



NCBI Pathogen Detection

From next generation sequencing reads to clusters of related bacterial pathogens based on SNPs

<https://www.ncbi.nlm.nih.gov/pathogens/>

National Center for Biotechnology Information • National Library of Medicine • National Institutes of Health • Department of Health and Human Services

Scope

Bacterial pathogens that originate from a food source are a serious concern in the US, with estimates by CDC that each year roughly 1 in 6 Americans (or 48 million people) get sick, 128,000 are hospitalized, and 3,000 die of foodborne diseases. Outbreaks are defined where two or more people become ill when ingesting the same food. If the bacterial pathogens that cause the illness are genetically related, it becomes easier to trace the source of the illness when samples from both patients and food source are compared. US Federal agencies involved in public health for foodborne illnesses have agreed to move forward on this by using whole genome sequencing. A pilot project was started in 2013 for *Listeria* with all federal agencies responsible for food safety (CDC, FDA, and USDA) sequencing all *Listeria* isolates, whether in clinical patients or identified during inspection of food products or processing facilities, and depositing the data in the public archives at NCBI. This project demonstrated that whole genome sequencing could be used to increase the number of clusters detected, with a decrease in median cluster size, and with more outbreaks solved [1]. The goal is to extend this to all 90,000 foodborne bacterial pathogens (*Campylobacter*, *Listeria*, *Escherichia coli* and *Shigella* (STECs), and *Salmonella*) that are collected in the US every year and sequenced in real time by the end of 2018.

The NCBI Pathogen Detection pipeline takes the incoming sequencing data and assembles, annotates, and clusters the genomes together to facilitate the analysis of genetically related strains to aid outbreak and traceback investigations. Phylogenetic trees are constructed for each SNP cluster using maximum compatibility [2]. Currently clusters are calculated daily for each organism group if new data are submitted and released publicly. The scope of the project has now expanded to include non-foodborne pathogens, especially those that are becoming increasingly antimicrobial resistant (AMR). The workflow uses a reference set of acquired resistance genes/proteins to report for each isolate the AMR genes that are encoded in the genome sequence. This database of antimicrobial resistance pathogens is part of the National Action Plan for Combating Antimicrobial-Resistant Bacteria by providing a resource for the research community on resistant organisms that have genomic sequence data [3].

Data Access

Several venues are available for accessing data from Pathogen Detection Pipeline:

1. The Pathogen Detection homepage (A) provides links both to the Isolates Browser and FTP results (B).
2. The FTP site (C) offers access to all of the phylogenetic trees and metadata tables. The tree files can be opened in NCBI Genome Workbench tool and explored locally. See the readme file for more information about the FTP.
3. The beta version of the Isolates Browser (D, linked off the Pathogen Detection home page) allows search and browse. In its display, each row is an assembled genome (either from SRA data by the Pathogen Detection Pipeline or genomes deposited in GenBank).

The table (E) on the main page specifies specific searches subcategorized by organism group and for the subset of isolates that have been newly deposited since the last cluster calculation.

The screenshot shows the NCBI Pathogen Detection homepage. At the top, there is a navigation bar with the NIH logo, 'U.S. National Library of Medicine', 'NCBI National Center for Biotechnology Information', and a 'Log in' link. The main heading is 'Pathogen Detection BETA' with the URL <https://www.ncbi.nlm.nih.gov/pathogens/>. Below the heading, a brief description states: 'NCBI Pathogen Detection integrates bacterial pathogen genomic sequences originating in food, environmental sources, and patients. It quickly clusters and identifies related sequences to uncover potential food contamination sources, helping public health scientists investigate foodborne disease outbreaks.' There are several callout boxes labeled A through E. Box A points to the 'Pathogen Detection' link in the navigation bar. Box B points to the 'Isolates Browser' link in the 'Data Resources' section. Box C points to the 'Submit' section, specifically the 'How to submit antibiotic resistance phenotypes' link. Box D points to the 'Find isolates now!' link. Box E points to the 'Explore the Data' table. Below the homepage, the 'Isolates Browser' interface is shown, featuring a search bar and a table of isolates. The table has columns for #, Organism Group, Strain, Serova, Isolate, Create D, Location, Isolation S, Isolation typ, Host, SNP cluster, Min-s, Min-d, BioSample, Assembly, K-mer group, AST phenotype, and AMR genotypes. Two rows are visible, both for *Elizabethkingia anophelis*.

