## **DNA and Protein Motif Pattern searches in FungiDB**

This tutorial uses https://beta.fungidb.org

## Example 1. Identify GATA motifs in *Aspergillus fumigatus* Af293 using DNA Motif Pattern Search

There are three major classes of eukaryotic zinc finger proteins. For example, Class II represents a group of proteins that includes transcription factors binding to GATA DNA motif. This motif occurs in the regulatory regions of the target genes and can be defined as [(A/T)GATA(A/G)] (A or T, followed by GATA, followed by A or G).

Using the VEuPathDB Regular Expression Tutorial, we can re-write the motif as follows: [AT]GATA[AG], where [] means "match any character contained in the brackets".

To deploy a search for GATA DNA motif in *Aspergillus fumigatus* Af293, follow these steps:

- Navigate to the Search for... and filter for motif,
- Select DNA Motif Pattern,
- Find *A. fumigatus* Af293 genome by filtering genomes for *fumig*,
- Enter [AT]GATA[AG] in the *Pattern* parameter and click *Get Answer*

S	earch for		Identify Ge • Organism	nomic Segments based on DNA Mot	f Pattern
C	motif	× 0	1 selected, fumig	out of 164	
	Genes		<ul> <li>Fungi</li> <li>Eurotiomycetes</li> <li>Aspergillus</li> </ul>		
	Sequence analysis		Seprijus fur Aspergilus Aspergilus Aspergilus Aspergilus Aspergilus	Inigiatus Turingiatus A1163 Turingiatus A2939 onovolumigiatus En Jacobi Carlo Managanus IBT 16806	
	Genomic Segments		aud tries private trie select all	l peret vir unen Jecor all	
	QDNA Motif Pattern				
i9465 G Strategy Genomic	enomic Segments from Step 1 Revise r. DNA Motif(9) Segment Results Genomic Locations				
	<ul><li>▲ 1 2 3 60 </li></ul>	Rows per page: 1000		Download Add to Basket	Add Columns
	Segment ID	LE Organism 🕄	🔷 Genomic Location 🝞 🔇	Motif 3	🔷 Length 😢
-	Chr1_A_fumigatus_Af293:1000359- 1000365:r	Aspergillus fumigatus Af293	Chr1_A_fumigatus_Af293:1,000,3591,000,365 (-)	CGCAGGGATAAACTACAGGG <mark>AGATAG</mark> TGCTACTTTCATGTCAGAGT	7
-	Chr1_A_fumigatus_Af293:1000521- 1000527:f	Aspergillus fumigatus Af293	Chr1_A_fumigatus_Af293:1,000,5211,000,527 (+)	TCAAGGTCACTTTACGTTGTTGATAGAAATAACCCCCAATTTTCTGT	7
-	Chr1_A_fumigatus_Af293:1001672- 1001678:r	Aspergillus fumigatus Af293	Chr1_A_fumigatus_Af293:1,001,6721,001,678 (-)	ACGATAGTCGCCGCTATCAGTGATAATGTCGCGGTTGCGAACTCTA	7
-	Chr1_A_fumigatus_Af293:1001679- 1001685:f	Aspergillus fumigatus Af293	Chr1_A_fumigatus_Af293:1,001,6791,001,685 (+)	CGCAACCGCGACATTATCACTGATAGCGGCGACTATCGTCGTGGAG	7
-	Chr1_A_fumigatus_Af293:1001866- 1001872:r	Aspergillus fumigatus Af293	Chr1_A_fumigatus_Af293:1,001,8661,001,872 (-)	GGTGGCAAGACACCAATATCTGATAGCACTGCTCGCCGTCCCGCTG	7

Further reading: https://mmbr.asm.org/content/70/3/583

## Example 2. Find fungal genes downstream of a regulatory DNA motif.

DNA motifs are important for the regulation of processes linked to host cell invasion, production of secondary metabolites, etc. The basic-helix-loop-helix (b-HLH) motif is involved in transcriptional regulation and cell-type determination. *S. cerevisiae* transcription factor PHO4 is one of the b-HLH motif proteins. It binds to the CACGTG DNA motif and positively regulates the acid phosphatase *PHO5*.

Identify all genes with upstream (600bp) CACGTG DNA motif in *S. cerevisiae* S288c and determine if PHO5 is one of those genes.

