

**CLIA-compliant Validation of
Whole Genome Sequencing (WGS)
for Clinical Microbiological Applications:
Experience of one State Public Health Reference Laboratory**

Varvara Kozyreva, PhD

Genotyping Unit Chief

Microbial Diseases Laboratory Program

California Department of Public Health

✉ varvara.kozyreva@cdph.ca.gov

ASM NGS 2018

September 24th 2018

Disclosure

I have no conflict of interest relevant to this presentation.



Clinical Laboratory Improvement Amendments (CLIA)- federal regulatory standards ensuring quality of clinical laboratory testing performed on samples from humans for the diagnostic purposes.



CLIA is implemented and enforced by the Centers for Medicare & Medicaid Services (CMS).

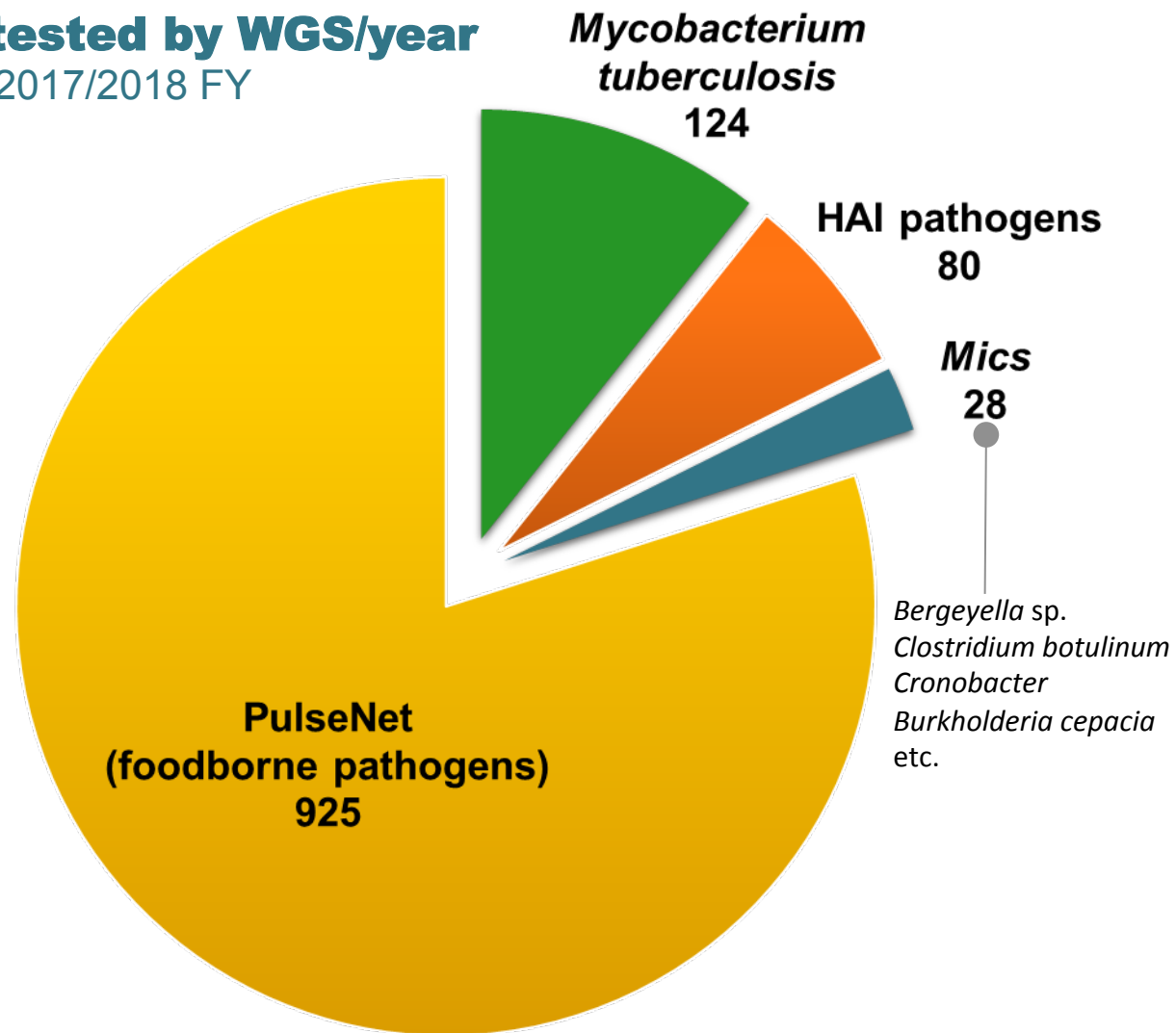
WGS Workload in the Microbial Diseases Laboratory of California Department of Public Health

Isolates tested by WGS/year
2017/2018 FY

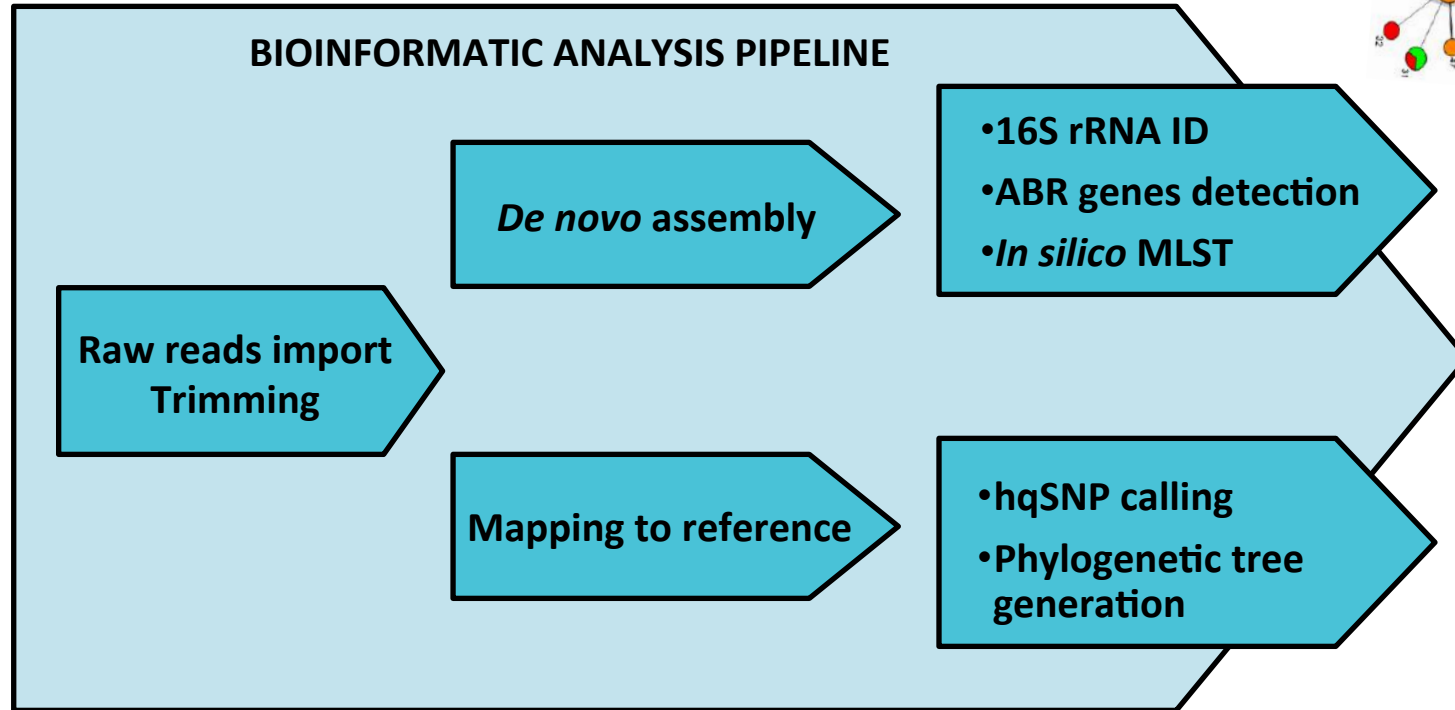
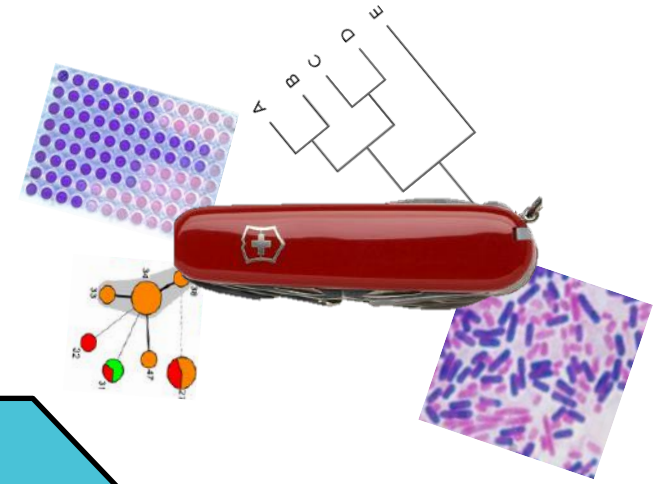


MDL

Microbial Diseases Laboratory
Pathogen Experts Keeping California Safe



Simple public health lab WGS toolkit



Report to
state epidemiologists
& outside clients

In-house sequencing:
MiSeq sequencer
Nextera XT library prep
500- or -600-cycle seq kits

... but dangerous if used incorrectly



WGS

- a laboratory-developed test (LDT)
- not approved/cleared by FDA

WGS validation in microbiological
PHL settings

Challenges:

- **Complexity** of technology with extensive computational analysis requirements
- Insufficient guidelines for validation of WGS for **microbiological applications**
 - Publications addressing NGS standardization and validation are mainly focused on cancer diagnostics, inherited human diseases, and pathology



- **CDC Nex-StoCT Workgroups I & II:** two publications on standardization of clinical NGS- *validation, QC, ref materials, PT* (Gargis et al. Nat Biotechnol 2012 & 2015)



- **College of American Pathologists (CAP):** *validation, QC, QA, ref materials* (CAP NGS Inspection Checklist, 2014; Aziz et al. Arch Pathol Lab Med, 2015; Jennings et al. J Mol Diag, 2017; Roy et al. J Mol Diag, 2018)



- **Clinical and Laboratory Standards Institute (CLSI):** MM09 - Nucleic Acid Sequencing Methods in Diagnostic Laboratory Medicine (2014) – *validation, QC, QA, PT, competency, ref materials*



- **American College for Medical Genetics and Genomics (ACMG):** Clinical laboratory standards for NGS- *Validation, QC, QA, ref material, PT.* (Rehm et al. Genet Med. 2013)



- Validation of NGS in **Microbial Forensics** (Budowle et al. Investigative Genetics, 2014)



- **Food and Drug Administration (FDA):** Draft Guidance for Infectious Disease NGS-Based Diagnostic Devices: ID, antimicrobial resistance, virulence (2016)- *for Industry and FDA staff*



- Validation of **Metagenomic NGS** for pathogen detection (Schlaberg et al. Arch Pathol Lab Med, 2017)

WGS

- a laboratory-developed test (LDT)
- not approved/cleared by FDA

WGS validation in microbiological PHL settings

Challenges:

- **Complexity** of technology with extensive computational analysis requirements
- Insufficient **guidelines** for validation of WGS for **microbiological applications**
 - Publications addressing NGS standardization and validation are mainly focused on cancer diagnostics, inherited human diseases, and pathology
- **No examples** of the validation of WGS-based tests used in PHL¹
- Poorly established **performance parameters**
- Lack of concept for CLIA **validation of epidemiological tests**
- Lack of established **reference materials**
 - ✧ Well-characterized
 - ✧ Relevant to public health
 - ✧ Diverse

Possible Reference Materials sources:

- **CDC PulseNet**- foodborne (PulseNet) organisms only (Isolates)
- **NIST**- National Institute of Standards and Technology – *S. enterica*, *S. aureus*, *P. aeruginosa*, *Clostridium sporogenes*² (DNA, not isolates)²
- **FDA-ARGOS**- Database for Reference Grade Microbial Sequences (*in silico* validation)³

¹ Examples of Validation of Metagenomic NGS for pathogen detection (Schlaberg et al. Arch Pathol Lab Med, 2017)

² <https://www.nist.gov/programs/projects/microbial-genomic-measurements>

³ <https://www.ncbi.nlm.nih.gov/bioproject/231221>

Main objectives of WGS validation in microbiological PHL:

- Establish the **performance specifications** of WGS used for public health applications
- Create a **validation set of microorganisms** which can be used for future validations of WGS platforms
- Determine the **optimal conditions** that will generate reliable, reproducible, and accurate results for the intended application
- Develop a **set of Quality Assurance (QA) & Quality Control (QC) measures** to ensure high quality and consistency of routine testing
- Establish **reporting language**

VALIDATION SET

10- *Enterobacteriaceae*
5- Gram-positive cocci isolates
5- Gram-negative non-fermenting bacterial isolates
9- *Mycobacterium tuberculosis*
5- representatives of miscellaneous species

Validated a test panel of bacterial isolates for WGS

- **Source:**
 - ATCC strains
 - Strains sequenced by the CDC
- **34 bacterial isolates**
- **19 species**
- **Genome sizes** 1.8 to 6.7 Mb
- Wide range of **GC content** 32.1-66.1%
- Can be used as **reference materials for future validation** of new sequencing platforms and tests, QC procedures, proficiency testing, and the independent assessment of test performance

Acinetobacter baumannii
Aeromonas hydrophilia
Bacteroides fragilis
Corynebacterium jeikeium
Enterobacter cloacae
Enterococcus faecalis
Escherichia coli
Haemophilus influenzae
Legionella pneumophila
Moraxella catarrhalis
Mycobacterium tuberculosis
Neisseria gonorrhoeae
Pseudomonas aeruginosa
Salmonella enterica ser Adelaide
Salmonella enterica ser Enteritidis
Salmonella enterica ser Infantis
Salmonella enterica ser Typhimurium
Salmonella enterica ser Worthington
Staphylococcus aureus
Staphylococcus epidermidis
Staphylococcus saprophyticus
Stenotrophomonas maltophilia
Streptococcus pneumoniae



- The **platform validation is not sufficient** to demonstrate that sequencing assay is able to detect the [disease-associated] sequence variation.
- A validation tailored **specifically to the application** (the specific implementation of the **bioinformatics pipeline**) should be undertaken.¹

- NGS tests are better approached through **method-based validation** than analyte-specific validation, if it is **impractical/impossible to obtain positive samples** of all targeted mutations.
- The **ability of a laboratory to validate**, perform, and interpret certain **types of genomic sequencing analysis** is not wholly dependent on the specific variants being tested, but rather on **ability of the lab to conduct that type of sequencing analysis in its totality**.²



It is **not** possible to validate all theoretically possible variants that can occur, therefore it is necessary to use a **combination** of a '**methods-based**' and '**analyte-specific**' validation approach for determining a test's analytic performance.³

wet & dry

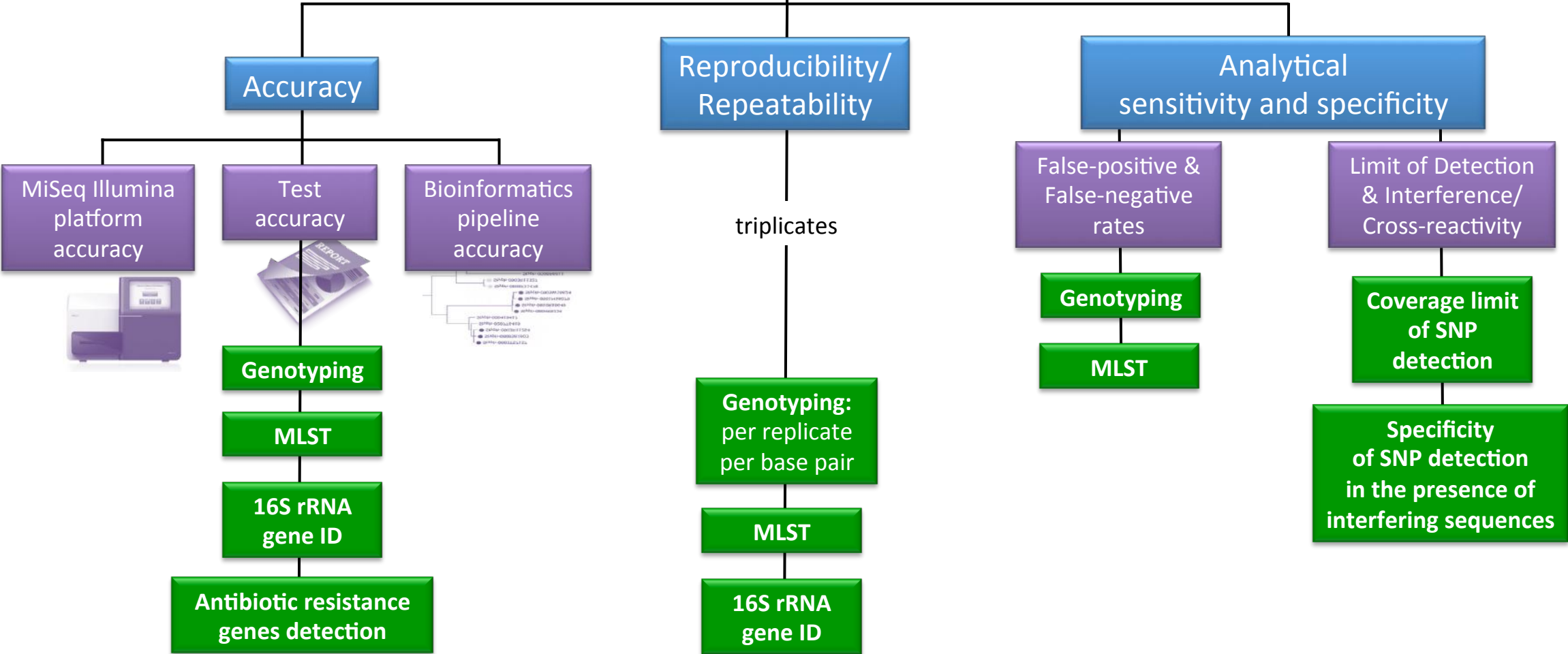
beginning-to-end

1. CLSI. MM09-A2, 2nd ed., 2014 CLSI, Wayne, PA
2. Schrijver I et al., J Mol Diagn. 2014 May;16(3):283-7.
3. Aziz N et al., Arch Pathol Lab Med. 2015 Apr;139(4):481-93.

