

# Build a Bioinformatics Analysis Platform and Apply it to Routine Analysis of Microbial Genomics and Comparative Genomics

**Hualin Liu**

Huazhong Agricultural University <https://orcid.org/0000-0002-3630-5522>

**Bingyue Xin**

Huazhong Agricultural University

**Jinshui Zheng** (✉ [jszheng@mail.hzau.edu.cn](mailto:jszheng@mail.hzau.edu.cn))

Huazhong Agricultural University <https://orcid.org/0000-0002-0147-5900>

**Hao Zhong**

Huazhong Agricultural University

**Yun Yu**

Huazhong Agricultural University

**Donghai Peng**

Huazhong Agricultural University

**Ming Sun** (✉ [m98sun@mail.hzau.edu.cn](mailto:m98sun@mail.hzau.edu.cn))

Huazhong Agricultural University <https://orcid.org/0000-0001-5465-5983>



---

## Method Article

**Keywords:** comparative genomics, COG annotation, phylogenetic orthology, phylogenetic analysis, variants calling, pan-genome, genome assembly, gene prediction, genome annotation, genome distance, Average Nucleotide Identity, antimicrobial and virulence genes mining, PGCGAP

**Posted Date:** March 12th, 2020

**DOI:** <https://doi.org/10.21203/rs.2.21224/v2>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

---

# Abstract

Genomics and comparative genomics have been increasingly used as routine methods for general microbiological researches. However, it is usually necessary to call several tools or even write some scripts to complete some simple analysis, which is complicated for most biological researchers. To simplify the operation process, especially for the convenience of microbiologists in the analysis, here we have developed PGCGAP, a comprehensive, malleable and easily-installed prokaryotic genomics and comparative genomics analysis pipeline, which implements genome assembly, gene prediction and annotation, average nucleotide identity (ANI) calculation, phylogenetic analysis, COG annotation, pan-genome analysis, inference of orthologous gene groups, variants calling and annotation and screening for antimicrobial and virulence genes. Although we have tried our best to simplify the installation and usage of PGCGAP, it may be difficult for non-bioinformatician users to master it. So, a protocol was created to help microbiologists without any experience in bioinformatics to establish their own bioinformatics platform and perform routine analysis. This protocol shows how to choose equipment, to install a Linux subsystem on a laptop with windows 10 system, to install PGCGAP and perform all analysis with an example dataset. The protocol requires a basic understanding of Linux, so an additional web page was written to help uninitiated users learn Linux and whole-genome sequencing (<https://github.com/liaochenlanruo/pgcgap/wiki/Learning-bioinformatics> or <http://bcam.hzau.edu.cn/linuxwgs.php>).

## Introduction

Genome sequencing has become a routine method for common microbiological studies as continuously decrease in the cost of genome sequencing. Various tools have been developed for genome analysis. But for general users, it takes time to install and learn to use various programs and prepare the related input files. Even for some simple purposes, users need to spend much effort to integrate several tools or even write some scripts. For example, when we need a core-genome-SNP based phylogenetic analysis for the isolates from same species, we should successively use the following tools, Bowtie2<sup>1</sup> or BWA<sup>2</sup> for reads mapping, Samtools<sup>3</sup> or GATK<sup>4</sup> for SNP calling, and FastTree<sup>5</sup> or RAxML<sup>6</sup> for phylogenetic tree construction. Therefore, a comprehensive, flexible and efficient pipeline for general analysis is urgently needed. We developed a prokaryotic genomics and comparative genomics analysis pipeline named PGCGAP to coordinate several genomic analysis software packages and in-house scripts to meet the various needs of microbiologists.

### Development of the protocol

PGCGAP was developed to facilitate the work of genomics and comparative genomics analysis of microbes. Considering the important role of basic bioinformatics in microbial research and most microbiologists lacking analysis skills, this protocol describes in detail the installation of Linux systems and demonstrates the software installation methods. Finally, we demonstrated all the usages of PGCGAP step by step through the example datasets.

### Applications of the protocol

PGCGAP can be used for (i) genome assembly, (ii) gene prediction and annotation, (iii) genome distance estimation, (iv) phylogenetic analysis, (v) COG annotation, (vi) pan-genome analysis, (vii) inference of orthologous gene groups, (viii) variants calling and annotation and (ix) screening for antimicrobial and virulence genes. It is worth noting that although the entire pipeline was developed for prokaryotes, some of the modules such as “Assemble”, “MASH”, “OrthoF”, “CoreTree”, “AntiRes” and “STREE” can also be used for the analysis of eukaryotic genomes. Besides, “VAR” applies to the analysis of any haploid genome.

### **Advantages and limitations of this pipeline**

PGCGAP is versatile, feature-rich, easy to install and use, and friendly to microbiologists and bioinformatics beginners. New features will continue to be added. But a Graphical User Interface (GUI) has not been developed.

### **Expertise required to implement the protocol**

Users need to be skilled in using computers, and it will be easier to master this protocol if they have some Linux skills. A webpage introducing the basics of Linux, usage of common commands, software installation, and the whole-genome sequencing technology was developed to help users get started with bioinformatics. Please visit <https://github.com/liaochenlanruo/pgcgap/wiki/Learning-bioinformatics> or <http://bcam.hzau.edu.cn/linuxwgs.php> for more information.

