# The COLLECTF external submission guide (v1.2)

COLLECTF, a play on words using the French *collectif* [collective] and the acronym for transcription factor [TF], is a database of prokaryotic transcription factor binding sites (TFBS). The main aim of COLLECTF is to provide high-quality, manually-curated information on the experimental evidence for transcription factor binding sites, and to map these onto reference bacterial genomes for ease of access and processing.

This document is intended to provide external submitters with a companion guide for the submission process. Additional information is available on the COLLECTF wiki. COLLECTF is accessible at http://collectf.umbc.edu.

## The COLLECTF submission process

### Data
COLLECTF uses data from only one type of source: published experimental evidence on transcription factor binding sites. COLLECTF distinguishes between two types of experimental support: evidence of binding (e.g. EMSA) and evidence of TF-mediated regulation (e.g. β-galactosidase assay). Identification of TF-binding sites through in silico means is recorded as part of the curation process, but not admitted as the single source of evidence for a TF-binding site. *Please do not submit data without some form of experimental evidence*.

### Step 0: Publication selection
The submission process starts with the submitter selecting a publication for curation. You must submit a publication for curation before starting a curation. You can submit several publications for curation.



### Step 1: Genome and TF information
Once a publication has been selected, the submitter must link the reported species (both for the sites and the transcription factor) to sequences present in the NCBI RefSeq database, by providing RefSeq accession number for a genome file (e.g. NC_005363.1; including version number)

and the protein corresponding to the TF (e.g. NP_970244; no version number). This is often a simple step, but can get more complex if the sequence for the exact strain used in your work is not available as an NCBI RefSeq record. Please try to identify a parental or related strain among those in NCBI RefSeq genomes. If there is no clear way to identify a surrogate genome in NCBI RefSeq, please use that of a common lab strain (most likely already in the database) for data submission.



In this screen you must first select the *TF* you are reporting on. If you are reporting that the TF acts in some multi-meric form, please indicate this too in the `TF structure` drop list. **COLLECTF** associates activation and repression information globally for each curation. This means that if your work describes both activated and repressed sites, you must submit two curations, selecting the appropriate `TF function` in each case and including, only the activated/repressed sites in each curation. If a site is both repressed and activated, it should be present in both curations. While this may sound tedious, the submission system will pre-populate all fields in your second (or third…) submission, facilitating the process enormously.

If the work you are reporting uses a strain different from the selected RefSeq genome/TF, please type/paste the original strain in the `Organism of origin...` and `Organism TF binding sites...` text fields. This allows us to keep track of the correspondence between reported and mapped strains.

## Step 2: Experimental methods

**COLLECTF** is all about gathering and validating information on experimentally-validated TF-binding sites. The next two steps in the curation process specifically target these two fundamental points. Step 2 requires that you report _all the techniques used in the paper to verify the TFBS_ that are being _reported in this submission_. In this step we also ask that you provide a _brief written summary_ of the process used to verify the submitted TFBS (not the overall experimental process, but just how the selected experimental techniques were combined to define reported TFBS)[1]. You can also indicate whether the paper contains promoter information (e.g. location of transcriptional start) or expression data (evidence of TF-mediated regulation), and provide external database accession numbers for expression data (e.g. GEO accession numbers), or details on whether the TF forms complex with other molecules in order to bind.



---

[1] For instance: _"Sites were first identified using a computer search, then binding was validated with EMSA. TF-mediated expression was confirmed with β-gal assays on w-t vs. tf- mutant"._ See curations in the database for examples.

## Step 3: Reported sites

In this step, you will enter the primary information for **COLLECTF**: binding sites reported in your work *using the techniques specified in Step 2*. Notice that if your work describes different groups of sites using different sets of techniques, *you should accordingly create multiple independent submissions*.

*Motif-associated*

TF-binding sites can be defined at two different levels. By definition, a TF-binding site is simply a (relatively short) stretch of DNA to which a transcription factor is shown to bind. Many TFs target known specific sequence patterns in the DNA that lead to aligned collections of sites (motifs), providing a much more concise definition of TF-binding site. In **COLLECTF** we refer to the latter as motif-associated and the former as non-motif associated. If you are confident that the sites you report conform to a known motif or you demonstrate that they do through experimental work (e.g. site-directed mutagenesis), you should check the `Reported sites are motif associated` checkbox.