10- Enterobacteriaceae
5- Gram-positive cocci isolates
5- Gram-negative non-fermenting bacterial isolates
9- *Mycobacterium tuberculosis*
5- representatives of miscellaneous species

Validation Set
34 bacterial isolates

WHOLE GENOME SEQUENCING VALIDATION IN PUBLIC HEALTH MICROBIOLOGY LAB SETTINGS

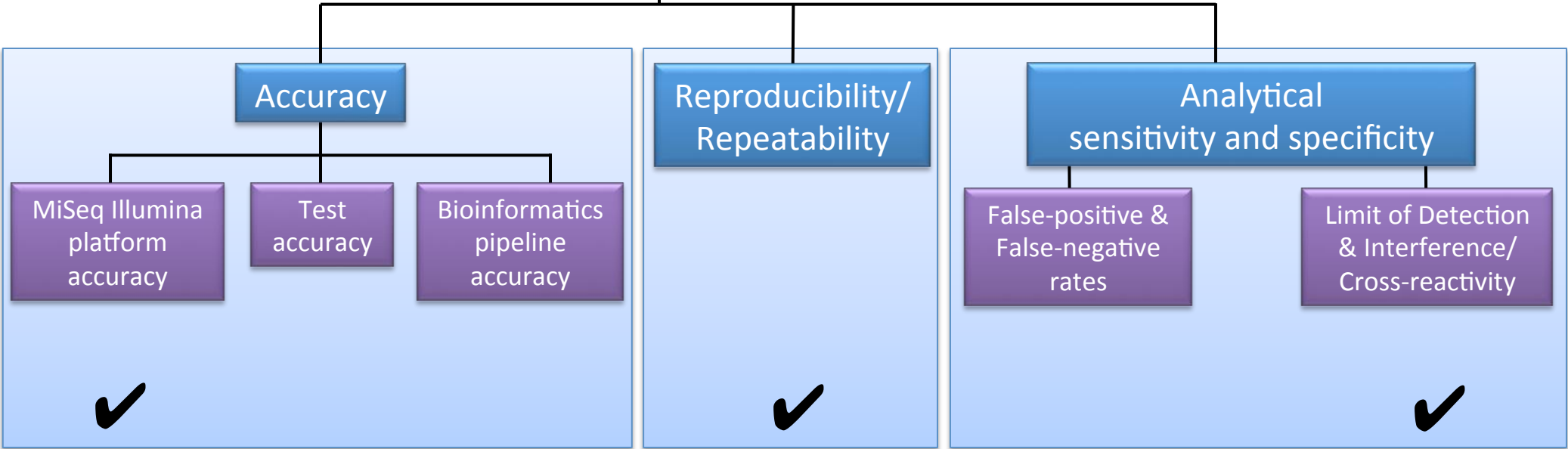


Base calls vs. WGS assays

- 10- *Enterobacteriaceae*
- 5- Gram-positive cocci isolates
- 5- Gram-negative non-fermenting bacterial isolates
- 9- *Mycobacterium tuberculosis*
- 5- representatives of miscellaneous species

Validation Set
34 bacterial isolates

WHOLE GENOME SEQUENCING VALIDATION IN PUBLIC HEALTH MICROBIOLOGY LAB SETTINGS



Correct base calls (SNPs)

Genotyping (topology)

MLST

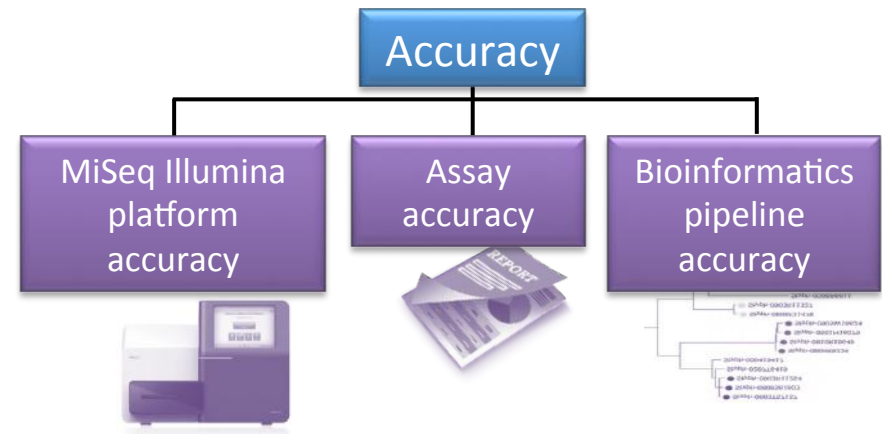
16S rRNA gene ID

Antibiotic resistance genes detection

- ✓
- ✓
- ✓
- ✓

- ✓
- ✓
- ✓

- ✓
- ✓

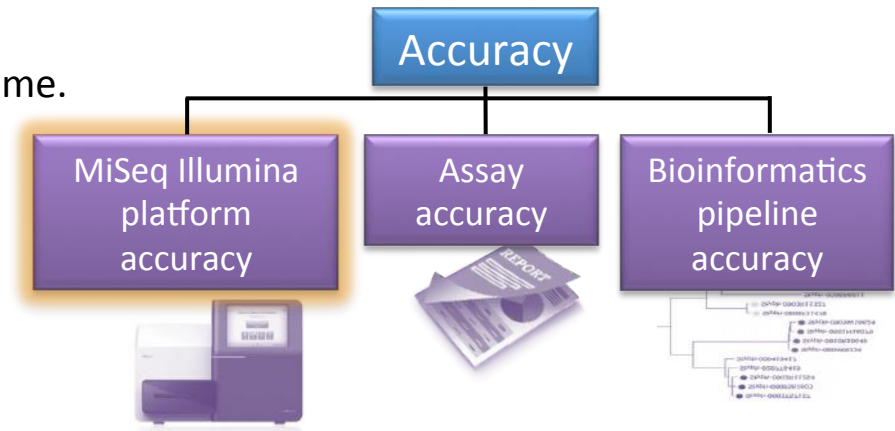


Accuracy of the platform -

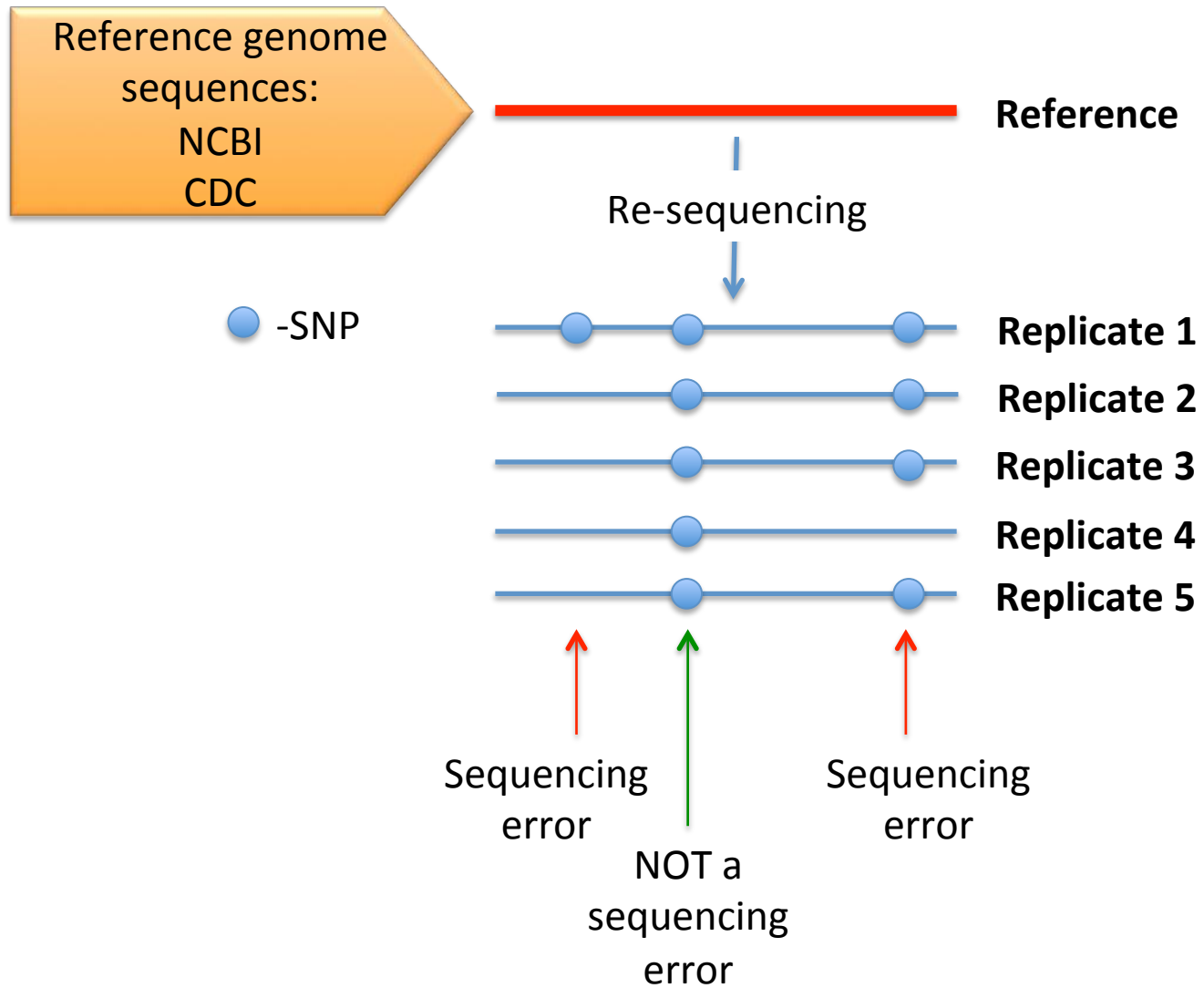
accuracy of the identification of individual base pairs (“base calling”) in the bacterial genome.

Determined by:

the proximity of agreement between base calling made by MiSeq sequencer (measured value) and NCBI/CDC reference complete genome sequence (the true value)

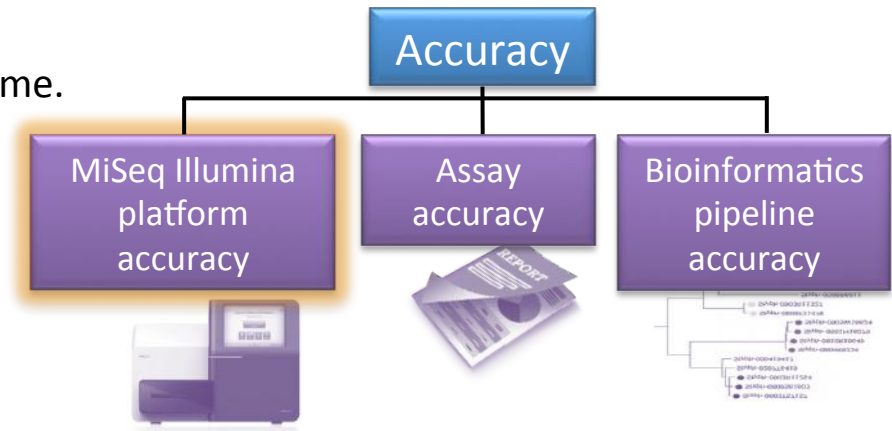


ID	Species	Total # of SNP difference with the reference
C1	<i>Escherichia coli</i>	5
C2	<i>Aeromonas hydrophilia</i>	1
C3	<i>Escherichia coli</i>	22
C4	<i>Enterobacter cloacae</i>	10
C5	<i>Staphylococcus aureus</i>	0
C6	<i>Salmonella ser Typhimurium</i>	12
C46	<i>Enterococcus faecalis</i>	3
C47	<i>Staphylococcus epidermidis</i>	184
C48	<i>Staphylococcus saprophyticus</i>	27
C51	<i>Stenotrophomonas maltophilia</i>	39
C52	<i>Legionella pneumophila</i>	2
C55	<i>Escherichia coli</i>	14
C72	<i>Escherichia coli</i> O121:H19	0
C73	<i>Salmonella ser Enteritidis</i>	0
C74	<i>Salmonella ser Infantis</i>	3
C75	<i>Salmonella ser Adelaide</i>	0
C76	<i>Salmonella ser Worthington</i>	1
C105	<i>Corynebacterium jeikeium</i>	4



Accuracy of the platform -

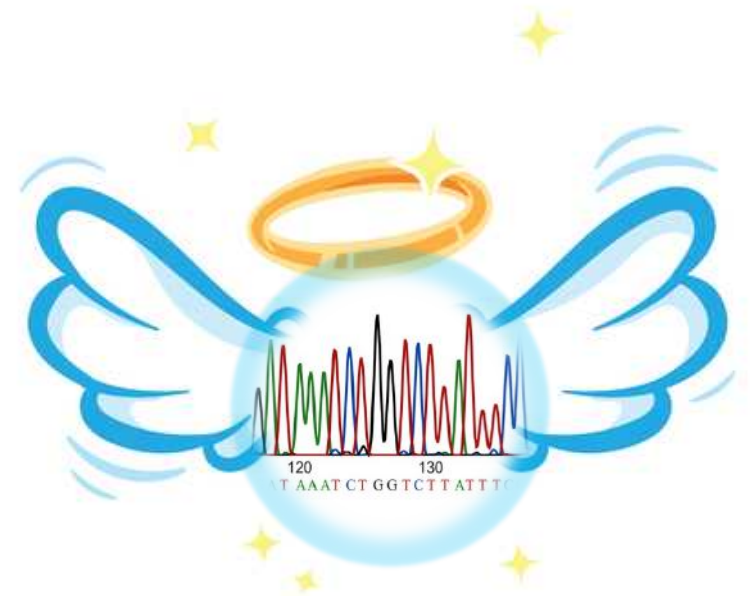
accuracy of the identification of individual base pairs (“base calling”) in the bacterial genome.



