

---

# MetaPathways

**BCB2**

**Sep 23, 2023**

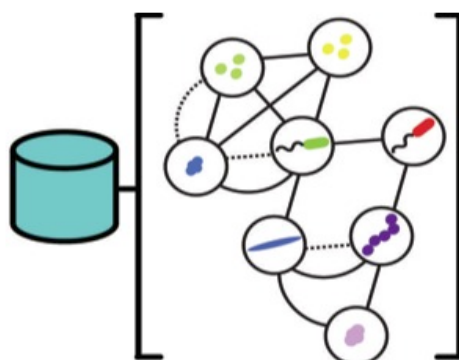


# DOCUMENTATION

<b>1</b>	<b>Overview</b>	<b>1</b>
1.1	Pipeline Overview . . . . .	1
1.2	Output Format . . . . .	3
1.3	Visualizing Output . . . . .	3
<b>2</b>	<b>Installation</b>	<b>5</b>
2.1	Container Install . . . . .	5
2.2	Installing with Pip and Conda . . . . .	5
2.3	Reference Sequences . . . . .	7
<b>3</b>	<b>Running MetaPathways</b>	<b>9</b>
3.1	Input . . . . .	9
3.2	Parameter File . . . . .	9
3.3	Run . . . . .	9
<b>4</b>	<b>Phandi GUI Overview</b>	<b>11</b>
<b>5</b>	<b>Indices and tables</b>	<b>13</b>
<b>6</b>	<b>Contact</b>	<b>15</b>
	<b>Bibliography</b>	<b>17</b>



## OVERVIEW



MetaPathways [CIT2002] is a meta'omic analysis pipeline for the annotation and analysis for environmental sequence information. MetaPathways include metagenomic or metatranscriptomic sequence data in one of several file formats (.fasta, .gff, or .gbk). The pipeline consists of five operational stages including

### 1.1 Pipeline Overview

MetaPathways is composed of five general stages, encompassing a number of analytical or data handling steps (**Figure 1**):

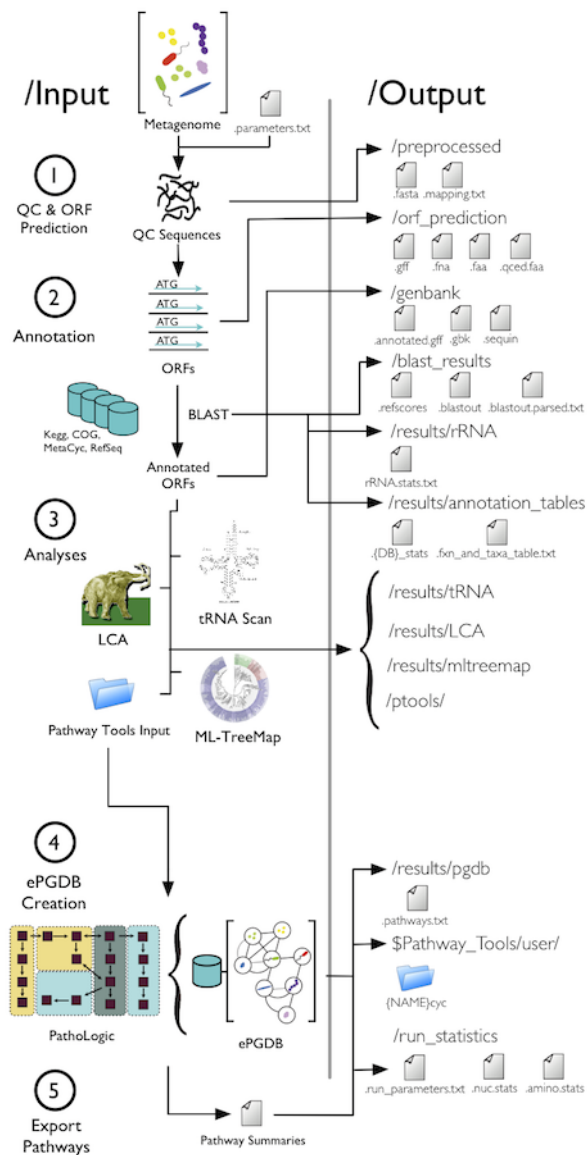
1. **QC and ORF Prediction:** Here MetaPathways performs basic quality control (QC) including removing duplicate sequences and sequence trimming. Open Reading Frame (ORF) prediction is then performed on the QC'ed sequences using Prodigal [PRODIGAL2010] or GeneMark [GeneMark12]. The final translated ORFs are now also trimmed according to a user-defined setting.
  - MetaPathways steps: *PREPROCESS INPUT*, *ORF PREDICTION*, and *FILTER AMINOS*
2. **Functional and Taxonomic Annotation:** Using seed-and-extend homology search algorithms (B)LAST [BLAST90], [LAST11], MetaPathways can be used to conduct searches against functional and taxonomic databases.
  - MetaPathways steps: *FUNC SEARCH*, *PARSE FUNC SEARCH*, *SCAN rRNA*, and *ANNOTATE ORFS*
3. **Analyses:** After sequence annotation, MetaPathways performs further taxonomic analyses including the **Lowest Common Ancestor (LCA)** algorithm [MEGAN07] and **tRNA Scan** [TRNASCAN97], and prepares detected annotations for environmental Pathway/Genome database (ePGDB) creation via Pathway Tools.
  - MetaPathways Steps: *PATHOLOGIC INPUT*, *CREATE ANNOT REPORTS*, and *COMPUTE RPKM*.

