

# PanGeT v.1.0

## User Manual

*Last modified on October, 2016*

**TABLE OF CONTENTS**

<b>1. Introduction.....</b>	<b>3</b>
<b>2. Dependencies for Linux.....</b>	<b>3</b>
<b>3. Dependencies for Windows.....</b>	<b>4</b>
<b>4. Dependencies for Mac.....</b>	<b>4</b>
<b>5. Running Procedure.....</b>	<b>4</b>
<b>5a. PanGeT_BLASTN.....</b>	<b>4</b>
<b>5b. PanGeT_BLASTP.....</b>	<b>6</b>
<b>6. Sample Results.....</b>	<b>8</b>
<b>7. Pathogenic and nonpathogenic strain comparison.....</b>	<b>13</b>

## 1. Introduction:

PanGeT is a user friendly pan-genome tool. Theoretically, PanGeT can work with a single query and 'N' number reference genomes given by the users. It differentiates the annotated CDS/protein based on the homology score or H-value cut-off scores ( $H \geq 0.42$ ). It computes the core, strain specific and dispensable/accessory genes based on the above threshold and computes the pan-genome. Further, the computed core and strain specific genes will be graphically displayed in the form of a 'Flower plot'. Link for the list of dispensable genes are also provided in the flower plot.

**Tested Operating Systems:** Ubuntu 11.04, Fedora core 14, Mac os X EI Captain, Windows 8

**Architecture support:** 64 bit

## 2. Dependencies for Linux

(1) Make sure that you have installed BLAST 2.4.0+ in your machine.

### Installation commands

-----  
rpm -ivh ncbi-blast-2.4.0+-2.x86\_64.rpm

OR

sudo apt-get install ncbi-blast+

(in case of Ubuntu)

-----  
(2) Make sure that your machine has latex installed.

Check it by typing command 'latex' in the terminal.

-----  
This is pdfTeXk, Version 3.141592-1.40.3 (Web2C 7.5.6).

-----  
If it is not installed in your machine. Install it manually.

(3) In case of Ubuntu type,

-----

```
sudo apt-get install texlive-full
```

-----

### **3. Dependencies for Windows**

Download the following packages and make the default installation

ActivePerl-5.16.3.1604-MSWin32-x64-298023.exe

basic-miktex-2.9.5105-x64.exe

ncbi-blast-2.4.0+-win64

### **4. Dependencies for Mac**

Download the following packages and make the default installation

Basic TeX

Website: (<https://tug.org/mactex/>)

ncbi-blast-2.4.0+

### **5. Running procedure**

#### **5a. PanGeT BLASTN mode**

(1) Create folder in the name of a genus (or any other name) that you want to analyze. For example: Salmonella, Ecoli, Brucella etc.,

(2) The '.fna' files and '.ptt' files of strains you want to compare in a genus should be kept inside the created Genus folder. Specify input and output folder paths, H-value, E-value cutoffs in the parameter file.

(Ex: /home/ubuntu/input\_files/**Salmonella**)

(Ex: C:\Users\username\input\_files\**Salmonella** in windows)

(3) Example : Parameter File for Linux / Mac OS

Input\_directory        /home/user/Documents/*Salmonella*  
Output\_directory       /home/user/Documents/Salmonella\_output  
Hvalue                0.42  
Evalue                1e-5  
Threads                4

(4) Example : Parameter File for windows users

Input\_directory        C:\Users\username\Documents\*Salmonella*  
Output\_directory       C:\Users\username\Documents\Salmonella\_result  
Hvalue                0.42  
Evalue                1e-5  
Threads                4

(5) Run the executable file PanGeT using the following command in the terminal

**`./PanGeT_BLASTN parameter.txt or perl PanGeT_blastn.pl parameter.txt`**

(6) An entire list of RefSeq Ids (Ex: NC\_000913) of the genus will be listed down to choose the reference strain

0. NC\_003197  
1. NC\_003198  
2. NC\_004631  
3. NC\_006511  
4. NC\_006905  
5. NC\_010067

6. NC\_011080
7. NC\_011083
8. NC\_011094
9. NC\_011147
10. NC\_011149

Enter the corresponding number entry to keep it as 'Query' strain

- (7) The Program will create two folders 'process\_genus\_name' and 'output\_genus\_name'. (eg. process\_Salmonella and output\_Salmonella).
- (8) Genus output folder (eg. output\_Salmonella) has Result.pdf, which contains the pan-genome flower plot with links, core genes output named CORE and many sub-folders with the name of strain specific RefSeq ids, which has different outputs, such as list of conserved genes, dispensable genes, unique genes and the total genes.
- (9) For further output details, please go through **Sample results** section in the manual.