ID	Species	Total # of SNP difference with the reference	# of sequencing errors (SNP is supported only by 4 or less validation replicates)
C1	<i>Escherichia coli</i>	5	0
C2	<i>Aeromonas hydrophilia</i>	1	0
C3	<i>Escherichia coli</i>	22	0
C4	<i>Enterobacter cloacae</i>	10	0
C5	<i>Staphylococcus aureus</i>	0	0
C6	<i>Salmonella ser Typhimurium</i>	12	0
C46	<i>Enterococcus faecalis</i>	3	0
C47	<i>Staphylococcus epidermidis</i>	184	2
C48	<i>Staphylococcus saprophyticus</i>	27	0
C51	<i>Stenotrophomonas maltophilia</i>	39	0
C52	<i>Legionella pneumophila</i>	2	0
C55	<i>Escherichia coli</i>	14	1
C72	<i>Escherichia coli</i> O121:H19	0	0
C73	<i>Salmonella ser Enteritidis</i>	0	0
C74	<i>Salmonella ser Infantis</i>	3	0
C75	<i>Salmonella ser Adelaide</i>	0	0
C76	<i>Salmonella ser Worthington</i>	1	0
C105	<i>Corynebacterium jeikeium</i>	4	0

Searching for the reference...

Is Sanger sequencing so infallible?



Clinical Chemistry 62:4
647-654 (2016)

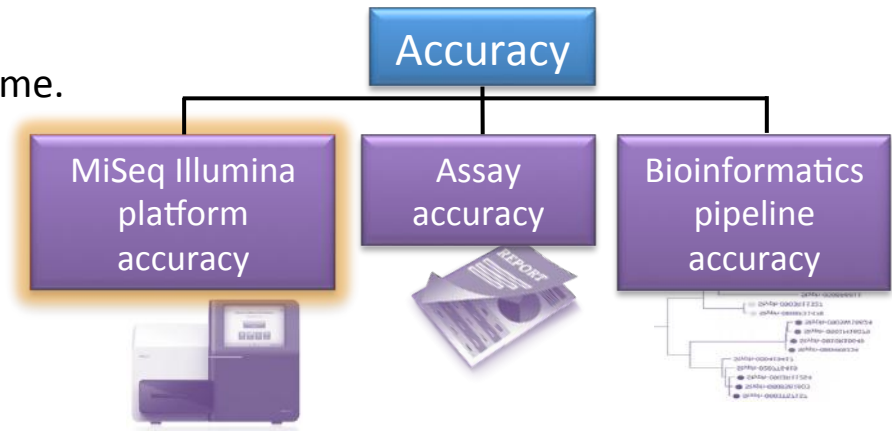
Molecular Diagnostics and Genetics

Systematic Evaluation of Sanger Validation of Next-Generation Sequencing Variants

Tyler F. Beck,¹ James C. Mullikin on behalf of the NISC Comparative Sequencing Program,^{1,2} and
Leslie G. Biesecker^{1*}

Accuracy of the platform -

accuracy of the identification of individual base pairs (“base calling”) in the bacterial genome.



ID	Species	Total # of SNP difference with the reference	# of sequencing errors (SNP is supported only by 4 or less validation replicates)
C1	<i>Escherichia coli</i>	5	0
C2	<i>Aeromonas hydrophilia</i>	1	0
C3	<i>Escherichia coli</i>	22	0
C4			
C5			
C6			
C7			
C8			
C9			
C10			
C11			
C12			
C13			
C14			
C15			
C16			
C17			
C18			
C19			
C20			
C21			
C22			
C23			
C24			
C25			
C26			
C27			
C28			
C29			
C30			
C31			
C32			
C33			
C34			
C35			
C36			
C37			
C38			
C39			
C40			
C41			
C42			
C43			
C44			
C45			
C46			
C47			
C48			
C49			
C50			
C51			
C52			
C53			
C54			
C55			
C56			
C57			
C58			
C59			
C60			
C61			
C62			
C63			
C64			
C65			
C66			
C67			
C68			
C69			
C70			
C71			
C72			
C73			
C74	<i>Salmonella ser Infantis</i>	3	0
C75	<i>Salmonella ser Adelaide</i>	0	0
C76	<i>Salmonella ser Worthington</i>	1	0
C105	<i>Corynebacterium jeikeium</i>	4	0

% agreement with reference =*

$$\frac{(Covered\ genome\ length) - (Total\ \#\ of\ SNP\ differing\ from\ reference)}{Covered\ genome\ length} \times 100\%$$

*Each nucleotide call - an independent test.

Accuracy

MiSeq Illumina platform accuracy

Assay accuracy

Bioinformatics pipeline accuracy

Accuracy of base calling against reference sequence

Data quality parameters affecting accuracy of base calling

Sequencing Run Q

- % of bases with quality score \geq Q30 for the run (%Q30)
- PhiX error rate
- Cluster density
- Cluster passing filter

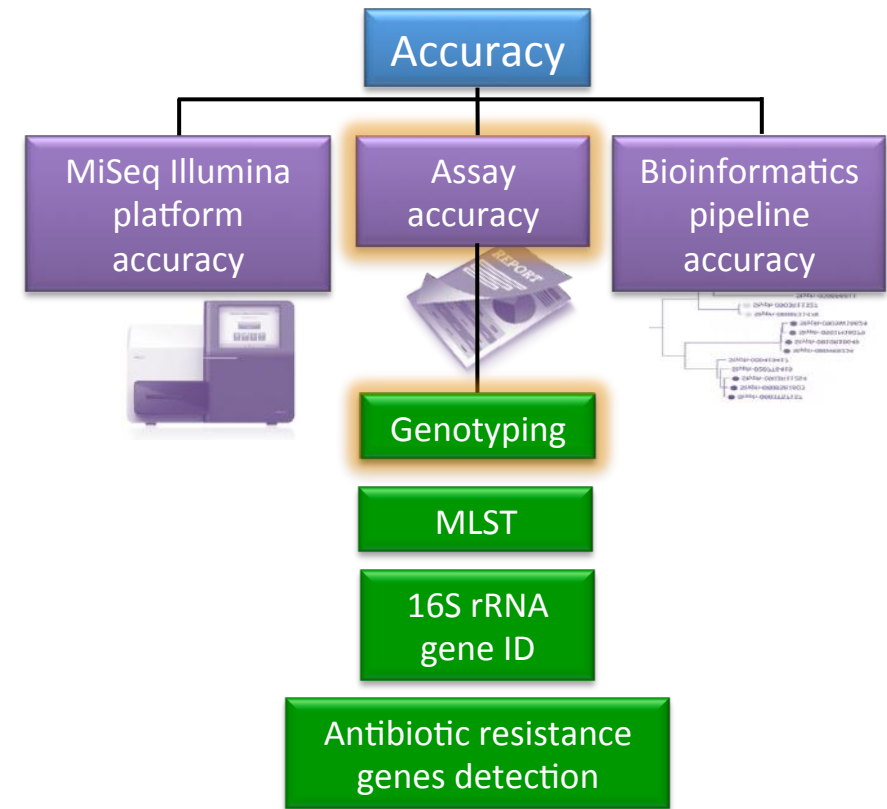
Raw Data Q

- Average depth of coverage
- The average read length with \geq Q30
- Minimum read length of the fragments which have \geq 75% of bases with Q30 score

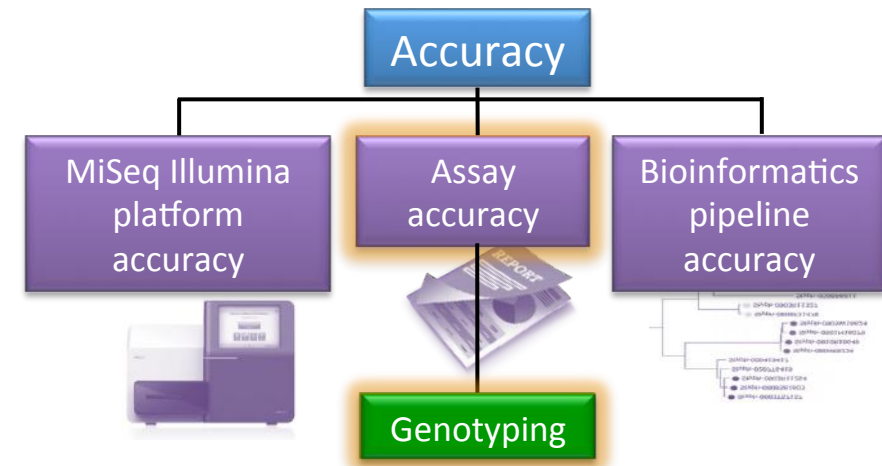
Secondary/Tertiary Data Analysis Q

- % genome covered after the mapping to reference
- Uniformity of coverage
- N50 for *de novo* assembled reads
- Number of assembled contigs

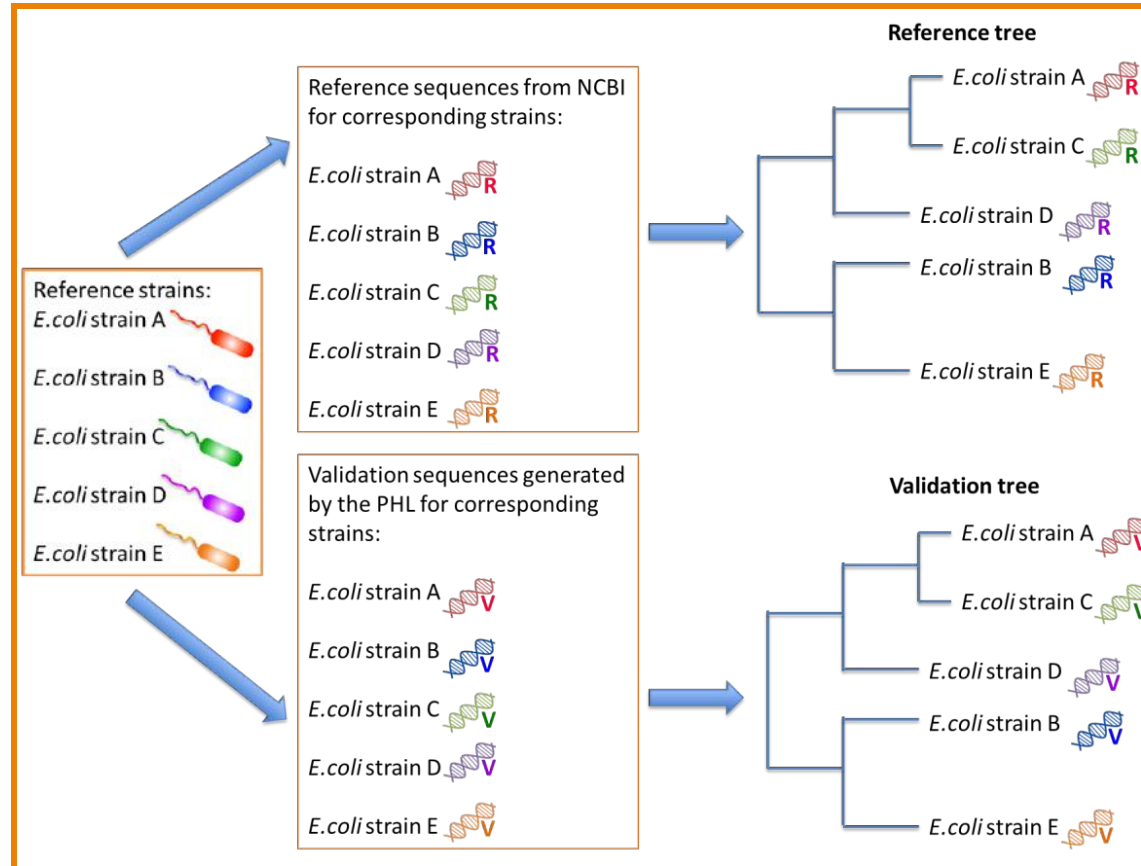
Assay accuracy - an agreement of the assay results for the validation sequences with the assay results for the reference sequences of the same strains.



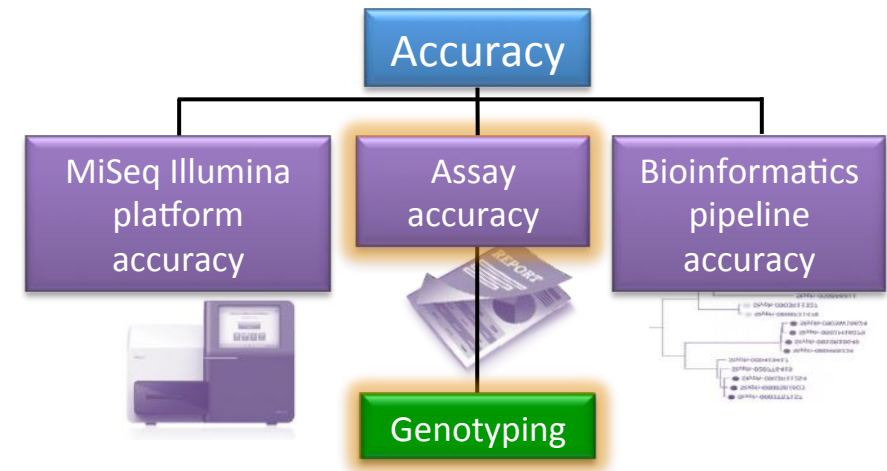
Assay accuracy - an agreement of the assay results for the validation sequences with the assay results for the reference sequences of the same strains.



Accuracy of genotyping assay:
congruence of phylogenetic trees built using reference sequences and validation sequences



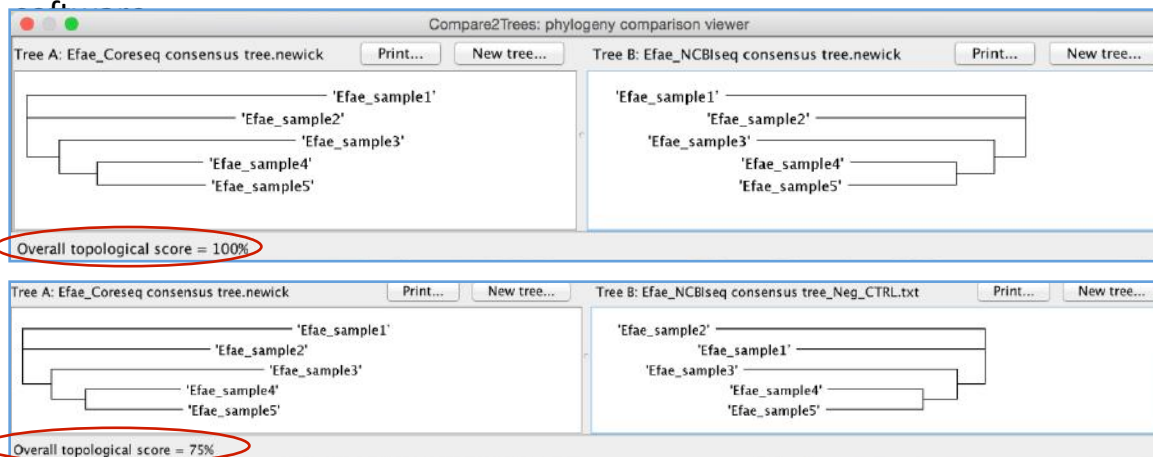
Assay accuracy - an agreement of the assay results for the validation sequences with the assay results for the reference sequences of the same strains.



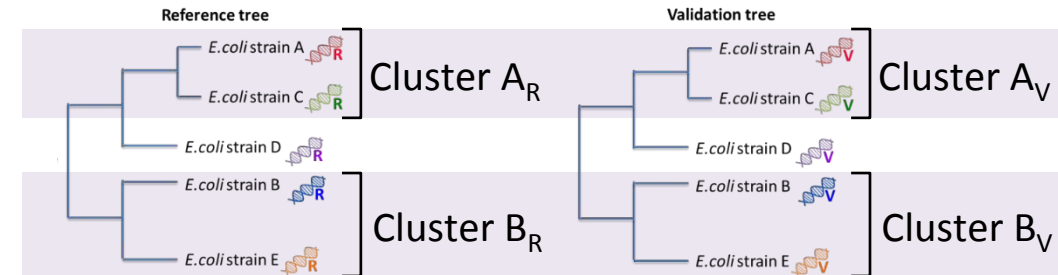
Accuracy of genotyping assay- congruence of phylogenetic trees built using reference sequences and validation sequences

Topological similarity between reference tree and validation tree

Percentage of topological similarity measured by Compare2Trees*

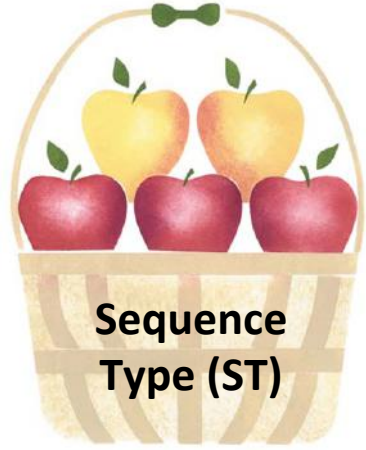


Comparison of clustering pattern of validation tree and reference tree



*Nye TM et al. Bioinformatics. 2006 Jan 1;22(1):117-9.