Species	New Isolates	Total Isolates
Salmonella enterica	53	66,525
E.coli and Shigella	134	27,957
Listeria monocytogenes	18	13,7
Campylobacter jejuni	8	8,983
Mycobacterium tuberculosis	1	5,325
Klebsiella pneumoniae	13	3,304
Acinetobacter baumannii	0	3,059
Neisseria	0	2,714
Pseudomonas aeruginosa	0	2,399
Enterobacter	4	830
Vibrio parahaemolyticus	0	830

Usage Examples

A. Isolation Browser: Basic navigation

From the Pathogen Detection homepage's table (p1, **E**), you can click on the number in the 'New Isolates' column for an organism group to locate new clusters for that group. You can search directly in the Isolate Browser using structured terms. For *E. coli* and *Shigella* group, the structured query terms are: `taxgroup_name:"E.coli and Shigella" AND new:1` (<http://bit.ly/ncbi-pathogen-1>). The search results will change every time new data is deposited. You can toggle open the Filters section (**A**, more on this in p 3) to narrow down isolates displayed. In the Isolate Browser, the **Min-same** and **Min-diff** columns (insert) report on the SNP differences for each isolate (**B**) based on Isolation type (clinical vs. environmental). Sorting on the **Min-diff** column (**C**) in ascending order brings isolates with the shortest SNP distance between a clinical and environmental samples to the top. This similarity is critically important for public health labs.

Health > Pathogen Detection > Isolates Browser

taxgroup_name:"E.coli and Shigella" Search

E.coli and Shigella PDS000000987.32 (631)
E.coli and Shigella PDS000000952.131 (543)
Show more clusters (up to 100)

Min-same Min-diff **A**

E.coli and Shigella Filters Expand All Columns Download

Page 48 of 1,403 20 View 941 - 960 of 28,051

#	Organism Group	Strain	Serova	Isolate	Create	Location	Isolation	Isolation type	Host	SNP cluster	Min-s	Min-d	BioSampl	Assembly	K-mer group	AST phenotypes	AMR genotypes
941	E.coli and Shigella	CDPHFI F15M01 2	O157:H7	PDT000191571	2017-03-07	USA:CA	Feces	environment		PDS00000095	2	3	SAMN0646		PG0000000004		blaEC
946	E.coli and Shigella	PN		PDT000202811	2017-04-14	USA		clinical		PDS00000431	2	3	SAMN0675		PG0000000004		blaEC
947	E.coli and Shigella	WAPHL A00003		PDT000205116	2017-04-25	USA:NV	Chocolate marquisse dessert	environment		PDS00000095	8	3	SAMN0435		PG0000000004		blaEC

Extra rows removed for brevity.

B. Use Case One: To identify potential clusters of interest for public health

We will use an environmental isolate PDT000205116.1 (row 947) in this exercise, which is 3 SNPs away from a clinical isolate (**D**). Clicking the isolate identifier opens the SNP Tree Viewer (**E**), which zooms to the particular branch of the phylogenetic tree where that isolate is located. Our environmental isolate is clustered with several clinical isolates (**F**). We can examine the details of a branch in the SNP distance tree by clicking a node (such as **G**).

Distance between isolates in the cluster: minimum=0 SNPs, maximum=131 SNPs, average=73.01 SNPs

Tools

Success

Distance between selected 12 isolates: minimum=0 SNPs, maximum=11 SNPs, average=4.45 SNPs **H**

Success

Nodes 793(15 selected) View port at (0,1485) of 123

Filters Columns Selected: 12 Download

With a branch selected, the summary at the top (**H**) updates to show the minimum, maximum, and avg. SNP distance for selected isolates. In this case, the avg. SNP distance of 4.45 indicates a very closely related set of 12 isolates. We can click "Selected ..." button (**I**) to bring the selected isolates to the top in the datatable below the tree to analyze the metadata associated with these isolates.