- Navigate to the *Search for* and filter for *motif*;
- Select DNA Motif Pattern, which is located in the Genomic Segments menu;
- Find S. cerevisiae genome by filtering genomes for cerev
- Enter CACGTG in the Pattern parameter and click Get Answer

parch for		Identi	fy Genomic Segments bas	ed on D	NA Motif Patte	ern		
		No	te: You may select up to 1 values for this parameter. 1 selected, out of 164					
motif	×	select all   cle cerev	ar unser presect only mean at all	0				
Genes	-	Ascon						
Sequence analysis		add these [ c] select al [ cir	Saccharomycetaceae     Saccharomycea cerevisiae     Saccharomycea cerevisiae     Saccharomycea cerevisiae S288c     art these   select only these     r all					
Genomic Segments		Ø Pattern	e Pattern					
Q DNA Motif Pattern		CACGTG						
My Search Strategi Opened (1) All (963) Public (37) Unnamed Search Strategy *	<b>es</b> Help							
				6	* 🖻 < 🖮   ×			
1,906 Segments								
1,906 Segments Revise this s	earch	-						
Step 1  1,906 Segments Revise this s  Genomic Segment Results Genomic Locations	earch	-						
1,906 Segments         P Add a step Step 1           1,906 Genomic Segments         Revise this s           Genomic Segment Results         Genomic Locations           4         1         2         Rows per	earch page: 1000 😨	3		Download Add to Be	asket Add Columns			
1,906 Segments Step 1 1,906 Genomic Segments Revise this s Genomic Segment Revise this s 4 1 2 b Rows per 5 Segment ID	earch page: 1000 👩	Genomic Location 🖉 3	Motif O	Download Add to Ba	sket Add Columns			
1,906 Segments     Produ a tategy       Brg 1     1,906 Genomic Segments     Revise this a       Cenomic Segment Results     Genomic Locations       Image:	earch page: 1000 📴 Lie Organism O Saccharomyces cerevisiae 5288c	Genomic Location	Molif  Molif  CATTAGCACTACCATGAATGCACGTGTCSCTGTCSTCATCA	Download Add to Ba	asket Add Columns			
1,906 Segments     P Add a step Step 1       1,906 Genomic Segments     Revise this is       Genomic Segment Results     Genomic Locations       I     2     Image: Revise this is       BX0069341958-1964.ft     BK0069341958-1964.ft	earch page: 1000 🝙 La Organism O Saccharomyces cerevisiae 5288c Saccharomyces cerevisiae 5288c	Cenomic Location  Control Cont	Motif O	Download Add to Ba CTGCT CTGATG	asket Add Columns			
1,906 Segments     P. Add a http://www.segments       Step 1     1,906 Genomic Segments       Revise this s     Genomic Segment Results       Genomic Segment Results     Genomic Lacations       Image: Segment ID     BK006934.1988-1964.1       BK006934.1988-1964.7     BK006934.197.1964.7       BK006934.107613-107619.f	earch pape: 1000 j corganism Saccharomyces cerevisiae 5286 Saccharomyces cerevisiae 5286 Saccharomyces cerevisiae 5286	© Genomic Location @ © BK006934.1,958.1,964.(+) BK006934.10,563.1,964.(-) BK006934.107,613.107,619.(+)	Motif CATTAGCACTACCATGAATGCACGTGTCGCTGTCCTCATCA ACCACTAGCACAGCGCACACCGTGCATTCATGGTAGTG ATTGCAAAGTAGTATTTTGTCACCGTGATTTGATCCAATTA	Download Add to Ba CTGCT CTAATG AAACT	asket Add Columns tength 7 7 7 7			

• Next, click on the + Add a step button



• Select Use Genomic Colocation... > A new Search > Genes

<del>&lt;</del>	Add a step to your search strategy 🛛 🥹	x
<b>Combine</b> with other Genomic Segments	Use the relative position of features on the genome between your existing step and the new step to identify features to keep in the final result.	
DNA Motif 7,806 Segments Step 1 Step 2	Choose which features to colocate. From  A new search An existing strategy My basket  expand all collapse all	7
Use Genomic Colocation to combine with other features	Filter the searches below T	

• Expand *Genes* options by clicking and select *Taxonomy,* and then *Organism*.



• Filter on *cerev* to select *S. cerevisiae* genome and click on the *Continue...* button

	Add	a step to your search strategy 🛛 🥹	
Organism			
	1 selected, out of 148	8	
add these   clear these   select on	ly these		
cerev		* 0	
Eunai			
Ascomycota			
Saccharomycetes			
Saccharomycetac	eae		
🗹 Saccharomyce	S		
Saccharom	yces cerevisiae S288c		
add these   clear these   select on select all   clear all	ly these		
		Continue	
		Give this search a name (optional)	

• Set up colocation tool to *Return each Gene from the new step* whose *upstream 600 bp region overlaps* the *exact* region of a Genomic Segment (which is the CACGTG motif) from the current step and is on either strand.

