# DNA Motifs and Genomic Colocation

**Identify a specific DNA motif and colocate this motif with genes:**
**For this exercise use http://microsporidiadb.org**

a. **Find all *Bam*HI restriction sites in all microsporidia genomic sequences available in MicrosporidiaDB.** The DNA motif search can be used to finds simple DNA motifs or complex motifs like transcription factor binding sites using regular expressions. A DNA restriction site can be defined by a DNA motif. For example the *Bam*H1 restriction site is GGATCC.



- How many times does the *Bam*HI site occur in the genomes you searched? Take a look at your results. Notice the Genomic location and the Motif columns.

b. **Find genes that have one of these *Bam*HI sites within 500 nucleotides upstream of their start.**
In 4a you found *Bam*HI sites anywhere in the genome. 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 *Bam*HI sites and genes. Add a "Genes by Organism" step to the motif search and select the "1 relative to 2, using genomic locations" option.

**1**

DNA Motif
30994 Segments

Step 1

Add Step

**Add Step**

| Run a new Search for | Genes | Text, IDs, Organism | Text (product name, notes, etc.) |
|---|---|---|---|
| Add contents of Basket | Genomic Segments | Genomic Position | Gene ID(s) |
| Add existing Strategy | ORFs | Gene Attributes | Organism |
| Filter by Weight | | Protein Attributes | User Comments |
| | | Protein Features | |
| | | Similarity/Pattern | |
| | | Transcript Expression | |
| | | Cellular Location | |

**2**

**Add Step**

## Add Step 2 : Organism

Organism ❓  select all | clear all | expand all | collapse all | reset to default

**3**

- ☑ Anncaliia
- ☑ Edhazardia
- ☑ Encephalitozoon
- ☑ Enterocytozoon
- ☑ Mitosporidium
- ☑ Nematocida
- ☑ Nosema
- ☑ Ordospora
- ☑ Spraguea
- ☑ Trachipleistophora
- ☑ Vavraia
- ☑ Vittaforma

select all | clear all | expand all | collapse all | reset to default

⊞ **Advanced Parameters**

### Combine Genomic Segments in Step 1 with Genes in Step 2:

- ○ 1 Intersect 2
- ○ 1 Minus 2
- ○ 1 Union 2
- ○ 2 Minus 1
- ● 1 **Relative to** 2 , using genomic colocation

Continue

**4**

**Add Step**

## Genomic Colocation ❓📹

Combine Step 1 and Step 2 using relative locations in the genome

*You had **30994 Genomic Segments** in your Strategy (Step 1).    Your new **Genes** search (Step 2) returned **67093 Genes**.*

"Return each [Gene from Step 2 ▾] whose **upstream region** [overlaps ▾] the **exact region** of a Genomic Segment in Step 1 and is on [either strand ▾] "

*(67093 Genes in Step )*

Region

Gene

- ○ Exact
- ● Upstream: [500] bp
- ○ Downstream: [1000] bp
- ○ Custom:
  - begin at: [start ▾] [- ▾] [500] bp
  - end at: [start ▾] [- ▾] [1] bp

*(30994 Genomic Segments in Step )*

Region

Genomic Segment

- ● Exact
- ○ Upstream: [1000] bp
- ○ Downstream: [1000] bp
- ○ Custom:
  - begin at: [start ▾] [+ ▾] [0] bp
  - end at: [stop ▾] [+ ▾] [0] bp

Submit

Close

- How did you modify the location relative to genes? How many genes did you get?



**"Return each** Gene from Step 2 ▾ **whose upstream region** overlaps ▾

(67093 Genes in Step )

Region

Gene

○ Exact
◉ Upstream: 500 bp
○ Downstream: 1000 bp

○ Custom:
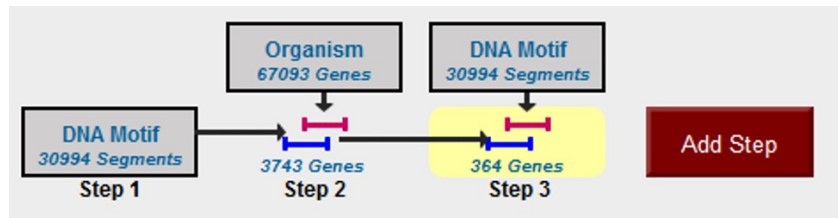  begin at: start ▾ - ▾ 500 bp
  end at: start ▾ - ▾ 1 bp

Organism
67093 Genes

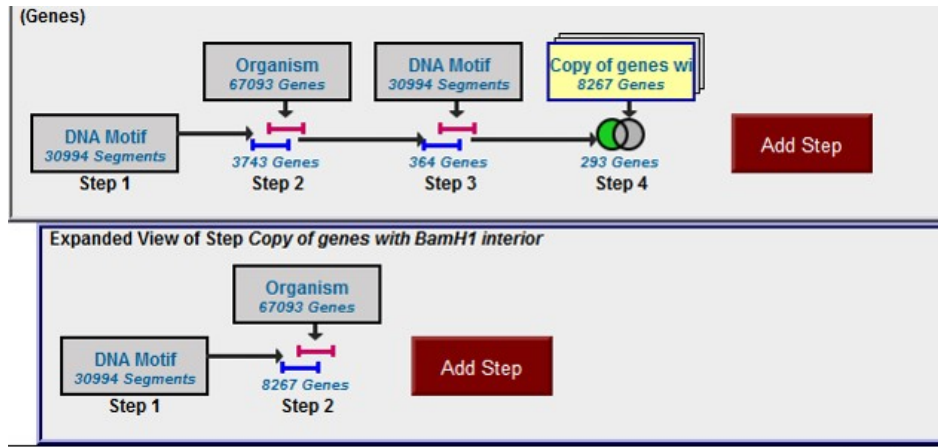DNA Motif
30994 Segments
Step 1

3743 Genes
Step 2

Add Step

c. **Using a similar sequence of steps as in 4b, define which of these genes also have a** *Bam***HI site in their 500 nucleotide downstream region.**

Hint: add a search for the BamH1 site and collocate that with the genes that have a BamH1 site upstream of their start sites.



Organism
67093 Genes

DNA Motif
30994 Segments

DNA Motif
30994 Segments
Step 1

3743 Genes
Step 2

364 Genes
Step 3

Add Step

**Step 1** = BamH1 anywhere in genome
**Step 2** = all genes
**2 (1+2)** = genes with BamH1 1000bp upstream
**Step 3** = BamH1 anywhere in genome
**3 (2+3)** = genes from 2 with BamH1 1000bp downstream

**d.** **Taking this a step further, define which of these genes do NOT contain a *Bam*HI site within them.** 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 *Bam*HI 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 display 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.