### **5b. PanGeT BLASTP mode**

- (1) Create folder in the name of a genus (or any other name) that you want to analyze. For example: Salmonella, Ecoli, Brucella etc.,
- (4) The '.faa' of strains you want to compare in a genus should be kept inside the created Genus folder. Specify input and output folder paths, H-value, E-value cutoffs in the parameter file.

(Ex: /home/ubuntu/input\_files/**Salmonella**)

(Ex: C:\Users\username\input\_files\**Salmonella** in windows)

- (3) Example : Parameter File for Linux / Mac OS

Input\_directory        /home/user/Documents/*Salmonella*

Output\_directory      /home/user/Documents/Salmonella\_output

Hvalue                0.42

Evalue                1e-5

Threads 4

(4) Example : Parameter File for windows users

Input\_directory C:\Users\username\Documents\Salmonella

Output\_directory C:\Users\username\Documents\Salmonella\_result

Hvalue 0.42

Evalue 1e-5

Threads 4

(5) Run the executable file PanGeT using the following command in the terminal

**./PanGeT\_BLASTP parameter.txt or perl PanGeT\_blastp.pl parameter.txt**

(6) An entire list of RefSeq Ids (Ex: NC\_000913) of the genus will be listed down to choose the reference strain

0. NC\_003197

1. NC\_003198

2. NC\_004631

3. NC\_006511

4. NC\_006905

5. NC\_010067

6. NC\_011080

7. NC\_011083

8. NC\_011094

9. NC\_011147

10. NC\_011149

Enter the corresponding number entry to keep it as 'Query' strain

(7) The Program will create two folders 'process\_genus\_name' and 'output\_genus\_name'. (eg. process\_Salmonella and output\_Salmonella).

(8) Genus output folder (eg. output\_Salmonella) has Result.pdf, which contains the pan-genome flower plot

with links, core genes output named CORE and many sub-folders with the name of strain specific RefSeq ids, which has different outputs, such as list of conserved genes, dispensable genes, unique genes and the total genes.

(9) For further output details go through **Sample results** section in the manual.

## 6. **Sample results:**

The '**Flower plot**' (**Result.pdf**) generated by '**PanGeT**' is an output, which describes the number of 'core genes' found within the entire 'Genus' or 'Species' while selecting the specific 'Query' and 'Reference' genomes. The flower plots will be opened in 'Adobe Reader' or 'Document viewer' utilities. The list of 'Dispensable genes' is shown through a link at the bottom right side of the flower plots.

### **Genus output folder**

Genus output folder (eg. output\_Salmonella) contains the following files

#### **Result.pdf**

Result.pdf has a flower plot which has hyper-links for the list of core, dispensable and strain specific genes (see example: Fig 1). When you click on the hyper-links, the core genes.html (see example: Fig 2), Strain specific genes.html (see example: Fig 3) will open in your default browser. Further, genes.html pages where users can download each core sequence through sequence link (see example: Fig 4) and also retrieve annotations from KEGG for every gene sequence (see example: Fig 5).