### **Overview of the procedure**

Twelve frequently used prokaryotic genomics and comparative genomics analysis processes were integrated into PGCGAP as different modules. Modules can be used separately or in different combinations for various purposes (Fig. 1). (i) “Assemble” performs genome assembly of Illumina reads, the third generation reads, and the hybrid reads using ABySS<sup>7</sup>, SPAdes<sup>8</sup>, Canu<sup>9</sup>, and Unicycler<sup>10</sup>, respectively. (ii) “Annotate” performs gene prediction and genome annotation by Prokka<sup>11</sup>. (iii) “ANI” computes Average Nucleotide Identity (ANI) between each genome pair by fastANI<sup>12</sup>. Three scripts “triangle2list.pl”, “get\_ANImatrix.pl” and “Plot\_ANIheatmap.R” have been developed here to generate the ANI matrix and plot the correlation matrix heat map (Supplementary Figure S1), respectively. (iv) “MASH” estimates genome and metagenome distance and similarity using MinHash<sup>13</sup>, and a heat map of genome similarity will be generated by two scripts “get\_Mash\_Matrix.pl” and “Plot\_MashHeatmap.R” (Supplementary Figure S2). (v) “Pan” calls Roary<sup>14</sup> to calculate the pan-genome. Two scripts “fmplot.py” and “plot\_3Dpie.R” were developed for result visualization (Supplementary Figure S3). A phylogenetic tree based on single-copy core proteins called by Roary<sup>17</sup> will be constructed (Supplementary Figure S4). (vi) COG (Clusters of Orthologous Group) annotation was conducted by module “pCOG”. Amino acid sequences of each genome were blasted against the COG database, and then all hits were mapped to the COG functional category by in-house scripts. R script “Plot\_COG.R” was written for result visualization (Supplementary Figure S5). Comparison and visualization of COG functional categories among different genomes can be done by Perl script “get\_flag\_relative\_abundances\_table.pl” and R script “Plot\_COG\_Abundance.R” (Supplementary Figure S6). (vii) “OrthoF” uses OrthoFinder<sup>15</sup> for phylogenetic orthology inference. Gene duplication events will be also predicted (Supplementary Figure S7). (viii) “CoreTree” was developed for genome-wide phylogenetic analysis based on the protein sequences or SNPs of

single-copy core genes. Firstly, CD-HIT<sup>16</sup> was used to rapidly generate protein clusters, and then the protein sequences of single-copy core genes were extracted by Perl scripts and aligned by MAFFT<sup>17</sup>. Secondly, on one hand, alignments of protein sequences were concatenated, and the phylogenetic tree was constructed by ModelTest-NG<sup>18</sup> and RAxML-NG<sup>19</sup> (Supplementary Figure S8). On the other hand, the protein sequence alignments were converted into corresponding codon alignments by PAL2NAL v14<sup>20</sup>. Then, the codon alignments were concatenated, and SNP-sites<sup>21</sup> was called to find SNP sites. Finally, ModelTest-NG<sup>18</sup> and RAxML-NG<sup>19</sup> were used to construct the SNPs phylogenetic tree (Supplementary Figure S9). (ix) “AntiRes” calls abricate<sup>22</sup> to screen for antimicrobial and virulence genes from contigs. (x) “VAR” performs genome-wide variants calling by mapping methods. Firstly, single-end or paired-end reads were mapped to a reference genome by BWA<sup>2</sup> after filtered by Sickle<sup>23</sup>. Secondly, variants calling and annotation were performed by FreeBayes<sup>24</sup> and snpEff<sup>25</sup>, respectively. Then, the whole genome SNP alignment and core SNP alignment were obtained by snippy-core<sup>26</sup>. Finally, Gubbins<sup>27</sup> was used to remove SNPs influenced by recombination events of the whole genome SNP alignment and construct a phylogenetic tree (Supplementary Figure S10). Furthermore, the phylogenetic tree of core SNP alignment was constructed by ModelTest-NG<sup>18</sup> and RAxML-NG<sup>19</sup> (Supplementary Figure S11). (xi) “STREE” Construct a phylogenetic tree based on multiple-FASTA sequences in one file. Firstly, the sequences were aligned by MUSCLE<sup>28</sup>, and then Gblocks<sup>29</sup> was used to get the conserved blocks of the aligned sequences. Finally, IQ-TREE<sup>30</sup> was used for phylogenomic inference (Supplementary Figure S12). (xii) “ACC” integrated other useful gadgets, and it now includes function “Assess” for filtering short sequences in the genome and assessing the status of the genome only.

## Experimental design

### Selection of reference genome format for variants calling

The reference genome can be files in FASTA format and GenBank format. If a GenBank file rather than a FASTA file was supplied as the reference, annotation information of the variants will be generated to show to the user which feature was affected by the variants.

### How to balance speed and assembly quality when assembling Illumina reads

According to our experience, ABySS<sup>7</sup> can complete genome assembly faster with fewer computer resources. While the assembly quality of Unicycler<sup>10</sup> and SPAdes<sup>8</sup> is better than ABySS<sup>7</sup>, but it occupies more computer resources and runs very slow. Therefore, we strongly recommend that users choose the auto mode for Illumina data assembly. PGCGAP calls ABySS<sup>7</sup> for Illumina reads assembly firstly. When N50 of the assembled genome is less than 50,000, it automatically calls Unicycler<sup>10</sup> and SPAdes<sup>8</sup> to try multiple parameters for another assembly.