B. Use Case One: To identify potential clusters of interest for public health (cont.)

In the sorted table (A), the clinical isolates (boxed) are from 2015, while the environmental isolates are from raw beef in 2016 and desserts in 2017. Any public health lab reviewing this data would need to make a determination based on additional metadata if further investigation is warranted.

Analysis Results FTP

Filters Columns Selected: 12 x Download

Page 1 of 28 20 View 1 - 20 of 543

#	Organism Gro	Strain	Serov	Isolate	Create	Location	Isolation	Isolation ty	Host	SNP cluster	Min-	Min-	BioSampl	Assembly	K-mer group	AST phen	AMR gen
1	E.coli and Shigella	WAPHL A00004		PDT00020511E	2017-04-25	USA:NV	Chocolate marquise dessert	environmen		PDS00000095	10	6	SAMN0435		PDG00000000		blaEC
2	E.coli and Shigella	WAPHL A00003		PDT00020511E	2017-04-25	USA:NV	Chocolate marquise dessert	environmen		PDS00000095	8	3	SAMN0435		PDG00000000		blaEC
3	E.coli and Shigella	PNUSA	E. coli O157:H	PDT00010109	2015-12-31	USA	Stool	clinical		PDS00000095	1	3	SAMN0435		PDG00000000		blaEC
11	E.coli and Shigella	PNUSA	E. coli O157:H	PDT00010110	2015-12-31	USA	Stool	clinical		PDS00000095	1	4	SAMN0435		PDG00000000		blaEC
12	E.coli and Shigella	FSIS15	O157	PDT00012741	2016-04-26	USA:ID	Product-Raw-Intact-Beef	environmen		PDS00000095	8	5	SAMN0490		PDG00000000		blaEC

Extra rows removed for brevity.

C. Use Case Two: To identify isolates encoding antimicrobial resistance genes

The Filter control (B) allows filtering by location, isolation source, sample collecting lab, host, whether the isolate has antimicrobial resistance genes or has antibiotic susceptibility test phenotypic data, a date range control, and by the scientific name of the isolate. In this example we set the organism group control above the Filter to 'Klebsiella pneumoniae,' check 'has AMR geno-

types,' and set date range control to 2016 (C, D, E). This setting combination selects all *K. pneumoniae* released in 2016 that encode an antimicrobial resistance gene, such as chromosomal encoded ampC beta lactamases. Click the up arrow (F) to close the Filter control.

Find one or more isolates ... Search

Klebsiella pneumoniae

Filters 2 x

Location: ATHENS (22), CHINA (46), DURBAN (20), GREECE (25), NETHERLANDS (26), OXFORD (176), SOUTH AFRICA (20), UNITED KINGDOM (417), USA (495), VIRGINIA (94)

Property: has AMR genotypes (1,398), has AST phenotypes (259)

Target Creation: Mon Dec 28 2015 09:52:24 GMT-0500 (Eastern Standard Time) -> Fri Dec 30 2016 02:54:00 GMT-0500 (Eastern Standard Time)

Extra filters removed for brevity.

You can customized the columns shown in the metadata table using the Column control (G). The example has 'Organism group' and 'serovar' column removed (H) since we are dealing only with *Klebsiella* isolates (through filter), and 'Collected by', 'Length' plus 'Contigs' add (I) to provide information on isolate collectors and the assembly quality. Clicking the "AST Phenotypes" column header sorts to the top those isolates tested for antibiotic susceptibility (J), many of which are multi-drug resistant. Clicking "Expand All" button (K) expands the rows in the table to reveal details (p4, A).