#### **CORE**

This file contains the core orthologue genes present in all the genomes found by reciprocal blast hits

#### **Strain Folders**

(Eg. NC\_003917)

It will create sub-folders with the name of strain specific RefSeq ids. Each folder contains

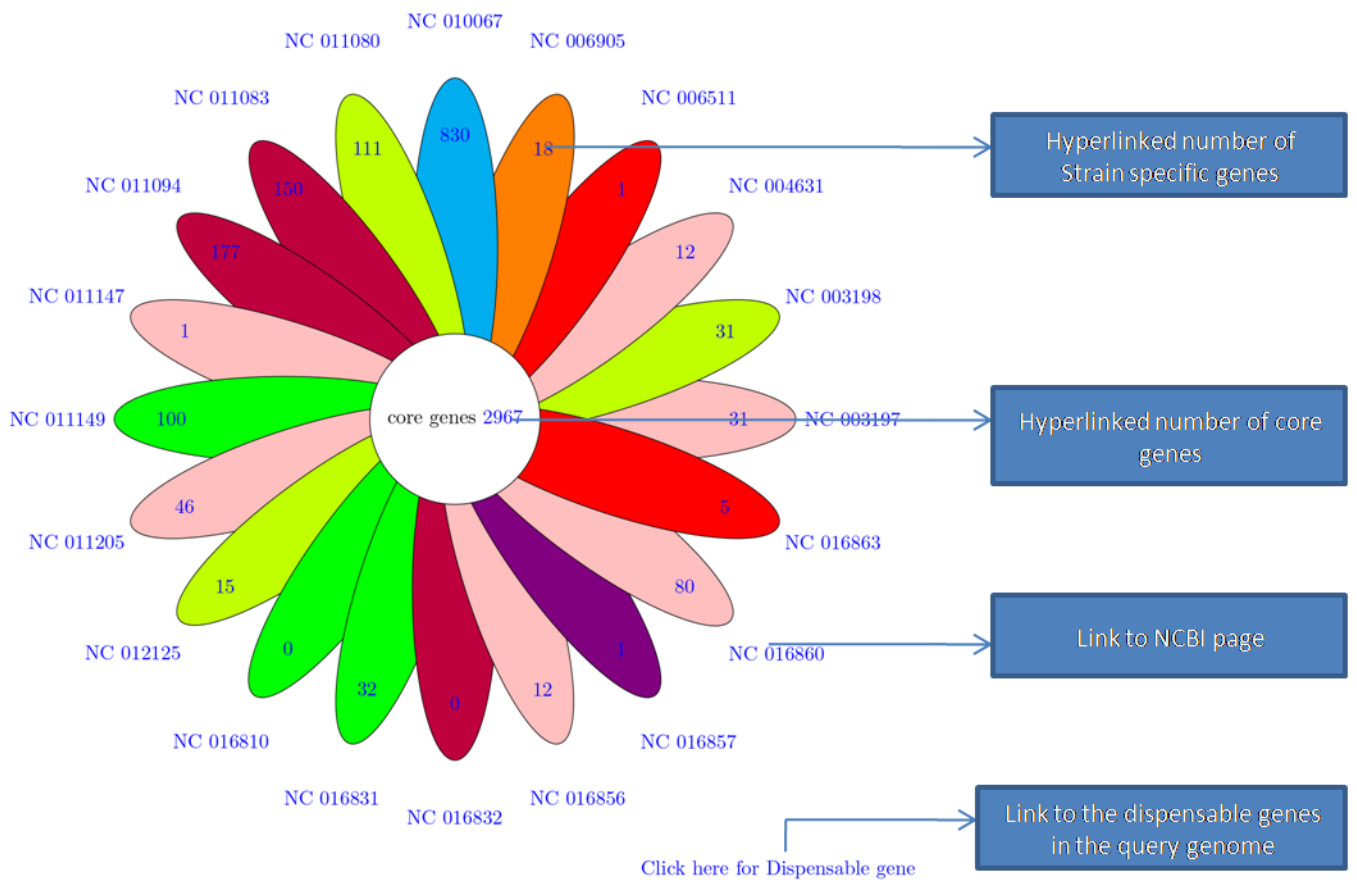
1. **Conserved genes files:** conserved with one genome to 'n' genome you have run
2. **Dispensable genes file:** contains genes present in at least 1 to n-1 genomes but not in all genomes.
3. **Conserved in all genomes files:** contains genes present in all genomes, but not reciprocal blast hits to every other genome.



4. **Unique genes file:** genes specific to the particular strain
5. **Total genes file:** All the genes with their H-value (Homology value) with other genome.

