



healthy all life long

Pathogen profiling in a clinical and public health context

The use of whole genome sequencing (WGS) for the
typing of bacterial pathogenic isolates

Capita selecta in Bioinformatics
17/02/2020

Dr. ir. Sigrid De Keersmaecker – Transversal activities in Applied Genomics



Overview

- Sciensano, the Belgian Scientific Institute of Public Health
- Context pathogen profiling
- Case study: The use of whole genome sequencing (WGS) for the typing of bacterial pathogenic isolates
- Bioinformatics @ Sciensano
- Constraints, including working in a quality system (ISO)



Who is Sciensano?

- Belgian federal research centre (established April 2018)
- Merger between  & 
Scientific Institute of
Public Health - Belgium
- CODA - CERVA
Veterinary
counterpart
- >1 century of scientific expertise in human and animal health
 - 700 staff members
- Support to policy, professionals and citizens through innovative research, analyses, monitoring activities and expert advice



.be

Our core: One health and beyond

Core: “One health” concept

Holistic view on health

Combining different research perspectives and disciplines to prevent, evaluate and mitigate health problems.

Unique contribution to health





Transversal activities in Applied Genomics

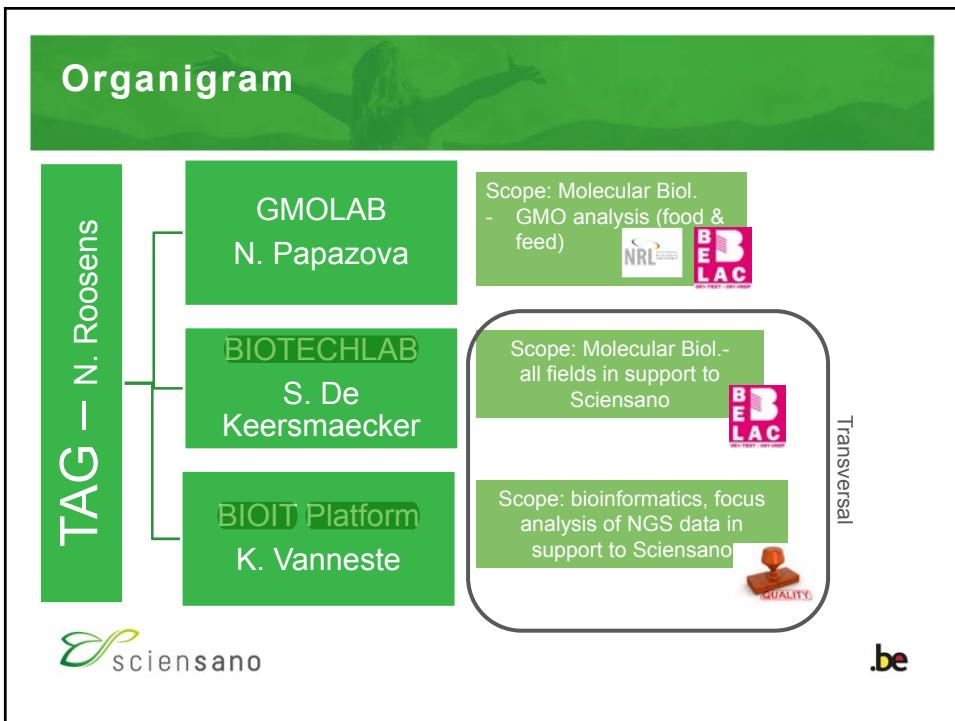
Objective? Transversal unit developing and implementing molecular biology and bioinformatics to perform both **routine** analyses, and scientific and technological research.

How? By setting up internal and external partnerships, we generate new knowledge and tools tailored to face and anticipate the current and future challenges affecting human and animal health.

Activities?

- in the domain of molecular detection, identification and characterisation of (pathogenic) micro-organisms (bacteria, viruses, and fungi), biotech organisms (GMO) and **human nucleic acid biomarkers**.
- high-tech equipment & expertise (e.g. real time & digital PCR, Sanger and NGS sequencers).
- bioinformatic tools & expertise





Sciensano in figures

700 staff members	10 patents in the field of public health	110 parliamentary questions processed annually	260 peer-reviewed articles published annually
More than 130 000 sample analyses conducted annually	30 surveys, studies and registries to monitor the evolution of public health	135 scientific projects financed by Belgian and international partners	More than 80 reference centres and labs

[.be](http://sciensano.be)

National reference centers & laboratories

Tasks: service, research and expertise for the Belgian authority

National reference centers (human) (NRC): RIZIV/INAMI – ECDC - WHO

```

graph TD
    ECDC[ECDC and WHO] <--> WIVISP[WIV-ISP: Coordinating center]
    MTAB[Medical-Technical Advisory Board (MTAB)] --> WIVISP
    RE[RIZIV/INAMI: Financial and legal support] <--> WIVISP
    RE <--> NRC[40 NRCS]
    EXP[Scientific experts: Evaluation of the NRC candidates] --> WIVISP
    CL[Clinical laboratories] --> NRC
    NS[National surveillance networks] --> NRC
    PHO[Public health officers] --> NRC
    
```

National reference laboratories (food) (NRL): FAVV/AFSCA - EFSA

[.be](http://sciensano.be)

Context pathogen profiling

Table 1 | Classes of determinants for pathogen profiling

Class of determinant	Data type	Uses
Pathogen identification	Presence of pathogen, genus and species-specific gene	Confirmation of identity of a pathogen
Virulence	Presence or absence of individual genes or mutants associated with virulence	Primary risk assessment or outcome prediction*
Transmissibility	Presence or absence of individual genes associated with transmissibility	Secondary risk assessment or outcome prediction*
Antimicrobial resistance	Presence or absence of individual genes or mutations associated with resistant phenotype	Treatment response prediction
Clonality	Genotypes and epidemiological data	Confirmation of epidemiological links or generation of hypotheses about relationships in the absence of epidemiological data ^a ; Tracking geographical and temporal spread of pathogens of public health importance
Clinical information	Patient's demographics and location, laboratory number	Unique identifier, temporal and geo-positioning

Sintchenko et al., 2007 *Nature Rev Microbiol*



Outbreak



Legionnaires' disease outbreak in Belgium affects 30, with 2 dead

A collage of three images: a close-up of a paper mill's wooden structure, a view of wind turbines and industrial buildings under a cloudy sky, and a view of a paper mill building with a tall chimney.

VRT NWS

1



150 leerlingen en leerkrachten
Spermalie Brugge ziek na besmetting
met salmonella

Tartaarsaus oorzaak salmonella-uitbraak in hotelschool Spermalie



Vers tartarsus is de bron van de salmonella-uitbraak in de Brugse hotelschool Spermelie begin september.

Op vrijdag 6 september en in de daaropvolgende dagen zijn in het hotel- en toerisme school spelen in Krugge een weekenderleiding en leerkrachten niet geweest. Uit laboratoriumonderzoek van steekplaagmuggen bleek al dat sporenlagen werden door een uitbraak van mückenpest. Daarop werden de steden gealarmeerd van de mogelijkheid dat in het schoolverblijf werden gescreven. Ook werd er een online enquête bij de leerlingen en leerkrachten opgestart om na te gaan wie in welk restaurant was gereisd had.

De resultaten van de voorbereidende en uitvoerende fase van de schoolrestauranthoreca-onderzoeken bleek al snel dat spermatische en een uitbraak van salmonella. Daarop werden de staten geanalyseerd die in het schoolrestaurant werden geserveerd. Ook werd er een calamiteiten en leerkrachten opgestart om na te gaan wie in welk restaurant

Case study:
The use of whole genome sequencing (WGS) for the typing of bacterial pathogenic isolates

SCIENTIFIC OPINION

eJ EFSA Journal

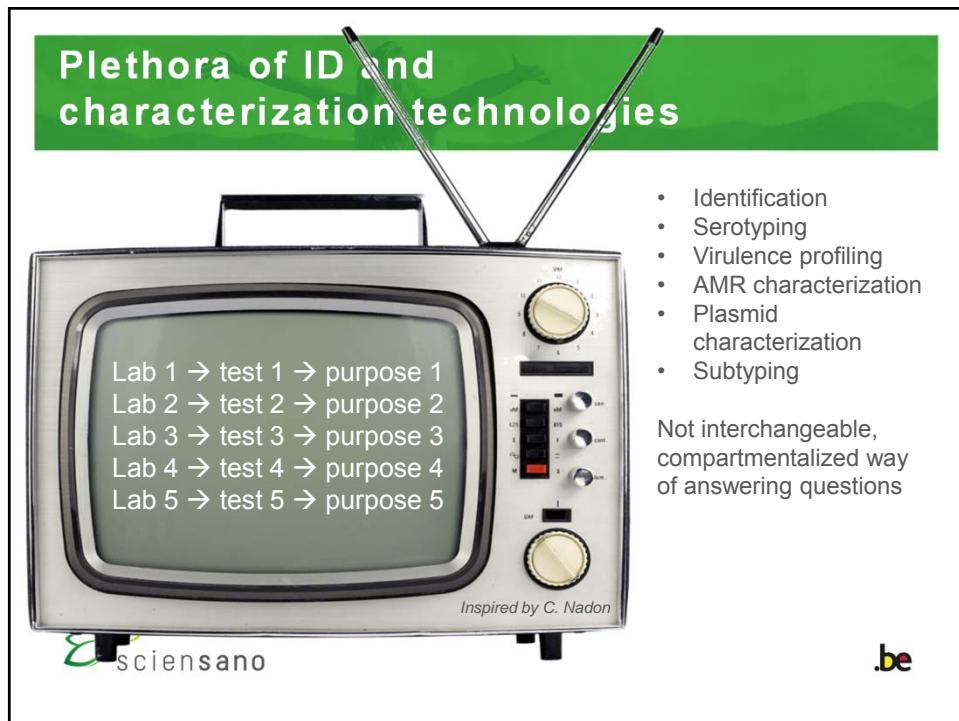
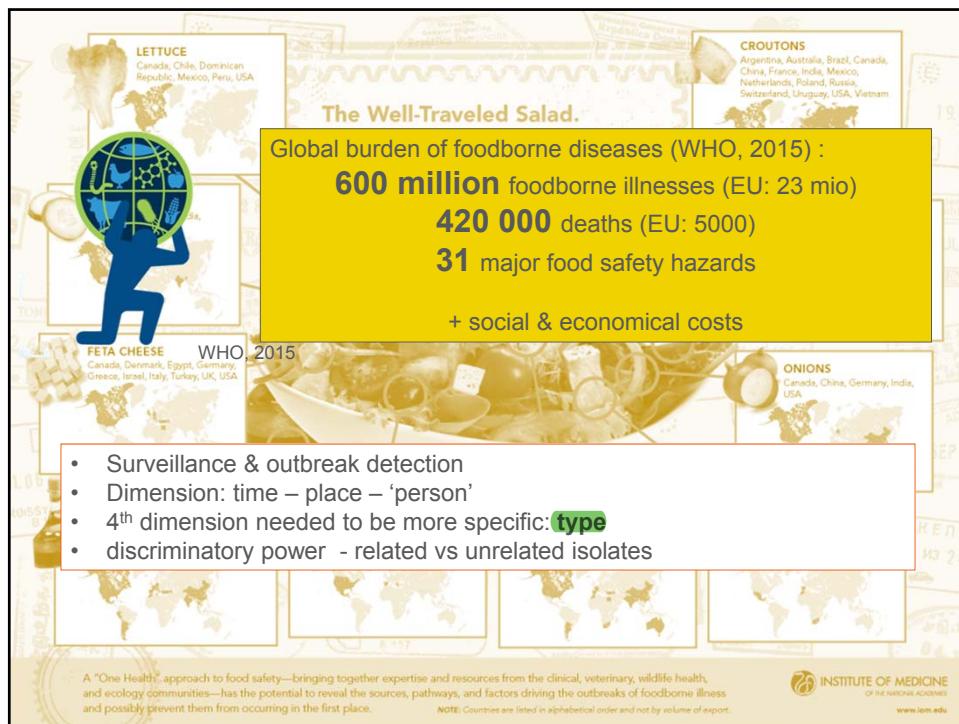
ADOPTED: 23 October 2019
doi: 10.2903/j.efsa.2019.5898