4. **ePGDB Creation:** MetaPathways then predicts [MetaCyc pathways](#) using the [Pathway Tools](#) software and its pathway prediction algorithm PathoLogic [KARP11], resulting in the creation of an environmental Pathway/Genome Database (ePGDB), an integrative data structure of sequences, genes, pathways, and literature annotations for integrative interpretation.

- MetaPathways Steps: *BUILD ePGDB*

5. **Pathway Export:** Here MetaCyc pathways or reactions are exported in a tabular format for downstream analysis. *As of the v2.5 release, MetaPathways will perform this step automatically.*

- MetaPathways Steps: *BUILD ePGDB*



## 1.2 Output Format

## 1.3 Visualizing Output





## INSTALLATION

MetaPathways supports installing the software using Conda and Pip in a 64-bit Linux environment, or from a container image that can be used with Docker or Singularity. If you do not have administrator (i.e., “root”) access to your computer, we recommend that users install MiniConda if they do not already have it set up. For users wanting to use MetaPathways in an academic grid computing environment, we recommend using the container image *via* Singularity. Below please find a description of how to install MetaPathways using the two supported options:

### 2.1 Container Install

Our container images are hosted at [Quay.io](https://quay.io/hallamlab/metapathways). The following commands assume that you are already familiar with installing and running Docker containers via the `docker` or `singularity` executables:

Using [Docker](#):

```
sudo docker pull quay.io/hallamlab/metapathways
```

Using [Singularity](#):

```
singularity build metapathways.sif docker://quay.io/hallamlab/metapathways:latest
```

More advanced container-related commands are available as Make targets in the `Makefile`.

### 2.2 Installing with Pip and Conda

#### 2.2.1 Summary

Assuming that you have all prerequisites satisfied, installing can be as simple as:

```
conda create --name metapathways python=3.10
conda activate metapathways
pip3 install git+https://bitbucket.org/BCB2/metapathways.git@dev#egg=MetaPathways
metapathways-install-deps.sh
```

Read on to learn the details.

### 2.2.2 Detailed Install

We currently offer a way to use Pip to install the MetaPathways Python package, along with using Conda to install all dependencies. We do not yet have a Conda package for MetaPathways. It is in the works for a future release.

For this to work, we assume that you have the following already set up in your command line environment:

- You have Python 3 (python3) and pip3 installed
- You have already installed Conda, and it is activated
- Development files for zlib, liblzma and libbz2 (required to install PySAM via pip)
- You have wget installed

If you are using a version of Linux that uses apt, and you have root access, then you can execute the following to get all of the dependencies except Conda:

```
sudo apt-get update -y
sudo apt-get install -y \
    python3 \
    python3-pip \
    zlib1g-dev \
    liblzma-dev \
    libbz2-dev \
    wget
```

#### Installing Python Package as Root

If you have root/administrator access, install the MetaPathways Python package using the following command:

```
pip3 install git+https://bitbucket.org/BCB2/metapathways.git@dev#egg=MetaPathways
```

#### Installing Python Package as an Unprivileged User

Use this form to install the package to the user's home directory:

```
pip3 install --user git+https://bitbucket.org/BCB2/metapathways.git@dev#egg=MetaPathways
```

Make sure to add \$HOME/.local/bin to your \$PATH environment variable. This will allow you to use the programs without having to type the full path each time.