Assay accuracy - an agreement of the assay result for validation sequences generated by the lab with the assay result for reference sequences of the same strains.



Sequence Type (ST)

vs.

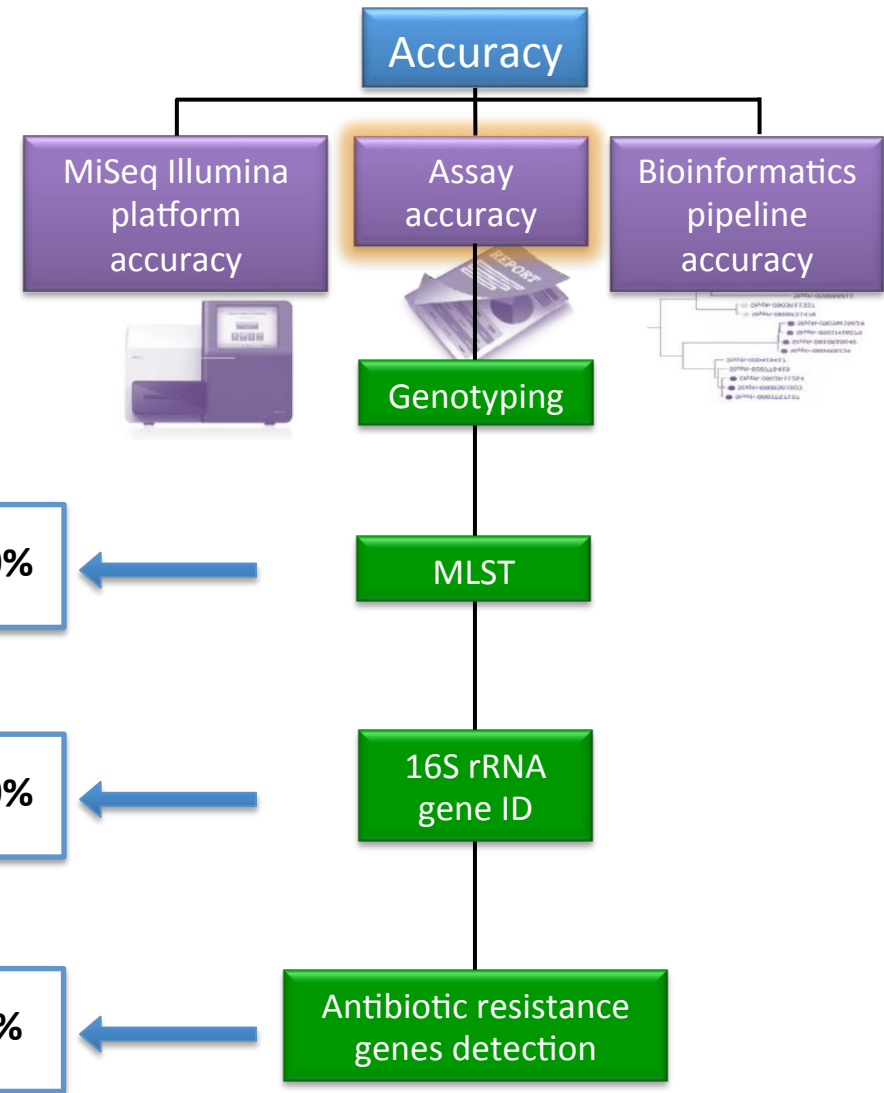


Allele

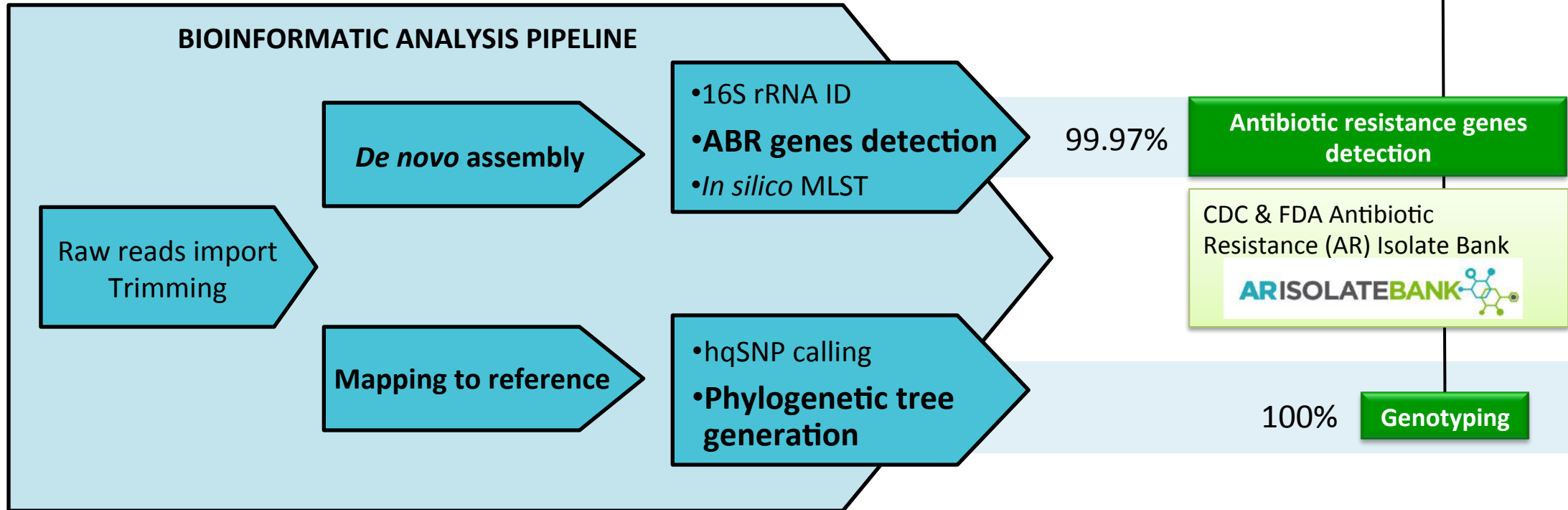
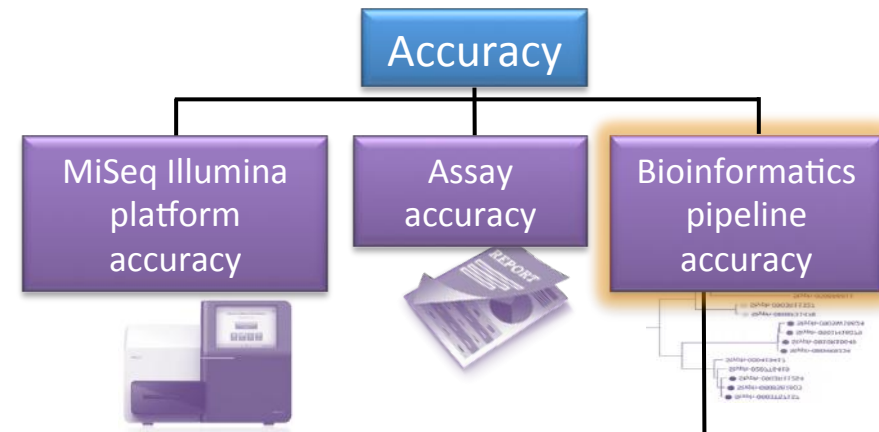
$$\frac{\text{\# of correctly called MLST alleles}}{\text{Total \# of MLST alleles tested}} \times 100\% = 100\%$$

$$\frac{\text{\# of correct IDs}}{\text{Total \# of IDs tested}} \times 100\% = 100\%$$

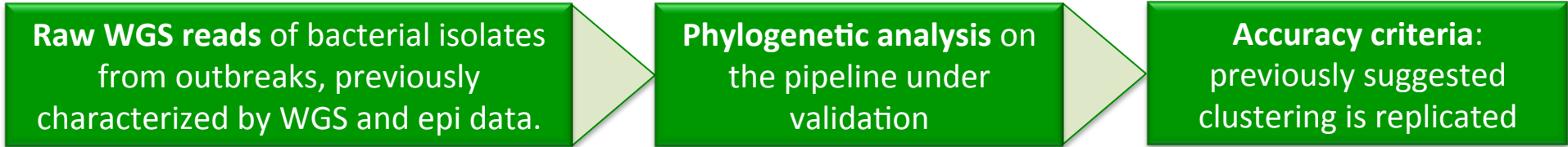
$$\frac{\text{\# of correctly detected ABR genes}}{\text{Total \# of ABR genes tested}} \times 100\% = 100\%$$



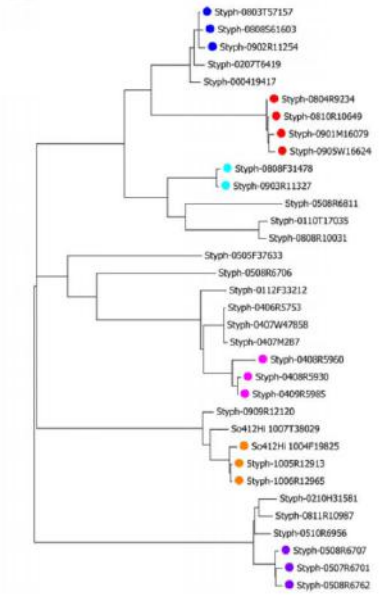
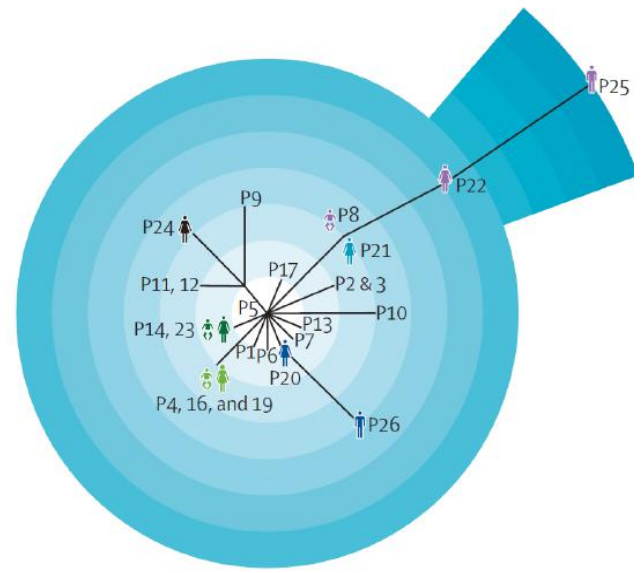
Bioinformatics pipeline accuracy



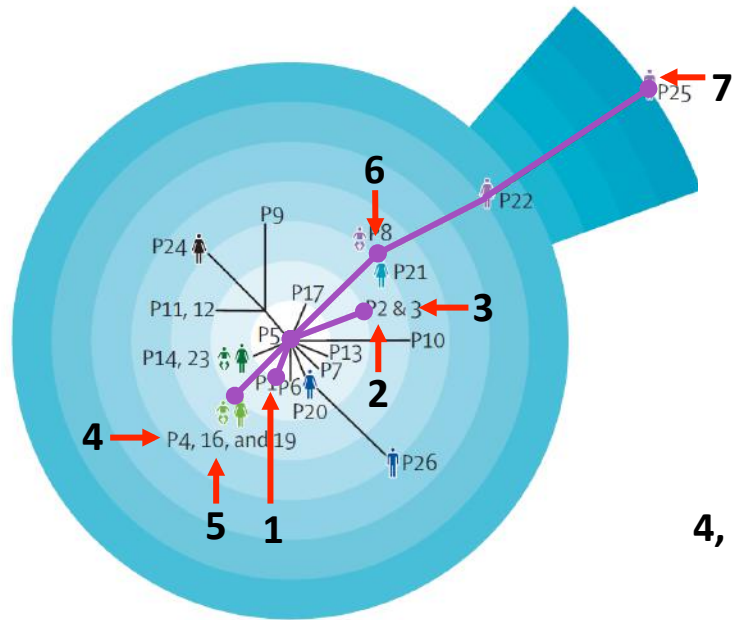
Phylogenetic bioinformatics pipeline accuracy:



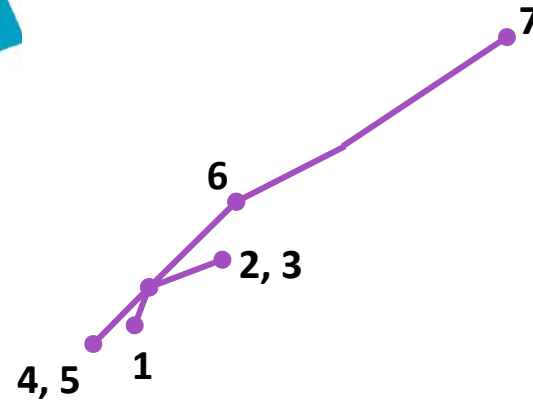
Study	Study 1. SR Harris et al. PMID: 23158674	Study 2. P Leekitcharoenphon et al. PMID: 24505344
Microorganism	Methicillin-resistant <i>Staphylococcus aureus</i>	<i>Salmonella enterica</i> serovar Typhimurium
Source of isolates	Human	Human
Number of isolates analyzed	7 outbreak isolates (1 outbreak cluster) + 2 epi unrelated isolates	9 outbreak isolates (4 outbreak clusters) + 2 epi unrelated isolates
Type of outbreak	Hospital-associated outbreak	Foodborne outbreaks



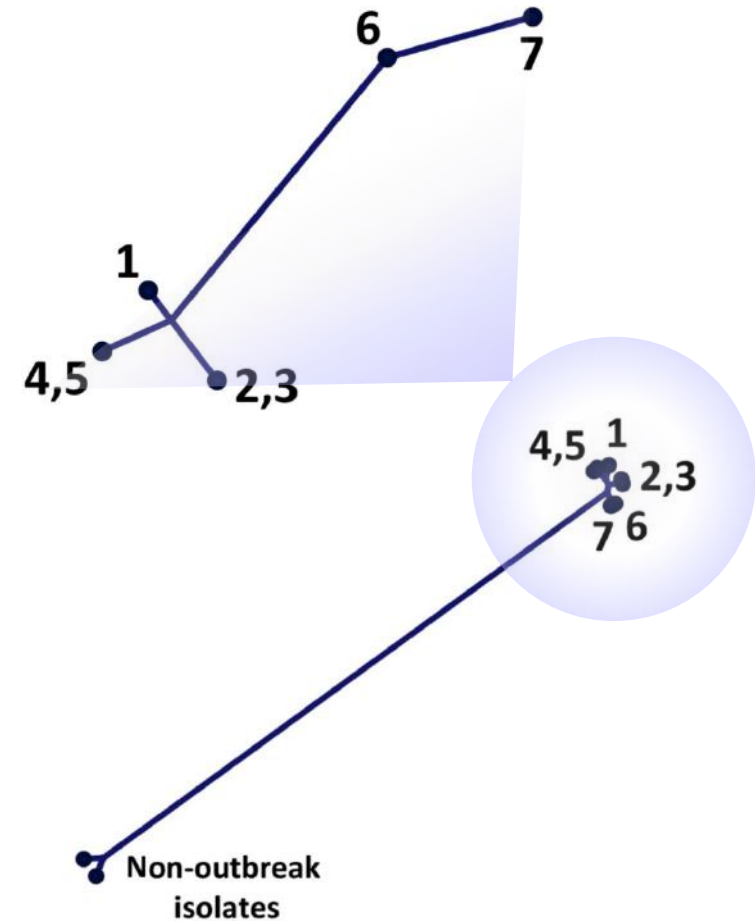
Study	Study 1. SR Harris et al. PMID: 23158674
Microorganism	Methicillin-resistant <i>Staphylococcus aureus</i>
Source of isolates	Human
Number of isolates analyzed	7 outbreak isolates (1 outbreak cluster) + 2 epi unrelated isolates
Type of outbreak	Hospital-associated outbreak