### Choice of a module to construct the phylogenetic tree of single-copy core proteins

Both “CoreTree” and “Pan” can be used to construct a phylogenetic tree of single-copy core proteins. That which module should be used depends on which type of input file the user has. “CoreTree” takes only the amino acid sequence files as inputs, while “Pan” needs both amino acid sequence files and Gff3 files.

## Choice of a module to calculate pairwise genome distance

Both “ANI” and “MASH” can calculate pairwise genome distance. “MASH” is more suitable for dealing with thousands of genomes because of its faster running speed. In addition to nucleotide sequences and assembled genomes, “MASH” can also take amino acid sequences and raw sequencing reads as inputs and can be used to calculate distances between metagenomic samples.

## Reagents

No reagents are required for this protocol. But example datasets were needed to practice the application of PGCGAP.

The example datasets used in this protocol can be downloaded at [http://bcam.hzau.edu.cn/PGCGAP/PGCGAP\\_Examples.tar.gz](http://bcam.hzau.edu.cn/PGCGAP/PGCGAP_Examples.tar.gz).

The output files of example datasets by PGCGAP can be downloaded at [http://bcam.hzau.edu.cn/PGCGAP/PGCGAP\\_Results.tar.gz](http://bcam.hzau.edu.cn/PGCGAP/PGCGAP_Results.tar.gz).

## Equipment

A laptop, desktop PC or server can be used to build a bioinformatics analysis platform, and the suggested hardware requirements are listed in **Table 1**. Slightly lower features are also allowed (CPU must have four logical processors, memory must be greater than 8 G), but the computing speed may decrease, and the capacity of the hard disk can be adjusted according to actual requirements.

## Procedure

### Building a analysis platform on Windows 10

Windows Subsystem for Linux (WSL) allows users to install Linux subsystems directly on Windows 10 system. It can easily run Linux commands and install Linux software, avoiding the installation of third-party virtual machine software. The advantage of WSL is that it makes better use of computer memory and does not require copying files between the host and the virtual machine.

### Configuration of WSL

Timing ~1 min

System requirements: Windows 10 Version 1709, Build 16299 or above, 64-bit systems.

1. Enable WSL: Open "Settings", click "Apps", then find and click "Programs and Features", click "Turn Windows features on or off", find "Windows Subsystem for Linux" and check the box, click OK and restart the computer (**Supplementary Video 1**).

### Install Linux

Timing ~59 min

2. Open the Microsoft Store, search Ubuntu, and choose to install Ubuntu 18.04 LTS. Follow the prompts to set up your username and password. Here we create an account with the username “bio” (**Supplementary Video 2**). After the installation is finished, we need to do some configuration on the system (**Supplementary Video 3**).

3. Enter the following command in the terminal to update the source:

```
$sudo apt-get update
```

4. Set the password for root.

```
$sudo passwd root
```

5. Change the content of `ptrace_scope` from 1 to 0.

```
$su root
```

```
#echo 0 > /proc/sys/kernel/yama/ptrace_scope
```

6. Switch to the ordinary user “bio”.

```
#su bio
```

7. Installation of Miniconda

(A) Installation of Miniconda on Linux

(i) Here, [Miniconda](#) will be installed, go to the official website, and select the installation file suitable for your system and python version. Here, Miniconda 3 will be installed (**Supplementary Video 4**).

```
$wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

(ii) Start installation

```
$bash Miniconda3-latest-Linux-x86_64.sh
```

Press Enter when prompted to visualize the license agreement, enter “yes” and press “Enter” to continue. Press “Enter” to confirm the default installation location. Miniconda was installed in the `miniconda3` directory under the user’s home directory. Type “yes” and press “Enter” to initialize `miniconda3`. Finally, run the command “`source ~/.bashrc`” in the terminal.

```
$source ~/.bashrc
```

(iii) Set up Bioconda channel. Add the channels by entering the following three commands in the terminal.

```
$conda config --add channels defaults
```

```
$conda config --add channels bioconda
```

```
$conda config --add channels conda-forge
```

(B) Install Miniconda on MacOS

(i) Installation of Miniconda3

```
$wget https://repo.anaconda.com/miniconda/Miniconda3-latest-MacOSX-x86_64.sh
```

```
$sh Miniconda3-latest-MacOSX-x86_64.sh
```

```
$source ~/.bash_profile
```

(ii) Add channels of Bioconda

```
$conda config --add channels defaults
```

```
$conda config --add channels bioconda
```

```
$conda config --add channels conda-forge
```

Installation of PGCGAP (**Supplementary Video 5**).

Timing ~34 min

8. Create a pgcgap environment for the installation of PGCGAP.

```
$conda create -n pgcgap python=3
```

9. Activate the pgcgap environment.

```
$conda activate pgcgap
```

10. Installation of PGCGAP.

```
$conda install pgcgap
```

11. Check if the dependent software packages were installed.

```
$pgcgap --check-external-programs
```

12. Set up the COG database.

```
$pgcgap --setup-COGdb
```

13. Exit the pgcgap environment.