*Coordinates and quantitative data*

Sites can be entered as sequences or using genome coordinates (if you already mapped them to the reference strain in your work), and the appropriate box should be checked. Sites should be entered one per line (FASTA format is also accepted for sequence entry). If you report quantitative data for sites (e.g. peak intensities, estimated $K_d$), please check the `Sites with quantitative data` checkbox. Quantitative data can then be appended with a tab/space after the sequence/coordinate entry. A brief description of its nature (method used and range of quantitative data) should be entered in the `Quantitative data format` textbox.

*ChIP data*

**COLLECTF** also captures ChIP peak data. When the `This paper reports ChIP data` box is checked, the submission process will request additional data.



The `Assay conditions` field should report the specific biological conditions on which the ChIP assay was carried out (e.g. iron deprivation). The `ChIP method notes` field should detail specific ChIP protocols (cross-linking method, sequencing, etc.) as reported in the *Materials and Methods* section of the manuscript. The `Supporting ChIP quantitative data` field accepts peak data with quantitative peak intensities and allows you to assign this quantitative data to reported motif associated sites. If this field is populated, the **COLLECTF** submission system will scan peak data for instances of reported motif associated sites (in the `Sites` field) and automatically associate the peak value to the motif associated site.
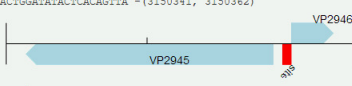
## Step 4: Verify sites (exact)

After you enter the sites, the **COLLECTF** submission system will download the specified genome sequence and search for entered sites. The sites are reported back to the curator specifying their location in the sequence and nearby genes. Gene annotation details can be accessed by hovering over any gene locus. You can use this information to verify that the sites identified by the **COLLECTF** submission system in the NCBI RefSeq genome sequence correspond to the sites you report in the paper.
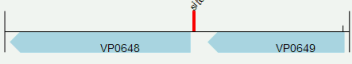


## Step 6: Verify sites (inexact)

In some cases, especially if using a sequence that is not an exact match to the reported strain, some sites may not be found using an exact search. In this case, the **COLLECTF** submission system will use the available evidence to construct a scoring matrix and search the genome for slightly inexact matches (up to two mismatches away from the reported site). These will be reported in the same way as exact matches and you will be asked to validate them in the same manner.

## Step 6: Verify quantitative data

If quantitative data has been entered for sites, the COLLECTF system will also display it along with sites in their genomic context for verification and editing.



## Step 7: Gene regulation

If the manuscript reports experimental evidence for TF-mediated regulation of target genes through TFBS, the COLLECTF submission system will ask you to specify, for each reported site, which genes have been shown to be regulated by the TF.



## Step 8: Curation information

The submission process ends with a final assessment of the curation. You will be asked whether the submission requires review (`Revision required`). Checking this option is indicated in several circumstances. For instance, it is quite possible that no appropriate sequence has been located in NCBI to perform a valid curation. In this case, the curation is marked for revision. The TFBS data is stored, but it will not be linked to a RefSeq sequence until a matching RefSeq record is posted.

You will also be asked whether the curation is considered valid for submission to NCBI. Curations will only be *considered for submission to NCBI if the sequence for the exact reported strain is available at NCBI or if a sequence matching the species of the reported strain is available and at least 90% of the sites you report have been located in the reference RefSeq record as exact matches*.

The system also requires that you specify if the `Curation for this paper is complete`. Do not check this box if, for instance, you must report *additional sites with different support* that need to be reported in a subsequent curation, you must report *additional sites in a different chromosome or for a different TF* or if you need to report *additional sites under a different mode of regulation* (e.g. repressed instead of activated). All these situations require multiple curations. The COLLECTF submission system will pre-populate fields to facilitate this process.

Finally, after you check `I want to submit this curation` and click `Next`, the system will ask that you verify the submission (a new window summarizing the submission will pop up for inspection).



Once a submission is completed, the data is uploaded to COLLECTF. The submission will be then reviewed by a COLLECTF curator and, if approved, tagged for submission to NCBI.