Original tree

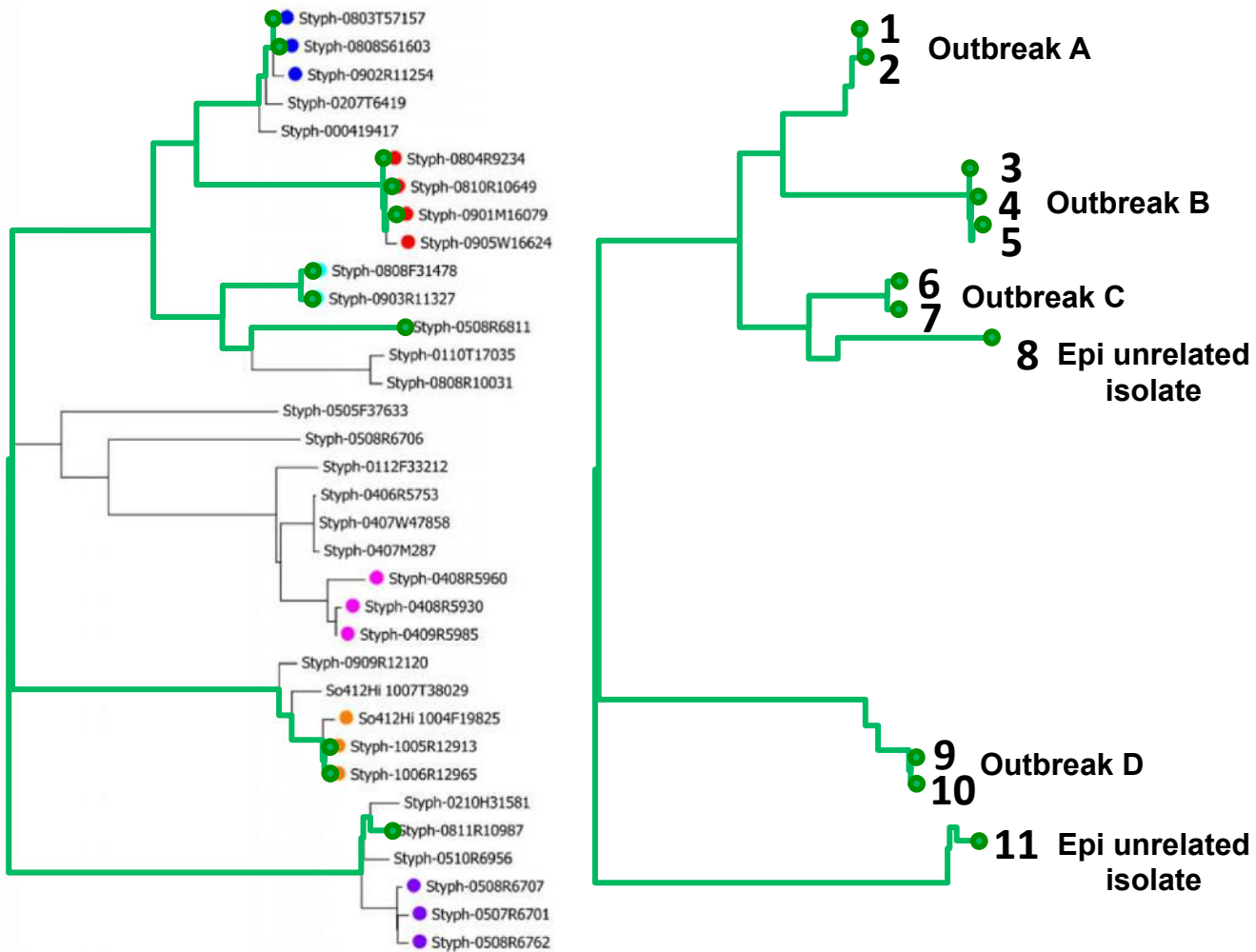


Validation tree

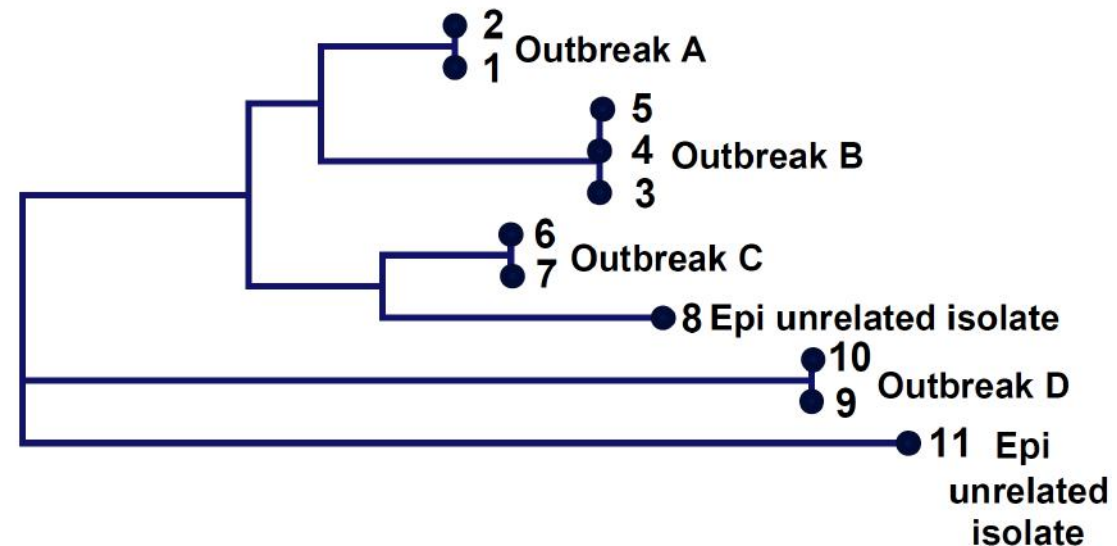


Study	Study 2. P Leekitcharoenphon et al. PMID: 24505344
Microorganism	<i>Salmonella enterica</i> serovar Typhimurium
Source of isolates	Human
Number of isolates analyzed	9 outbreak isolates (4 outbreak clusters) + 2 epi unrelated isolates
Type of outbreak	Foodborne outbreak

Original tree



Validation Tree



Laboratory Investigation of *Salmonella enterica* serovar Poona Outbreak in California: Comparison of Pulsed-Field Gel Electrophoresis (PFGE) and Whole Genome Sequencing (WGS) Results

[Varvara K. Kozyreva](#), [John Crandall](#), [Ashley Sabol](#), [Alyssa Poe](#), [Peng Zhang](#), [Jeniffer Concepción-Acevedo](#),* [Morgan N. Schroeder](#), [Darlene Wagner](#), [Jeffrey Higa](#), [Eija Trees](#), and [Vishnu Chaturvedi](#)

Between-laboratory comparison of genotyping pipelines

Microbial Diseases Lab CA

PFGE pattern color coding

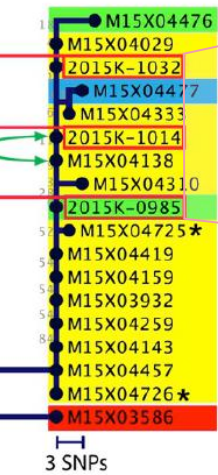
- JL6X01.0018
- JL6X01.0375
- JL6X01.0778
- JL6X01.0776

* Cucumber isolates

120-123 SNPs difference

0 SNP

The same isolate sequenced by 2 labs- 0 SNP



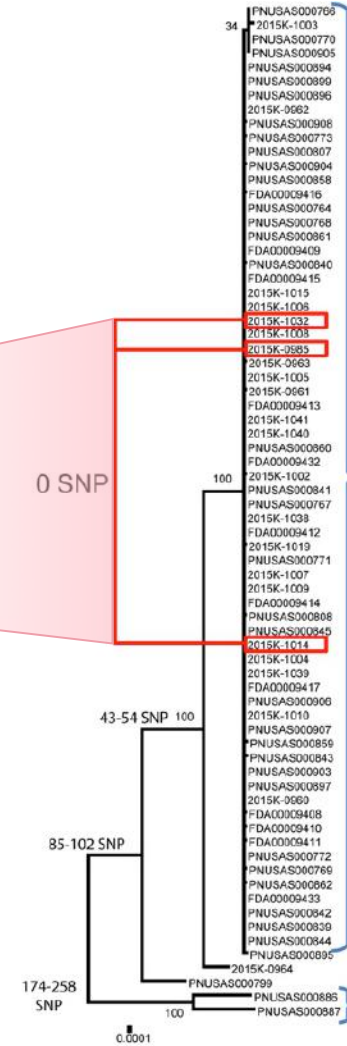
3 SNPs

0-5 SNPs difference

VS.

0 SNP

CDC



0-4 SNPs

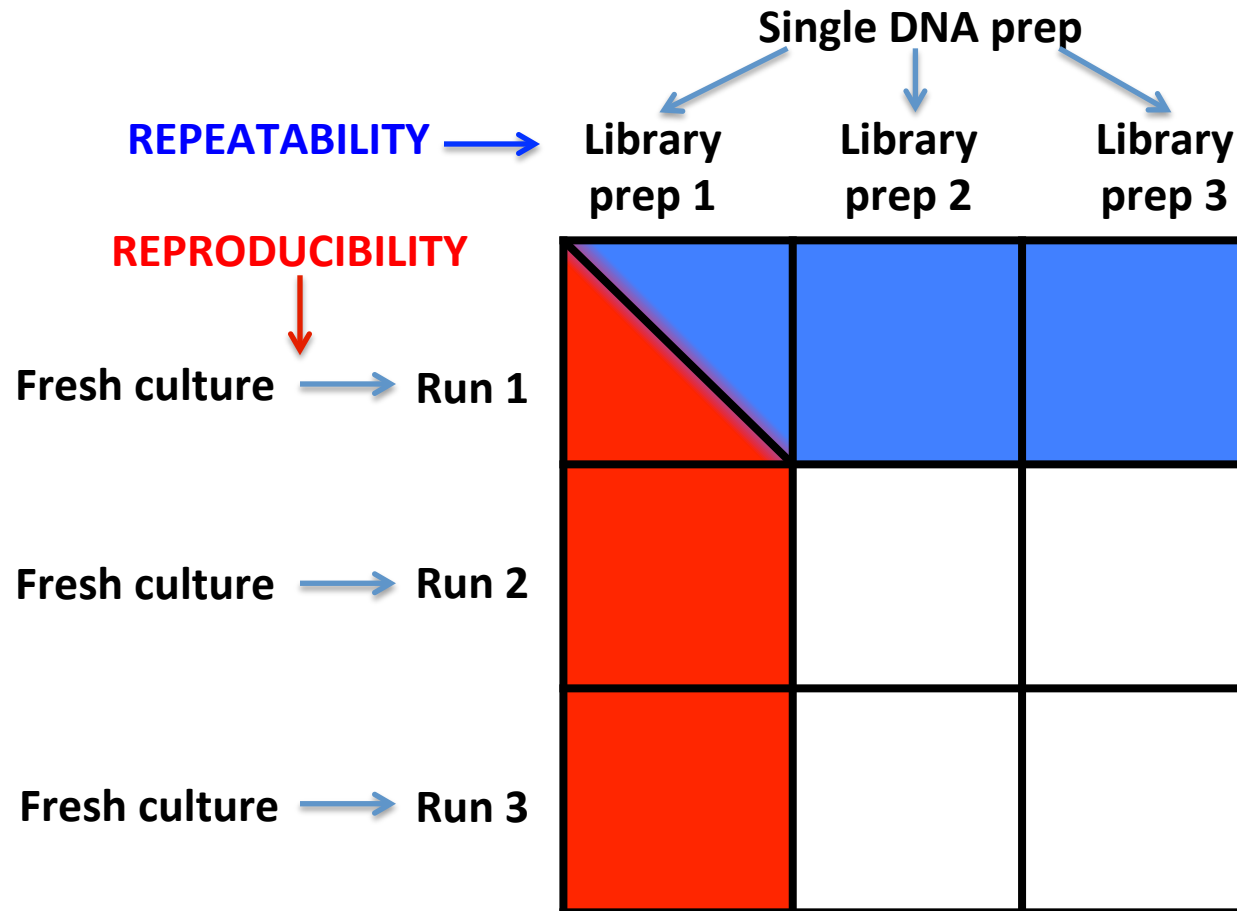
174-258 SNP

95 SNP

0.0001

Repeatability (=precision within run=inter-assay agreement) –
The concordance of the assay results and quality metrics obtained for a sample tested multiple times within the same sequencing run.

Reproducibility (=precision between runs=intra-assay agreement) –
the consistency of the assay results and quality metrics for the same sample sequenced on different occasions by different operators.



Repeatability & Reproducibility

triplicates

Genotyping

MLST

16S rRNA
gene ID

Consistency of
quality metrics

Repeatability and reproducibility of SNP calling—
between- and within-run precision of single nucleotide variant detection.

Repeatability & Reproducibility

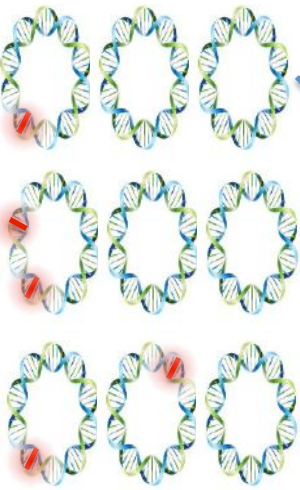
Evaluation of absolute precision
per replicate

The whole genome of a single replicate = an independent test

Evaluation of precision relative to the genome size (**per base pair**)

Each nucleotide call in the WGS = an independent test (a percentage of the number of bp called in the genome)

Rep 1 Rep 2 Rep 3

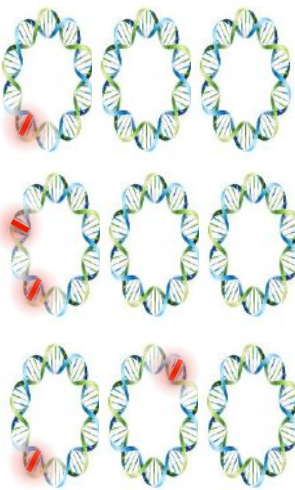


66.7% agreement
between replicates

33.3% agreement
between replicates

bp sequenced in agreement
bp sequenced in total

Rep 1 Rep 2 Rep 3



$$\frac{(5,000,000 \times 3) - 1}{5,000,000 \times 3}$$

$$\frac{(5,000,000 \times 3) - 2}{5,000,000 \times 3}$$

$$\frac{(5,000,000 \times 3) - 2}{5,000,000 \times 3}$$

Repeatability and reproducibility of SNP calling— between- and within-run precision of single nucleotide variant detection.