Klebsiella pneumoniae

Filters 2 Expand All Columns Download

Page 1 of 70 20 View 1 - 20 of 1,390

#	Str	Isolate	Creat	Locati	Isolati	Isolation	Host	SNP clus	Mir	Mir	BioSar	Asse	K-mer gro	AST phenotyp	AMR genoty	Length	Contigs	Collected
1		PDT000130	2016-05-16	USA: Minnes	Urine	clinical	Homo sapiens	PDS000000	38	38	SAMN0		PDG000000	Resistant (16) Intermediate (1) Susceptible (5) Other (7) Expand All	aac(3) aac(6)-lb-cr aph(3)-lb Show all 16 genes	5,662,625	127	MN State Health Department
2		PDT000130	2016-05-16	USA: Indiana	Sputum	clinical	Homo sapiens	PDS000000	8	n/a	SAMN0		PDG000000	Resistant (15) Intermediate (2) Susceptible (5) Other (7) Expand All	blaCTX-M-15 blaOXA-181 blaSHV-11 Show all 6 genes	5,537,627	144	IN State Health Department
3		PDT000130	2016-05-16	USA: Illinois	Urine	clinical	Homo sapiens	PDS000000	40	n/a	SAMN0		PDG000000	Resistant (20) Intermediate (1) Susceptible (1) Other (7) Expand All	aac(6)-lb-cr blaCTX-M-15 blaOXA-1 Show all 10 genes	5,559,588	40	IL State Health Department

Column controls: Organism Group, Strain, Serovar, Isolate, Create Date, Location, Isolation Source, Isolation type, Host, SNP cluster, Level, NSI, Length, Contigs, Method, BioProject, Collected by, Collection Date, Host Disease, Lat/Lon.

AST Phenotyp: Resistant (16), Intermediate (1), Susceptible (5), Other (7), Expand All

AMR genoty: aac(3), aac(6)-lb-cr, aph(3)-lb, Show all 16 genes

Length: 5,662,625

Contigs: 127

Collected: MN State Health Department

Expand All button (K)

Column control (G)

Removed columns (H): Organism Group, Strain, Serovar

Added columns (I): Length, Contigs, Collected by

Sorted by (J): AST phenotyp

C. Use Case Two: To identify isolates encoding antimicrobial resistance genes (cont.)

The expanded example shown (A) is a multi-drug resistant *Klebsiella* that has been phenotypically tested and shown to be resistant to multiple drugs. Additionally, the genome sequence shows that this encodes multiple antimicrobial genes, some of which may confer resistance to those drugs. Note, the presence of an antimicrobial resistance genes does not guarantee resistance to a specific drug in a clinical setting and care must be taken in interpretation. This information is provided to aid researchers and other experts in resistance and bacterial pathogens.

#	Strain	Isolate	Creat	Locati	Isolati	Isolator	Host	SNP clus	Mir	Mir	BioSar	Asse	K-mer gro	AST phenotyp	AMR genoty	Length
1	PDT00013	2016-05-16	USA: Minnes	Urine	clinical	Homo sapiens	PDS00000	38	38	SAMNO		PDG00000	Resistant (16) Intermediate (1) Susceptible (5) Other (7) Expand All	aac(3) aac(6)-Ib-cr aph(3)-Ib Show all 16 genes	5,662,625	

AST phenotypes

Resistant (16)
amikacin
amoxicillin-clavulanic acid
ampicillin-sulbactam
ampicillin
aztreonam
cefazolin
cefepime
cefotaxime-clavulanic acid
cefotaxime
ceftriaxone
ciprofloxacin
ertapenem
gentamicin
levofloxacin
tetracycline
tobramycin

Intermediate (1)
imipenem

Susceptible (5)
cefotaxim
doripenem
meropenem
tigecycline
trimethoprim-sulfamethoxazole

Other (7)
cefazidime-clavulanic acid
cefazidime
chloramphenicol
colistin
minocycline
piperacillin-tazobactam
polymyxin B