```
$conda deactivate
```

**Step by Step examples**

Timing ~2.3 d

The usage and parameters of PGCGAP can be viewed by typing “pgcgap -h” in the terminal. Next, we show how to run all the modules of PGCGAP through a dataset.

14. Download and decompress the example dataset.

```
$wget http://bcam.hzau.edu.cn/PGCGAP/PGCGAP_Examples.tar.gz
```

```
$tar -zxvf PGCGAP_Examples.tar.gz
```

In this example, the working directory locates at the H drive. All hard disks in Windows were mounted in the “/mnt” directory of Ubuntu Linux. The “PGCGAP\_Examples/Reads/Illumina” directory contains six Illumina Hiseq paired-end reads of *Escherichia coli*, the “PGCGAP\_Examples/Reads/Oxford” directory contains the Oxford Nanopore reads of *Escherichia coli* K12, and the “PGCGAP\_Examples/Reads/PacBio” directory contains the Pacific Biosciences released P6-C4 chemistry reads of *Escherichia coli* K12.

“PGCGAP\_Examples/Reads/MG1655” is the GenBank format file of *E.coli* K-12 *substr.* MG1655, used as the reference genome. The “PGCGAP\_Examples/Reads/Hybrid” directory contains two short reads files and one long reads file of the same strain. “PGCGAP\_Examples/Other\_inputs/ proteins.fas” contains 18 protein sequences of MFS transporter from several bacteria species.

15. Activate the pgcgap environment.

```
$conda activate pgcgap
```

16. Example 1: Genome assembly with Illumina reads.

Paired-end reads of six strains in the directory “Reads/Illumina/” were used as inputs. In the dataset, the naming format of the genome is “strain\_1.fastq.gz” and “strain\_2.fastq.gz”. The string after the strain name is “\_1.fastq.gz”, and its length is 11, so “-suffix\_len” was set to 11. Users can choose “abyss”, “spades” and “auto” for the genome assembly. The assembly speed with “abyss” is faster, and the assembly quality with “spades” is better. To take into account the speed and quality of assembly, we suggest using “auto” mode for assembly. “-filter\_length” was set here to remove sequences shorter than 200 bp from the assembled genomes.

```
$pgcgap --Assemble --platform illumina --assembler abyss --filter_length 200 --ReadsPath Reads/Illumina --reads1 _1.fastq.gz --reads2 _2.fastq.gz --kmmer 81 --threads 4 --suffix_len 11
```

```
$pgcgap --Assemble --platform illumina --assembler spades --filter_length 200 --ReadsPath Reads/Illumina --reads1 _1.fastq.gz --reads2 _2.fastq.gz --threads 4 --suffix_len 11
```

```
$pgcgap --Assemble --platform illumina --assembler auto --filter_length 200 --ReadsPath Reads/Illumina --reads1 _1.fastq.gz --reads2 _2.fastq.gz --kmmer 81 --threads 4 --suffix_len 11
```

New directories and documents were generated after the program is finished. The assembly results for each genome are in the “Results/Assemblies/Illumina” directory, and all scaffolds of the strains were stored in



“Results/Assemblies/Scaf/Illumina”. “\*.filtered.fas” is the genome with short sequences removed. “\*.prefilter.stats” describes the status of the genome before filtering, and “\*.filtered.stats” describes the status of the genome after short sequences filtering. While “abyss” was chosen as the assembler, users are advised to check the assembly stats file (such as Results/Assemblies/Illumina/SRR9620252\_assembly/SRR9620252-stats.tab) of each genome to ensure that the value of N50 is greater than 50,000 bp. The file “scaf.list” under the working directory contains the absolute path of all genomes.

#### 17. Example 2: Oxford reads assembly.

Oxford nanopore only produces one reads file (“Reads/Oxford/oxford.fasta”), so only the parameter of “–reads1” needs to be set, here the value is “.fasta”. “–genomeSize” is the estimated genome size, and users can check the genome size of similar strains in the NCBI database for reference. The parameter was set to “4.8m” here. The suffix of the reads file here is “.fasta” and its length is 6, so “–suffix\_len” was set to 6.

```
$pgcgap --Assemble --platform oxford --ReadsPath Reads/Oxford --reads1 .fasta --genomeSize 4.8m --threads 4 --suffix_len 6 --filter_length 200
```

The results are stored in the “Results/Assemblies/Oxford” directory and the “Results/Assemblies/Scaf/Oxford” directory. The former contains all intermediate files and genome files, the latter contains only the assembled genome.

#### 18. Example 3: PacBio reads assembly.

PacBio also produces only one reads file too (“Reads/PacBio/pacbio.fastq”), the parameter settings are similar to Oxford. The strain name is “pacbio” with the suffix “.fastq” and the suffix length is 6, so “–suffix\_len” was set to 6.

```
$pgcgap --Assemble --platform pacbio --ReadsPath Reads/PacBio --reads1 .fastq --genomeSize 4.8m --threads 4 --suffix_len 6 --filter_length 200
```

The results are stored in the “Results/Assemblies/PacBio” directory and in the “Results/Assemblies/Scaf/PacBio” directory. The former contains all intermediate files and genome files, the latter containing only the assembled genome.

#### 19. Example 4: hybrid assembly of short reads and long reads.

Paired-end short reads and long reads in the directory “Reads/Hybrid/” were used as inputs. Illumina reads and long reads must be from the same isolates.