Repeatability & Reproducibility

Evaluation of absolute precision
per replicate

Evaluation of precision relative to the
genome size (**per base pair**)

The whole genome of a single
replicate = an independent test

Each nucleotide call in the WGS
= an independent test (a percentage of
the number of bp called in the genome)

Isolate ID		C1	C2	C3	C4	C5	C6	C46	C47	C48	C49	C50	C51	C52	C53	C54	C55	C72	C73	C74	C75	C76	C103	C104	C105	C106	C56	C57	C58	C59	C61	C65	C67	C68	C69
Total # of SNP difference for replicates	within-run	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	between-run	0	0	0	0	0	0	0	2	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Precision per replicate:

Within-run precision =
 $(101/102) \times 100\%$
 = 99.02%

Between-run precision =
 $(99/102) \times 100\%$
 = 97.05%

99.99998% 99.99999% 99.99999%

99.99997%

... of covered genome (sequenced in triplicate)

Precision per base pair:

Within-run precision (avg) =
 99.9999997%

Between-run precision (avg) =
 99.999998%

Repeatability & Reproducibility

triplicates

Genotyping

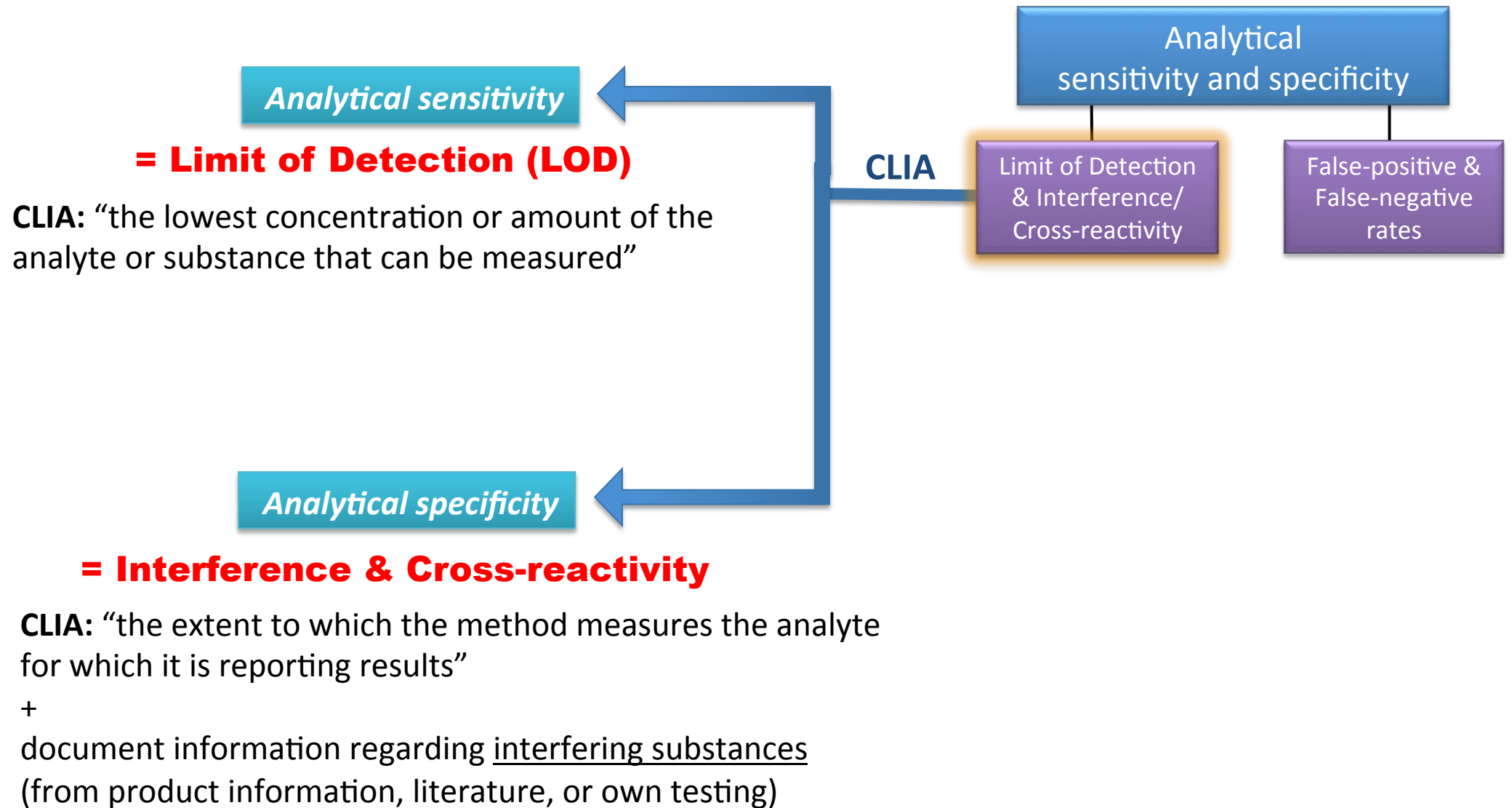
MLST

16S rRNA
gene ID

Between-run	$\frac{\text{\# of alleles in agreement}}{\text{Total \# of alleles tested}} \times 100\% = \mathbf{100\%}$
Within-run	

Between-run	$\frac{\text{\# of isolates' IDs in agreement}}{\text{Total \# of IDs tested}} \times 100\% = \mathbf{100\%}$
Within-run	





Analytical sensitivity of WGS

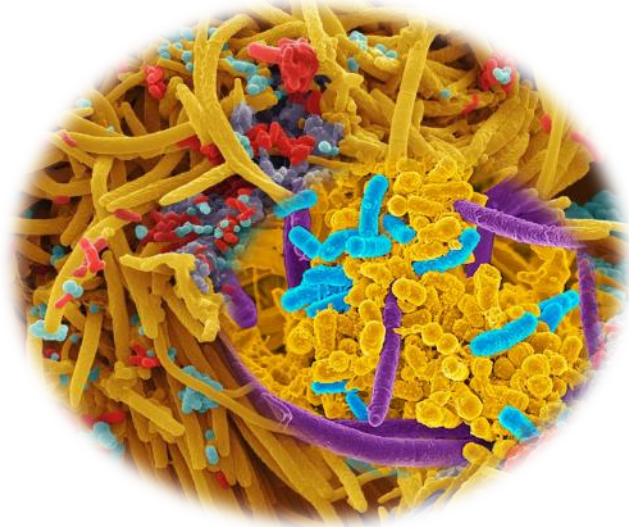
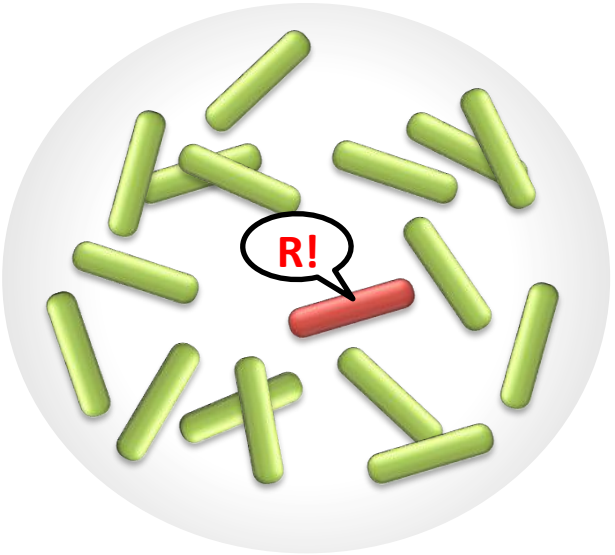
Analytical sensitivity and specificity

Limit of Detection & Interference/ Cross-reactivity



LOD for heteroresistance detection in *Mycobacterium tuberculosis*

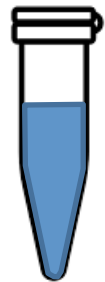
LOD for metagenomic pathogen detection



Analytical specificity of WGS

- the ability of an assay to detect only the intended target in the presence of potentially cross-reacting nucleotide sequences

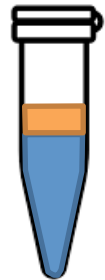
Example of modeled interference/cross-reactivity from contaminating DNA:



Pure genome of *E. coli*

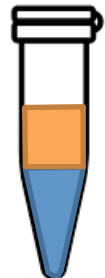
Number of SNPs
differing w/reference

22



Contamination
25% (*in silico*)
w/ *Salmonella enterica*

22



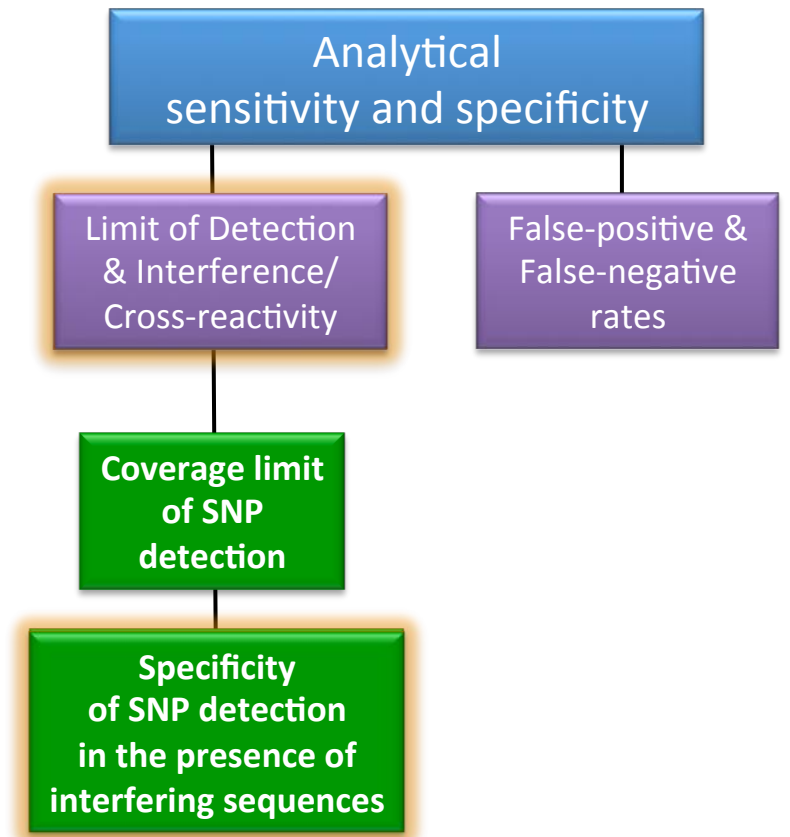
Contamination
50% (*in silico*)
w/ *Salmonella enterica*

24

1 SNP missed
1 non-specific SNP detected



68% reads mapped
59.6% reads in pairs



Analytical specificity of WGS

Samples		C3 <i>E.coli</i>	C3 <i>E.coli</i> + C75 <i>Salmonella</i> <i>enterica</i>	C3 <i>E.coli</i> + C1 <i>E.coli</i>	C3 <i>E.coli</i> + C54 <i>Acinetobacter</i> <i>baumannii</i>	C3 <i>E.coli</i> + C57 <i>Mycobacterium</i> <i>tuberculosis</i>	C3 <i>E.coli</i> + C5 <i>Staphylococcus</i> <i>aureus</i>
Mapping quality metrics	% of Contamination reads	no contamination	50%	50%	50%	50%	50%
	% of Mapped reads	99%	68%	91%	51%	67%	49%
	% of Not mapped reads	1%	32%	9%	49%	33%	51%
	% of Reads in pairs	90%	59.60%	81%	48.30%	60%	46.60%
	% of Broken paired reads	9%	8%	10%	3%	6%	3%
	% of reference covered	99%	99%	99%	99%	99%	99%
SNP calling specificity	SNPs between sequence and reference	22	24	22	40	22	36
	Number of FN SNPs	NA	1	0	0	0	0
	Number of FP SNPs	NA	1	0	18	0	14

Analytical sensitivity and specificity

Limit of Detection & Interference/
Cross-reactivity

False-positive &
False-negative rates

Coverage limit
of SNP
detection

Specificity
of SNP detection
in the presence of
interfering sequences

$$\text{Sensitivity} = \frac{TP}{TP + FN} \times 100\%$$

$$\text{Specificity} = \frac{TN}{TN + FP} \times 100\%$$

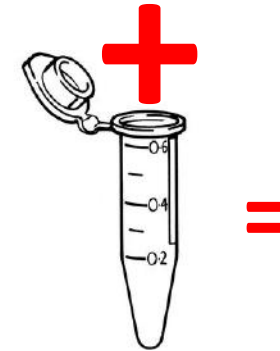
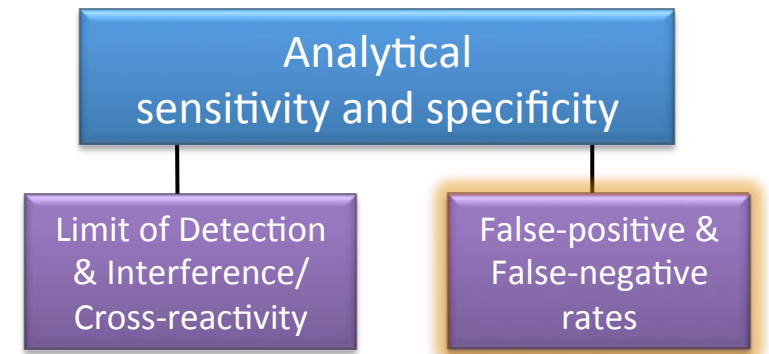
Often used in
DIAGNOSTIC (or clinical)
 sensitivity and specificity calculations

DIAGNOSTIC SENSITIVITY & SPECIFICITY-
 ability to predict the disease or condition in a person:

vs.

Analytical Sensitivity & Specificity
 as false-positive & false-negative rates of
 WGS assays

TP- True positive results
 TN- True negative results
 FP- False positive
 FN- False negative



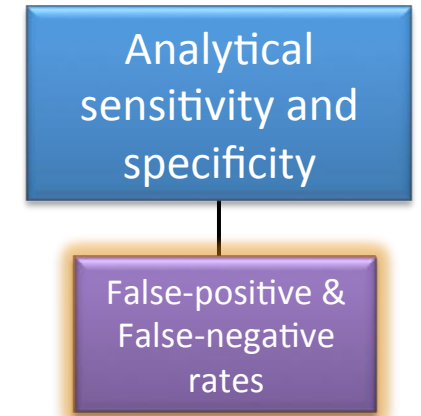
Analytical Sensitivity & Specificity as false-positive & false-negative rates of WGS assays

Analytical sensitivity- the likelihood that a WGS assay will detect sequence variation when it is present (this value reflects a false negative rate of the test).

$$\text{Analytical sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \times 100\%$$

Analytical specificity- the probability that a WGS assay will not detect sequence variations when none are present (this value reflects a test's false positive rate).

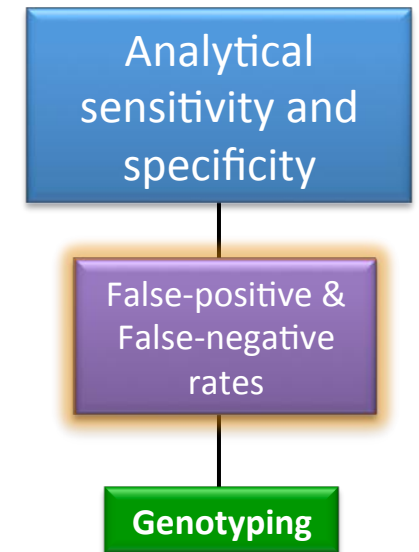
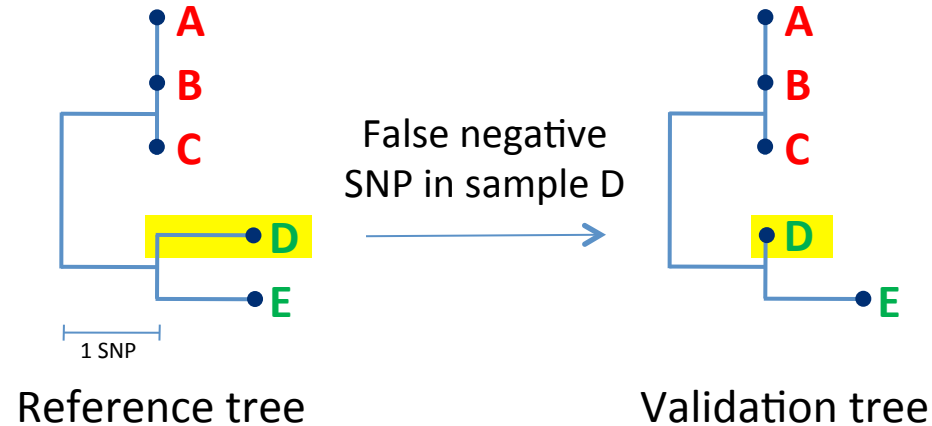
$$\text{Analytical specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \times 100\%$$



TP- True positive results
TN- True negative results
FP- False positive
FN- False negative

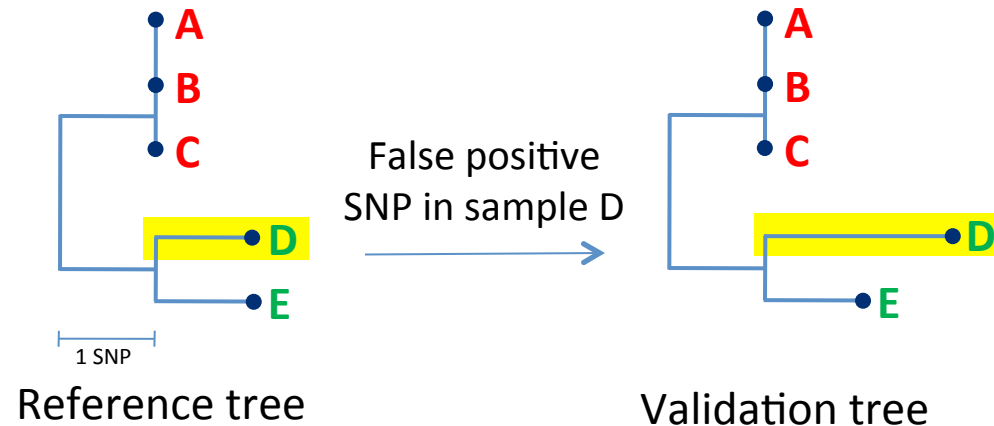
Analytical sensitivity in genotyping assay is the likelihood that all the SNPs differing between the isolates will be detected.

