# JGI Microbial Single Cell Program
## Single Cell Data Decontamination

**Author: Scott Clingenpeel**

Despite our best efforts, it is likely that there are some contigs in your single cell genome(s) that are from contaminant organisms.  Common contaminants that are known to be in the reagents we purchase are *Delftia*, *Pseudomonas*, and *Ralstonia*.  Other contaminants that we commonly see are *Propionibacterium* and *Lactobacillus*.  In addition, there may be contaminants from your particular sample in the form of free DNA that made it into the well along with your single cell.  Although we do an automated screen of your data for the known common contaminants and this information is provided to you in the JGI Single-cell Assembly QC report, this data is not removed because this could result in the removal of legitimate, highly conserved genes from your genome.  Thus, it is necessary for you to do a manual screening of your data to remove contaminant sequences.  While there are no clear rules on the identification and removal of contamination (i.e. phage or horizontal gene transfer may be difficult to discriminate from contamination), we would like to provide some recommendations and guidance.

**JGI warrants a 6-month review period for single cell genome(s) sequenced at the JGI, during which you can curate and clean the data set(s).  After 6 months, the data will be released to the public.**

When you first log into the IMG/MER system there are two ways to find your genomes.

Under the Find Genomes tab you can search for any genome in IMG by a variety of criteria.

Your own genomes should be under this link.



Click on the name of the genome that you want in order to bring up the genome overview page.

A good place to start with contamination screening is to look at any ribosomal RNA sequences in your genome.  When you pull up the overview of your genome and scroll down, you will see some genome statistics.



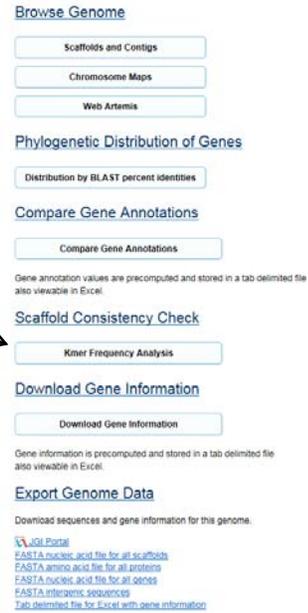Clicking here will bring up all your rRNA genes



If the length of the gene is too short then it will not be phylogenetically informative.

Clicking the links in this column will allow you to retrieve the sequence of the gene.

This column tells you which scaffold each gene is located on.

BLAST these rRNA sequences to see if they come from your target genome or from a contaminant.  In this example, rRNA genes from scaffolds 185, 302, and 408 match Rhizobium genes, which are the target organism.  Scaffolds 168, 271, and 356 contain rRNA genes that match Caulobacter and scaffolds 241 and 426 contain rRNA genes that match Betaproteobacteria.  These five scaffolds should be removed as contaminants.  Note the scaffold numbers that contain rRNA genes as this will be important in the next step.

Now go back to the genome overview and scroll down further than before.  Click on the button for Kmer Frequency Analysis.

This number is the minimum size of scaffold to include in the analysis.

This number indicates that a point will be plotted for every 500bp in the scaffold.

Begin by sticking with the defaults.  The larger scaffolds have more statistical power, which will produce a more defined cloud of points.  Also, it is easier to get a feel for the data with the few large scaffolds than if you included all of the data.  Later you will want to rerun this analysis with a smaller fragment window to include all your scaffolds in the screen.

Clicking the Generate button will produce a Kmer plot.  There is both a 2D and a 3D view.  We find the 3D view to be the most useful.

You can click and drag on the image to rotate it in three dimensions.  First, look at the percent of variation explained by each principal component.  If the percentages are all small (<~5%) then you have a very clean genome and the outliers are unlikely to be a problem.

Below is a fairly contaminated genome.  Most points are in a large mass which is our target genome, but there is a distinct cloud of contaminant scaffolds to the left of the main cloud.  By clicking on any of the points in the plot it will open a separate window of that scaffold.



Note the Orange scaffold.  This one starts in the main cloud, extends into the contaminant zone, and then returns to the main cloud.  Upon examining this scaffold, we find that the region that extends out from the main cloud contains rRNA genes that match the target organism.  Ribosomal RNA genes often contain a different GC content from the rest of the genome and thus will plot outside the main cloud of your target genome.  Scaffolds that extend from the main genome cloud can also contain other interesting features.

