J. Rodriguez, Michael Lim, et al. 2016. "Next-Generation Sequencing in Neuropathologic Diagnosis of Infections of the Nervous System." *Neurol Neuroimmunol Neuroinflamm* 3 (4). ter Science,; Biostatistics (S.L.S.), Johns Hopkins University, Baltimore, MD.: e251. doi:10.1212/NXI.0000000000000251.