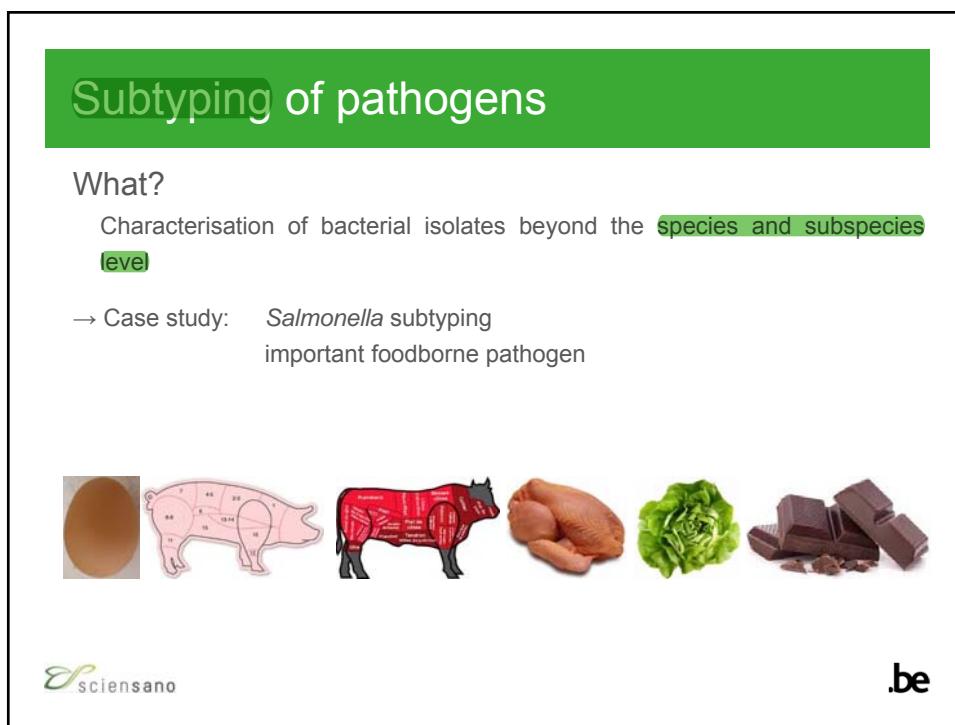
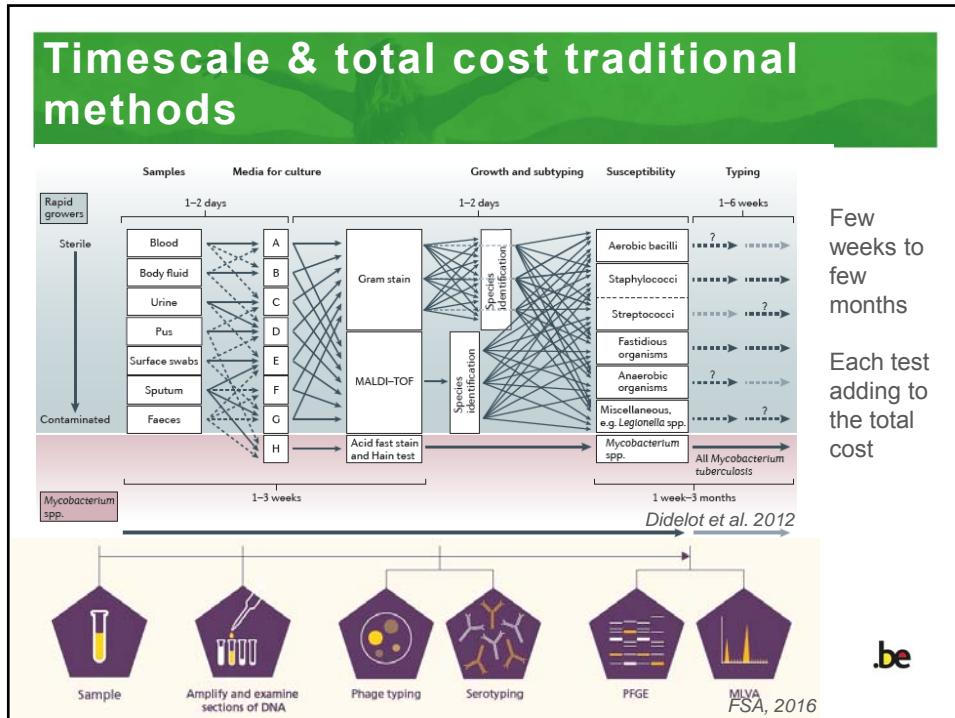
Whole genome sequencing and metagenomics for outbreak investigation, source attribution and risk assessment of food-borne microorganisms

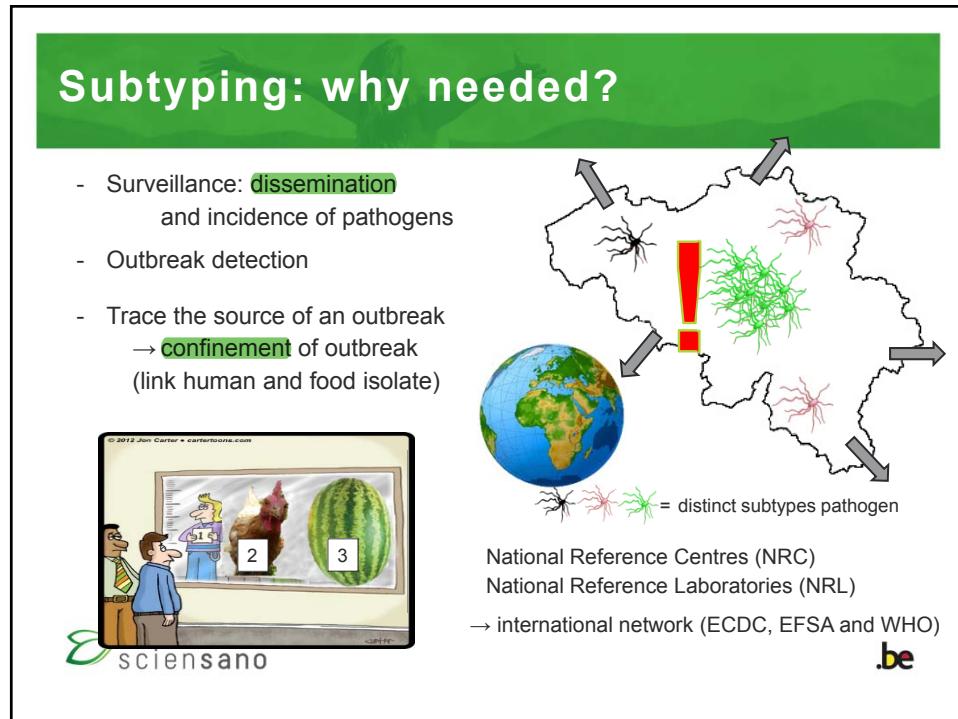
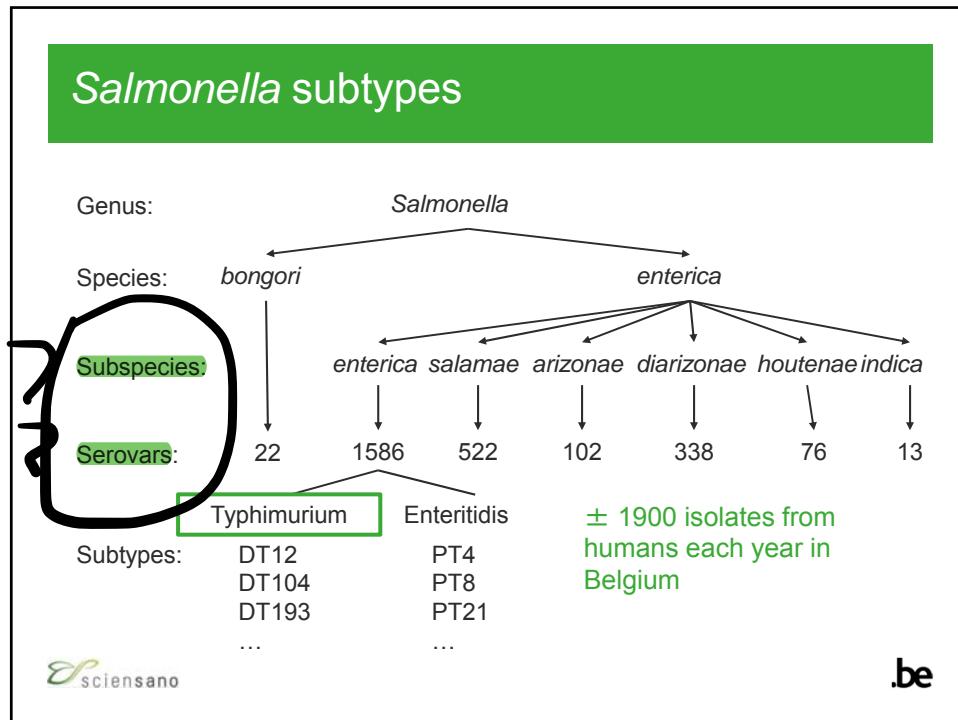
EFSA Panel on Biological Hazards (EFSA BIOHAZ Panel),
Kostas Koutsoumanis, Ana Allende, Avelino Alvarez-Ordóñez, Declan Bolton, Sara Bover-Ordóñez,
Marianne Chemaly, Robert Davies, Alessandro De Cesari, Friederike Hilbert, Roland Lindqvist,
Maarten Nauta, Luisa Peixe, Giuseppe Ru, Marion Simmons, Panagiotis Skandamis,
Elisabetta Suffredini, Claire Jenkins, Burkhard Malorny, Ana Sofia Ribeiro Duarte,
Mia Torpdahl, Maria Teresa da Silva Felicio, Beatriz Guerra, Mirko Rossi and Lieve Herman

.be









Ideal global public health molecular tool

- Robust
- Easily implemented and standardised
- Sufficiently discriminative
- Objective data → appropriate interpretation and transfer between laboratories, reproducibility & portability
- Rapid
- Inexpensive
- Universally applicable



.be

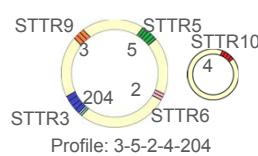
Classical subtyping methods for *Salmonella*

Phage typing



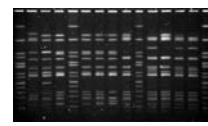
- + inexpensive
- + discriminative
- **high level of expertise required**
- interpretation of lysis patterns is **subjective**
- NT and RDNC strains

MLVA



- + rapid
- + profiles easily compared between laboratories
- relatively expensive
- only available for small number of pathogens
- data analysis with commercial software
- too discriminative for *Salmonella* Typhimurium

PFGE

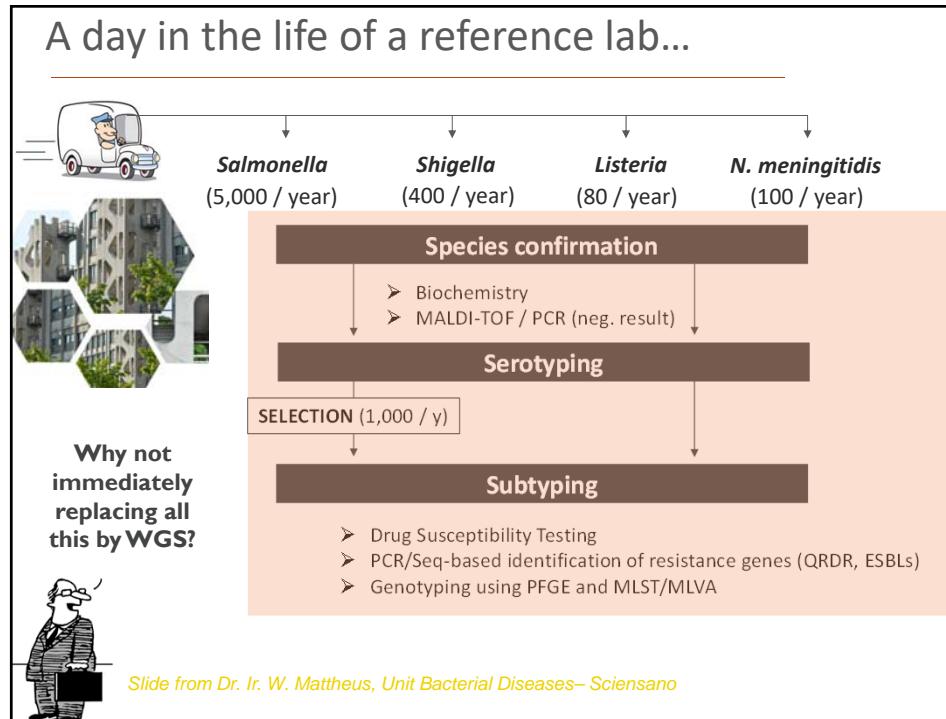
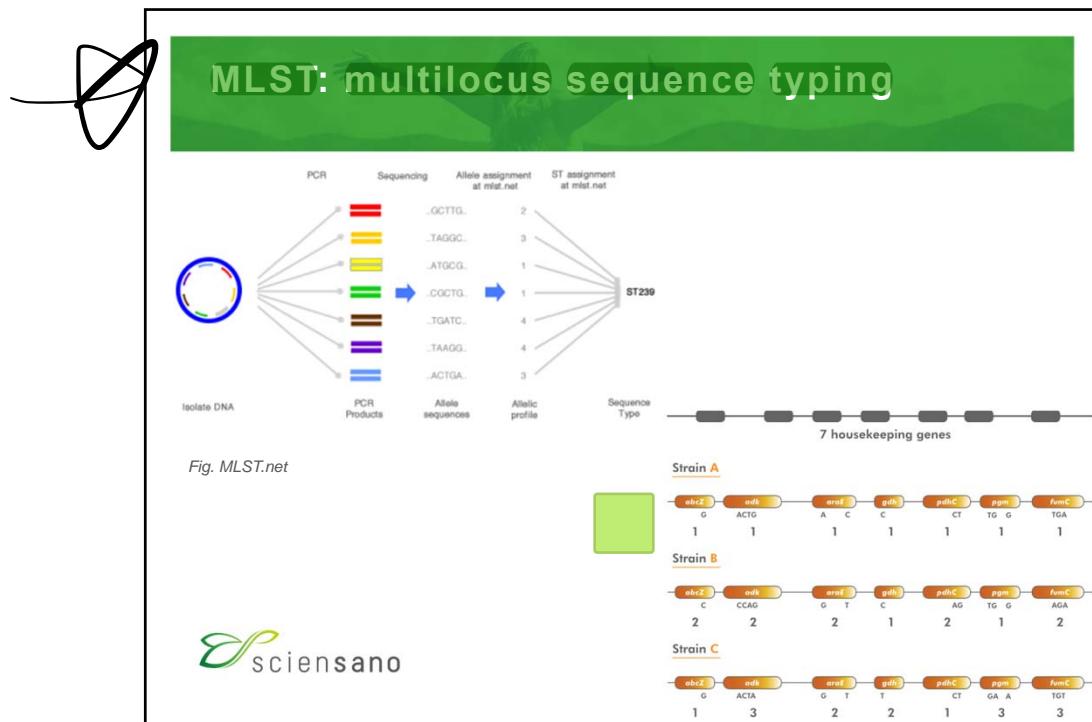


- + medium discriminative
- + public database available
- relatively slow
- limited reproducibility
- data analysis with commercial software



+ AMR, (q)PCR virulence genes....