This red scaffold has points in the main cloud but extends well out.

Clicking on the points in this scaffold opens a separate window with more detail on the scaffold shown below.



This scaffold is from a Rhizobium single cell and it contains 30 predicted genes.  By doing a BLASTx search on each of the genes we found that the ends have high matches to proteins from various Rhizobium species.  However, the genes in the middle, which caused the scaffold to stick out from the main cloud in the plot, have best matches to phage proteins.  This cell appears to be infected with a lysogenic phage.

When a scaffold is suspicious based on the plot you need to verify if it is a contaminant by doing BLAST searches.  Left-clicking a point on the plot will open a separate window for the scaffold.  Right clicking a point will add that scaffold to your scaffold cart.  By holding shift while you click and drag you can select multiple scaffolds and then right click to add them to the scaffold cart.

From the scaffold cart select a scaffold.



Click here to go to the scaffold.

Click here to get a list of genes
on the scaffold.  From there
you can select each gene to get
its sequence for BLASTing.

Click here to get graphical
depiction of the scaffold
(see below).

The genes are color coded according to their COG category (scroll down for key).

Mouse over the genes to get a summary of their annotation.

Click on a gene to open the Gene Detail page (see below).

Click here to get the nucleotide sequence of the gene in fasta format.

Click here to get the amino acid sequence of the gene in fasta format.

If you scroll down to the bottom of the page there is a link to BLAST your gene's sequence.

Click here to do a
blastp on your gene.

If you don't want to BLAST the genes individually, then go the Scaffold Detail page and click to get a list
of genes from that scaffold.

Select the genes you want and
then add them to the gene cart.

From the Upload & Export & Save tab you can export all your genes as nucleotide or amino acid sequences in fasta format.

Because of limitations with the databases, you should BLAST multiple genes from each scaffold and take the consensus taxonomy for all of them before deciding whether a particular scaffold is contaminant or not.

Although we primarily rely on the Kmer plot to identify suspicious scaffolds, there are other tools that you may find useful.  First go back to the genome overview and scroll down

Click this button to get a phylogenetic distribution of your genes based on BLAST.

The Phylogenetic Distribution of Genes allows to assess the phylogenetic composition of a genome sample based on the distribution of best BLAST hits of protein-coding genes in the dataset. The phylogenetic distribution can be projected onto the families in a phylum/class (click on phylum/class name), and then further onto species in a family. For a reference genome within a species, the genome genes can be viewed using the Protein Recruitment Plot or the Reference Genome Context Viewer.

| Distribution of Best Blast Hists | COG Functional Category Statistics | COG Pathway Statistics |

### Distribution of Best Blast Hits (Gene count)

Domains(D): *=Microbiome,
B=Bacteria, A=Archaea, E=Eukarya, P=Plasmids, G=GFragment, V=Viruses.

> ● **hint:** *Hit genome count is in brackets ( ).*
> *Histogram is a count of best hits within the phylum / class*
> *at 30%, 60%, and 90% BLAST identities.*
> *Unassigned are the remainder of genes less than the percent identity cutoff, or that are not best hits at*
> *the cutoff, or have no hits.*