#### Conda-Based Setup

Once you have installed the Python package, you will have the following executables either in the system Python install path, or in ~/.local/bin, so be sure to add those paths to your \$PATH environment variable.

```
MetaPathways
metapathways-install-deps.sh
metapathways-data-install.sh
metacount
fastal
fastdb
```

Execute metapathways-install-deps.sh to install pipeline dependencies using Conda.

## 2.3 Reference Sequences

### 2.3.1 Summary

Assuming that you have MetaPathways installed, installing the reference DB can be as simple as:

```
metapathways-data-install.sh /media/ref-db-dir stage_fast_full
```

Read on for detailed instructions.

### 2.3.2 Details

MetaPathways relies on reference databases of sequences to assign functional and taxonomic annotations to the user's sequences. The reference databases, and the index files for each database, take up a significant amount of disk storage. See below for an anecdotal example.

You cannot install these large reference databases within the container, though. You should have a directory on a disk with plenty of capacity, and use Docker's and Singularity's bind options to mount that external directory within the container. Here's an example using Singularity:

```
singularity shell --bind /mnt/sandbox/user:/data docker://quay.io/hallamlab/
↪metapathways:latest
```

The above example binds the host operating system's `/mnt/sandbox/user` directory within the running container as `/data`.

*Warning:* Circa 2021-10, using a beefy computer with many cores and plenty of RAM, performing the staging of the full Blast databases may take an hour, and staging the full set of FAST databases will take around *24 hours*. The Blast `refseq_protein` databases take up ~90 GB of disk capacity, while the FAST `refseq_protein` database takes up ~375 GB. The combination of other staged databases (including both Blast and FAST versions) consumes an additional ~20 GB. Please make sure you have adequate disk capacity before starting the database staging.

We use Snakemake to automate the staging of reference databases needed by MetaPathways. We have installed Snakemake via Conda. If you are using the Docker container, then Conda is already initialized. If you are using the container via Singularity, you must first initialize Conda as follows (note the space between the period character, and the first slash character):

```
. /opt/conda/etc/profile.d/conda.sh
```

Now we can run the `metapathways-data-install.sh` script

```
metapathways-data-install.sh /media/ref-db-dir stage_fast_lite
```

... where `/media/ref-db-dir` is the reference database installation directory (make sure this directory has adequate capacity for the data to be installed).

Above we issued the `stage_fast_lite` command to Snakemake, as an example that runs quickly. There are actually four options for staging the data:

- All databases, indexed for use with Blast: `stage_blast_full`
- All databases except RefSeq Proteome, indexed for use with Blast: `stage_blast_lite`
- All databases, indexed for use with FAST: `stage_fast_full`
- All databases except RefSeq Proteome, indexed for use with FAST: `stage_fast_lite`

So, first decide whether you want to use Blast or FAST, and then decide whether you have the disk space and the install time to install the NCBI RefSeq Proteome reference database. FAST runs faster than Blast, with comparable sensitivity. And the RefSeq Proteome is currently required for MetaPathways to accurately annotate contigs taxonomically. Thus, we recommend running `stage_fast_full`, if you have the disk storage and the time to let it run.

## RUNNING METAPATHWAYS

### 3.1 Input

MetaPathways inputs are fasta files provided in an input folder. The file names must end with a *.fasta* or *.fas*. These fasta files contains the contigs or DNA sequences from assembling.

### 3.2 Parameter File

The parameter file must indicate the setting for any MetaPathways run. An example paramter file can be downloaded as

```
$ wget https://github.com/kishori82/MetaPathways_Python.3.0/raw/kmk-develop/data/text/  
↪template_param.txt
```

Below we describe the settings in the parameter file.