.be



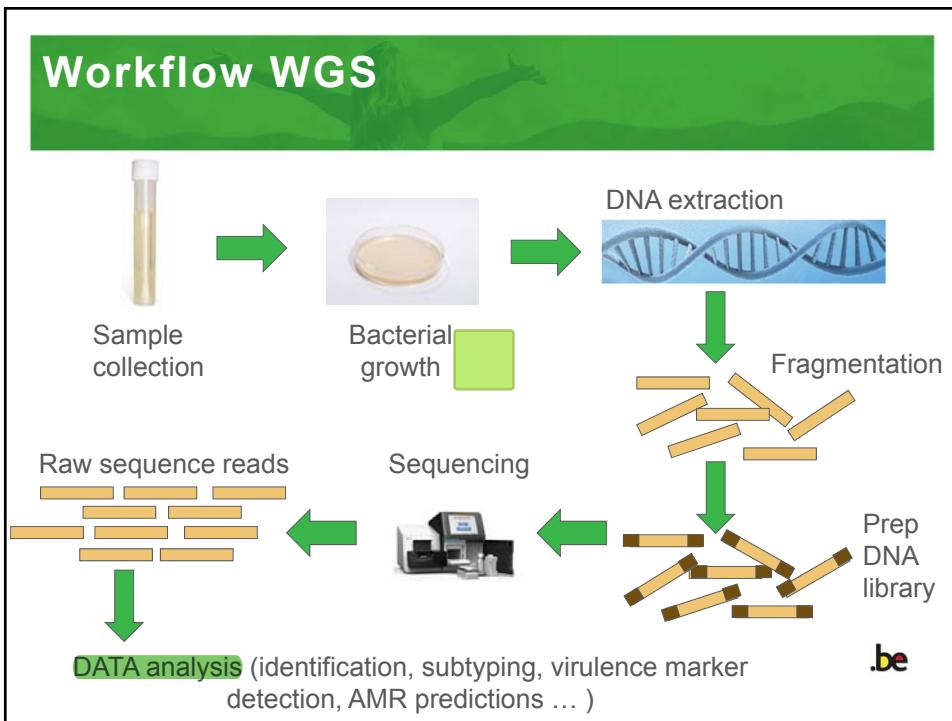
Whole genome sequencing (WGS)

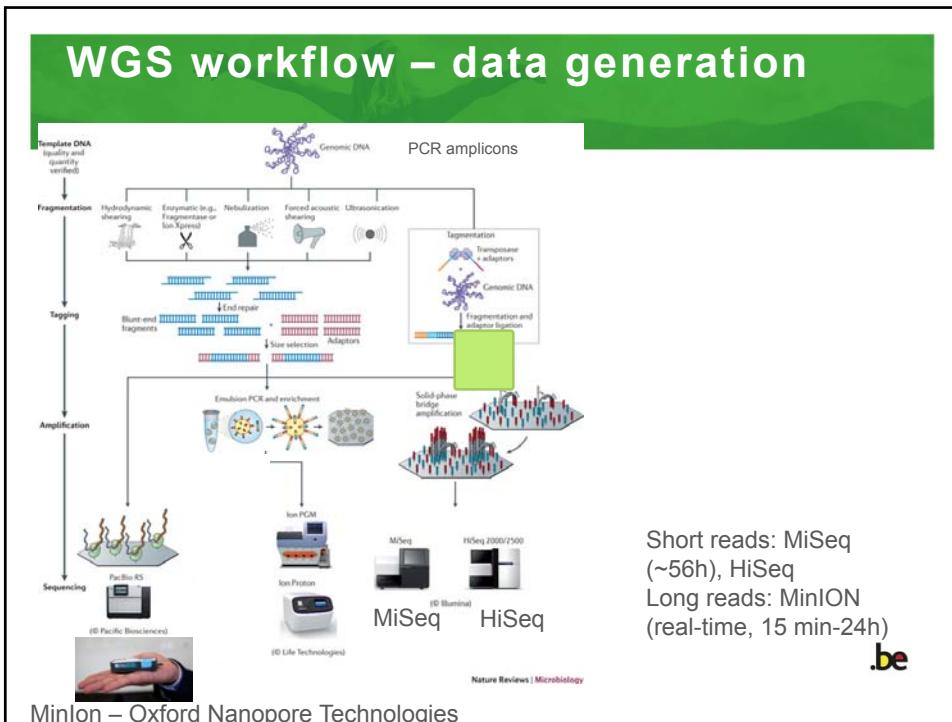
- fast
- conceptual simple workflow
- ONE universal method (i.e. species independent)
- high resolution (maximal strain discrimination)
- multiplex
- no *a priori* knowledge
- more and more affordable
- 1 test to answer >1 question at a time:
 - Identification
 - Define lineage & surveillance
 - Predict resistance
 - Serotyping
 - Virulence profiling
 - AMR characterization
 - Plasmid characterization
 - Subtyping – outbreak investigation




« One method to rule them all... »

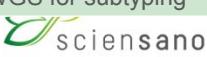


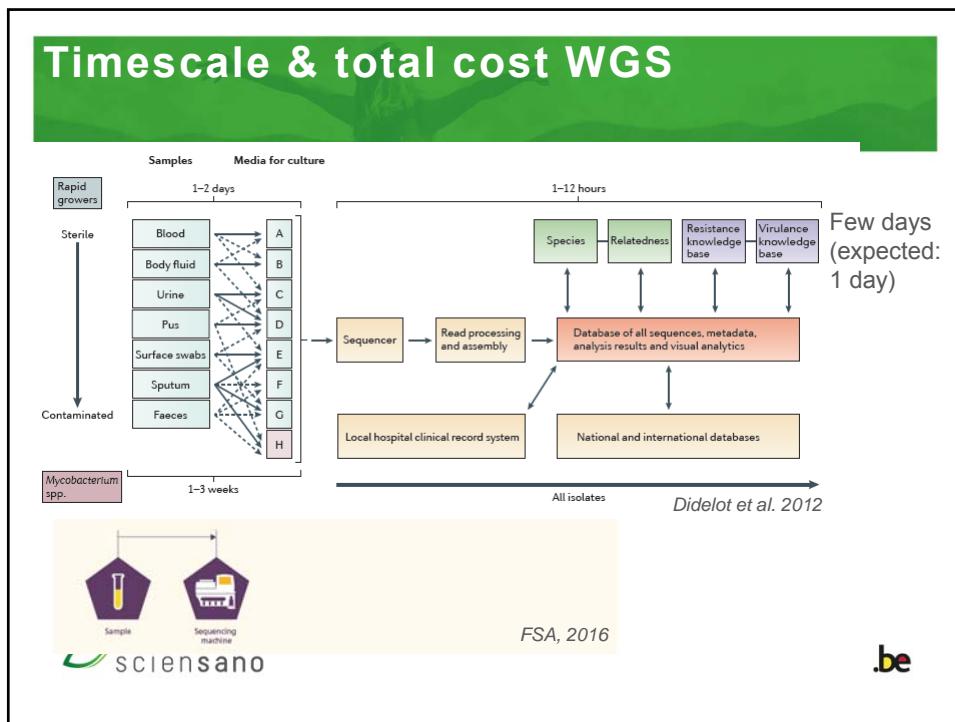




In silico typing Backward comparability

'Classical' technique	Via WGS?	Comment
PCR-based	yes	multiplex
Antimicrobial susceptibility testing	yes	Likely to be predicted from WGS with sufficient accuracy, depending on mechanism
Serotyping	yes	Likely to be predicted from WGS with sufficient accuracy
Multiple-locus VNTR analysis (MLVA)	no	Cannot yet be derived from routine WGS data, potentially in future
Multilocus sequence typing (MLST)	yes	Can be derived unambiguously from WGS data
Pulse-field gel electrophoresis (PFGE)	no	Cannot be derived from routine WGS, to be determined if possible in the future
WGS for subtyping	yes	SNP-based for high resolution

 .be



Cost and timeliness typing methods

Pathogen	Typing method	Weight	Countries	Samples per year	Consumables per isolate (EUR)	Operator time per isolate (h)	Total time per isolate (d)
<i>Salmonella</i>	Isolation, confirmation and pure culture	100%	5-7	18315-18465	1-4-15	0.1-0.5-1.5	1-1.5-4
	Serotyping agglutination*	100%	6-8	21430-21500	5-25-55	0.00-0.5-1	<1-2-17
	MLVA Typhimurium and Enteridis	67%	6-8	2395-2430	7-17.5-30	0.67-0.3-3	1-1.5-2
	PFGE XbaI	33%	6-8	1165-1435	20-25-42	0.36-0.9-1	2-3.5-6
	Resistance testing	0%	5-6	460-470	1.5-9-20	0.17-0.5-1	1-2-3
Total!	N/A	N/A	17-49-104	0.7-2-4.8	3.3-5.7-24.3		
<i>Listeria monocytogenes</i>	Isolation, confirmation and pure culture	100%	4-5	700-775	1.5-3-20	0.1-0.69-3	1-2-6
	Serotyping agglutination	33%	4-5	143-148	10-14-40	0.25-0.49-1	<1-1-2
	Serotyping PCR	67%	5-6	290-360	1.2-6.5-35	0.31-1.5	0.09-1-2
	PFGE ApaI+AspI	100%	7-8	440-515	24-30-60	1-2-2.8	3-5.6
Total!	N/A	N/A	30-42-117	1.4-3.5-8.5	4.4-8.5-14		
<i>Campylobacter</i>	Isolation, confirmation and pure culture	100%	4-5	1015-1165	1.5-3-20	0.1-0.69-3	1-2-6
	Species typing	0%	4-6	2000-2150	1-5-30	0.07-0.5-2	<1-1-1
	MLST	100%	3-4	650	25-61.6-250	3-3.5-5	2-3-5
	Total!	N/A	N/A	27-65-270	3.1-4.2-8	3-5-11	
<i>E. coli</i>	Isolation, confirmation and pure culture	100%	6-7	2960-3150	1-6-15	0.1-0.75-3	1-2-3
	vibL+vib2-hse gene detection (not subtyping)	100%	7-8	4105-4275	7-20-40	0.7-1-2	0.5-2
	O-group typing agglutination	100%	6-8	2285-2375	5-10-85	0.2-1-2.5	1-2-7
	PFGE XbaI	100%	7-8	765-800	20-25-50	0.36-0.9-1.27	3-4-12
	ESBL production	0%	2-4	303-403	10-13.5-15	0.25-0.5-2	1-1.5-2
	Total!	N/A	N/A	33-61-210	1.4-3.7-8.8	5.5-9-24	
WGS*	Isolation, confirmation and pure culture	100%	19-24	N/A	1-35	0.1-0.5-3	1-2-3
	DNA Extraction	100%	10-19	N/A	1.5-5-10	0.5-0.67-3	0.25-1-2
	Library preparation and sequencing	100%	6-8	N/A	48.3-80-151	0.09-0.48-1	2-3-12
	Data analysis	100%	-	-	-	-	-
	Total!	N/A	N/A	51-91-196	0.7-1.7-7	3.25-6-17	

Price: assumption that full batch

WGS: no standard yet

Equipment costs not included (min 17 euro per isolate for WGS)

Data storage costs not included (min 2.7 euro/isolate)

E. coli & *Campylobacter*: less costly

Listeria: cost idem

Salmonella: higher cost, depending on throughput

ECDC – Expert opinion on the introduction of next-generation typing methods for food- and waterborne diseases in the EU and EEA, 2015

Logos

- sciensano
- .be

Economic impact of using WGS for foodborne outbreaks

An economic analysis of salmonella detection in fresh produce, poultry, and eggs using whole genome sequencing technology in Canada

Sonali Jain^a, Kakali Mukhopadhyay^{b,*}, Paul J. Thomasson^a

^aDepartment of Agricultural Economics, McGill University, Canada
^bGokhale Institute of Politics and Economics, Pune 411004, India

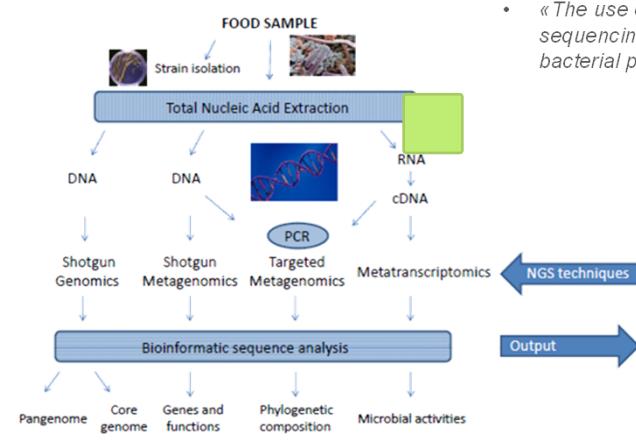
Canada	
Incidence of illness	47082 per year
Direct healthcare costs (hospitalization, physician, traditional lab costs PFGE, medication)	\$6.43 mio (\$98.7 mio value of statistical life)
Indirect costs (productivity loss)	\$21.13 mio
Federal costs	\$161.44 mio
Total illness cost traditional methods	\$287.78 mio
Total net benefit use WGS	\$90.25 mio

Belgium:
24.81 mio euro ?
(Salmonellosis only)

.be



NGS for food microbiology

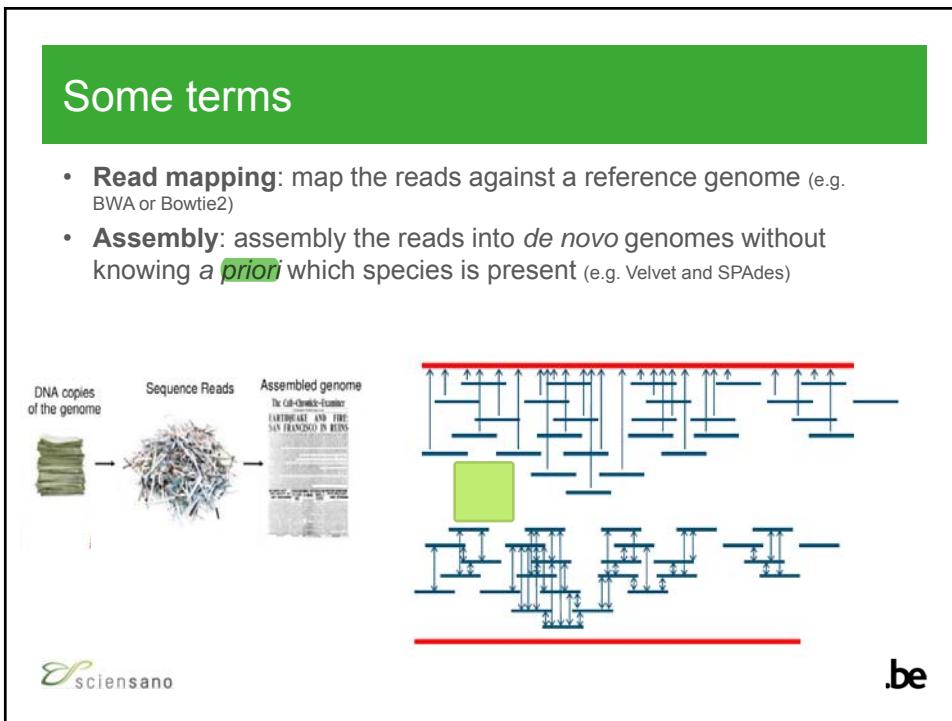
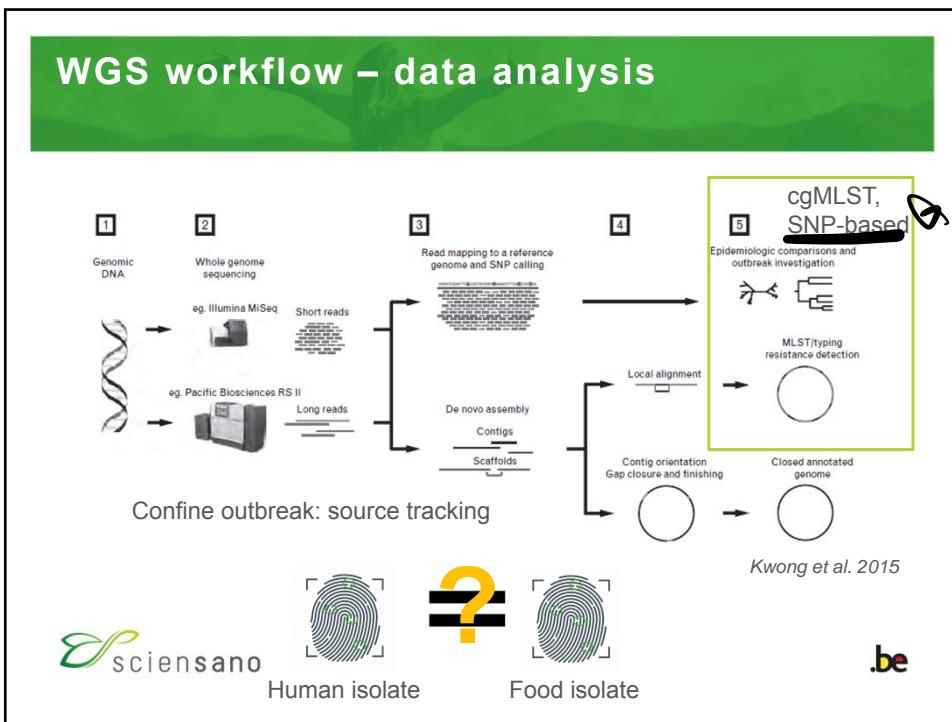