```
$pgcgap --Assemble --platform hybrid --ReadsPath Reads/Hybrid --short1 short_reads_1.fastq.gz --short2 short_reads_2.fastq.gz --long long_reads_high_depth.fastq.gz --threads 4
```

The results are stored in the “Results/Assemblies/Hybrid” directory, and the final assembly was named as “assembly.fasta”.

#### 20. Example 5: Gene prediction and annotation.

Here, the assembly results of Illumina reads were taken as inputs (“Results/Assemblies/Scaf/Illumina/\*.fa”). The suffix of the genome is “-8.fa”. When running the program, the value of the “-Scaf\_suffix” parameter cannot be quoted. Here, -8.fa should not be quoted.

```
$pgcgap -Annotate -scafPath Results/Assemblies/Scaf/Illumina -Scaf_suffix -8.fa -genus Escherichia -species “Escherichia coli” -codon 11 -threads 4
```

The generated files are stored in the “Results/Annotations” directory, and files in the directories “Results/Annotations/AAs”, “Results/Annotations/CDs” and “Results/Annotations/GFF” will be used for subsequent analysis.

#### 21. Example 6: Constructing the single-copy core proteins tree and core SNPs tree.

The phylogenetic trees of single-copy core proteins and single-copy core genes SNPs will be constructed using the six *E. coli* genomes sequenced by Illumina as datasets. The input files are the amino acid sequence files (“Results/Annotations/AAs/\*.faa”) and the nucleotide sequence files (“Results/Annotations/CDs/\*.ffn”) obtained by the genome annotation. Amino acid files and nucleotide files must be suffixed with “.faa” and “.ffn”, respectively. The “.faa” and “.ffn” files of the same strain should have the same prefix name. The name of protein IDs and gene IDs in the Amino acids file and nucleotide file should be started with the strain name. The Prokka<sup>14</sup> software was suggested to generate the input files.

```
$pgcgap -CoreTree -CDsPath Results/Annotations/CDs -AAsPath Results/Annotations/AAs -codon 11 -strain_num 6 -threads 4
```

The result files are stored in the “Results/CoreTrees” directory. “ALL.core.protein.\*.support” and “ALL.core.snp.\*.support” are the phylogenetic tree files of the single-copy core proteins and the core SNPs constructed with the best-fit model of evolution, respectively. Users can import these two files into MEGA<sup>31</sup> or iTOL<sup>32</sup> to view the topology.

#### 22. Example 7: Constructing the single-copy core protein tree only.

If the “-CDsPath” was set to “NO”, the nucleotide files will not be needed, and the phylogenetic tree of core SNPs will not be constructed too.

```
$pgcgap -CoreTree -CDsPath NO -AAsPath Results/Annotations/AAs -codon 11 -strain_num 6 -threads 4
```

#### 23. Example 8: pan-genome analysis and phylogenetic tree construction.

GFF3 files (With “.gff” as the suffix) of each strain placed into a directory (“Results/Annotations/GFF/\*.gff”). They must contain the nucleotide sequence at the end of the file. Protein sequence files (one per species) in FASTA format under another directory were also needed (“Results/Annotations/AAs/\*.faa”) if the parameter “-PanTree” was provided for constructing a phylogenetic tree. It should be noted that the “\*.gff” file and the “\*.faa” file must correspond. We strongly recommend using Prokka<sup>14</sup> to generate the files. If the “-Annotate” function was run first, the files will be generated automatically.

```
$pgcgap -Pan -codon 11 -strain_num 6 -threads 4 -GffPath Results/Annotations/GFF -PanTree -AAsPath Results/Annotations/AAs
```

The results are stored in the “Results/PanGenome” directory. A spreadsheet named “gene\_presence\_absence.csv” lists each gene and which sample is presented in. Users can take the gene\_presence\_absence.csv file and a traits file to conduct pan-genome wide association studies with Scoary<sup>33</sup> software. At the same time, some visual results (“\*.pdf”) are also outputted. “Results/PanGenome/Core/Roary.core.protein.BIC.AIC.AICc.HIVW+I+G4+F.raxml.support” is the phylogenetic tree constructed based on the single-copy core proteins called by Roary<sup>14</sup> software. “HIVW+I+G4+F” represent the best-fit model of evolution for the protein alignments according to AIC<sup>34</sup>, AICc, and BIC<sup>35</sup> statistical criteria. If the parameters “-PanTree” and “-AAsPath” were not provided, the phylogenetic tree will not be constructed.

#### 24. Example 9: Inference of orthologous gene groups.

The input files are also the amino acid sequence files suffixed with “.faa” (“Results/Annotations/AAs/\*.faa”).

```
$pgcgap -OrthoF -threads 4 -AAsPath Results/Annotations/AAs
```

The resulting files are placed in the “Results/OrthoFinder/Results\_orthoF” directory.

#### 25. Example 10: Compute whole-genome Average Nucleotide Identity.

The input file named “scaf.list” contains the absolute path of each genome, one per line. If the “-Assemble” function was run first, the list file will be generated automatically. The value of the parameter “-Scaf\_suffix” depends on the actual situation, here is “-8.fa”.

