## 7.1 Identification of specific DNA motifs.
   Note: For this exercise use http://microsporidiadb.org

a.  Find all BamHI restriction sites in all microsporidia genomic sequences available in MicrosporidiaDB.  Note: you can use the DNA motif search to find complex motifs like transcription factor binding sites using regular expressions.

   Hint:  BamHI = GGATCC and the DNA motif search is under the heading "Genomic Segments".



b.  How many times does the BamHI site occur in the genomes you searched? Take a look at your results; notice the Genomic location and the Motif columns.

## 7.2 Find genes that have one of these BamHI sites within 500 nucleotides upstream of their start.

In the section 7.1 you found BamHI sites, but now you are looking for genes that have one of these sites located within 500 nucleotides upstream of their start.

Hint: You can achieve this by running a genomic collocation search that defines the genomic relationship between the BamHI sites and genes. Add a "Genes by Organism" step to the motif search and select the "1 relative to 2, using genomic locations" option.

How did you modify the location relative to genes?
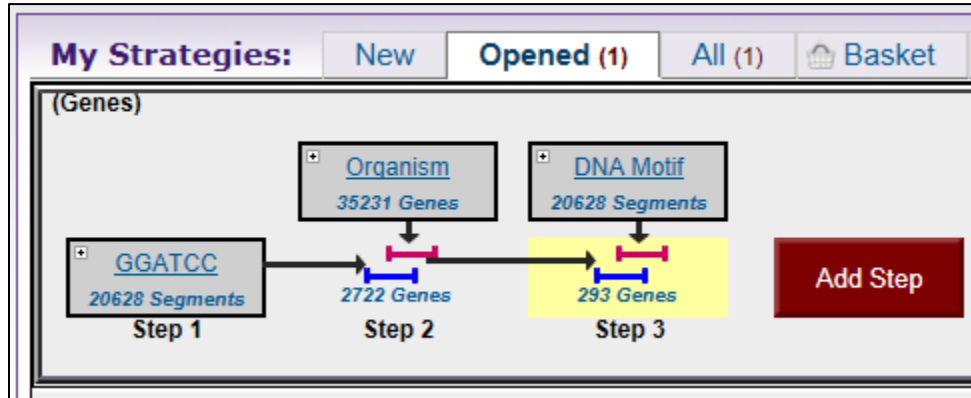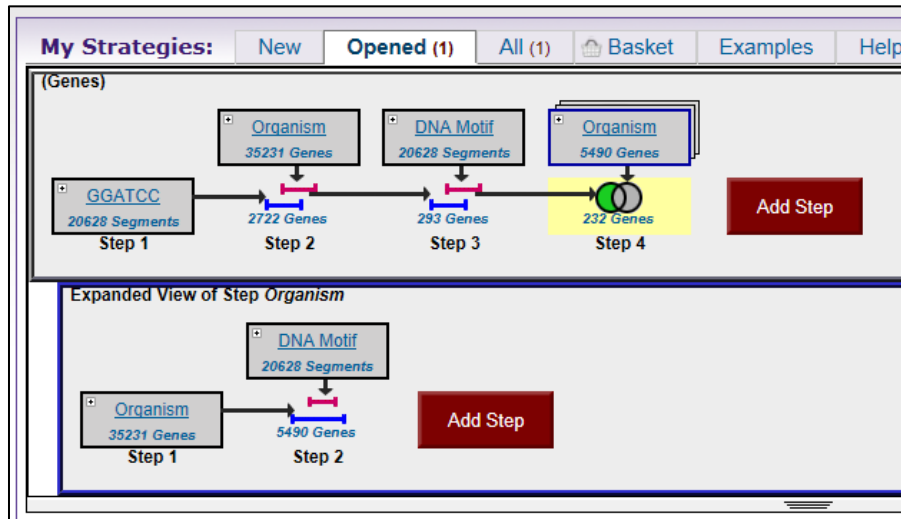


How many genes did you get?

**7.3 Using a similar sequence of steps as in part 7.2, define which of these genes also have a BamHI site in their 500 nucleotide downstream region.**

Hint: after you click on add step you will have to select DNA motif search and select the genomic collocation option.



**7.4** Taking this a step further, define which of these genes do **NOT** contain a BamHI site within them.

Hint: you will have to use a nested strategy.



Look at your results. Do they make sense? Confirm your results by looking at one of the genes in Gbrowse and showing BamHI restriction sites.

**Note:** you can add a column to any result table that allows you to go directly to GBrowse at the genomic coordinates of any ID in your result list. Click on the Add Columns button.

**Note:** you can configure restriction sites by clicking on the configure button in GBrowse and selecting the restriction sites you would like to display. To view restriction sites, the "Restriction Sites" data track must be turned on. Go to the "Select Tracks" page and click "Restriction Sites" under the "Analysis" section.