- «The use of whole genome sequencing (WGS) for the typing of bacterial pathogenic isolates »
- WGS vs NGS
- Isolate vs population
- Bacterial vs virus
- Pathogenic

Mayo et al. 2014



.be



Some terms (2)

- **Coverage**: average number of times the data is covered in the genome

$$C = N \cdot L / G$$

N = number of reads

L = read length

G = genome size (target or assembly)

Example:

N = 5 million

L = 100 bp

G = 5 Mb

$$C = 5 \cdot 100 / 5 = 100x$$

On average, 100 reads covers each position in the genome

- **Depth**: number of reads that cover a particular nucleotide in each position in the genome

Reads/site = depth

coverage

depth

Deriving phylogeny

Some considerations

- Assess whether 2 or more isolates originate from an epidemiologically relevant **common ancestor (source)** (dendograms)
- Entire genome available for analysis: not straightforward analysis & **interpretation** for epidemiologists
- Amount of diversity generated in one unit of time depends on the **rate of reproduction and the evolutionary pressure** that organisms are under (environment specific) -> highly unlikely that the number of SNPs difference found between 2 isolates can be reliably related to the number of days that have passed since they diverged from a common ancestor
- Distantly related strains may still be epidemiologically linked (several contamination foci along the same food chain)

Deriving phylogeny

Two approaches

SNP-based typing/relatedness

- Best estimate/approximation for the true evolutionary distance -> very high accuracy at the highest possible resolution (**SNP**)
 - Novel technique (databases, standards schemes etc. lacking)

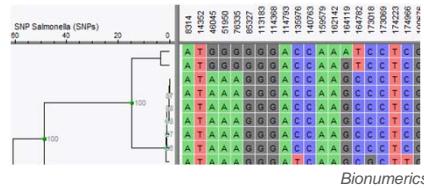
(wg)MLST-based typing/relatedness

- Lower resolution but very **robust** approach
 - MLST established technique for which WGS can be used -> expand to **wgMLST** to include core and/or pan genome (databases, standards schemes etc. however lacking – under development)



.be

Phylogeny based on SNPs



- Occur with a relatively well understood frequency
 - Derived from alignment with reference genome, also reference-free methods available (kmer-based), but less performant
 - Difference with sequencing error/mapping error -> filtering out false SNPs (e.g. consecutive SNPs, SNPs at end of read)
 - Computationally demanding
 - Most rewarding for **closely** related isolates (low proportion of false SNPs): highly accurate, highest possible resolution (e.g. after cgMLST)
 - Different algorithms/pipelines possible – different outcome?



.be

SNP analysis - examples

CSI Phylogeny 1.4 (Call SNPs & Infer Phylogeny)

CSI Phylogeny calls SNPs, filters the SNPs, does site validation and infers a phylogeny based on the concatenated alignment of the high quality* SNPs.
The service is having some issues with files compressed. Please submit all files uncompressed.

Input data

Upload reference genome (fasta format)
Note: Reference genome must not be compressed.

 Include reference in final phylogeny.

Select min. depth at SNP positions
10x

Select min. relative depth at SNP positions
10 %

Select minimum distance between SNPs (prune)
10 bp

Select min. SNP quality
30

Select min. read mapping quality
25

Select min. Z-score
1.96

Ignore heterozygous SNPs

Comment (to yourself)
This comment will appear unfiltered on your output page. It has no effect on the analysis.

Use filtered FastTree (more accurate)

Upload read files and/or assembled genomes (fasta or fastq format)
Note: If you upload a file, please make sure it is a plain text file. If you upload a compressed file, please decompress it first.
If you get an "Access forbidden" error, try this: Make sure the user of the web address is logged in and not just this. Fix it by clicking here.

File list

Name	Size	Progress	Status
<input type="button" value="Upload"/>	<input type="button" value="Remove"/>		

<https://cge.cbs.dtu.dk/services/CSIPhylogeny/>

PeerJ

An evaluation of alternative methods for constructing phylogenies from whole genome sequence data: a case study with *Salmonella*
James B. Pettengill, Yan Luo, Steven Davis, Yi Chen, Nariol Gonzalez-Escalona, Andrea Ottesen, Hugh Rand, Marc W. Allard and Errol Strain
Center for Food Safety & Applied Nutrition, U.S. Food & Drug Administration, College Park, MD, USA

.be

SNP - interpretation

- SNP Distance**

How close are the isolates? No single threshold for all species/types – need to incorporate epidemiological data

Rough, conservative guides @ FDA

 - Inclusion: <=20 SNPs match, virtually identical (examples outbreaks – at most 1 SNP apart)
 - Inconclusive: 20-100 SNPs
 - Exclusion: > 100 SNPs exclude

remark: Outbreak – 2 isolates 100 SNPs different, might not be considered related if other clustered isolates were differing by 10 SNPs. However if 90 of the 100 SNPs co-located within a small genomic segment suggestive of a recombination event, still represent evidence of transmission
- Bootstrapping**

Do the isolates form a unique cluster w/ $\geq 95\%$ support?

Is the cluster distinct from other isolates in the tree?

 **.be**

Quality trimming

SNP – different workflows, different results

PLOS ONE | <https://doi.org/10.1371/journal.pone.0198052>

RESEARCH ARTICLE
Comparison of SNP-based subtyping workflows for bacterial isolates using WGS data, applied to *Salmonella enterica* serotype Typhimurium and serotype 1,4,[5],12:i:-

Assia Saltykova^{1,2}, Véronique Wuyts¹, Wesley Mattheus³, Sophie Bertrand³, Nancy H. C. Roosens¹, Kathleen Marchal^{2,4,5*}, Sigrid C. J. De Keersmaecker¹

PhD
Assia Saltykova

Phylogeny based on allele differences (wgMLST - cgMLST)

a) Assignment of unique allele identifiers

Locus:	A	B	C	D	E	F
Isolate 1:	A9	B7	C5	D3	-	F1
Isolate 2:	A9	B7	C6	D3	E3	F9
Isolate 3:	A6	B7	C6	D3	E3	F9

Pan genome
Core genome
Accessory genome
Seven gene

b) Distance between isolates according to different schemes

Isolate pair	Seven gene	Core genome	Pan genome
1-2	0	1	3
1-3	1	2	4
2-3	1	1	1

c) Assignment of sequence types

Isolate	Seven gene	Core genome	With sublevels (AB / CD)
1	23 (A9-B7)	115 (A9-B7-C5-D3)	23-9 (A9-B7 / C5-D3)
2	23 (A9-B7)	267 (A9-B7-C6-D3)	23-15 (A9-B7 / C6-D3)
3	47 (A6-B7)	349 (A6-B7-C6-D3)	47-93 (A6-B7 / C6-D3)

Based on read mapping (better results) or assembly
Only coding regions (excludes a large amount of potentially informative genomic information, such as phages, insertion sequences, and other mobile genetic elements that may indicate direct transmission)
Allele sequences identical or different, regardless of the number of actual evolutionary events
Missing loci if bad quality WGS data
Reduced resolution, but more robust
Database needed
Pan/core/accessory genome

ECDC – Expert opinion on the introduction of next-generation typing methods for food- and waterborne diseases in the EU and EEA, 2015

.be

MLST, cgMLST, wgMLST

The diagram shows the relationship between different types of genes in a species population:

- Housekeeping genes for MLST & eMLST**: Shared by all strains.
- Core (c) genes ('present in all strains in a species')**: Shared by all strains.
- Pan-genome (wg) ('all genes in the whole population of a species')**: All genes present in the population.
- Serotyping genes**: Used for strain differentiation.
- Genes for genus/species/subspecies identification**: Used for taxonomic classification.
- Virulence genes**: Involved in pathogenesis.
- Antimicrobial resistance genes**: Involved in drug resistance.

R. Lindsey, CDC

MLST analysis - examples

Center for Genomic Epidemiology

MLST 2.0 (Multi-Locus Sequence Typing)

Select MLST configuration: [Enterobacteriaceae](#)

Please note that for four organisms, two or three different MLST schemes are available:

- Acinetobacter baumannii (Acinetobacter baumannii #1) [\[edit\]](#)
- Escherichia coli (Escherichia coli #1) [\[edit\]](#)
- Escherichia coli (Escherichia coli #2) [\[edit\]](#)
- Escherichia coli (Escherichia coli #3) [\[edit\]](#)
- Leptospiral serovars. Icterohaemorrhagiae, Pomona, Sehmi, Sehmi mumpslike K2 (modified) [\[edit\]](#)
- Leptospiral serovars. Pomona, Sehmi K2, Leptospiral K3 [\[edit\]](#)

Select type of data input: [Raw sequencing data](#) [Profile data](#)

Only data from one single isolate should be uploaded. If new sequencing reads are uploaded, KMA will be used for mapping. KMA supports the following sequencing platforms: Illumina, Ion Torrent, Roche 454, SOLiD, Galaxy NGS, and Pacific.

Comments: [View comments](#)

Please note that "Associated GenomeConfig" should be selected if you have already assembled your short sequencing reads into one continuous genome or into several contigs. In a different case, the type of short sequence reads were used to produce the genomeconfig.

Name	Size	Progress	Status
Inouye et al. <i>Genomic Medicine</i> 2014, 6:90 http://genomemedicine.com/content/6/11/90			

SOFTWARE **Open Access**

SRST2: Rapid genomic surveillance for public health and hospital microbiology labs

Michael Inouye^{1,2}, Harriet Dashnow^{3,4}, Lesley-Ann Raven¹, Mark B Schultz¹, Bernard J Pope^{5,6}, Takehiro Tomita^{2,6}, Justin Zobel⁶ and Kathryn E Holt^{1*}

Nomenclature

- Short, yet still informative human readable code for isolate
- Same code = share properties
- Causal relationship between vehicle and human cases
- Efficient (international) communication between stakeholders
- Stability of system



 sciensano

.be

wgMLST nomenclature

c) Assignment of sequence types

Isolate	Seven gene	Core genome	With sublevels (AB / CD)
1	23 (A9-B7)	115 (A9-B7-C5-D3)	23-9 (A9-B7 / C5-D3)
2	23 (A9-B7)	267 (A9-B7-C6-D3)	23-15 (A9-B7 / C6-D3)
3	47 (A6-B7)	349 (A6-B7-C6-D3)	47-93 (A6-B7 / C6-D3)

ECDC, 2015

- Fixed set of loci for a particular species or genus, agreed upon by all parties
- Core genome (cgMLST), pan genome (wgMLST) – different for each species
- Unique alleles for these loci -> unique identifier in global nomenclature database accessible to all (allele nomenclature)
- Unique identifier to each combination of alleles observed (schemes) (sequence type – 7-gene MLST; sub-sequence types, covering the next N most stable loci per top level ST) (strain nomenclature)
- Allele calling algorithm (extract allele sequences from a new sequence for all defined loci, match against known allele sequences) & strain classification algorithm (determine type based on allele identifiers)
- No phylogenetic relationships (e.g. 23-25 and 47-93 with 1 SNP in 1 of the 7 HKG -> other place in tree)

;

MLST databases

a) wgMLST schemes

Scheme	Loci
PanGenome	A-B-C-D-E-F
CoreGenome	A-B-C-D
SevenGene	A-B
Sublevels	A-B / C-D

wgMLST locus definitions

RefGenome	Locus	Location
Ref1	A	50001-51123
Ref1	B	250001-251320
Ref1	C	300001-301566
Ref1	D	321285-322284

wgMLST allele sequences

Locus	Allele	Sequence
A	A1	ATGCATTAT...
A	A2	ATGCAATAG...
B	B1	ATGGGCCTG...
B	B2	ATGGGCCTG...

wgMLST sequence-type nomenclature

Scheme	SeqType	Alleles
CoreGenome	115	A9-B7-C5-D3
CoreGenome	267	A9-B7-C6-D3
SevenGene	23	A9-B7
Sublevels	23-15	A9-B7 / C6-D3

e.g. PubMLST

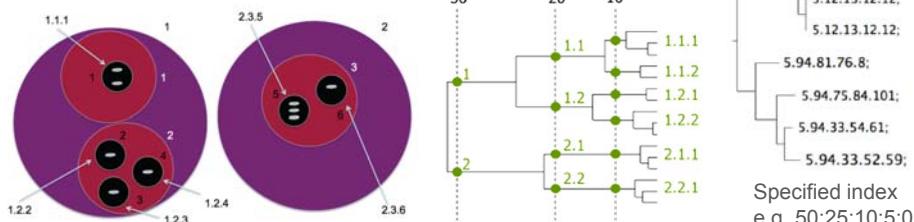
ECDC, 2015

 .be

SNP nomenclature (taxonomical)

- Phylogenetic relationships
- Hierarchical description of where a particular isolate fits in a dendrogram that spans all isolates
- SNP address (PHE)

Ashton ... Dallman et al, 2015, PHE



Height cut-offs of 50, 20 and 10 correspond to the average number of SNPs or allele differences

ECDC, 2015

Specified index
e.g. 50:25:10:5:0

Each number represents the cluster membership at each descending SNP distance threshold. The resultant 'SNP address' provides an isolate level nomenclature where two isolates with the same SNP addresses have 0 SNP differences.