```
$pgcgap -ANI -threads 4 -queryL scaf.list -refL scaf.list -ANIO Results/ANI/ANIs -Scaf_suffix -8.fa
```

The results are stored in the “Results/ANI” directory. The file “ANI” contains comparison information of genome pairs. The document is composed of five columns, each of which represents ANI value, count of bidirectional fragment mappings and total query fragments, respectively. A heat map file “ANI\_matrix.pdf” was generated.

#### 26. Example 11: Genome similarity estimation using MinHash

It takes genome files (complete or draft) in a directory as inputs (Default: Results/Assemblies/Scaf/Illumina).

```
$pgcgap -MASH -scafPath Results/Assemblies/Scaf/Illumina -Scaf_suffix -8.fa
```

The results are stored in the “Results/MASH” directory. The file “MASH” shows pairwise distance between pair genomes and each column represents Reference-ID, Query-ID, Mash-distance, P-value and Matching-hashes, respectively. A heat map file “MASH\_matrix.pdf” was generated.

#### 27. Example 12: COG annotation.

The input files are also the amino acid sequence files suffixed with “.faa” (“Results/Annotations/AAs/\*.faa”).

```
$pgcgap -pCOG -threads 4 -strain_num 6 -AAsPath Results/Annotations/AAs
```

The results are stored in the "Results/COG" directory. The super COG table of each strain ("\*.2Scog.table") and its plot ("\*.2Scog.table.pdf") will be generated. "All\_flags\_relative\_abundances.table" is a table containing the relative abundance of each flag for all strains, and "All\_flags\_relative\_abundances.pdf" is the corresponding visualization result.

28. Example 13: Variants calling, and phylogenetic tree construction based on a reference genome.

The six genomes sequenced by Illumina were chosen as datasets ("Reads/Illumina/\*.gz"). Taking *Escherichia coli* K-12 *substr.* MG1655 as the reference genome and the reference file "MG1655.gbff" in the GenBank format is stored in the "Reads" directory. The absolute path of the reference genome (here is "/mnt/h/PGCGAP\_Examples/Reads/MG1655.gbff") is required to run the program.

```
$pgcgap -VAR -threads 4 -refgbk /mnt/h/PGCGAP_Examples/Reads/MG1655.gbff -ReadsPath Reads/Illumina -reads1 _1.fastq.gz -reads2 _2.fastq.gz -suffix_len 11 -strain_num 6 -qualtype sanger
```

The resulting files are stored in the "Results/Variants" directory, where the "Core" directory contains the core SNPs of all strains and their phylogenetic tree.

29. Example 14: Screening of contigs for antimicrobial and virulence genes

It takes genome files (complete or draft) in a directory as inputs (Default: Results/Assembles/Scaf/Illumina).

```
$pgcgap -AntiRes -scafPath Results/Assembles/Scaf/Illumina -Scaf_suffix -8.fa -threads 4 -db ncbi -identity 75 -coverage 50
```

The resulting files are stored in the "Results/AntiRes" directory. "\*.tab" files are screening results of each strain, and the "summary.txt" file contains a matrix of gene presence/absence for all strains.

30. Example 15: Perform all functions for paired-end reads.

Only the reads file and reference file should be provided. For the sake of flexibility, the "VAR" function needs to be added extra.

```
$pgcgap -All -platform illumina -ReadsPath Reads/Illumina -reads1 _1.fastq.gz -reads2 _2.fastq.gz -suffix_len 11 -kmmer 81 -genus Escherichia -species "Escherichia coli" -codon 11 -strain_num 6 -threads 4 -VAR -refgbk /mnt/h/PGCGAP_Examples/Reads/MG1655.gbff -qualtype sanger -PanTree
```

31. Example 16: Filter short sequences in the genome and assess the status of the genome.

"Assess" takes the assembled genomes as inputs. Firstly, it assesses the stats of the genome; Secondly, the sequences shorter than "-filter\_length" were deleted from the genome; Finally, the stats of the filtered genome was assessed. The results files will be stored in the same directory as the inputs.

```
$pgcgap -ACC -Assess -scafPath Results/Assembles/Scaf/Illumina -Scaf_suffix -8.fa -filter_length 200
```

32. Example 17: Construct a phylogenetic tree based on multiple-FASTA sequences in one file.

“STREE” takes the file containing multiple-FASTA sequences as input. The results files will be stored in “Results/STREE”. The file “proteins.fas.aln.gb.treefile” contains the final phylogenetic tree.

