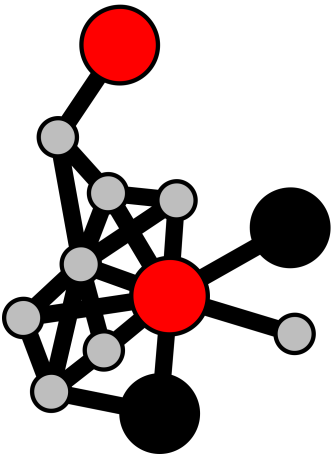
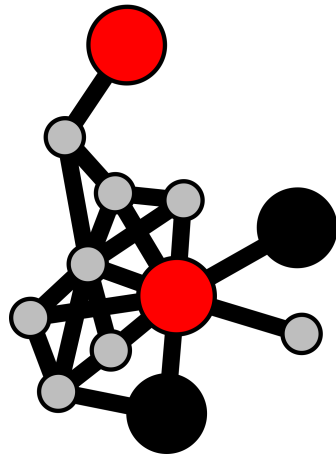
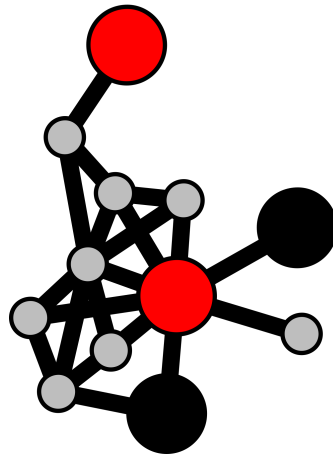
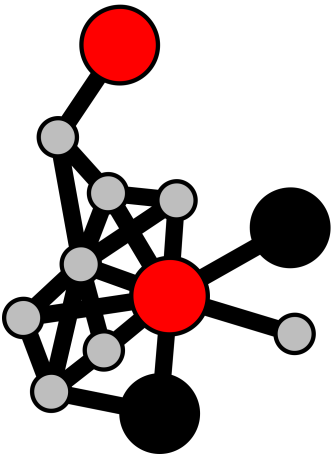
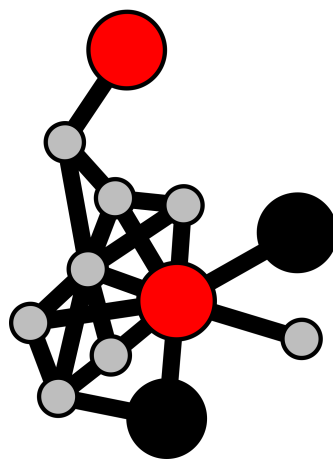
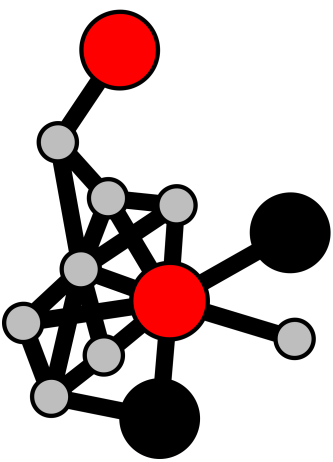