.be

Snapper DB (PHE)

(A) Illustration of SNP difference between a reference sequence (top row) and a set of isolate sequences (remaining row). (B) Single linkage clustering of SNP differences into 0, 5 and 10 SNP levels for a set of isolates. (C) Examples of SNP addresses based on the seven descending SNP thresholds; 250, 100, 50, 25, 10, 5, 0

A

B

C

SNP threshold	250	100	50	25	10	5	0
Isolate 1	1	1	1	158	199	222	243
Isolate 2	1	1	1	158	199	222	256
Isolate 3	1	2	2	35	60	125	160

Bioinformatics, Volume 34, Issue 17, 01 September 2018, Pages 3028–3029, <https://doi.org/10.1093/bioinformatics/bty212>
The content of this slide may be subject to copyright: please see the slide notes for details.

OXFORD
UNIVERSITY PRESS

Bioinformatics, 34(17), 2018, 3028–3029
doi: 10.1093/bioinformatics/bty212
Advance Access Publication Date: 5 April 2018
Applications Note

OXFORD

Genome analysis

SnapperDB: a database solution for routine sequencing analysis of bacterial isolates

Timothy Dallman^{1,*}, Philip Ashton^{1,2}, Ulf Schafer³, Aleksey Jironkin³, Anais Painset¹, Sharif Shaaban⁴, Hassan Hartman¹, Richard Myers³, Anthony Underwood³, Claire Jenkins¹ and Kathie Grant¹

¹Gastrointestinal Bacteria Reference Unit, National Infections Service, Public Health England, London, NW9 5EQ, UK, ²Oxford University Clinical Research Unit, Ho Chi Minh City, Vietnam, ³Infectious Disease Informatics, National Infections Service, Public Health England, London, NW9 5EQ, UK and ⁴Roslin Institute, University of Edinburgh, Scotland, EH25 9RG, UK

MLST vs SNP		
	SNP	MLST
Epidemiological concordance	High	High
Stable nomenclature	(No)	Yes
Reference characterization: identification, serotyping, virulence & resistance markers	No	Yes
Speed	Slow SNP calling, slow analysis	Slow allele calling, fast analysis
Local computing requirements	Medium-High	Low
Local bioinformatics expertise	Yes	No
Reference used to perform analysis	Most often sequence of closely related annotated strain (reference-free also possible)	Allele database
Requires curation	No	(Yes)


frontiers
 in Microbiology

ORIGINAL RESEARCH
 published: 18 December 2019
 doi: 10.3389/fmicb.2019.02897

PhD

Assia Saltykova

OPEN ACCESS

Edited by:

David W. Ussery,

Detailed Evaluation of Data Analysis Tools for Subtyping of Bacterial Isolates Based on Whole Genome Sequencing: *Neisseria meningitidis* as a Proof of Concept

Assia Saltykova^{1,2}, Wesley Mattheus³, Sophie Bertrand³, Nancy H. C. Roosens¹, Kathleen Marchal^{2,4} and Sigrid C. J. De Keersmaecker^{1*}

¹ Transversal Activities in Applied Genomics, Sciensano, Brussels, Belgium, ² IDlab, IMEC, Department of Information Technology, Ghent University, Ghent, Belgium, ³ Belgian National Reference Centre for Neisseria, Human Bacterial Diseases, Sciensano, Brussels, Belgium, ⁴ Department of Plant Biotechnology and Bioinformatics, VIB, Ghent University, Ghent, Belgium

	Geno-by-gene	Assembly-based (pairing)	SNP-based (assembly)	SNP-based, kmer	SNP-based, mapping
Workflows:	QIIME ¹	Pairwise ²	Pairwise ²	Pairwise ²	SNAP ³ (pairwise) SNAP ³ (kmer) SNAP ³ (mapping) Repatnayake ⁴ (pairwise) Repatnayake ⁴ (kmer) Cai ⁵ (pairwise) Cai ⁵ (kmer)
I. Pre-processing					Readmapping (pairwise) (reference)
		Alignment-Pairwise (reference-based)	Alignment-Pairwise (reference-based)		
II. Detection of variants, alleles, or analysis of genomic regions	Alele ⁶ qBlast ⁷	Genomic region detection (Pairwise)	SNP detection (Pairwise)	SNP detection (SNP-SNP)	SNP detection (SNAP ³) SNP filtering (reference) Repatnayake ⁴ (reference) Cai ⁵ (reference)
III. Recombination filtering				Res. 80; Gubbins ⁸ Gubbins ⁸ (diff.)	Res. 60; Gubbins ⁸ Gubbins ⁸ (diff.)

Phenotype prediction

- Antimicrobial resistance
- Virulence
- Serotype
- Complexity of prediction dependent on cellular process involved (single gene for antimicrobial resistance vs. Multiple virulence genes in *E. coli* or AMR through mutations or *Salmonella*)
- Associated phenotypic data used to 'train' predictive engine



.be

Mining WGS for AMR

Repositories

- Webserver
 - ResFinder
 - Comprehensive Antibiotic Resistance Database (CARD)
- Offline
 - ARG-ANNOT database
- Combination (to be installed locally), e.g. SSTAR, a Stand-Alone Easy-To-Use Antimicrobial Resistance Gene Predictor

AMR efforts e.g. at NCBI

- Database of sequenced isolates with standardized AMR metadata (i.e. accept antibiograms)

[http://www.ncbi.nlm.nih.gov/biosample/?term=antibiogram\[filter\]](http://www.ncbi.nlm.nih.gov/biosample/?term=antibiogram[filter])
- Stable, up-to-date database of AMR genes with standardized nomenclature
- Implement & validate tools for identifying AMR genes in new isolates



.be

Webserver - ResFinder

Center for Genomic Epidemiology

ResFinder 3.2

ResFinder identifies acquired antimicrobial resistance genes and/or chromosomal mutations in total or partial sequenced isolates of bacteria.

ResFinder consists of two programs, ResFinder by identifying acquired genes, and PointFinder by identifying chromosomal mutations. Software and databases are available online.

ResFinder software: 3.2 (2020-02-06)
 ResFinder database: (2020-02-11)
 PointFinder software: 3.2 (2020-02-23)
 PointFinder database: (2019-07-02)

Chromosomal mutations

Resistance caused by mutations

Select species: **E. coli** (highlighted in blue)

- Enterococcus faecalis
- Enterococcus faecium
- H. pylori
- Klebsiella
- M. tuberculosis

Show unknown mutations
 Show only known mutations (highlighted in blue)
 Show all mutations, known and unknown

Select threshold for %ID: 90 %

Select minimum length: 60 %

Acquired antimicrobial resistance genes

Select Antimicrobial configuration

Select multiple items, with Ctrl-Click (or Cmd-Click on Mac) - by default all databases are selected

- Aminoglycoside
- Beta-lactam
- Collistin
- Fluoroquinolone
- Fosfomycin
- Fusidic Acid

Select threshold for %ID: 90 %

Select minimum length: 60 %

Assembled Genome/Contigs*

- 454 - single end reads
- 454 - paired end reads
- Illumina - single end reads
- Illumina - paired end reads** (highlighted in blue)
- Ion Torrent
- SOLID - single end reads
- SOLID - paired end reads
- SOLID - mate pair reads

ResFinder results

ResFinder-2.1 Server - Results

Aminoglycoside						
Resistance gene	%Identity	Query/HSP length	Contig	Position in contig	Predicted phenotype	Accession number
<i>aadD</i>	99.74	771 / 771	scaffold28_size1986	574..1344	Aminoglycoside resistance Alternate name: ant(4')-Ia and <i>aadD2</i>	AF181950

Beta-lactam						
Resistance gene	%Identity	Query/HSP length	Contig	Position in contig	Predicted phenotype	Accession number
<i>b₁T_{EM-116}</i>	100.00	861 / 861	scaffold25_size2701	605..1465	Beta-lactam resistance	AY425988

No resistance genes found.

No resistance genes found.

No resistance genes found.

Sciensano .be

ResFinder – database updates

Database Updates (Acquired antimicrobial resistance)

[The ResFinder database download site](#)

- 01-Oct-2019 Add tet(X3) and tet(X4) to tetracycline db
- 28-Apr-2019 Add mcr-9.1 to colistin db
- 22-Feb-2019 Updates of fusidic acid db
- 20-Feb-2019 Updates of blaCARB and blaFRI genes in beta-lactam db
- 28-Nov-2018 Updates and a few format changes in several db
- 10-Sep-2018 colistin db was updated to mirror the official mcr nomenclature (JAC doi:10.1093/jac/dky262)
- 10-Aug-2018 blaIMP-64 sequence was corrected in beta-lactam db
- 03-Aug-2018 Genes encoding NDM, VIM, KPC, IMP, TEM, SHV, CTX-M, CMY, ACC, DHA in beta-lactam db were revised and updated
- 04-Jul-2018 mcr genes in beta-lactam db were revised and updated
- 04-Jul-2018 New genes added in quinolone, macrolide, oxazolidinone, phenicol and sulphonamide db
- 04-Jul-2018 Glycopeptid db was restructured
- 04-Jul-2018 All databases were formatted and corrected for duplicates
- 21-Mar-2018 Major updates tetracycline db
- 21-Mar-2018 mcr-4_2_MG026621, mcr-5_2_MG384740 and mcr-7_1_MG267386 added to colistin database
- 19-Feb-2018 Major updates and corrections in quinolone and colistin database
- 30-Nov-2017 General updates and corrections
- 14-Sep-2017 Fosfomycin database was updated with 20 genes
- 25-Aug-2017 Colistin database was updated with the genes mcr-4 and mcr-5
- 25-Aug-2017 Beta-lactam database was updated with the gene blaOXA-427
- 25-Aug-2017 Beta-lactam database entries for blaSHV-5, blaSHV-12 and blaSHV-129 were corrected
- 25-Aug-2017 Sulphonamide database entries for sul3 were corrected
- 04-Jul-2017 Colistin database was updated with the gene mcr-3
- 03-Jul-2017 ResFinder database updated
- 17-Feb-2017 62 new genes in beta-lactam, oxazolidinone and quinolone databases. Removed duplicates in phenicol and oxazolidinone data
- 02-Nov-2016 Beta-lactam database was updated with the gene blaBRO-1 [beta-lactamase](#)
- 02-Nov-2016 Beta-lactam database was updated with the gene blaBRO-2 [beta-lactamase](#)
- 02-Nov-2016 Beta-lactam database was updated with the gene blaBKC-1 [beta-lactamase](#)
- 01-Sep-2016 Colistin database was updated with the gene mcr-1.2. [colistin](#)
- 08-Jul-2016 Colistin database was updated with the gene mcr-2. [colistin](#)
- 21-Jun-2016 Default threshold (%ID) changed to 90%
- 02-Jun-2016 blaVCC-1 was added to [beta-lactamase](#) and 13 new variant of oprA to [oxazolidinone](#) and [phenicol](#)

Limitations of using WGS for predicting resistance

• Technical

- Difficulty in closing plasmids
- Exclusion of repeat regions with short reads
- Missing IS that might interrupt regulatory loci
- Flawed phenotyping
- Incorrect breakpoint

• Scientific (incomplete understanding of resistance)

- Chromosomal mutations that alter gene expression leading to:
 - Modified expression of OMPs
 - Modified expression of efflux pumps
 - Changes to the LPS
- Possibility to predict susceptibility in the presence of an undiscovered determinant
- Genetic instability (e.g. loss of plasmids)
- Heteroresistance (e.g. Staph, MTB)
- Resistance from loss of function in pro-drug activators (any nt possible)

Virulence

Center for Genomic Epidemiology

VirulenceFinder 2.0

Software version (2020-02-06)
Database version (2019-10-12)

Select species
Listeria
Escherichia coli
Enterococcus

Select threshold for %ID
90 %

Select minimum length
60 %

Select type of your reads
Only data from one sample include should be uploaded. If you have several samples, please upload them in separate files. Roche 454, SOLiD, 4x Assembled or Draft Genome/Contig* (fasta)

Choose File(s)
Name
@Upload @Remove

Identification of acquired virulence genes

Center for Genomic Epidemiology

PathogenFinder 1.1

View the [version history](#) of this server.

Choose the phylum or class of your organism:
choose 'All' if you want to use the model created using all bacteria
Automatic Model Selection

Sequencing Platform
Select the sequencing platform used to generate the uploaded reads. (Note: Select 'Assembled Genome' if you are uploading preassembled reads)
Proteome

Choose File(s)
Name
@Upload @Remove

Prediction of a bacteria's pathogenicity towards human hosts



Plasmids

Center for Genomic Epidemiology

PlasmidFinder 2.1

Software version 2.1 (2020-02-06)
Database version (2019-01-02)

Select Database
List of databases available for this tool

Select threshold for minimum % identity
90 %

Select minimum % coverage
60 %

Select type of your reads
Only data from one sample include should be uploaded. If more sequencing reads are uploaded KMA will be used for mapping. KMA supports the following sequencing platforms: Illumina, Roche 454, SOLiD, Oxford Nanopore, and Pacific. Assembled or Draft Genome/Contig* (fasta)

Phases note that "Assembled Genome/Contig" should be selected if you have already assembled your short sequencing reads into one continuous genome or into several contiguous genomes. A contig is defined as one or several contigs in one FASTA file (one entry per contig). It is irrelevant which type of short sequence reads were used to produce the genome/contig.