```
$pgcgap -STREE -seqfile Other_inputs/proteins.fas -seqtype p -bsnum 500
```

## Troubleshooting

Troubleshooting advice can be found in **Table 2**.

## Time Taken

The following marked time was tested in the WSL on the laptop, the features of the laptop are as follows: i7-4710MQ CPU (with 4 cores and 8 logical processors), 16 GB DDR3L RAM, 240 G SSD, 1 T HDD. All commands called 4 threads.

Step 1, configuration of WSL, 1 min.

Step 2, installation of Linux, 43 min.

Step 3-6, the configuration of Linux, 10 min.

Step 7, installation of Miniconda, 6 min.

Step 8-13, installation of PGCGAP, 34 min.

Step 14, download and decompress the example dataset, 11 min.

Step 15, activate the pgcgap environment, 8 s.

Step 16, Illumina reads assembly by abyss, spades and auto, 43 min, 12.8 h, and 5.7 h, respectively.

Step 17, Oxford reads assembly, 1.6 h.

Step 18, PacBio reads assembly, 54 min.

Step 19, hybrid assembly of short reads and long reads, 18 min

Step 20, gene prediction and annotation, 1 h.

Step 21, constructing the single-copy core protein tree and core SNPs tree, 2.6 h.

Step 22, constructing the single-copy core protein tree only, 2.4 h.

Step 23, pan-genome analysis and phylogenetic tree constructing, 3 h.

Step 24, inference of orthologous gene groups, 51 min.

Step 25, compute whole-genome Average Nucleotide Identity, 17 s.

Step 26, genome similarity estimation using MinHash, 1 min.

Step 27, COG annotation, 20.3 h.

Step 28, variant calling, and phylogenetic tree construction based on the reference genome, 11.8 h.

Step 29, Screening of contigs for antimicrobial and virulence genes, 30 s.

Step 30, Perform all functions for paired-end reads, 1.9 d.

Step 31, Filter short sequences in the genome and assess the status of the genome, 15 s.

Step 32, Construct a phylogenetic tree based on multiple-FASTA sequences in one file, 3 h.

## Anticipated Results

## References

- 1 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357 (2012).
- 2 Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv:1303.3997v2 [q-bio.GN]* (2013).
- 3 Heng Li *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)* **25**, 2078-2079 (2009).
- 4 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297-1303 (2010).
- 5 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 6 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
- 7 Jackman, S. D. *et al.* ABySS 2.0: resource-efficient assembly of large genomes using a Bloom filter. *Genome Res.* **27**, 768-777 (2017).
- 8 Nurk, S. *et al.* Assembling Genomes and Mini-metagenomes from Highly Chimeric Reads. 158-170 (Springer Berlin Heidelberg).
- 9 Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* **27**, 722-736 (2017).

- 10 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads. *PLoS Comp. Biol.* **13**, e1005595 (2017).
- 11 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-2069 (2014).
- 12 Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* **9**, 5114-5114 (2018).
- 13 Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biology* **17**, 132 (2016).
- 14 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691-3693 (2015).
- 15 Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biology* **20**, 238 (2019).
- 16 Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659 (2006).
- 17 Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059-3066 (2002).
- 18 Darriba, D. *et al.* ModelTest-NG: A New and Scalable Tool for the Selection of DNA and Protein Evolutionary Models. *Mol. Biol. Evol.* **37**, 291-294 (2019).
- 19 Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453-4455 (2019).