AMR genotypes

aac(3)
aac(6)-Ib-cr
aph(3)-Ib
aph(6)-Ib
blaCTX-M-15
blaOXA-1
blaOXA-48
blaSHV-28
blaTEM-1
fosA
mph(A)
oxqA
oxqB
sul2
tet(A)

Show fewer genes

D. Use Case Three: To find isolates that encode specific resistance genes

The colistin resistance gene *mcr* has been found to occur on mobile elements and encodes resistance to one of the last line drugs that can be used. Isolates that also encode carbapenemases are of significant interest to the community. We can use structured search terms to find isolates that encode *mcr* and any allele of the KPC family of beta lactamases:

[AMR_genotypes:mcr*](#) AND [AMR_genotypes:blaKPC*](#)

Currently, only three isolates that encode both genes are in the database (B). The 'Collected by' and 'Collection date' columns (C, D) show that they were collected by different groups in 2014 and 2016. The genomes of all three isolates are available in GenBank. The expanded view for the first isolate shows that it encodes *mcr*-1.2.

#	Isolate	Creat	Locati	Isolati	Isolator	Host	SNP clus	Mir	Mir	Assembl	K-mer grc	AMR genot	Collected by	Collection date	Species
1	PDT00013	2016-07-06	Italy: Pisa	Rectal swab	clinical	Homo sapiens	PDS00000	25	26	GCA_001	PDG00000	aac(3)-II aac(6)-Ib aadA2 aadA5 blaKPC-3 blaOXA blaSHV-11 blaTEM-1 dfrA17 fosA mcr-1.2 oxqA oxqB qacEdelta1 sul1 Show fewer genes	Carlo Tascini	2014	<i>Klebsiella pneumoniae</i>
2	PDT00020	2017-04-18	Brazil: Vitoria	urine	clinical	Homo sapiens		n/a	n/a	GCA_002	PDG00000	aac(3)-IIa aac(6)-Ib aac(6)-Ib-cr Show all 18 genes	Santa Casa da Misericord de Vitoria	2016-09	<i>Klebsiella pneumoniae</i>
3	PDT00018	2017-02-23	Brazil: Porto Alegre	rectal swab	clinical	Homo sapiens		n/a	n/a	GCA_002	PDG00000	aph(3)-Ib aph(3)-IIa aph(6)-Ib Show all 14 genes	LABRESIS	2014	<i>Escherichia coli</i>

<http://bit.ly/ncbi-pathogen-3>

Extra columns removed for brevity

Future Work

The NCBI Pathogen Detection team is working on a number of enhancements to the interface including: 1) genomic locations of resistance genes, 2) ability to search using specific MIC values.

Submissions and Contact

If you wish to submit data to the NCBI Pathogen Detection pipeline, please review the submission page:

<https://www.ncbi.nlm.nih.gov/pathogens/submit/>

You can address questions on this resource to the NCBI Pathogen Detection team directly:

pd-help@ncbi.nlm.nih.gov

References

- Allard et al. (2016). The Practical value of Food Pathogen Traceability through BUILDING a Whole-Genome Sequencing Network and database. *J Clin Microbiol.* 2016 Mar 23. pii: JCM.00081-16.
- Cherry J. (2017). A practical exact maximum compatibility algorithm for reconstruction of recent evolutionary history. *BMC Bioinformatics.* 2017 Feb 23;18(1):127. doi: 10.1186/s12859-017-1520-4.
- Jackson et al. (2016). Implementation of Nationwide Real-time Whole-genome Sequencing to Enhance Listeriosis Outbreak Detection and Investigation. *Clin Infect Dis.* 2016 Apr 18. pii: ciw242.
- National Action Plan for Combating Antimicrobial-Resistant Bacteria. <https://www.whitehouse.gov/blog/2015/03/27/our-plan-combat-and-prevent-antibiotic-resistant-bacteria>
https://obamawhitehouse.archives.gov/sites/default/files/docs/national_action_plan_for_combating_antibiotic-resistant_bacteria.pdf