Choose File(s)
Name
Size Progress Status
@Upload @Remove

Detecting & classifying variants of known plasmids based on similarity with plasmids present in plasmid database – not able to find new plasmids

plasmidSPAdes: Assembling Plasmids from Whole Genome Sequencing Data

Extended abstract

Dmitry Antipov¹; Nolan Hartwick², Max Shen³, Mikhail Raiko², Alla Lapidus¹ and Pavel A. Pevzner^{1,2}

¹Center for Algorithmic Biotechnology, Institute for Translational Biomedicine, St. Petersburg State University, St. Petersburg, Russia
²Department of Computer Science and Engineering, University of California, San Diego, USA
³Bioinformatics and Systems Biology Program, Massachusetts Institute of Technology, USA

- <http://spades.bioinf.spbau.ru/plasmidSPAd>

 .be

Long read sequencing



 sciensano

.be

Short read versus long read sequencing

Input DNA 

Short reads 

Adapted from Peona et al 2018

Long reads 


Distance Brussels-Midi station – Sciensano
(STEC genome)
5.500.000 bp

 Autobus
(long reads)

 Smartphone
(short reads)
150 bp - 300 bp

Putting things in perspective



Taken from W. Timp

Building a genome is like putting together a puzzle

- Larger (blurry) pieces = easier puzzle
- Small pieces = hard to put together, can't figure out « blue sky »
- Small + large pieces = clearly and easily to put together

.be

Hybrid assembly

NGS strategy for plasmid reconstruction: combining short and long reads

For the reconstruction of plasmids with NGS the Illumina (MiSeq) and Nanopore (MinION) technologies will be used:

MiSeq:

- Short reads (25-250 bp)
- + High accuracy (99.9%)
- + Affordable
- + Standardized protocols

MinION:

- + Long reads (>8000 bp)
- Low accuracy (75-90%)
- Expensive (~1000 euro for 1 flowcell)
- Protocols constantly evolving and improving
- Stringent requirements for input DNA

Title
Combining short and long read sequencing to characterize antimicrobial resistance genes on plasmids applied to an unauthorized genetically modified *Bacillus*

Authors
Bas Berbers^{1,2}, Assia Saltykova^{1,2}, Cristina Garcia-Graells³, Patrick Philipp⁴, Fabrice Arella⁴, Kathleen Marchal^{2,5}, Raf Winand¹, Kevin Vanneste¹, Nancy H. C. Roosens^{1,f}, Sigrid C.J. De Keersmaecker^{1,e}

¹Transversal activities in Applied Genomics, Sciensano, Brussels, Belgium; ²Department of Information Technology, IDLab, Ghent University, IMEC, Ghent, Belgium; ³Foodborne Pathogens, Sciensano, Brussels, Belgium; ⁴Service Commun des Laboratoires, Illkirch-Graffenstaden, France; ⁵Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent, Belgium.

A. MiSeq D. Hybrid
B. MinION E. PhD Bas Berbers
C. G.

(User-friendly) bioinformatics (1/3)

- Commercial packages
 - Bionumerics
 - CLC Genomics Workbench
 - SeqSphere
 - ...
 - **Advantages**
 - Very powerful
 - Graphical user interface -> more user-friendly
 - **Disadvantages**
 - Licenses -> expensive
 - Black box
 - Limitations in applications
 - Might be complicated to learn how to use

The screenshot shows the homepage of the Center for Genomic Epidemiology. At the top, there is a green header bar with a logo on the right. Below it is a red banner with the text "Center for Genomic Epidemiology" and the ".be" logo. The main content area has a white background with a dark grey navigation bar at the top containing links for Home, Organization, Project, Services, and Contact. On the left, there is a sidebar with sections for "Services" (listing Pipeline, Phenotyping, and Phylogenetics) and "News" (with two recent articles). The central part of the page features a world map with various callout boxes. A vertical bar on the left is labeled "CLINIC SIDE" and "SERVICE SIDE". The bottom of the page includes the "sciensano" logo and the ".be" logo.

(User-friendly) bioinformatics (2/3)

- Webservices
 - e.g. Resfinder
 - Integrated pipeline (BAP): <https://cge.cbs.dtu.dk/services/CGEpipeline-1.1/>
 - PubMLST
- **Advantages**
 - Very powerful
 - Graphical user interface -> more user-friendly
 - Easy to use
- **Disadvantages**
 - Black box
 - Limitations in applications
 - Security? Privacy?
 - Accessibility/Availability at all times?

The slide has a green header bar with the title "(User-friendly) bioinformatics (2/3)". The main content area contains a bulleted list under the heading "Webservices". It includes a link to the CGE pipeline. The slide also lists advantages and disadvantages of such services. At the bottom, there is the "sciensano" logo and the ".be" logo.

(User-friendly) bioinformatics (3/3)

- « in-house » pipelines/databases, to be installed/implemented locally
 - Open source tools (e.g. SRST) (less expensive)
 - **Linux** -> use Galaxy as interface to use Linux-based tools in a command-line free way, centralized resource (browser/network (nothing locally), multiple simultaneous users, scalable)
 - Push on the button pipelines, customized to pathogen
 - **Advantages in-house**
 - Both entire system and individual processes known and controlled (regulated)
 - Ability to provide a quick response (accessibility)
 - No need to send data to external servers (privacy, security)
 - In-house controlled and tracked databases and code (traceability)
 - No dependencies on external providers (availability)
 - Can be tailored to specific needs (flexibility)
 - **Disadvantages**
 - Requires more development & maintenance compared to commercial applications



.be

PLOS ONE

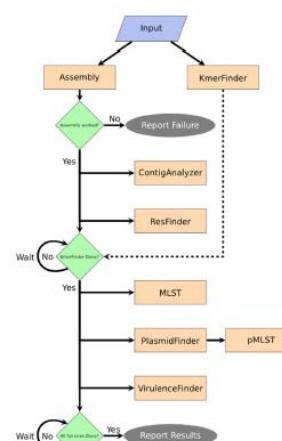
PLOS ONE | DOI:10.1371/journal.pone.0157718 June 21, 2016

A Bacterial Analysis Platform: An Integrated System for Analysing Bacterial Whole Genome Sequencing Data for Clinical Diagnostics and Surveillance

Martin Christen Frelund Thomsen^{1*}, Johanne Ahrenfeldt¹, Jose Luis Bellod Cisneros¹, Vanessa Jurz², Mette Voldby Larsen¹, Henrik Hasman², Frank Møller Aarestrup², Ole Lund¹

¹ Department of Systems Biology, Technical University of Denmark, Kemitorvet Building 208, 2800 Kgs. Lyngby, Denmark, ² National Food Institute, Technical University of Denmark, Sætlets Plads Building 221, 2800 Kgs. Lyngby, Denmark

<https://cge.cbs.dtu.dk/services/CGEpipeline-1.1/>



.be

Combine tools to make pipelines

Tweak parameters

Access using browser (simultaneous usage possible)

Standard tools

Custom tools

Workflows
Galaxy @ BIO-IT Sciensano – slide K. Vanneste

https://galaxy.wib.iop.be/workflow/editor?id=2815a520f5604ce#

Optimized ‘push-on-button’ pipelines

Centralized computational infrastructure

User-friendly access

Trade-off between quality and speed (outbreak vs. surveillance)

Automatically updated and traceable databases

Validated & optimized parameters

Galaxy @ BIO-IT Sciensano – slide K. Vanneste

**Added value of WGS – routine NRL/NRC
Salmonella Enteritidis outbreak**

9/12/2015 Whole Genome Sequence Analysis of *Salmonella Enteritidis* PT4 Outbreaks from a National Reference Laboratory's Viewpoint – PLOS Currents O...

PLOS CURRENTS OUTBREAKS

Whole Genome Sequence Analysis of *Salmonella Enteritidis* PT4 Outbreaks from a National Reference Laboratory's Viewpoint

September 11, 2015 · Research

Citation

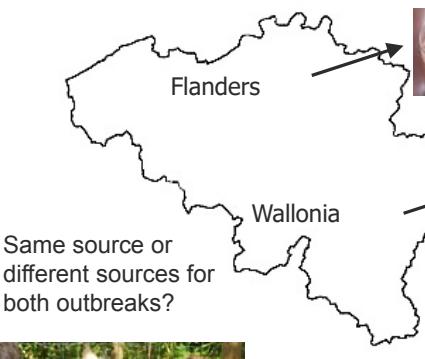
Wuyts V, Denayer S, Roosens NH, Mattheus W, Bertrand S, Marchal K, Dierick K, De Keersmaecker SC. Whole Genome Sequence Analysis of *Salmonella Enteritidis* PT4 Outbreaks from a National Reference Laboratory's Viewpoint. PLOS Currents Outbreaks. 2015 Sep 11 . Edition 1. doi: 10.1371/currents.outbreaks.aa5372d90826e6cb0136ff66bb7a62fc.



PhD Véronique Wuyts - SalMoType

**April – May 2014:
2 foodborne outbreaks**



Same source or different sources for both outbreaks?

 Social event: ± 220 guests
Food: catering service
Disease onset: April 23rd, 2014
No. cases: 45
No. hospitalised cases: 5

 Social event: ± 300 guests
Food: bbq by volunteers
Disease onset: May 1st, 2014
No. cases: ± 40
Hospitalised cases: yes

 Non-commercial eggs

 Human cases linked to egg-containing food samples



Microbiological investigation

Outbreak Flanders:

- NRL-FBO: 2 chocolate mousse samples tested positive for *Salmonella*
→ prepared with non-commercial eggs
- NRCSS: 11 human isolates received
→ 2 isolates randomly selected for further testing



Outbreak Wallonia:

- NRL-FBO: raw egg tested positive for *Salmonella*
→ non-commercial eggs used to prepare tiramisu
- NRCSS: 1 human isolate received



All outbreak isolates:

serotyped as *Salmonella enterica* subsp. *enterica* serovar Enteritidis



.be

Microbiological investigation

Outbreak	Isolate	Origin	Phage type	MLVA	Antimicrobial resistance
Flanders	S14FP01640	Chocolate mousse ^a	PT4	3-10-5-4-1	Colistin
Flanders	S14FP01642	Chocolate mousse ^a	PT4	3-10-5-4-1	Colistin
Flanders	S14BD01605	Human	PT4	3-10-5-4-1	Colistin
Flanders	S14BD01672	Human	PT4a	3-10-5-4-1	Colistin – ampicillin
Wallonia	S14FP01877	Raw egg ^b	PT4	3-10-5-4-1	Colistin
Wallonia	S14BD01753	Human	PT4	3-10-5-4-1	Colistin

^a Prepared with non-commercial eggs from private laying hens

^b Non-commercial egg from private laying hen, used to prepare tiramisu



.be

WGS analysis

WGS:

- Illumina HiSeq 2000
- 100 bp paired-end reads
- 40-plex

→ Raw FASTQ reads



WGS data analysis on a Windows 7 platform

- CLC Genomics Workbench:
 - quality trim of raw FASTQ reads
 - *de novo* assembly → contigs of chromosome and plasmids
 - read mapping on *S. Enteritidis* P125109 chromosome (NC_011294)
 - + *de novo* assembly of unmapped reads → contigs of plasmids
- Tools on server of the Center for Genomic Epidemiology (DTU):
 - MLST – ResFinder – PlasmidFinder
 - BRIG (BLAST Ring Image Generator)



.be

Center for Genomic Epidemiology: MLST – ResFinder - PlasmidFinder

Outbreak	Isolate	Origin	MLST ^a	ResFinder	PlasmidFinder
Flanders	S14FP01460	Chocolate mousse	ST-11	-	pSLA5
Flanders	S14FP01642	Chocolate mousse	ST-11	-	pSLA5
Flanders	S14BD01605	Human	ST-11	-	pSLA5
Flanders	S14BD01672	Human	ST-11	<i>bla</i> _{TEM-1B}	pSLA5 – pSD107
Wallonia	S14FP01877	Raw egg	ST-11	-	pSLA5
Wallonia	S14BD01753	Human	ST-11	-	pSLA5

^a *aroC-dnaN-hemD-hisD-purE-sucA-thrA*

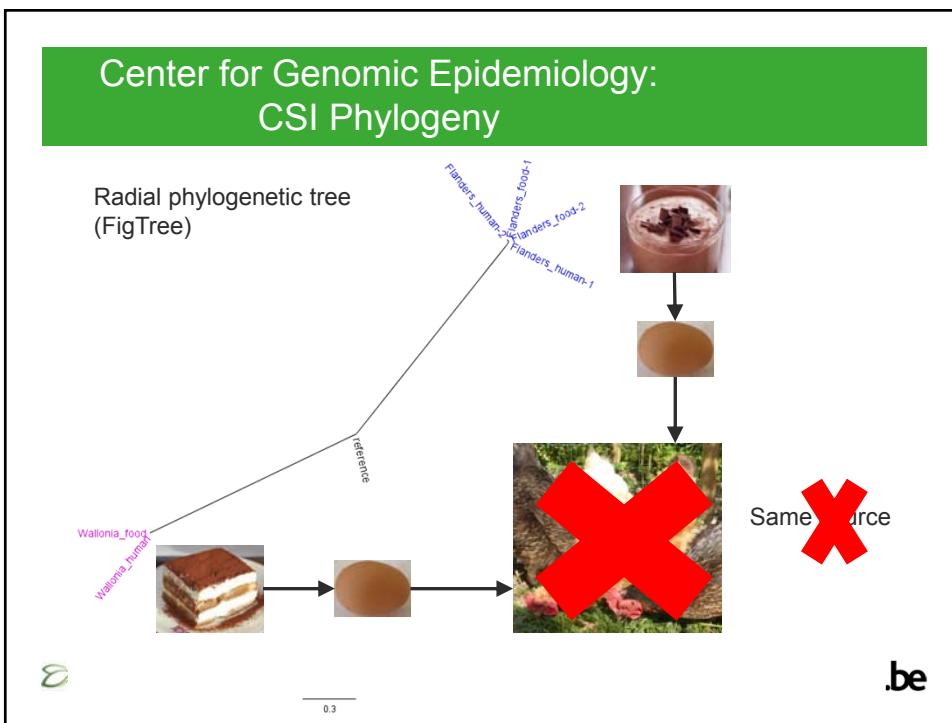
1 analysis → similar results as traditional methods