Add a step to your search strategy o								
"Return each Gene from the new step       whose upstream region         Region       Gene         Exact       Output         Upstream:       600         begin at:       start         custom:       begin at:         start       Constream:         begin at:       start         custom:       0         begin at:       0         begin at:       0         custom:       0         begin at:       0         custom:       0         c	overlaps       © the         exact region       of a Genomic Segment from the current step and is on         either strand       © "         Region       Genomic Segment         © Exact       Upstream: 1000 bp         Downstream:       1000 bp         Custom:       begin at:         start © + © 0 bp       bp         end at:       stop © + © 0 bp							
~	Run Step							

• Click Run Step

Organism     6,350 Genes       DNA Motif     + Add a step       1,906 Segments     49 Genes       Stan 1     Stan 2										
449 Genes (409 ortholog groups) Organism Filter select all   cear all   expand all   collapse all	Gene Res	Genor Rows per page	me View Ana e: 1000 📀	alyze Results		Down	load Add to Bask	et Add Colum	nns	
Hide zero counts           Search organisms         Q           ▶ □ Fungi         449           ▶ □ Oomycetes         Q	<u>_</u>		Transcript ID	12 Organism 😯 😣	Genomic Location (7 (Transcript)	Product Description I I III	⇔ Match Count     S     S	Region 8		
select all   clear all   expand all   collapse all	-	YHL038C	YHL038C- t26_1	Saccharomyces cerevisiae S288c	BK006934:2361725509(-)	Cbp2p	2	25510 - 26109 (	(-)	
Hide Organ	-	YHL036W	YHL036W- t26_1	Saccharomyces cerevisiae S288c	BK006934:2624127881(+)	Мир3р	2	25641 - 26240 (	(+)	

• To check if *PHO5* is one of the 449 genes returned, click on the **+** *Add a step* button, choose to *Combine with other Genes*, and *intersect* with *Text* search:

	Add a step to your search strategy 🛛 🧕
Combine with other Genes	Choose <i>how</i> to combine with other Genes
Organism 6,350 Ganes	( <u>()</u> 2 INTERSECT 3 <u>()</u> 2 UNION 3 <u>()</u> 2 MINUS 3 <u>()</u> 3 MINUS 2
449 Genee Step 2 Step 3	Ochoose which Genes to combine. From
Transform into related records	A new search     An existing strategy     My basket
Organism 6,350 Genes	(text x) 0
449 Genes Step 2 Step 3	Cext 
Use Genomic Colocation to	Find genes with a text search against their product name, notes, GO, EC, Domains, NRDB, or metabolic pathways.

• Search for a text term **PHO5** in **Gene name or symbol**:

	Add a step to your search strategy 💡
Cerev	<b>3</b> (x
E Fungi	
Ascomycota	
Saccharomycetes	
Saccharomycetales	
Saccharomycetaceae	
Saccharomyces	
Saccharomyces	cerevisiae
🕑 Saccharomy	ces cerevisiae S288c
add these   clear these   select only these	
select all prear all	
Text term (use * as wildc	ard)
Text term (use * as wildc	ard)
Text term (use * as wildc PHOS Fields	ard)
Text term (use * as wildc PH05 Fields	ard)
Text term (use * as wildc PH05 Fields C descriptions E descriptions and numbers	ard)
Text term (use * as wildc PH05 Fields Alternate product descriptions E0 descriptions and numbers Ephotopes from IEDB	ard)
Text term (use * as wildc PHOS Fields C descriptions EC descriptions and numbers Eptopes from IEDB External links	ard)
Text term (use * as wildc PH05 Fields C descriptions and numbers E descriptions and numbers Ext descriptions from IED8 External links External links	ard)
Text term (use * as wildc PH05 Fields Alternate product descriptions EC descriptions and numbers Eptiopes from IED8 External links Gene ID Gene ID Gene nor or symbol	ard)

		-										
pened (1) All (963) Publ	ic (37)	Help										
Unnamed Search Strategy *	1											
DNA Mott 1,905 Segments Step 1 Gene (1 ortholog groups)	Text 1 Gene 1 Gene Step 3	++	Add a ste	P							<	1
Organism Filter select all   clear all   expand all   collapse o Hide zero counts	dl		R	ows per page:	1000 ᅌ			Download	Add to Bas	iket	Add Col	umns
Search organisms	۵ 🔞			🌲 Gene ID	🜲 Transcript ID	Product Description ? S	🌲 Genomic Location (Gene) 💡 🕻	)	🜲 Gene N	ame or	Symbol	0
	1	←	-	YBR093C	YBR093C-t26_1	acid phosphatase PH05	BK006936:429,548430,951(-)		PH05			
<ul> <li>Fungi</li> <li>Oomycetes</li> </ul>	Comparetes     O     Select all [collapse all     Devanta l] collapse all     Eventa l] collapse											