**OUTPUT**

Core,Dispensable and Strain specific genes plot



\* click on the numbers to see the unique gene sequences  
 \* click on the NCBI-IDs to go to NCBI genome page

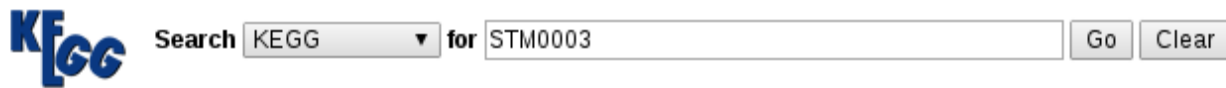
**Fig 1: Core, Dispensable and strain specific genes plot .**

CORE GENE OF NC_003197					
LOCATION:PID	STRAND	GENE_ID	SYNONYM	COG_CLASS	ANNOTATION: SEQUENCE
>337..2799:STM0002	+	thrA	<a href="#">STM0002</a>	COG0527E	bifunctional aspartokinase I/homoserine dehydrogenase
>2801..3730:STM0003	+	thrB	<a href="#">STM0003</a>	COG0083E	homoserine kinase: <a href="#">sequence</a>
>3734..5020:STM0004	+	thrC	<a href="#">STM0004</a>	COG0498E	threonine synthase: <a href="#">sequence</a>
>5114..5887:STM0005	-	-	<a href="#">STM0005</a>	COG30225	hypothetical protein: <a href="#">sequence</a>
>5966..7396:STM0006	-	yaaJ	<a href="#">STM0006</a>	COG1115E	alanine/glycine transport protein: <a href="#">sequence</a>
>7665..8618:STM0007	+	talB	<a href="#">STM0007</a>	COG01766	transaldolase B: <a href="#">sequence</a>
>8729..9319:STM0008	+	mogA	<a href="#">STM0008</a>	COG0521H	molybdochetalase: <a href="#">sequence</a>
>9376..9942:STM0009	-	-	<a href="#">STM0009</a>	COG15845	hypothetical protein: <a href="#">sequence</a>
>11593..13509:STM0012	+	dnaK	<a href="#">STM0012</a>	COG04430	chaperone protein DnaK: <a href="#">sequence</a>
>13595..14734:STM0013	+	dnaJ	<a href="#">STM0013</a>	COG04840	chaperone protein DnaJ: <a href="#">sequence</a>
>15014..15961:STM0014	+	-	<a href="#">STM0014</a>	COG0583K	LysR family transcriptional regulator: <a href="#">sequence</a>
>16088..16432:STM0015	+	-	<a href="#">STM0015</a>	-	bacteriophage protein: <a href="#">sequence</a>
>16493..17026:STM0016	-	-	<a href="#">STM0016</a>	COG3926R	hypothetical protein: <a href="#">sequence</a>
>17043..17486:STM0017	-	-	<a href="#">STM0017</a>	COG3710K	hypothetical protein: <a href="#">sequence</a>
>17867..19966:STM0018	+	-	<a href="#">STM0018</a>	COG3325G	exochitinase: <a href="#">sequence</a>
>23335..24039:STM0020	+	-	<a href="#">STM0020</a>	COG0664T	cytoplasmic protein: <a href="#">sequence</a>
>24469..25011:STM0021	+	bcfA	<a href="#">STM0021</a>	COG3539NU	fimbrial subunit: <a href="#">sequence</a>
>25112..25798:STM0022	+	bcfB	<a href="#">STM0022</a>	COG3121NU	fimbrial chaperone: <a href="#">sequence</a>
>25803..28424:STM0023	+	bcfC	<a href="#">STM0023</a>	COG3188NU	fimbrial usher: <a href="#">sequence</a>
>28425..29432:STM0024	+	bcfD	<a href="#">STM0024</a>	COG3539NU	fimbrial subunit: <a href="#">sequence</a>
>29433..29978:STM0025	+	bcfE	<a href="#">STM0025</a>	COG3539NU	fimbrial subunit: <a href="#">sequence</a>
>29994..30512:STM0026	+	bcfF	<a href="#">STM0026</a>	COG3539NU	fimbrial subunit: <a href="#">sequence</a>
>30478..31209:STM0027	+	bcfG	<a href="#">STM0027</a>	COG3121NU	fimbrial chaperone: <a href="#">sequence</a>
>31274..32119:STM0028	+	bcfH	<a href="#">STM0028</a>	COG16510	thiol-disulfide isomerase: <a href="#">sequence</a>
>32545..32994:STM0029	-	-	<a href="#">STM0029</a>	COG3710K	transcriptional regulator: <a href="#">sequence</a>
>46190..47356:STM0039	+	nhaA	<a href="#">STM0039</a>	COG3004P	Na(+)/H(+) antiporter NhaA: <a href="#">sequence</a>
>47418..48317:STM0040	+	nhaR	<a href="#">STM0040</a>	COG0583K	transcriptional activator NhaR: <a href="#">sequence</a>
>52280..52543:STM0043	-	rpsT	<a href="#">STM0043</a>	COG0268J	30S ribosomal protein S20: <a href="#">sequence</a>
>52649..52864:STM0044	+	yaaY	<a href="#">STM0044</a>	-	cytoplasmic protein: <a href="#">sequence</a>