Functional annotation lesson



Targeted HMM-based annotation:
FeGenie introduction and tutorial





Artwork: Nancy Merino

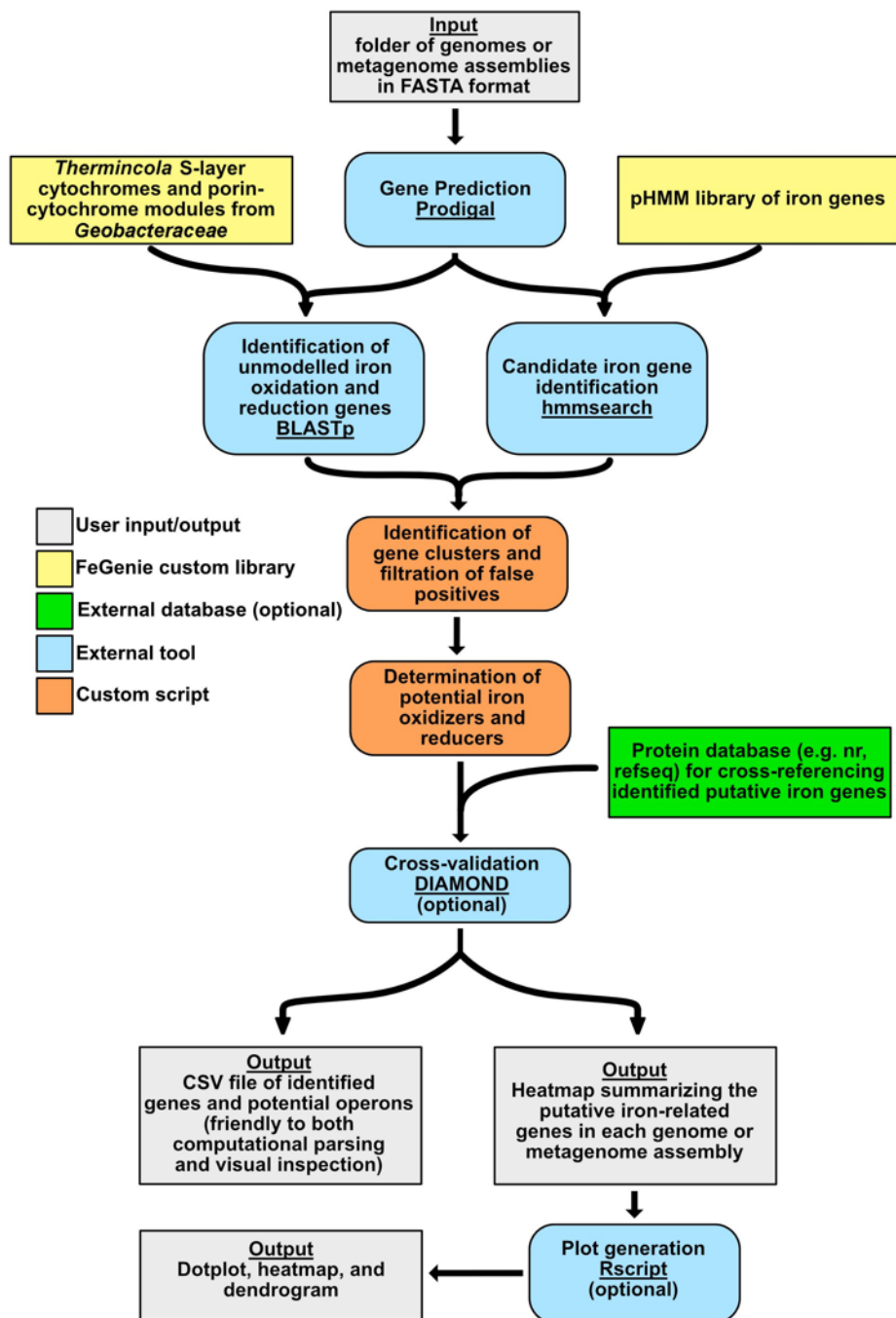
What is FeGenie

1. A collection of HMMs, based on a comprehensive set of genetic markers related to iron acquisition/scavenging, transport, efflux, storage, as well as iron redox (dissimilatory reduction and oxidation)
2. A bioinformatics software that uses these HMMs to profile genomic datasets (in FASTA format)
3. Uses rules for filtering out potential false positives (operon structure, genetic co-occurrence, bit scores)