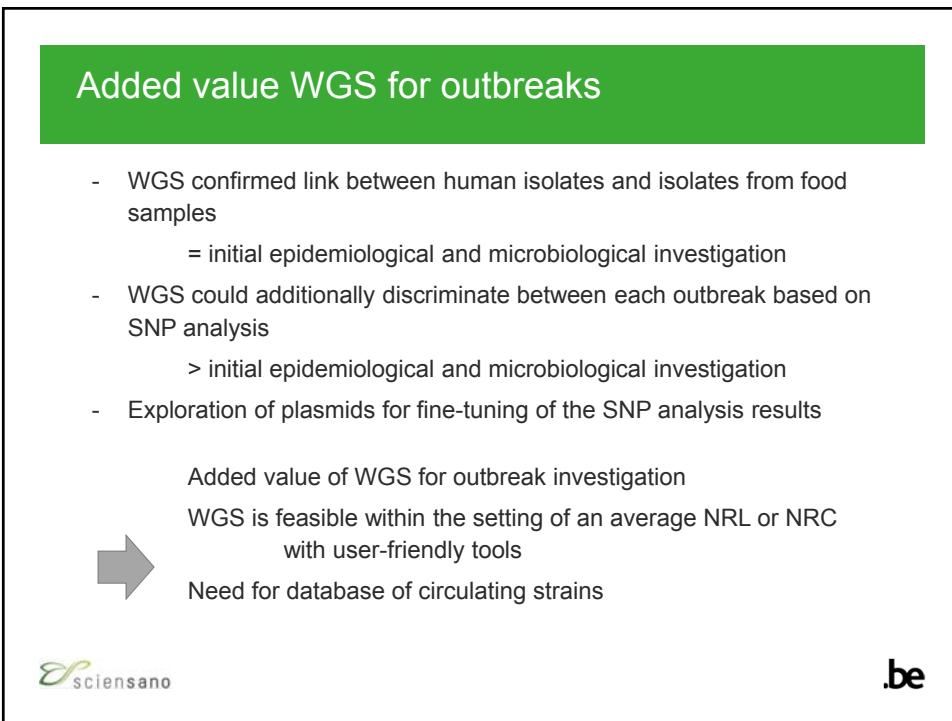
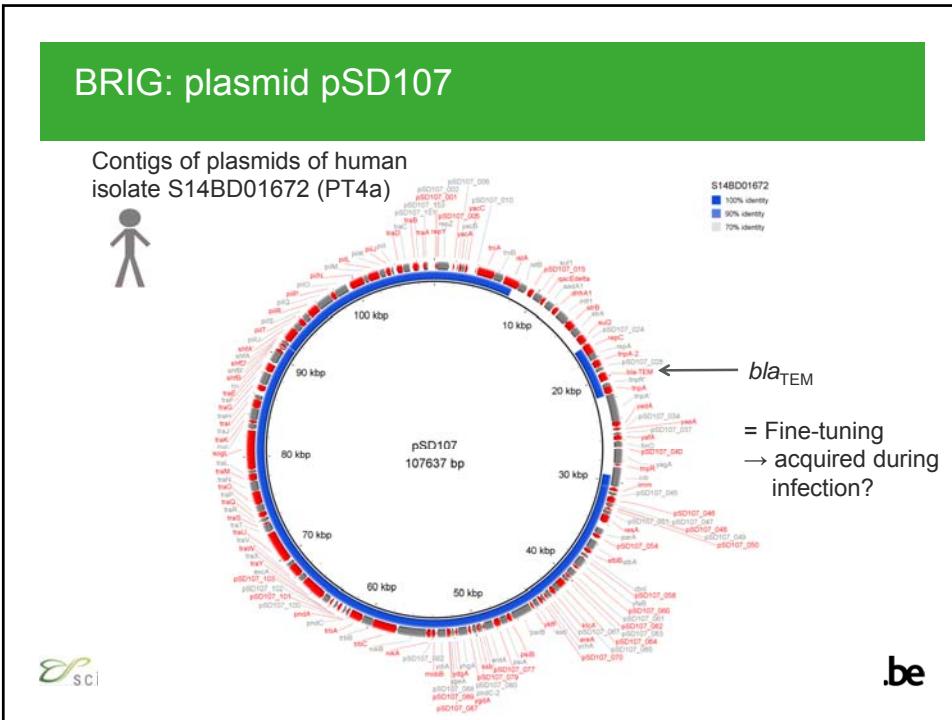


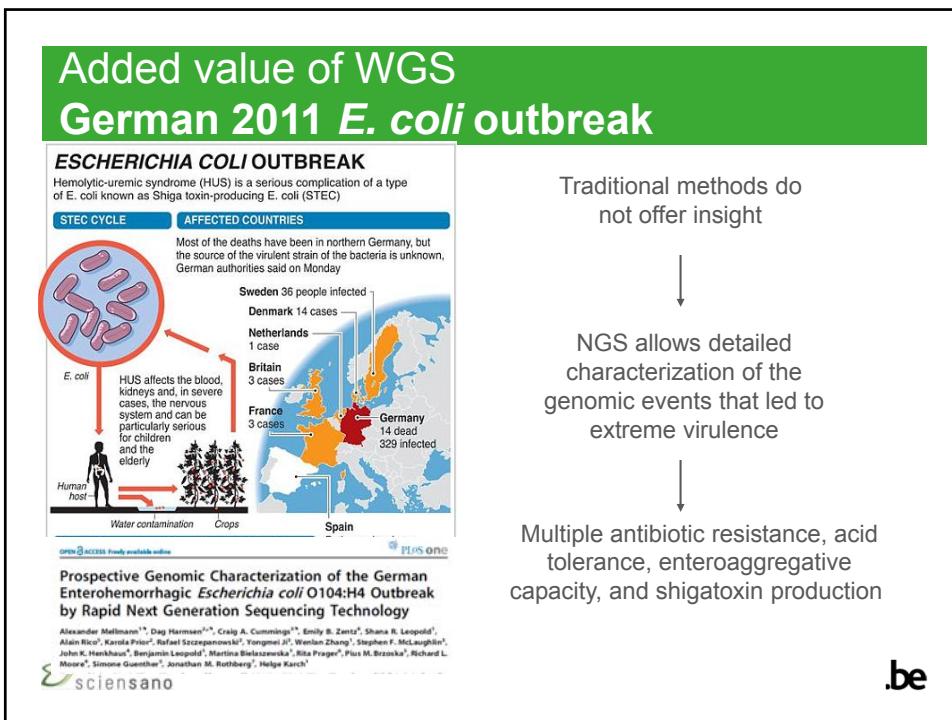
.be

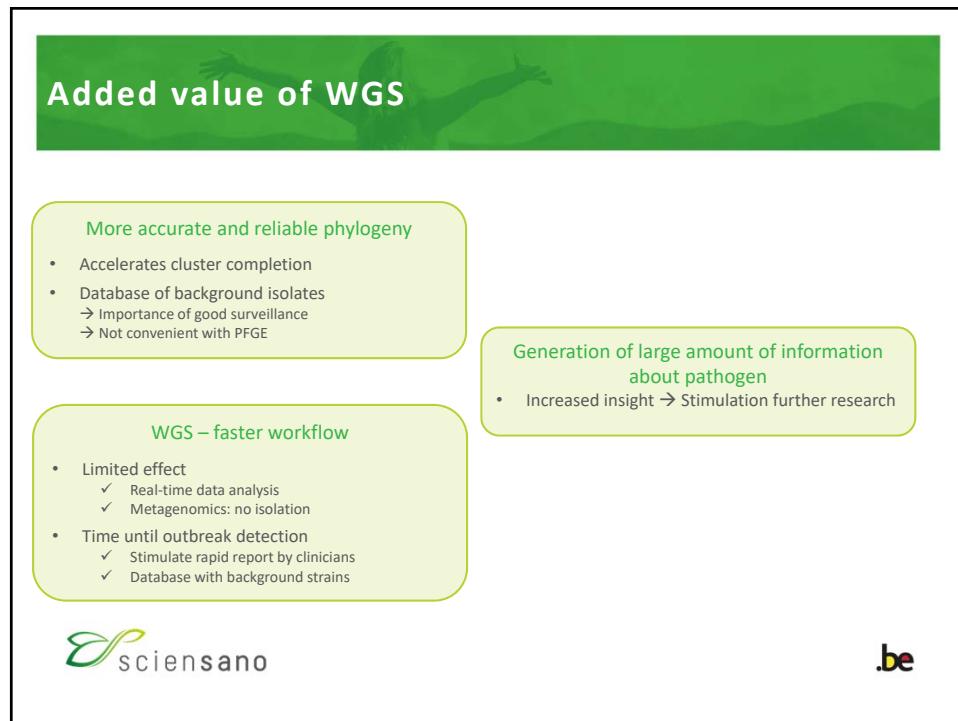
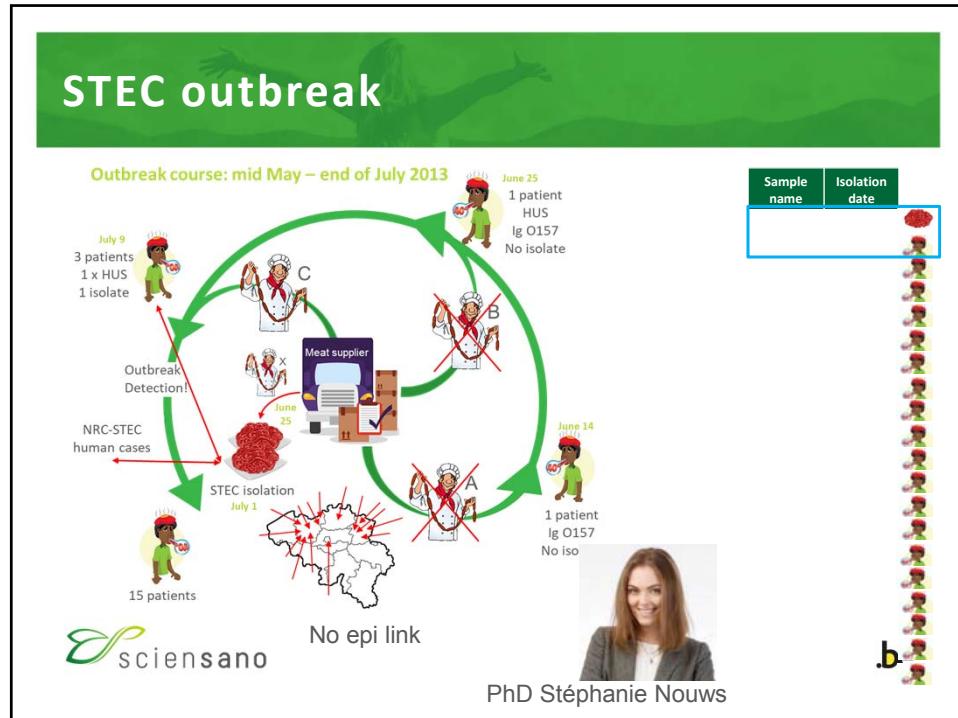
		Flanders				Wallonia		
		Human	Food	Human	Food	Human	Food	Reference P125109
Flanders	Human	0	0	1	1	52	52	46
		0	0	1	1	52	52	46
Flanders	Food	1	1	0	2	53	53	47
		1	1	2	0	51	51	45
Wallonia	Human	52	52	53	51	0	0	44
		52	52	53	51	0	0	44
Reference P125109		46	46	47	45	44	44	0

.be









Added value of WGS – routine NRL/NRC *Neisseria meningitidis*

Méningite sur le campus de Dijon, une campagne de vaccination à la rentrée

Après la mort de deux étudiantes, victimes d'une méningite à l'université de Dijon, une campagne de vaccination devrait être lancée dès le 3 janvier.

LE MONDE | 27.12.2016 à 11h45 • Mis à jour le 27.12.2016 à 12h03

- 2 cases in France (deaths)
- Student of University of Dijon died in Saint Pierre hospital in Brussels because of meningitis (20/12/2016)
- Isolate received @ NRC Sciensano (S. Bertrand) (23/12/2016)
- Classical testing: invasive strain W (or W 135) cfr. cases in France (23/12/2016)
- WGS started 03/01/2017 @ TAG
- WGS data available 9/01/2017
- Bioinformatic pipeline 'Neisseria': same cluster W:P1-5,2:F1-1:cc11 (09/01/2017)
- Phylogenetic analysis by French authorities: confirmation isolate belongs to same clone UK/South America, as the other isolates in Dijon (10/01/2017)



.be

WGS data storage requirements

Pathogen	Genome size (Mb)	Average coverage	Raw read FASTQ file size, uncompressed (MB)	Fully assembled genome FASTA file size, uncompressed (MB)
<i>Salmonella</i> spp.	5.1	50	510	5.1
<i>Listeria monocytogenes</i>	2.9	50	290	2.9
<i>E. coli</i>	4.6–5.4	50	460–540	4.6–5.4
<i>Campylobacter</i> spp.	1.6	50	160	1.6

RawReadsFileSize (MB, uncompressed) = GenomSize (Mb) x Coverage x 2
Size file fully assembled genome (MB) = size of genome in Mb

- Zipping of FASTQ files (2-3x reduction)
- Need to keep raw reads for accreditation purposes, litigation, reanalysis
- Back-up of data
- Sharing of data – bandwidth (secure FTP sites, file sharing services)

ECDC – Expert opinion on the introduction of next-generation typing methods for food- and waterborne diseases in the EU and EEA, 2015



.be

Global data sharing

Benefits

- ↳ Improves food safety sciences
 - ↳ Efficient use of resources (prevent duplication)
 - ↳ Enables trend analysis of pathogen evolution & rapid response
 - ↳ Mitigates health, social, economic impacts

Challenges

- ? Inequities in capacity & resources
 - ? Data ownership & metadata access
 - ? Accountability & transparency
 - ? Validation, standards & quality (ISO working group ongoing)
 - ? Fear of being 'scooped'
 - ? Technical trade barriers
 - ? Privacy law

FAO & WHO Technical background paper 2016: Applications of whole genome sequencing in food safety management



Added value open data sharing



Facility #2, Peppers



Genome Trakr



WGS resource considerations

« WGS is too advanced, too difficult and too expensive to implement »



PERCEPTION

« WGS is the perfect solution to all our problems »