**Fig 2: core genes output**

Sequences of the core genes retrieved from NC\_003197 genome through hyperlinked petals of the flower plot. The list of core genes identified from the PanGeT analysis is displayed with their genome location, strand orientation, Gene\_ID, synonym, COG functional class, annotation and link to the corresponding sequence, as shown above. The synonyms were given with hyperlinks to the KEGG GENES database to get more information.





Database: KEGG - Search term: STM0003

## KEGG GENES

stm:STM0003

thrB; homoserine kinase (EC:2.7.1.39); K00872 homoserine kinase [EC:2.7.1.39]

Fig 5: KEGG Genes database entry for 'STM0003' gene entry.

Gene Name	NC_003197	NC_003198	NC_011094	NC_000905	NC_011147	NC_010067	NC_011149	NC_020307	NC_010102	NC_011205	NC_021818	NC_011080	NC_011274	NC_021820
16763398	0.92	STY0008		0.97	SC0008	0.93	SARI_02083		0.97	SPAB_00009		0.92	SNL1254_A0008	
0.97	SEHA_C0008	0.97	SESA_A0008	0.92	SSPA0007	0.97	SeAg_B0008	0.97	SeD_A0008	0.97	SeD_A0008	0.97	SG0008	0.97
0.97	SEH0007	0.97	SPC_0008	0.90	-	0.97	CFSAN001992_11000	0.97	CFSAN002050_06515	0.97	CFSAN002050_06515	1.00	SE451236_06050	1.00
0.93	SEEB0189_19350													
16763400	0.99	STY0010		0.99	SC0010	0.00	-	0.99	SPAB_00011	0.99	SeD_A0010	0.99	SNL1254_A0010	
0.99	SeHA_C0010	0.99	SESA_A0010	0.99	SSPA0009	0.99	SeAg_B0010	0.99	SeD_A0010	0.99	SeD_A0010	0.99	SG0011	0.99
0.99	SEH0009	0.99	SPC_0010	0.91	SBG_0010	0.99	CFSAN001992_10990	0.99	CFSAN002050_06525	0.99	CFSAN002050_06525	1.00	SE451236_06060	1.00
0.98	SEEB0189_19340													
16763401	0.99	STY0011		0.99	SC0011	0.00	-	0.99	SPAB_00012	0.99	SeD_A0011	0.99	SNL1254_A0011	
0.99	SeHA_C0011	0.99	SESA_A0011	0.99	SSPA0010	0.99	SeAg_B0011	0.99	SeD_A0011	0.99	SeD_A0011	0.99	SG0011	0.99
0.99	SEH0010	0.99	SPC_0011	0.00	-	0.98	CFSAN001992_10985	0.99	CFSAN002050_06530	0.99	CFSAN002050_06530	1.00	SE451236_06065	1.00
0.99	SEEB0189_19335													
16763409	0.98	-	0.98	SC0018	0.27	-	0.99	SPAB_00023	0.99	SeD_A0020	0.99	SNL1254_A0021		
0.99	SeHA_C0021	0.98	SESA_A0021	0.98	SSPA0017	0.99	SeAg_B0022	0.99	SeD_A0020	0.99	SeD_A0020	0.99	SG0034	0.99
0.99	SEH0018	0.99	SPC_0021	0.85	SBG_0020	0.98	CFSAN001992_10940	0.98	CFSAN002050_06570	0.98	CFSAN002050_06570	1.00	SE451236_06105	1.00
0.99	SEEB0189_19290													