| Filter | | Select All | | Clear All | | Show All Phyla | | | | | | | |

| Select | D | Phylum/Class | No. Of Genomes | No. Of Hits 30% | % Hits 30% | Histogram 30% | No. Of Hits 60% | % Hits 60% | Histogram 60% | No. Of Hits 90% | % Hits 90% | Histogram 90% |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☑ | B | Acidobacteria | 7 (1) | 1 (1) | 0.08% | I | | | I | | | I |
| ☑ | B | Actinobacteria | 280 (1) | 1 (1) | 0.08% | I | | | I | | | I |
| ☑ | B | Cyanobacteria | 64 (2) | 2 (2) | 0.17% | I | | | I | | | I |
| ☑ | B | Bacilli | 568 (30) | 27 (15) | 2.23% | ▮ | 54 (19) | 4.47% | ▮ | 925 (17) | 76.51% | ▬▬▬▬▬ |
| ☑ | B | Clostridia | 231 (2) | 2 (2) | 0.17% | I | | | I | | | I |
| ☑ | B | Planctomycetes | 11 (2) | 2 (2) | 0.17% | I | | | I | | | I |
| ☑ | B | Alphaproteobacteria | 255 (17) | 13 (12) | 1.08% | ▮ | 8 (6) | 0.66% | I | 1 (1) | 0.08% | I |
| ☑ | B | Betaproteobacteria | 179 (2) | 4 (2) | 0.33% | I | | | I | | | I |
| ☑ | B | Deltaproteobacteria | 57 (1) | 1 (1) | 0.08% | I | | | I | | | I |
| ☑ | B | Gammaproteobacteria | 636 (8) | 6 (6) | 0.50% | I | 1 (1) | 0.08% | I | 1 (1) | 0.08% | I |
| ☑ | B | Spirochaetes | 60 (6) | 20 (6) | 1.65% | ▮ | 6 (1) | 0.50% | I | | | I |
| ☑ | B | Thermi | 19 (1) | 1 (1) | 0.08% | I | | | I | | | I |
| ☑ | B | Thermotogae | 14 (1) | 1 (1) | 0.08% | I | | | I | | | I |
| ☑ | P | Bacilli | 339 (8) | 2 (2) | 0.17% | I | 3 (3) | 0.25% | I | 3 (3) | 0.25% | I |
| ☑ | V | dsDNA viruses, no RNA stage | 902 (3) | 1 (1) | 0.08% | I | | | I | 3 (2) | 0.25% | I |
| ☑ | - | Unassigned | - | 120 | 9.93% | ▬▬ | 204 | 16.87% | ▬▬▬ | 276 | 22.83% | ▬▬ |

| Filter | | Select All | | Clear All | | Show All Phyla |

We mainly focus on hits that are 60% or 90% identity.  This particular genome is in the Bacilli and it is good to see that most of the hits are either Bacilli or Unassigned.  However, you can see that there are a few hits for Alpha- and Gammaproteobacteria, Spirochetes, and dsDNA viruses.  Clicking on the blue number in the "No. Of Hits" column will take you to a list of the genes that hit that phylum.

Another useful tool shows a histogram of the GC content of the scaffolds.  First, go back to the genome overview.  Near the top where you found the rRNA genes there is the number of scaffolds you have.



Clicking here will bring
up all your scaffolds



Select all your scaffolds
and then add them to
the Scaffold Cart.

Select all your scaffolds and then go to the Histogram tab.



Select GC Content from drop down menu then click Show Histogram.

For a clean genome, you should have a single peak and all scaffolds should be within ~10% to either side of the center. For the above genome I would investigate the four scaffolds with the lowest GC content and the 2 with the highest. Note that the histogram always gives you 10 bars/bins no matter how wide the spread in GC content is.

Once you have identified which scaffolds are contaminants, you need to remove them from your dataset. The IMG system does not have a direct way to remove data so the process involves selecting all the scaffolds that you want to keep and re-uploading them to IMG to replace the contaminated dataset.

One thing that can help with this process are scaffold sets. First, go to the Scaffold Cart.

Select the scaffolds you want in the set and then go to the Upload & Export & Save tab.

Name the set and then press the Save Selected to Workspace button.

You can examine the sets that you have created by going to Scaffold Sets under Workspace in the My IMG tab.



From here you can view the scaffolds in the sets that you have created and add them to the scaffold cart.

You can easily get the contents of a scaffold set by selecting the set and going to the Import & Export tab.



Here you can export the contents of your scaffold set as a fasta file that you save to your computer.

Once you have a file of just the clean scaffolds from your genome, you need to re-upload the data. This will overwrite the existing data so if you want to save the sequences from your contaminant scaffolds for future use you will have to save a fasta file of those scaffolds on your own computer.



Select Submit Data Set from the Companion Systems tab from anywhere in IMG/MER, or click the Data Submission Site button on the home page.

Select Metagenome Submissions.

Note the IMG Taxon OID number for your genome.

Scroll to the bottom of this page.



Click the Submit Genome Dataset to IMG/M ER button.

Find your genome's project.



Select your genome's project.





Enter the Taxon OID then click
on the Submit sequence file tab.

Enter your fasta file and change the gene caller to Isolate Genome Gene Calling.

Enter a Locus tag prefix then hit the submit button.

Congratulations!  Once this is loaded into IMG you will have a single cell genome to analyze that is free of contamination sequences.