### 3.3 Run

As an illustration we donwload a small input file *testsample1.fasta* in a folder named *mp\_input* and we want the output in a folder names *mp\_output*

```
$ mkdir mp_input  
$ cd mp_input  
$ wget https://github.com/kishori82/MetaPathways_Python.3.0/raw/kmk-develop/data/  
↪testdata/testsample1.fasta  
$ cd ..
```

Now we kick off a run as

```
$ MetaPathways --input mp_input --output mp_output -p template_param.txt -d ~/  
↪MetaPathways_DBs/
```



## PHANDI GUI OVERVIEW

RUNSTATS   KEGG   COG   MetaCyc   SEED   CAZY			
Include ALL, OR, X, or, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000			
Number of sequences in input file BEFORE QC (nucleotide)	1381		
avg length (bp)	2002		
avg length (bp)	13443		
avg length (bp)	4292118		
total base pairs (bp)	71732421		
Number of sequences AFTER QC (nucleotide)	1381		
avg length (bp)	2002		
avg length (bp)	48400		
avg length (bp)	4292118		
total base pairs (bp)	60882787		
Number of translated ORFs BEFORE QC (nucleotide)	61997		
avg length (bp)	138		
avg length (bp)	322		
avg length (bp)	14858		
total base pairs (bp)	10886115		
Number of translated ORFs AFTER QC (nucleotide)	60448		
avg length (bp)	60		
avg length (bp)	329		
avg length (bp)	14858		
total base pairs (bp)	10916097		
Number of hits from metaCyc-v6 (2011-07-03) (AA)	22027		
Number of hits from CAZY_2014_06_04 (AA)	4025		
Number of hits from NCBI_GAG_protocols (AA)	23770		
Number of hits from COG_2013-12-27 (AA)	47501		
Number of hits from KEGG_2012-10-18 (AA)	188715		
Number of hits from unref-2014-01-16 (AA)	238030		
Number of hits from unref-2014-01-16 (AA)	222495		
Annotation: meeting user defined thresholds from CAZY_2014_06_04	1188		
Annotation: meeting user defined thresholds from COG_2013-12-27	24880		
Annotation: meeting user defined thresholds from KEGG_2012-10-18	16107		
Annotation: meeting user defined thresholds from NCBI_GAG_protocols	7147		
Annotation: meeting user defined thresholds from metaCyc-v6 (2011-07-03)	4733		
Annotation: meeting user defined thresholds from unref-2014-01-16	70514		
Annotation: meeting user defined thresholds from unref-2014-01-16	68822		
Total Protein Annotations	41123		
ORF hits meeting user defined thresholds from CRENSONAL_2012-11-08	9		
ORF hits meeting user defined thresholds from COG_2013-12-27	10		
ORF hits meeting user defined thresholds from KEGG_2012-10-18	10		

MetaPathways (Phandi) GUI viewers is a stand-alone desktop tool for inspecting and exporting the large amount of outputs produced by the MetaPathways pipeline.





## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`



## CONTACT

contact



## BIBLIOGRAPHY

- [CIT2002] K. M. Konwar, N. W. Hanson, A. P. Pagé, S. J. Hallam, MetaPathways: a modular pipeline for constructing pathway/genome databases from environmental sequence information. *BMC Bioinformatics* 14, 202 (2013) <http://www.biomedcentral.com/1471-2105/14/202>
- [PRODIGAL2010] D. Hyatt et al., Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 11, 119 (2010).
- [GeneMark12] D. Hyatt, P. F. LoCascio, L. J. Hauser, E. C. Uberbacher, Gene and translation initiation site prediction in metagenomic sequences. *Bioinformatics* 28, 2223–2230 (2012).
- [BLAST90] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment search tool. *J Mol Biol* 215, 403–410 (1990).
- [LAST11] S. M. Kieľbasa, R. Wan, K. Sato, P. Horton, M. C. Frith, Adaptive seeds tame genomic sequence comparison. *Genome Res* 21, 487–493 (2011).
- [MEGAN07] D. H. Huson, A. F. Auch, J. Qi, S. C. Schuster, MEGAN analysis of metagenomic data. *Genome Res* 17, 377–386 (2007).
- [TRNASCAN97] T. M. Lowe, S. R. Eddy, tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Research* 25, 0955–0964 (1997).
- [KARP11] P. D. Karp, M. Latendresse, R. Caspi, The pathway tools pathway prediction algorithm. *Stand Genomic Sci* 5, 424–429 (2011).