39546288	0.99	STY0034		0.99	SC0027	0.92	SARI_02963	1.00	SPAB_00037	0.99	SeD_A0031	0.99	SNL1254_A0031	
1.00	SEH0028	0.99	SESA_A0031	0.99	SSPA0026	0.99	-	1.00	-	1.00	-	1.00	SG0031	1.00
1.00	SEEB0189_19240													
16763419	0.99	STY0035		0.61	SC0028	0.88	SARI_02962	0.61	SPAB_00038	0.99	SeD_A0031	0.99	SNL1254_A0032	
0.61	SeHA_C0033	0.96	SESA_A0032	0.99	SSPA0027	0.99	SeAg_B0033	0.61	SeD_A0031	0.99	SeD_A0031	0.99	SG0032	0.99
0.99	SEH0029	0.61	SPC_0032	0.00	-	0.00	-	0.61	CFSAN002050_06625	0.61	CFSAN002050_06625	1.00	SE451236_06160	1.00
0.99	SEEB0189_19235													
16763420	0.98	STY0036		0.00	-	0.89	SARI_02961	0.00	-	0.99	SeD_A0032	0.99	SNL1254_A0033	
0.00	-	0.98	SESA_A0033	0.00	-	0.98	SeAg_B0034	0.00	-	0.99	SeD_A0032	0.99	SG0033	0.99
0.99	SEH0030	0.00	-	0.00	-	0.00	-	0.00	-	0.99	SeD_A0032	0.99	SE451236_06165	1.00
0.99	SEEB0189_19230													
16763421	0.95	-	0.96	SSPA0029	0.85	SARI_02960	0.00	-	0.99	SeD_A0033	0.99	SNL1254_A0034		
0.00	-	0.96	SESA_A0034	0.00	-	0.97	SeAg_B0035	0.00	-	0.99	SeD_A0033	0.99	SG0034	0.99
0.99	SEH0031	0.00	-	0.00	-	0.00	-	0.00	-	0.99	SeD_A0033	0.99	SE451236_06170	1.00
0.99	SEEB0189_19225													
16763422	0.99	STY0039		0.00	-	0.90	SARI_02959	0.00	-	0.99	SeD_A0035	0.99	SNL1254_A0036	
0.00	-	0.99	SESA_A0036	0.00	-	0.99	SeAg_B0037	0.00	-	0.99	SeD_A0035	0.99	SG0035	0.99
0.99	SEH0032	0.00	-	0.00	-	0.00	-	0.00	-	0.99	SeD_A0035	0.99	SE451236_06175	1.00
0.99	SEEB0189_19220													
16763423	0.99	STY0040		0.00	-	0.93	SARI_02958	0.00	-	0.99	SeD_A0036	0.99	SNL1254_A0037	
0.00	-	0.99	SESA_A0037	0.00	-	0.99	SeAg_B0038	0.00	-	0.99	SeD_A0036	0.99	SG0036	0.99
0.99	SEH0033	0.00	-	0.81	SBG_0033	0.00	-	0.00	-	1.00	SeD_A0036	1.00	SE451236_05180	1.00

Fig 6 : List of Dispensable genes output

The list of 'Dispensable genes' identified by PanGeT is shown above with their respective homology scores in the reference genomes. The above list may be imported into MS-EXCEL to get more information.

## 7. Pathogenic and non pathogenic strains comparison

If the user want to find the genes which are conserved only in few pathogens and absent in non pathogens, user can run the program '*Find\_genes\_conserved\_few\_while\_comparing\_many.pl*'. Instructions for running is given inside the program file.