False negative result - missed sequence variations, e.g.:



Analytical specificity in genotyping assay is the likelihood that variation between the isolates (SNP) will not be detected when none are present.

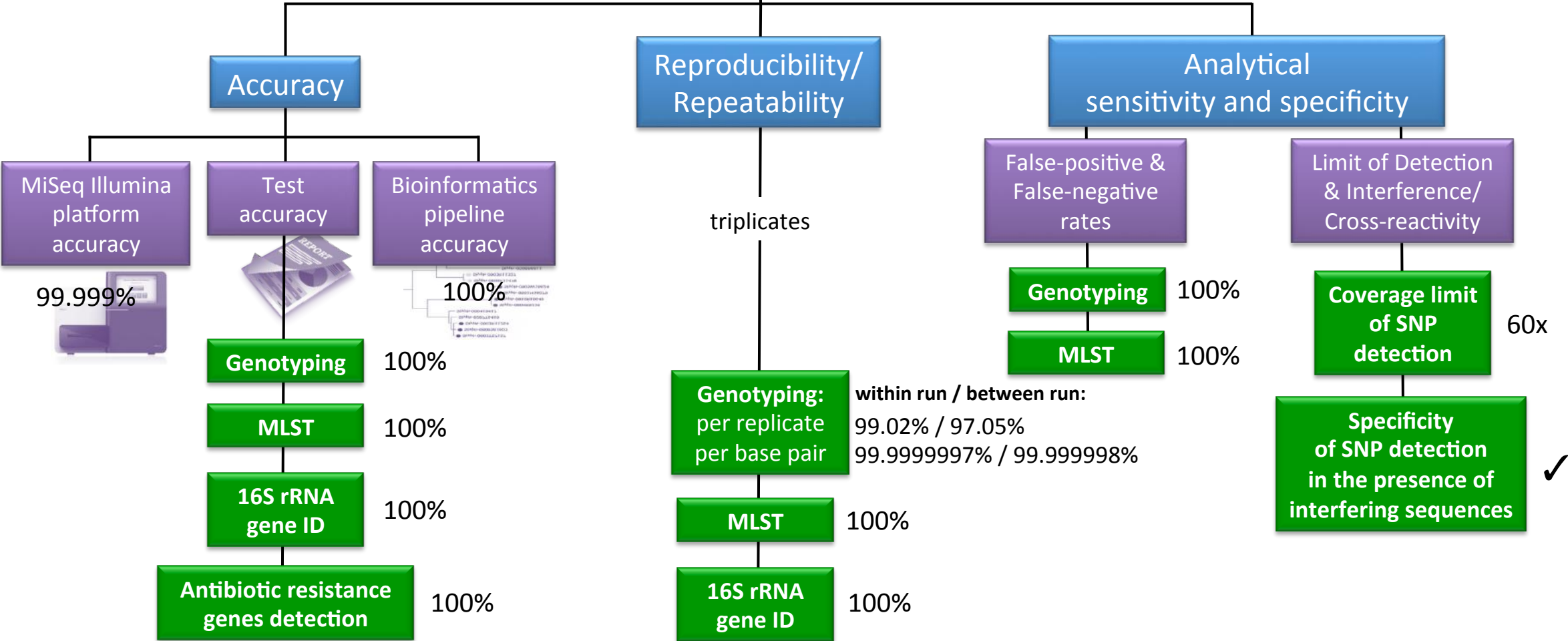
False positive result – false sequence variations (introduced either during library prep/sequencing, or data analysis), e.g.:



10- Enterobacteriaceae
5- Gram-positive cocci isolates
5- Gram-negative non-fermenting bacterial isolates
9- Mycobacterium tuberculosis
5- representatives of miscellaneous species

Validation Set
34 bacterial isolates

WHOLE GENOME SEQUENCING VALIDATION IN PUBLIC HEALTH MICROBIOLOGY LAB SETTINGS



WGS implementation in CLIA settings

ASSAY DEVELOPMENT

Wet-Lab / Dry-Lab

- Workflow development & optimization
- Empirically determined optimal assay conditions
- Documentation: SOP, worksheets, QC logs, calibration sheets, etc.
- Intended use, Acceptance and rejection criteria

ASSAY VALIDATION

Platform – Bioinformatics - Assay

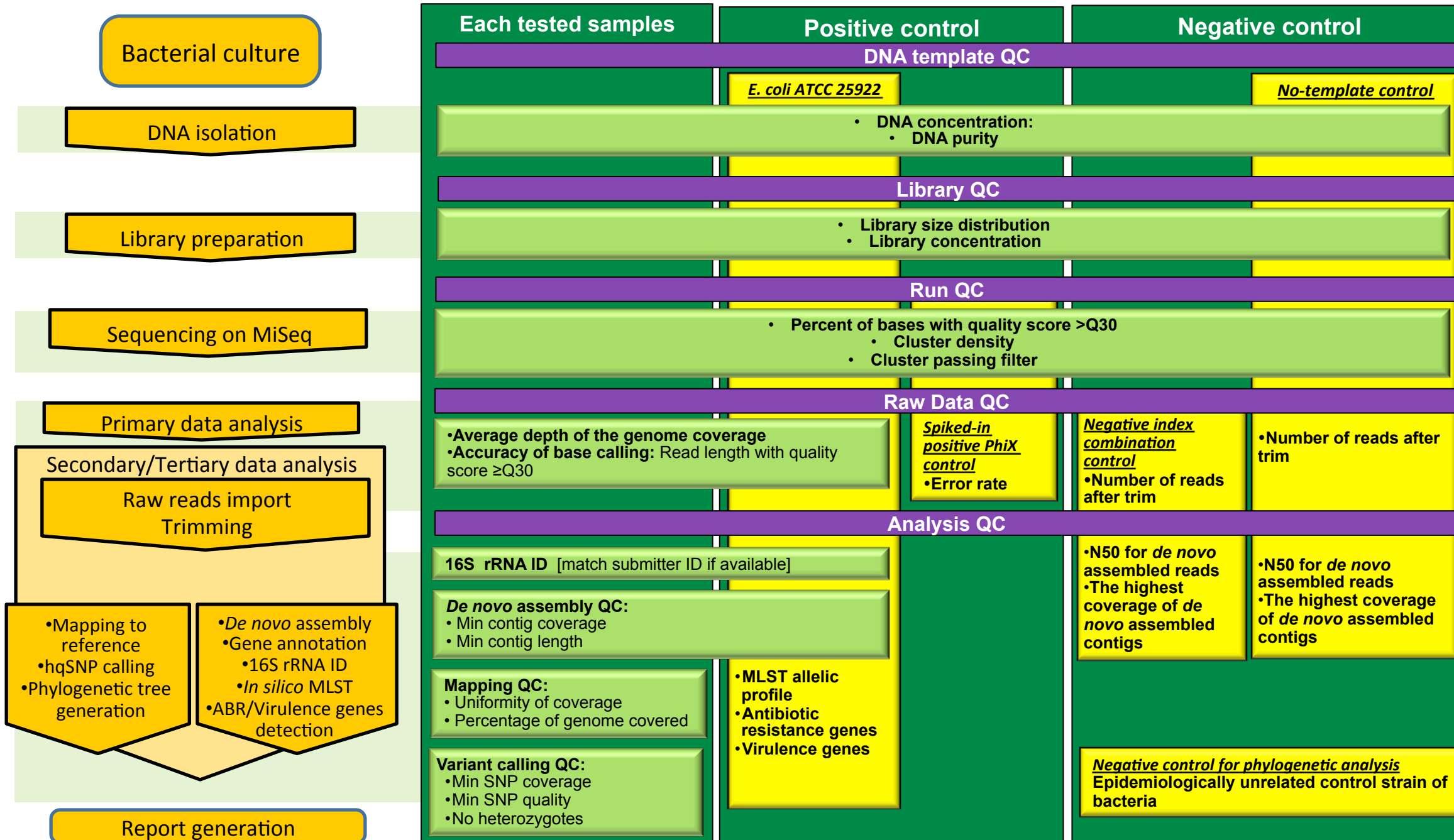
- Validation Plan
- Ability to sequence targets of interest
- Accurate results of SNP detection / genotype / ID / resistance genes / etc
- Determine QC metrics and checkpoints

QUALITY MANAGEMENT

- Documented quality management plan
- Routine QC of all reagents
- Equipment preventive maintenance and calibration
- Proficiency testing (PT)
- Competency assessment
- Re-validation/confirmation of the performance for assay modifications
- Instrument correlation
- Corrective actions

**QUALITY
WGS**

WGS QUALITY CONTROL SCHEME



REPORTING LANGUAGE

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Antibiotic resistance genes found:

Type of antibiotic resistance	Aminoglycoside					
	strB		strA		armA	
Resistance gene	strB		strA		armA	
Gene accession number in NCBI	M96392		M96392		AY220558	
Isolate ID	% ID	Query/aligned length	% ID	Query/aligned length	% ID	Query/aligned length
C239	100	837 / 837	100	804 / 804	100	774 / 774
C240	100	837 / 837	100	804 / 804	100	774 / 774
C241	100	837 / 837	100	804 / 804	100	774 / 774
C242	100	837 / 837	100	804 / 804	100	774 / 774
C243	100	837 / 837	100	804 / 804	100	774 / 774
C244	100	837 / 837	100	804 / 804	100	774 / 774
C245	100	837 / 837	100	804 / 804	100	774 / 774

Type of antibiotic resistance	Beta-Lactam	
Resistance gene	blaNDM-1	
Gene accession number in NCBI	FN398876	
Isolate ID		

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

The Whole Genome Sequencing of 7 *Klebsiella pneumoniae* isolates was performed by MDL Core laboratory using Illumina sequencing chemistry, 900bp x 2 paired-end reads, on Illumina MiSeq sequencer. Phylogenetic analysis, 16S rRNA identification, *in silico* MLST typing, and antibiotic resistance genes detection results are included in this report. Please find the report summary on page 6.

Isolate Details

MDL Core Lab Sample ID	Submitter	Submitter's ID number	Country	Date collected	Date received	Time of Day of Isolation	Case or Outbreak	Date of Reporting
C239	MDL PH90	M13401222	Alameda (Highland Hosp.)	7/30/2013	11/9/2013	11/9/2013	12/8/2013	12/8/2013
C240	MDL PH90	M13401223	Alameda (Highland Hosp.)	2/6/2014	12/9/2013	12/9/2013	12/8/2013	12/8/2013
C241	MDL PH90	M13401224	Alameda (Highland Hosp.)	6/28/2013	12/9/2013	12/9/2013	12/8/2013	12/8/2013
C242	MDL PH90	M13401225	Alameda (Highland Hosp.)	10/2/2013	12/9/2013	12/9/2013	12/8/2013	12/8/2013
C243	MDL PH90	M13401226	Alameda (Highland Hosp.)	10/3/2013	12/9/2013	12/9/2013	12/8/2013	12/8/2013
C244	MDL PH90	M13401227	Alameda (Highland Hosp.)	7/27/2013	12/9/2013	12/9/2013	12/8/2013	12/8/2013
C245	MDL PH90	M13401228	Alameda (Highland Hosp.)					

Species identity confirmation

- Species identity was confirmed
- Species identity is different from original identification. New ID _____
- Species identity is different from percent identity: _____

List genes used for identification and percent identity:

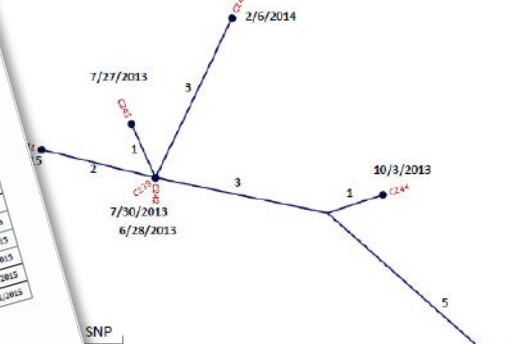
- 16S rRNA 99-100% identity *Klebsiella pneumoniae*
- cpn60 _____
- Other: _____

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Results of genotyping based on high quality (hg) Single Nucleotide Polymorphisms (SNPs)

Phylogenetic tree

- Algorithm used for phylogenetic tree building:
- Neighbor Joining
 - Maximum Likelihood
 - Other: Specify: _____



above the branches designate SNP difference between closest relatives are labeled by the corresponding isolates ID numbers.

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Summary

1. All *Klebsiella pneumoniae* isolates clustered together based on high quality (hg) Single Nucleotide Polymorphisms (SNPs) and differ by 0-8 SNPs. 1-18.4 substitutions/genome/year [1, 2], suggesting that the samples isolated 2 years apart could be related. All isolates belong to the same sequence type ST-147. *bla*_{NDM-1} gene does not encode a carbapenemase and is known to mediate aminoglycoside (streptomycin) resistance [3]. All isolates were found to be susceptible to all other antibiotics tested.

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Distance matrix (pairwise comparison):

	C239	C240	C241	C242	C243	C244	C245
C239	1	1	2	3	4	5	6
C240	2	1	3	4	5	6	7
C241	3	3	1	2	3	4	5
C242	4	4	2	1	2	3	4
C243	5	5	3	2	1	2	3
C244	6	6	4	3	2	1	2
C245	7	7	5	4	3	2	1

More similar (blue) / More different (red)

Explanation: Value in intersection shows the number of sites difference between two isolates.

In silico MLST results:

Isolate ID	MLST type (sequence type, ST)
C239	ST-147
C240	ST-147
C241	ST-147
C242	ST-147
C243	ST-147
C244	ST-147
C245	ST-147

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Date of analysis: 02/18/2015
Report status: Preliminary Final

Microbial Disease Laboratory / Core laboratory/ Whole Genome Sequencing report

Disclaimer: This report was generated using the Microbial Disease Laboratory's (MDL) Whole Genome Sequencing (WGS) pipeline. The results are provided based on a laboratory developed test (LDT) using Illumina MiSeq Sequencer. Additional investigation is necessary to confirm the results.

Isolate and genome wide analysis (including phylogenetic tree, group of isolates, etc.) were performed using the Microbial Disease Laboratory's (MDL) WGS pipeline. The results are provided based on a laboratory developed test (LDT) using Illumina MiSeq Sequencer. Additional investigation is necessary to confirm the results.