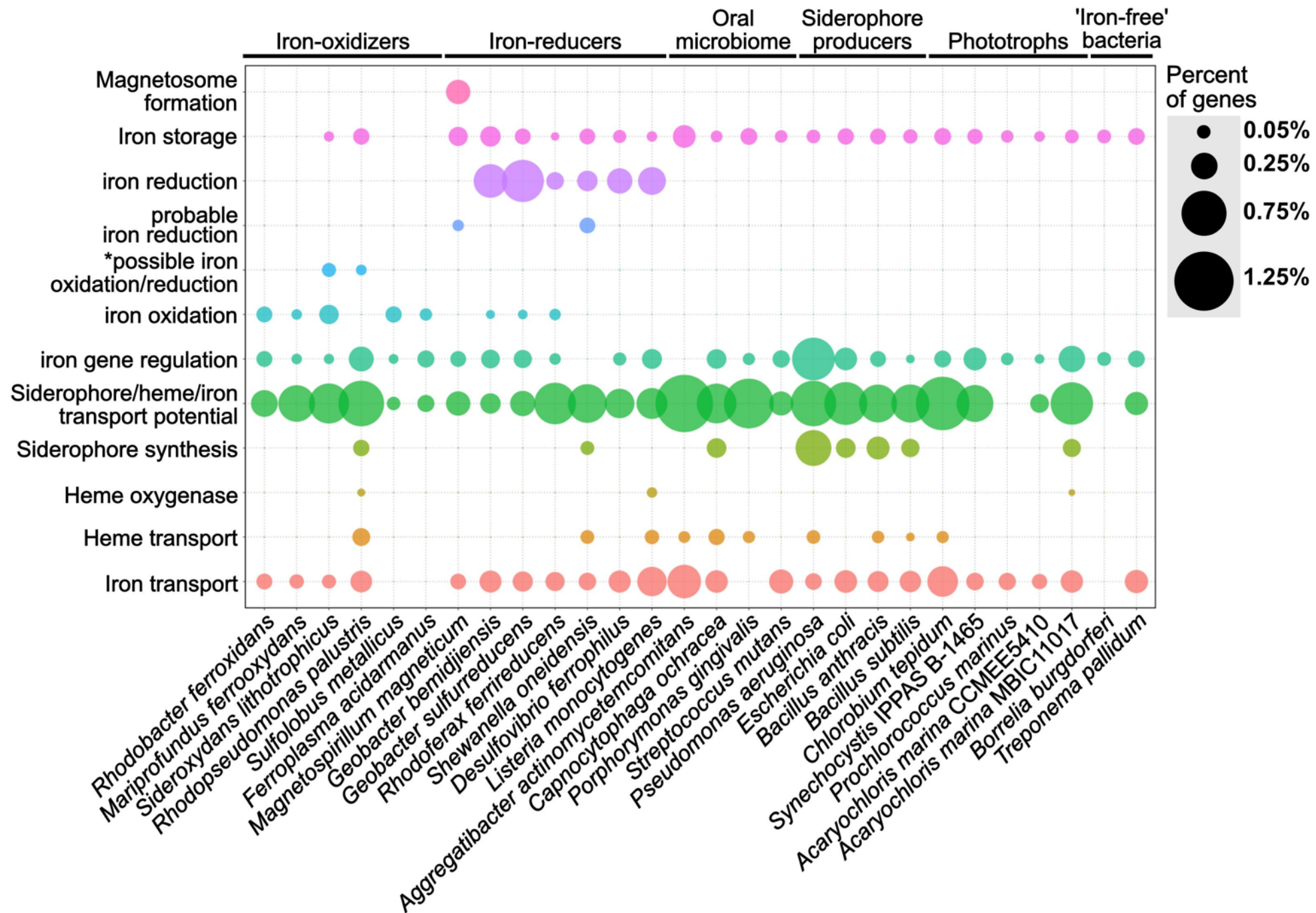
Category	Function	Protein Families
Iron acquisition	Iron(II)/(III) transport	Efe <u>UOB</u> ¹ , FbpABC ² , SfuABC ³ , YfuABC ⁴ , FeoAB(C) ⁵ , FutA ¹⁶ , FutA2 ⁶ , FutB ⁶ , FutC ⁶ , YfeABCD ⁷
	Heme oxygenase	ChuS ⁸ , ChuZ ⁹ , HemO ^{10,11} , PigA ^{10,11} , Hem <u>RSTUV</u> ¹² , HmoB ¹³ , HmuO ¹⁴ , HugZ ¹⁵ , HupZ ¹⁶ , Isd-LmHde ¹⁷ , IsdG ¹⁸ , IsdI ¹⁹ , MhuD ²⁰ , PhuS ²¹ (in PhuRSTUVW)
	Heme transport	Has <u>RADE</u> (B)F ²² , HmuR <u>STUV</u> ²² , HmuY ²³ , HmuY ²³ , HutZ ²⁴ , Hxu <u>CBA</u> ²⁵ , IsdX1 ²⁶ , IsdX2 ²⁶ , Phu <u>RSTUVW</u> ²¹ , Rv0203 ²⁷
	Transferrin/Lactoferrin	Tbp <u>AB</u> (Lbp <u>AB</u>) ²⁸ , Sst <u>ABCD</u> ²⁹
	Siderophore synthesis	Acs <u>ABCDE</u> F ³⁰ , AmoA ³¹ , AngR ³² , Asb <u>ABCDE</u> F ³³ , DhbACEBF ³⁴ , entD- fepA -fes-entF-fepEC <u>GDB</u> -entCEBA-ybdA ³⁵ , IroD in Iro <u>NBCDE</u> ³⁶ , lucABCD ^{37,38} , lutA ^{37,38} , MbtIJABCDEFGHI ³⁵ , LbtA ³⁹ (in LbtUABC), PchABCDEFGHI ³⁵ , PvdQAPMNOFEDJIHLGS ⁴⁰ , PvsABCDE ⁴¹ , VenB ⁴² , Vab genes in VabR-fur-vabGA-fur-VabCEBSFH-fur-ftvA-vabD ⁴³ , Vib genes in VibB-vibEC-vibA-vibH-viuP <u>DGC</u> -vibD and viuAB-vibF ^{44–46} , RhbABCDEF-rhrA-rhtA ⁴⁷