- 20 Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic Acids Res.* **34**, W609-612 (2006).
- 21 Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb Genom* **2**, e000056 (2016).
- 22 Seemann, T. Abricate, Github <https://github.com/tseemann/abricate>.
- 23 Joshi NA & JN, F. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files (Version 1.33) [Software]. Available at <https://github.com/najoshi/sickle>. (2011).
- 24 Garrison E & Marth G. Haplotype-based variant detection from short-read sequencing. *arXiv preprint arXiv:1207.3907 [q-bio.GN]* (2012).
- 25 Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80-92 (2012).

- 26 Seemann, T. Snippy: Rapid haploid variant calling and core genome alignment. Available at <https://github.com/tseemann/snippy>. (2014).
- 27 Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
- 28 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792-1797 (2004).
- 29 Castresana, J. Selection of Conserved Blocks from Multiple Alignments for Their Use in Phylogenetic Analysis. *Mol. Biol. Evol.* **17**, 540-552 (2000).
- 30 Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* (2020).
- 31 Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary Genetics Analysis across Computing Platforms. *Mol. Biol. Evol.* **35**, 1547-1549 (2018).
- 32 Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res.* **44**, W242-W245 (2016).
- 33 Brynildsrud, O., Bohlin, J., Scheffer, L. & Eldholm, V. Rapid scoring of genes in microbial pan-genome-wide association studies with Scoary. *Genome Biology* **17**, 238 (2016).
- 34 Akaike, H. A new look at the statistical model identification. *IEEE Transactions on Automatic Control* **19**, 716-723 (1974).
- 35 Schwarz, G. Estimating the Dimension of a Model. *Ann. Statist.* **6**, 461-464 (1978).

## Acknowledgements

This work was made possible through funding from the National Key R&D Program of China (2017YFD0201201), National Natural Science Foundation of China (31670085, 31970003 and 31770003), and China 948 Program of Ministry of Agriculture (2016-X21).

## Figures



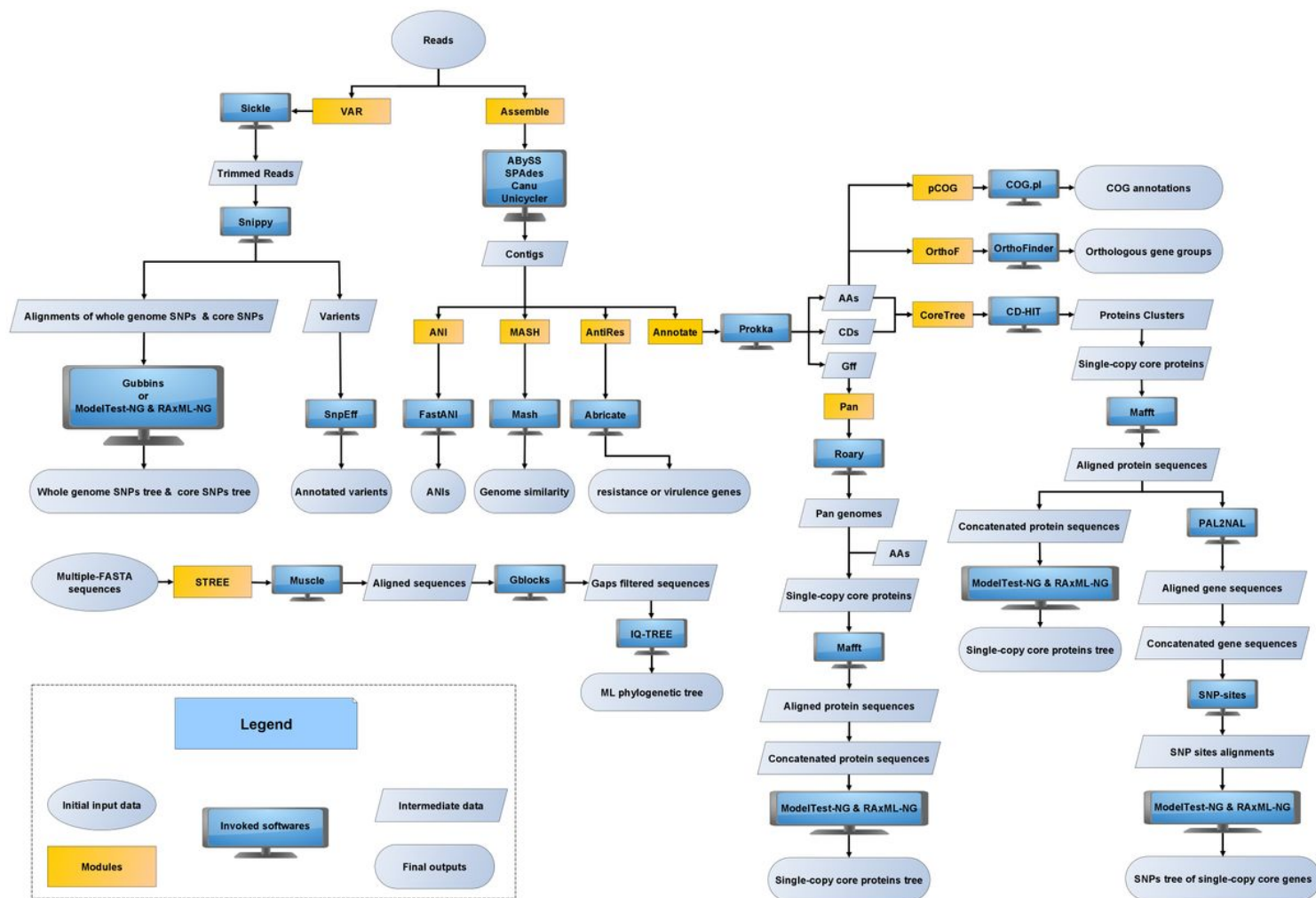


Figure 1

Overview of the PGCGAP workflow

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryFigureS11.AphylogenetictreeofthecoreSNPalignmentgeneratedbymoduleVAR.jpg](#)
- [SupplementaryVideo2InstallUbuntu.mp4](#)
- [SupplementaryFigureS3.Plotsofscriptsfplot.pyandplot3Dpie.R.jpg](#)
- [SupplementaryVideo1EnableWSL.mp4](#)
- [SupplementaryFigureS2.ThecorrelationmatrixheatmapdrawnbythemoduleMASHofPGCGAP.jpg](#)
- [SupplementaryFigureS1.ThecorrelationmatrixheatmapdrawnbythemoduleANIOFPGCGAP.jpg](#)
- [SupplementaryFigureS5.ApicturedescribesthefrequencyofeachflagforthestrainSRR9620252.jpg](#)
- [SupplementaryVideo3ConfigureUbuntu.mp4](#)
- [SupplementaryFigureS4.AphylogenetictreeofsinglecopycoreproteinscalledbymodulePan.jpg](#)
- [SupplementaryFigureS8.PhylogenetictreeofsinglecopycoreproteinsgeneratedbymoduleCoreTree.jpg](#)

- [SupplementaryFigureS9.SNPsphylogenetictreeofSinglecopycoregenesgeneratedbymoduleCoreTre.jpg](#)
- [SupplementaryFigureS10.AphylogenetictreeofthewholegenomeSNPalignmentgeneratedbymoduleVAR.jpg](#)
- [SupplementaryVideo4InstallBioconda.mp4](#)
- [Table2.Troubleshooting.docx](#)
- [SupplementaryFigureS6.Aheatmapdepictstherelativeabundanceofeachflagforallstrains.jpg](#)
- [Table1.Suggestedhardwarerequirementsforabioinformaticsanalysisplatform.docx](#)
- [SupplementaryVideo5InstallPGCGAP.mp4](#)
- [SupplementaryFigureS12.AphylogenetictreeoftheMFStransporterogeneratedbymoduleSTREE.jpg](#)
- [SupplementaryFigureS7.ArootedspeciestreeforthespeciesbeinganalyzedinferredbymoduleOrthoF.jpg](#)