**CLIA
INSPECTION**



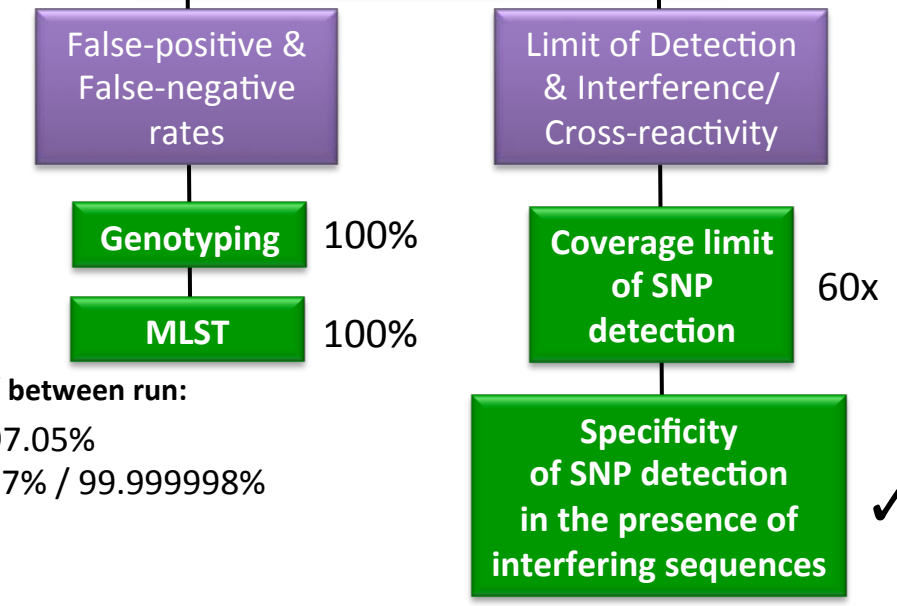
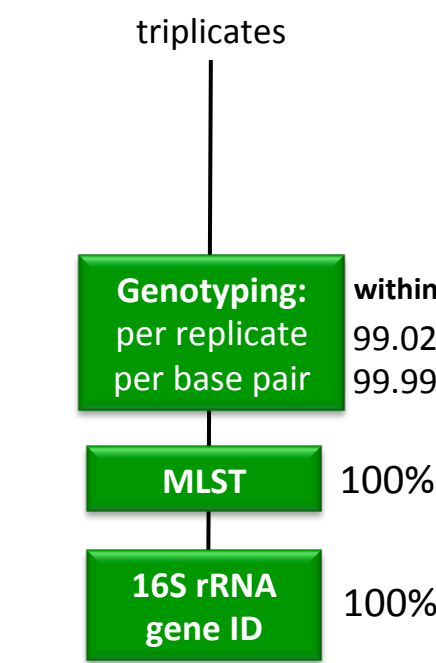
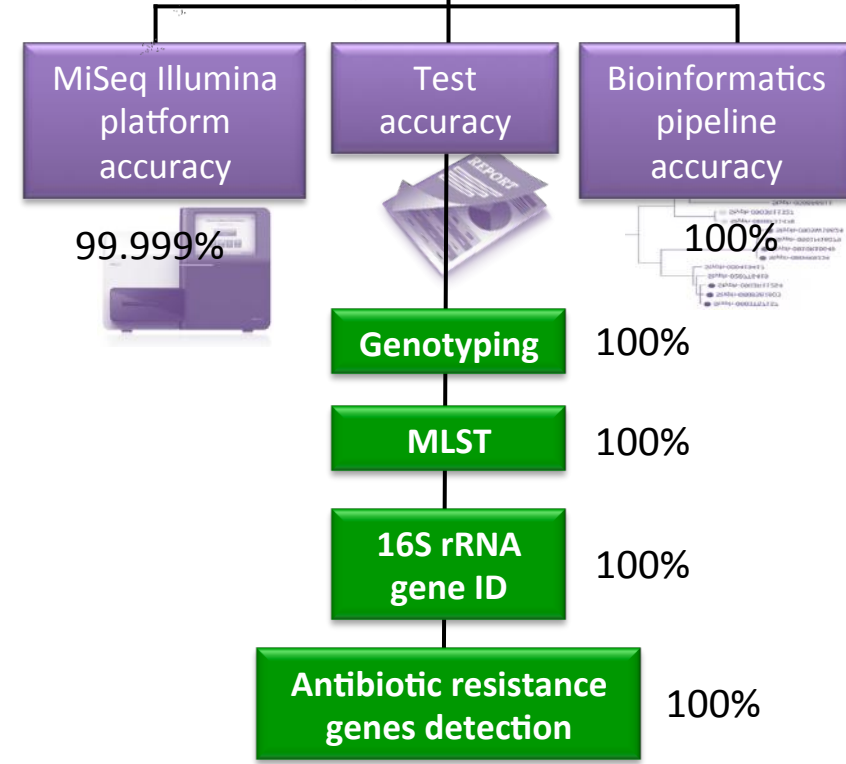
- 10- *Enterobacteriaceae*
- 5- Gram-positive cocci isolates
- 5- Gram-negative non-fermenting bacterial isolates
- 9- *Mycobacterium tuberculosis*
- 5- representatives of miscellaneous species

Validation Set
34 bacterial isolates

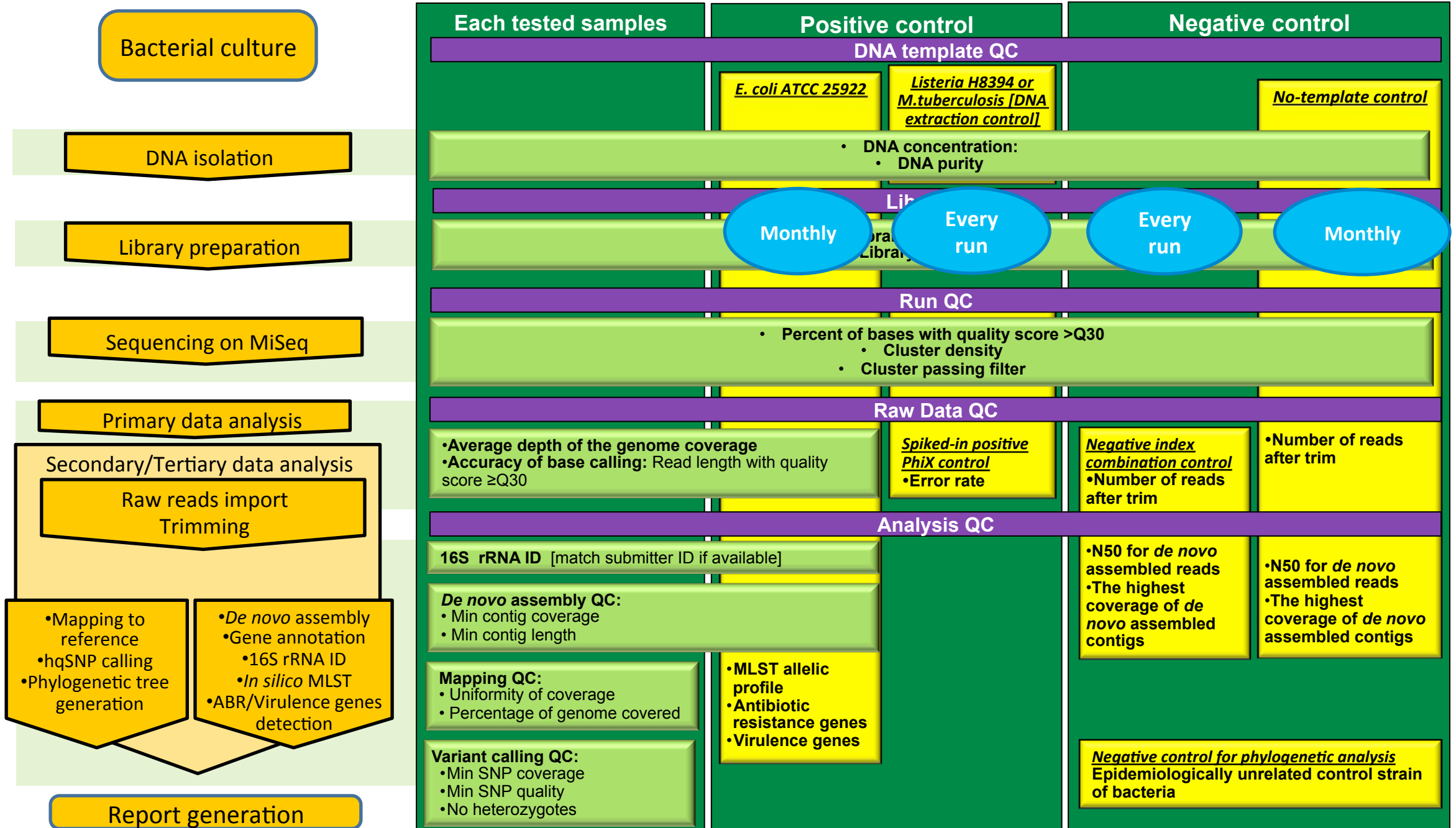
WHOLE GENOME SEQUENCING VALIDATION IN PUBLIC HEALTH MICROBIOLOGY LAB SETTINGS

Reproducibility/ Repeatability

Analytical sensitivity and specificity



WGS QUALITY CONTROL SCHEME



Bacterial culture

DNA isolation

Library preparation

Sequencing on MiSeq

Primary data analysis

Secondary/Tertiary data analysis

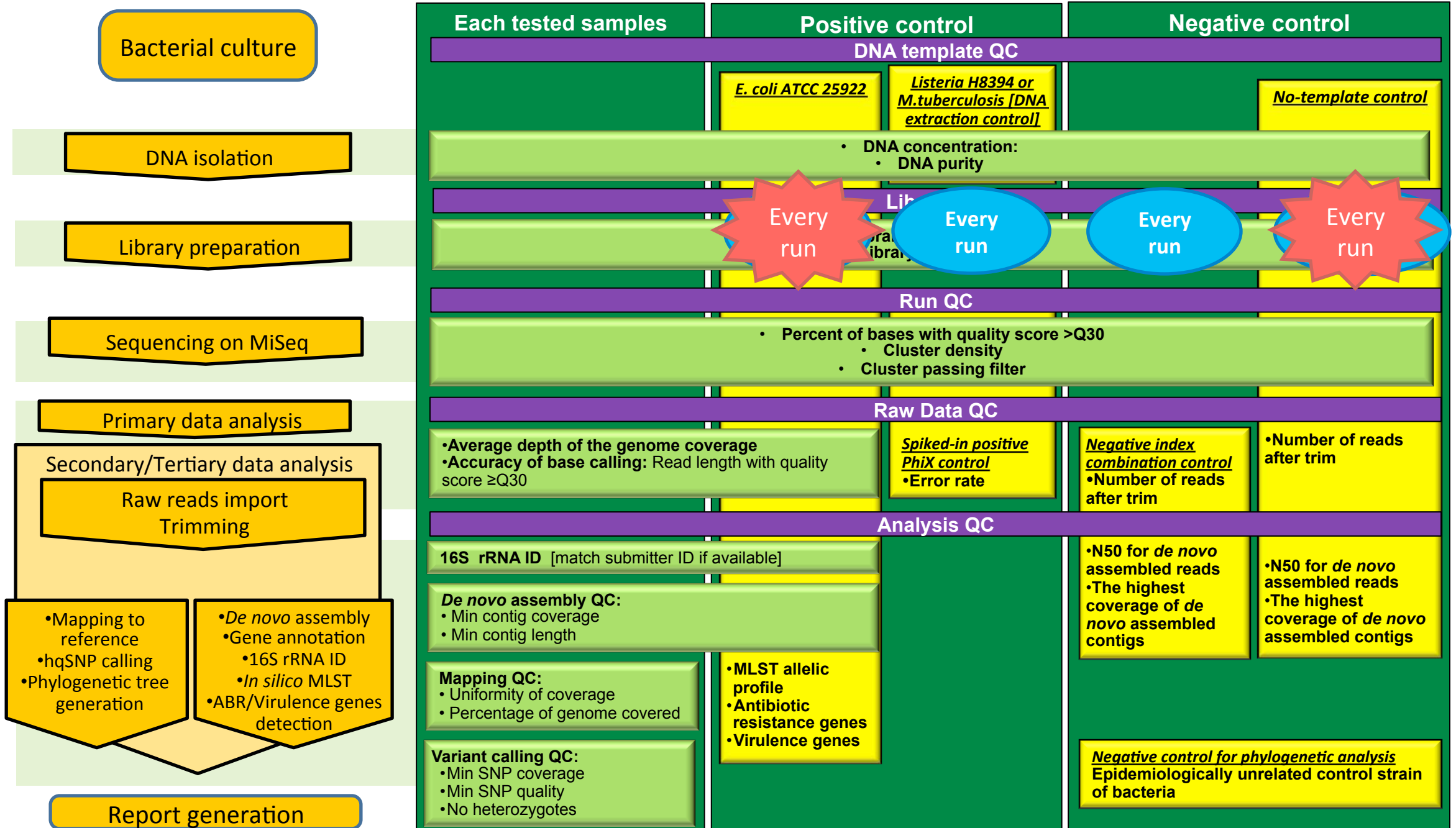
Raw reads import
Trimming

- Mapping to reference
- hqSNP calling
- Phylogenetic tree generation

- De novo* assembly
- Gene annotation
- 16S rRNA ID
- In silico* MLST
- ABR/Virulence genes detection

Report generation

WGS QUALITY CONTROL SCHEME



Bacterial culture

DNA isolation

Library preparation

Sequencing on MiSeq

Primary data analysis

Secondary/Tertiary data analysis

Raw reads import
Trimming

- Mapping to reference
- hqSNP calling
- Phylogenetic tree generation

- De novo assembly
- Gene annotation
- 16S rRNA ID
- In silico MLST
- ABR/Virulence genes detection

Report generation

Re-Validation / Confirmation of the Performance for the modified components of the WGS pipeline

- ✦ An amendment of Illumina MiSeq v.2 chemistry sequencing kits
- ✦ A new processing algorithm for highly-contagious pathogens
- ✦ Added virulence genes analysis
- ✦ Added Kmer species identification analysis
- ✦ Developed a custom script to create an automated bioinformatics pipeline
- ✦ ... Change is the only constant ...

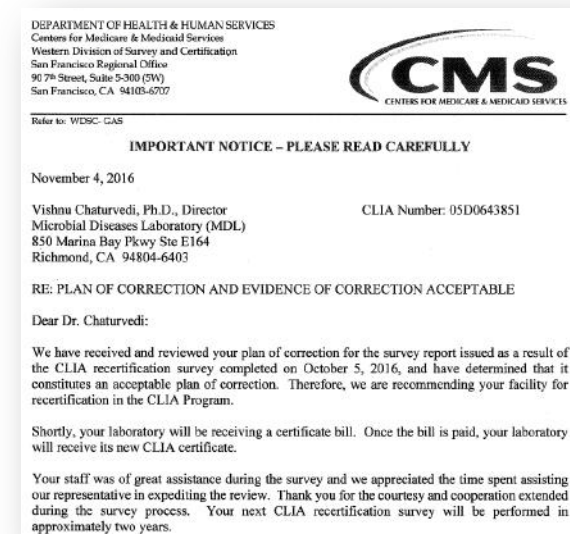


More validations to come...



Summary

- ✓ **Established the performance specifications** for WGS in the application to public health microbiology in accordance with CLIA guidelines for the LDTs
- ✓ **Developed quality assurance (QA) and quality control (QC) measures** for WGS.
- ✓ **Developed a validation set** of pathogenic bacteria for further validations of new WGS instruments/platforms, QC purposes, and multi-laboratory comparisons.
- ✓ **Developed a modular template for the validation** and re-validation of 'wet bench' and 'dry bench' components of WGS.
- ✓ **Designed laboratory reports** for end users with or without WGS expertise.
- ✓ **Successfully passed CLIA inspection!**
- ✓ Validation report is available as a supplementary material to the published manuscript in Journal of Clinical Microbiology, 2017 Aug; 55(8):2502-2520 (PMID: 28592550).



x 2

ACKNOWLEDGEMENTS

MDL Laboratory Director (former)

Dr. Vishnu Chaturvedi

Core laboratory

Chau-Linda Truong (former member)

Dr. Rituparna Mukhopadhyay

John Crandall

Dr. Matthew Sylvester

Dr. Zhirong Li

Food- and Waterborne Diseases Section

Dr. Stephanie Abromaitis

Greg Inami

Yue Qing Zhao

Francine Arroyo

Beverly Kaneko

Bacterial Diseases Section

Dr. Peng Zhang

Margot Graves (former member)

Robin Hogue (former member)

Mycobacteriology and Mycology Section

Dr. Ed Desmond

Terry Weber

Grace Lin

High Risk Pathogens Section

Dr. Jennifer Kyle

Mahtab Shahkarami

Alyssa Poe

Immunodiagnosics Section

Joseph Morgan

Frank Ni (former member)

Linda Sae-Jang

Compliance

Gillian Edwards

Collaborators outside of CDPH:

Dr. Alexander Greninger, UCSF/ University of Washington

Dr. Eija Trees, PulseNet Next Generation Subtyping Methods Unit, CDC

Dr. Heather Carleton, Enteric Diseases Laboratory Branch, CDC



**KEEP
CALM
AND
VALIDATE**