Iron Gene regulation	Transcriptional regulation	BesA ⁴⁸ , Cbr <u>BCD</u> ⁴⁹ , TonB-ExbB-ExbD ⁵⁰ , Fat <u>ABCD</u> ⁵¹ , FecI <u>RABCDE</u> ⁵² , FeuABC-yusV ⁵³ , Fhu <u>ACDB</u> ^{54–56} , FhuF ^{54–56} , Fpt <u>ABCX</u> ⁵⁷ , Fpu <u>AB</u> ⁵⁸ , FpuC ⁵⁸ , FpuD ⁵⁸ , Fpv <u>IR-FpvA</u> -FpvGHJKDEF ⁵⁹ , FvtA in VabR-fur-vabGA-fur-VabCEBSFH-fur-ftvA-vabD ⁴³ , HatCD <u>B</u> ³⁷ , IroNBCDE ³⁶ , LbtUABC ³⁹ , PirA ⁶⁰ , PiuA ⁶⁰ , Pvu <u>ABCDE</u> ⁴¹ , Viu genes in VibB-vibEC-vibA-vibH-viuP <u>DGC</u> -vibD and viuAB-vibF ^{44–46} , YfiZ -yfhA ⁶¹ , YfiY ⁶¹ , YqjH ⁶² , ybdA and Fep genes in entD- fepA -fes-entF-fepEC <u>GDB</u> -entCEBA-ybdA ³⁵
		DtxR ⁶³ , FecR (in FecIABCDE) ⁵² , FeoC in FeoAB(C) ⁵ , Fur ⁶⁴ , IdeR ⁶⁵ , Yqji ⁶² , RhrA in RhbABCDEF-rhrA-rhtA ⁴⁷
Iron oxidation and reduction	Iron oxidation	Cyc1 ^{66,67} , Cyc2 ^{66,67,68} , FoxABC ⁶⁹ , FoxEYZ ⁷⁰ , Sulfocyanin ⁷¹ , Pio <u>ABC</u> ⁷²
	Probable iron oxidation and possible iron reduction	Mto <u>AB</u> ⁷³ , Cyc2 (cluster 3)
	Dissimilatory iron reduction	CymA ⁷⁴ , Mtr <u>CAB</u> ⁷⁵ , OmcF ⁷⁶ , OmcS ⁷⁶ , OmcZ ⁷⁶ , FmnA-dmkA-fmnB-pplA-ndh2-eetAB-dmkB ⁷⁷ , DFE_0448-0451, DFE_0461-0465 ⁷⁸
	Probable iron reduction	MtrCB, MtrAB, MtoAB-MtrC
Iron storage	Iron storage	Bfr ⁷⁹ , DpsA ⁸⁰ , Ftn ⁸¹
Magnetosome-related	Magnetosome formation	MamABEKLMOPIQ ^{82,83} (Note: These genes are found in all known magnetotactic microorganisms, except for <i>mamL</i> which is found in magnetite-producing magnetotactic microorganisms ⁸¹)