Further reading: https://www.pnas.org/content/99/26/16893

## Example 3. Identify CaaX box proteins using Protein Motif Search

CaaX box proteins undergo several modifications essential to their targeting. These eukaryotic proteins contain a motif where cysteine is followed by the two aliphatic amino acids, followed by an amino acid that confers substrate specificity (e.g. proteins with C-terminal cysteine-aliphatic-aliphatic-M/S/Q/A/C/L/E amino acids, where M/S/Q/A/C/L/E confer substrate specificity).

Using the information below and the VEuPathDB Regular Expression help, define CaaX motif and identify corresponding gene in *S. cerevisiae*.

- 1. Aliphatic amino acids are non-polar and hydrophobic that are mostly found within protein molecules (e.g. isoleucine, leucine, valine), while alanine and glycine may be found either inside or outside a protein molecule.
- 2. Amino acids codes:

Cysteine – C	Isoleucine – I	Leucine – L	Valine – V
Alanine – A	Glycine – G	Serine – S	Methionine – M
Glutamine – Q	Glutamic acid - E		

**3.** Use \$ to match only at the end of the string or C-terminus

The C-terminal CaaX (cysteine-aliphatic-aliphatic-M/S/Q/A/C/L/E) motif can be rewritten as: C[ILVAG][ILVAG][MSQACLE]\$

C - Cysteine

[ILVAG] - matches any aliphatic amino acid (isoleucine, leucine, valine, alanine, glycine) [MSQACLE] - matches any amino acid that confers substrate specificity

\$ - match only at the end of the string or C-terminus

Deploy a search for CaaX protein motif in *S. cerevisiae* S288c:

- Navigate to the Search for ... and filter for motif;
- Select Protein Motif Pattern;
- Find *S. cerevisiae* genome by filtering genomes for *cerev*;
- Enter C[ILVAG][ILVAG][MSQACLE]\$ in the *Pattern* parameter and click *Get Answer*

Search for		Identify G	enes base	ed on Pro	otein Motif Pattern
motif ×	0	Pattern			
Genes		C[AGILV][AGILV][MSQA	CLE]\$		
Sequence analysis		Organism			
Q Protein Motif Pattern		1 selecte	ed, out of 148		
Genomic Segments		cere		×	
Q DNA Motif Pattern		Saccharomycetes Saccharomyce Saccharom Saccharom Saccharom	es yces cerevisiae romyces cerevisiae S288c		
		add these   clear select	these   select only these all   clear all		
My Strategies: New Opened (5) All (107	A Backet Public Strategies				
Hide search strategy panel	Basket Fublic Strategies				
Prot Mult 17 Genes Step 1					Rename Duplicate Sive As Share Delete
17 Genes from Step 1 Revise Strategy: Prot Motif(10) □ ▼ Click on a number in this table to limit/filter your result Gene Results Genome View Analyze Results					
Rows per page: 1000					Download Add to Basket Add Columns
Gene □ ID □ ID □ ID □ ID □ ID □ Organism	Genomic Location (Gene)	Product Description <b>2 8</b>	⊕ Match Locations     ⊗	Match 3	🗘 Motif 🕄
WBL061C YBL061C- t26_1 S. cerevisiae S288c	BK006936:105,316107,406(-)	Skt5p	(693-696)		SKKPTSLKNKKDKQGKKKKDCVIM
YDL009C YDL009C- t26_1 S. cerevisiae S288c	BK006938:432,925433,248(-)	hypothetical protein	(104-107) 1	I	AQPRIYHHSRLVMILKVSLECAVS
→ YDR461W YDR461W- t26_1 S. cerevisiae S288c	BK006938:1,385,1761,385,286(+	) mating pheromone a	(33-36) 1	1	MQPSTATAAPKEKTSSEKKDNYIIKGVFWDPAC)
→ YGL082W YGL082W- t26_1 S. cerevisiae S288c	BK006941:355,827356,972(+)	hypothetical protein	(378-381)	I	KFLPFNGSNKEKKRDKLKKNCVIM

Further reading: https://ec.asm.org/content/5/9/1560