*Bolded and underlined HMMs are derived from Pfam or TIGRFAMs databases. Other HMMs were created by using select sequences. See **Supplementary Table S1** for more information, including the corresponding Pfam or TIGRFAMs families and the sequences used to create the HMMs.* ¹Miethke et al., 2013, ²Adhikari et al., 1996, ³Angerer et al., 1990, ⁴Gong et al., 2001, ⁵Lau et al., 2016, ⁶Katoh et al., 2001, ⁷Bearden et al., 1998, ⁸Suits et al., 2006, ⁹Zhang et al., 2011, ¹⁰Friedman et al., 2003, ¹¹Friedman et al., 2004, ¹²Schneider et al., 2006, ¹³Park et al., 2012, ¹⁴Matsui et al., 2005, ¹⁵Hu et al., 2011, ¹⁶Sachla et al., 2016, ¹⁷Duong et al., 2014, ¹⁸Reniere et al., 2010, ¹⁹Skaar et al., 2004, ²⁰Graves et al., 2014, ²¹Ochsner et al., 2000, ²²Tong and Guo, 2009, ²³Wójtowicz et al., 2009, ²⁴Liu X. et al., 2012, ²⁵Morton et al., 2007, ²⁶Honsa et al., 2014, ²⁷Tullius et al., 2011, ²⁸Gray-Owen et al., 1995, ²⁹Morrissey et al., 2000, ³⁰Carroll and Moore, 2018, ³¹Barghouthi et al., 1991, ³²Wertheimer et al., 1999, ³³Oves-Costales et al., 2007, ³⁴May et al., 2001, ³⁵Crosa and Walsh, 2002, ³⁶Hantke et al., 2003, ³⁷Suzuki et al., 2006, ³⁸Martinez et al., 1994, ³⁹Cianciotto, 2015, ⁴⁰Lamont and Martin, 2003, ⁴¹Tanabe et al., 2003, ⁴²Tan et al., 2014, ⁴³Balado et al., 2008, ⁴⁴Wyckoff et al., 2001, ⁴⁵Keating et al., 2000, ⁴⁶Wyckoff et al., 1999, ⁴⁷Lynch et al., 2001, ⁴⁸Miethke et al., 2006, ⁴⁹Mahé et al., 1995, ⁵⁰Garcia-Herrero et al., 2007, ⁵¹Lemos et al., 2010, ⁵²Braun, 2003, ⁵³Peuckert et al., 2011, ⁵⁴Köster and Braun, 1989, ⁵⁵Coulton et al., 1987, ⁵⁶Braun et al., 2002, ⁵⁷Youard et al., 2011, ⁵⁸Dixon et al., 2012, ⁵⁹Brillet et al., 2012, ⁶⁰Moynie et al., 2017, ⁶¹Ollinger et al., 2006, ⁶²Wang et al., 2011, ⁶³Guedon and Helmann, 2003, ⁶⁴Escolar et al., 1998, ⁶⁵Rodriguez et al., 2002, ⁶⁶Castelle et al., 2008, ⁶⁷Barco et al., 2015, ⁶⁸McAllister et al., 2019, ⁶⁹Bathe and Norris, 2007, ⁷⁰Croal et al., 2007, ⁷¹Ilbert and Bonnefoy, 2013, ⁷²Liu J. et al., 2012, ⁷³Jiao and Newman, 2007, ⁷⁴Castelle et al., 2015, ⁷⁵Pitts et al., 2003, ⁷⁶Santos et al., 2015, ⁷⁷Light et al., 2018, ⁷⁸Deng et al., 2018, ⁷⁹Grossman et al., 1992, ⁸⁰Grant et al., 1998, ⁸¹Andrews, 1998, ⁸²Uebe and Schuler, 2016, ⁸³Kolinko et al., 2016.

Bolded and underlined genes represent HMMs taken from Pfam and TIGRFAMs



General workflow involves a combination of custom scripts and other bioinformatics software for ORF-prediction and orthology searching



Tutorial

🔗 Dependencies

- Python (version 3.6 or higher)
- Diamond (version 0.9.22.123) (only necessary if you are doing the cross-validation against a reference database)
- BLAST (version 2.7.1+)
- HMMER (version 3.2.1)
- Prodigal (version 2.6.3)
- R (version 3.5.1)
- Rscript

Absolutely required dependencies

Absolutely required dependencies

Installation

```
git clone https://github.com/Arkadiy-Garber/FeGenie.git
cd FeGenie
bash setup_noconda.sh
```

Put \$PATH to FeGenie into your bash profile (.bash_profile or .profile in your home directory):

PATH="/Users/arkadiygarber/topic-functional-annotation/FeGenie-tutorial/FeGenie:\$PATH"

Easy Installation (if you have Conda installed)

```
git clone https://github.com/Arkadiy-Garber/FeGenie.git
cd FeGenie
bash setup.sh
conda activate fegenie
FeGenie.py -h
```


Usage

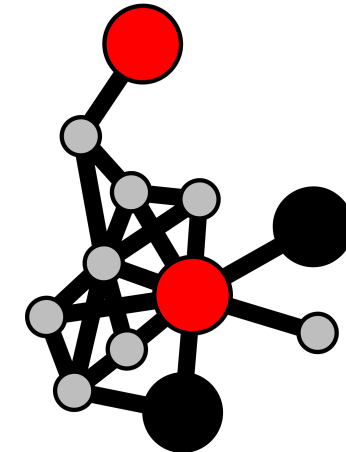
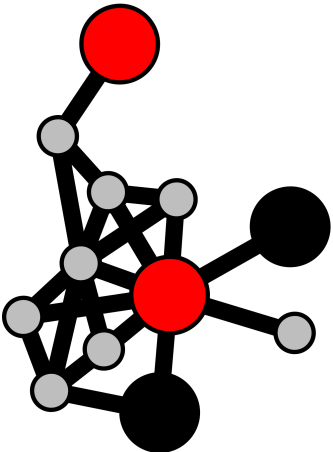
```
-h, --help          show this help message and exit
-bin_dir BIN_DIR    directory of bins
-bin_ext BIN_EXT     extension for bins (do not include the period)
-d D                maximum distance between genes to be considered in a
                    genomic 'cluster'. This number should be an integer and
                    should reflect the maximum number of genes in between
                    putative iron-related genes identified by the HMM
                    database (default=5)

-ref REF            path to a reference protein database, which must be in
                    FASTA format
-out OUT            name output directory (default=fegenie_out)
-inflation INFLATION inflation factor for final gene category counts
                    (default=1000)
-t T                number of threads to use for DIAMOND BLAST and
                    HMMSEARCH (default=1, max=16)
-bams BAMS          a tab-delimited file with two columns: first column
                    has the genome or metagenome file names; second column
                    has the corresponding BAM file (provide full path to
                    the BAM file). Use this option if you have genomes
                    that each have different BAM files associated with
                    them. If you have a set of bins from a single
                    metagenome sample and, thus, have only one BAM file,
                    then use the '-bam' option. BAM files are only
                    required if you would like to create a heatmap that
                    summarizes the abundance of a certain gene that is
                    based on read coverage, rather than gene counts.
-bam BAM            BAM file. This option is only required if you would
                    like to create a heatmap that summarizes the abundance
                    of a certain gene that is based on read coverage,
                    rather than gene counts. If you have more than one BAM
                    file corresponding to different genomes that you are
                    providing, please use the '-bams' argument to provide
                    a tab-delimited file that denotes which BAM file (or
                    files) belongs with which genome

--gbk [GBK]         include this flag if your bins are in Genbank format
--orfs [ORFS]       include this flag if you are providing bins as open-
                    reading frames or genes in FASTA amino-acid format
--meta [META]       include this flag if the provided contigs are from
                    metagenomic/metatranscriptomic assemblies
--norm [NORM]       include this flag if you would like the gene counts
                    for each iron gene category to be normalized to the
                    number of predicted ORFs in each genome or metagenome.
                    Without normalization, FeGenie will create a heatmap-
                    compatible CSV output with raw gene counts. With
                    normalization, FeGenie will create a heatmap-
                    compatible with 'normalized gene abundances'

--makeplots [MAKEPLOTS] include this flag if you would like FeGenie to make
                    some figures from your data?. To take advantage of
                    this part of the pipeline, you will need to have
                    Rscript installed. It is a way for R to be called
                    directly from the command line. Please be sure to
                    install all the required R packages as instructed in
                    the FeGenie Wiki: https://github.com/Arkadiy-
                    Garber/FeGenie/wiki/Installation. If you see error or
                    warning messages associated with Rscript, you can
                    still expect to see the main output (CSV files) from
                    FeGenie.
```

Customize your iron-gene annotation pipeline



Jupyter Binder tutorial