WGS – benefits & drawbacks

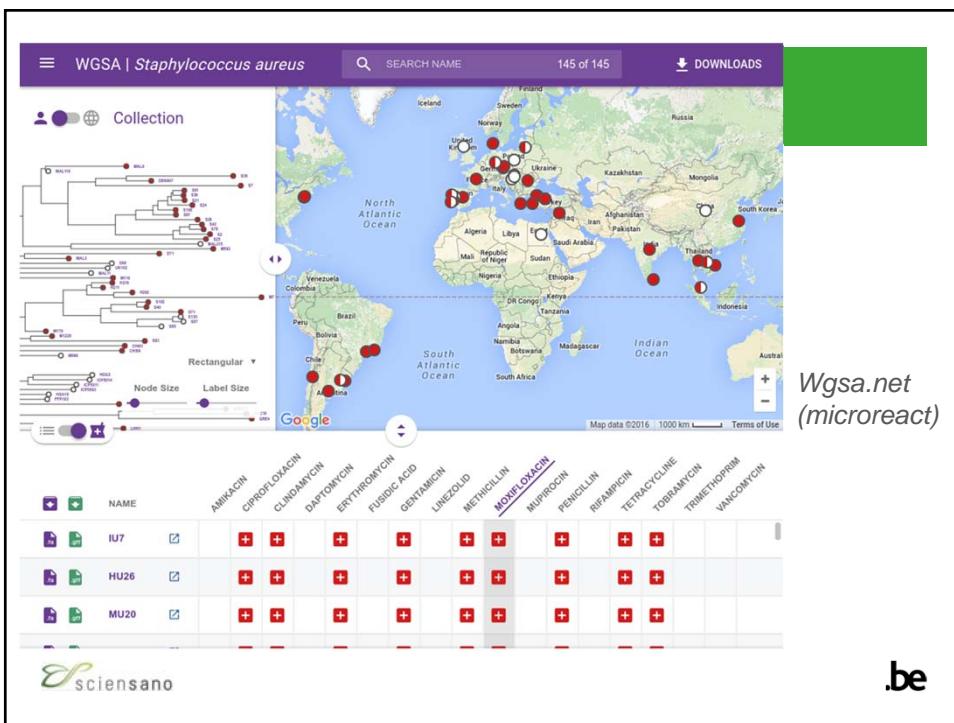
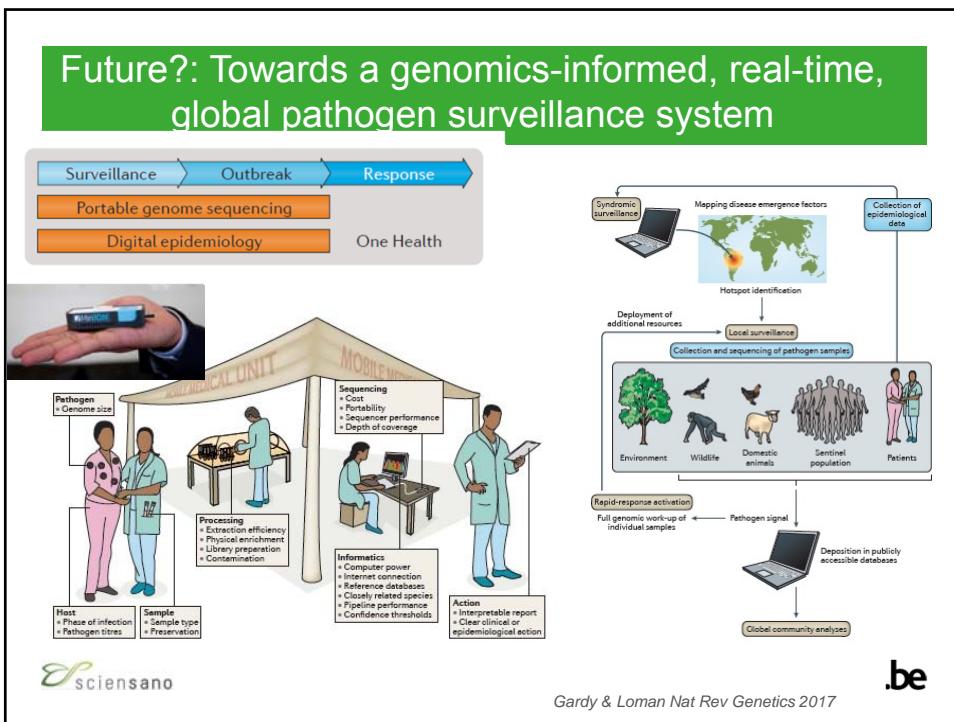
Performance (specificity/sensitivity)

- + Confidence in clusters and links isolates
- + Variety of outputs based on ONE « simple » laboratory procedure
- + Flexibility/universality – not pathogen specific
- + Speed, amenable to reanalysis
- + Ease of sharing (basic common language, data repositories)
- No phenotype
- Education & experience needed for analysis (bioinformatics) and interpretation
- Willingness to (global) data sharing
- Standardization needed

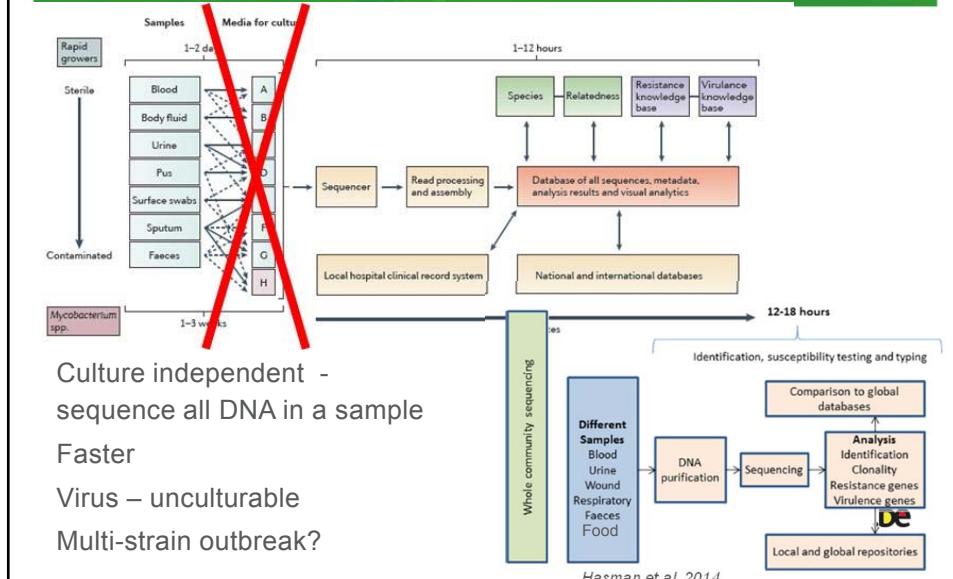
Cost

- + Cost effective alternative to classical typing if replacing several classical methods (not complementary) & adequate volumes of isolates; economic cost
- + Flexibility – same person can analyse different pathogens
- Expensive to establish (equipment, education, computational) – but can be shared across institutes





Isolates (WGS) vs. metagenomics



Isolates (WGS) vs. Metagenomics (2)

- **16S rRNA PCR & sequencing**
 - Taxonomically and phylogenetically informative marker
 - Only to the genus/species level
 - Bias PCR (primers not universal), artificial amplicons after assembly
 - Low sensitivity if insufficient pathogen DNA
 - Affected by presence of contaminating DNA from other bacterial species
 - No functional/biological information
- **Shotgun metagenomics**
 - Taxonomical & biological information
 - Filtering out unwanted DNA in post-sequencing analysis
 - Higher sensitivity as organisms can be identified from different segments of DNA
- **Considerations**
 - Dead & alive organisms
 - Massive sequence data linked to e.g. food processing facility (pathogen?)
 - Human reads? Privacy law

Future: metagenomics?
Towards isolation-free outbreak investigation

Metagenomics: no isolation... 

Human  Food 

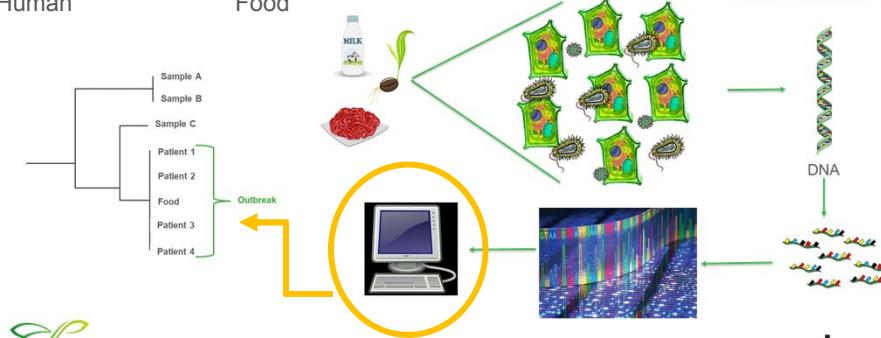
PhD Florence Buytaerts 

Sample A
Sample B
Sample C
Patient 1
Patient 2
Patient 3
Patient 4
Food

Outbreak

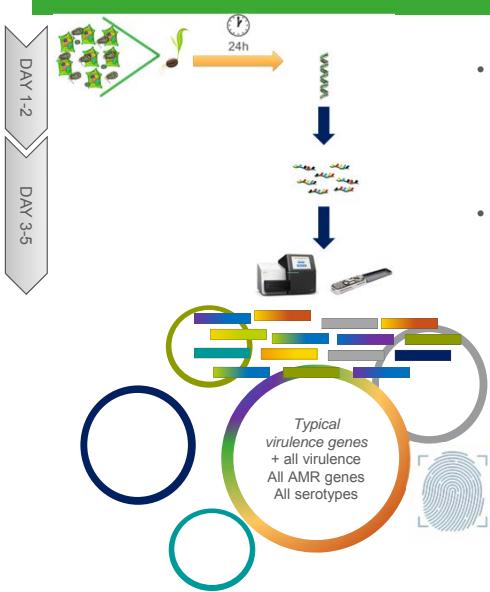
sciensano

Shotgun metagenomics approach 



DNA

Metagenomics-based outbreak investigation



- Information:**
 - All species present
 - No isolation needed -> shorter analysis time
- Complexity:**
 - all DNA present in sample sequenced
 - (Short) reads: big puzzle: what belongs to which genome?
 - Reconstruct genomes present
 - Phylogeny of reconstructed genomes (SNP-level): link human & food isolate

Typical virulence genes + all virulence All AMR genes All serotypes

.be

Metagenomics & food safety

AEM
Journal of AEM.org

Application of Metagenomic Sequencing to Food Safety: Detection of Shiga Toxin-Producing *Escherichia coli* on Fresh Bagged Spinach

Susan R. Leonard, Mark K. Mammel, David W. Lacher, Christopher A. Elkins
Division of Molecular Biology, Center for Food Safety and Applied Nutrition, U.S. Food and Drug Administration, Laurel, Maryland, USA

December 2015 | Volume 81 | Number 23 | Applied and Environmental Microbiology | aem.asm.org | 8183

TABLE 1 Results of spiking spinach with STEC Sakai^a

Sampling time and spike (CFU)	E. coli phylogroup E (%) ^b	Mapped reads (%) ^c	Coverage (fold) ^d
Preenrichment			
10	0.003	0.01	0.04
1,000	0.003	0.05	0.11
10,000	0.119	0.03	0.10
100,000	0.356	0.13	0.41
1,000,000	4.251	0.84	3.16
5-h enrichment			
10 ^e	0.008	0.08	0.24
1,000 ^f	0.066	0.12	0.39
10,000 ^g	7.418	3.30	13.3
100,000	30.449	17.09	70.9
23-h enrichment			
10 ^h	66.954	49.73	184
1,000 ⁱ	84.330	61.77	252

^aSamples of 100 g spinach were spiked with STEC Sakai.
^bPercentage of the bacterial population belonging to E. coli phylogroup E.
^cPercentage of total reads mapping to the STEC Sakai genome.
^dAverage chromosome coverage. The values were normalized to 10,000,000 total reads per sample. The coverage for the two plasmids was greater than that of the chromosome for all samples.
^eAverage of triplicate biological samples.
^fAverage of duplicate biological samples.

TABLE 2 Number of reads mapping to virulence genes and serotyping loci for spinach samples spiked with STEC Sakai

Sampling time and spike (CFU)	No. of reads mapping to the following virulence gene or serotyping locus ^j :
	stx _{1a} stx _{1b} γ-cote ehxA waz wzy fltC
Preenrichment	
10	0 0 0 0 0 0 0
1,000	0 0 0 0 0 0 0
10,000	1 1 1 1 0 0 0
100,000	2 0 7 4 3 0 2
5-h enrichment	
10 ^k	0 0 0 0 0 1 0
1,000 ^l	0 2 4 5 0 0 0
10,000 ^m	57 65 329 231 36 22 43
100,000	373 275 1,676 1,382 194 126 232
23-h enrichment	
10 ⁿ	1,225 1,361 3,467 4,013 988 850 978
1,000 ^o	1,715 1,647 4,773 4,529 1,419 1,110 1,298

^jNumber of reads mapping to virulence genes and serotyping loci for 100-g spinach samples spiked with STEC Sakai. For comparison, the numbers are normalized to 10,000,000 total reads per sample.
^kAverage of triplicate biological samples.
^lAverage of duplicate biological samples.

^mAverage of triplicate biological samples.
^oAverage of duplicate biological samples.

Shotgun metagenomics
LOD: 10 CFU/100g if enriched 8h
Virulence genes & assembled genome
.be

Strain-level metagenomic analysis of bacterial pathogens in enriched food samples

Future work: regulation to be changed – still isolate needed...

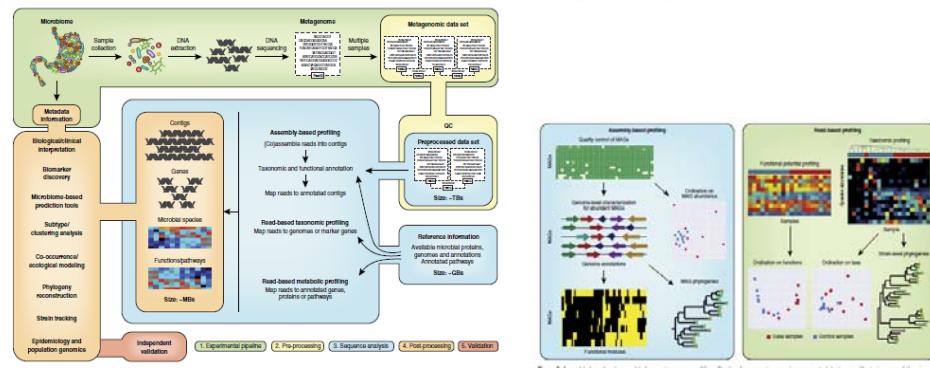
sciensano

PhD Assia Saltykova PhD Florence Buytaers

.be

Shotgun metagenomics, from sampling to analysis

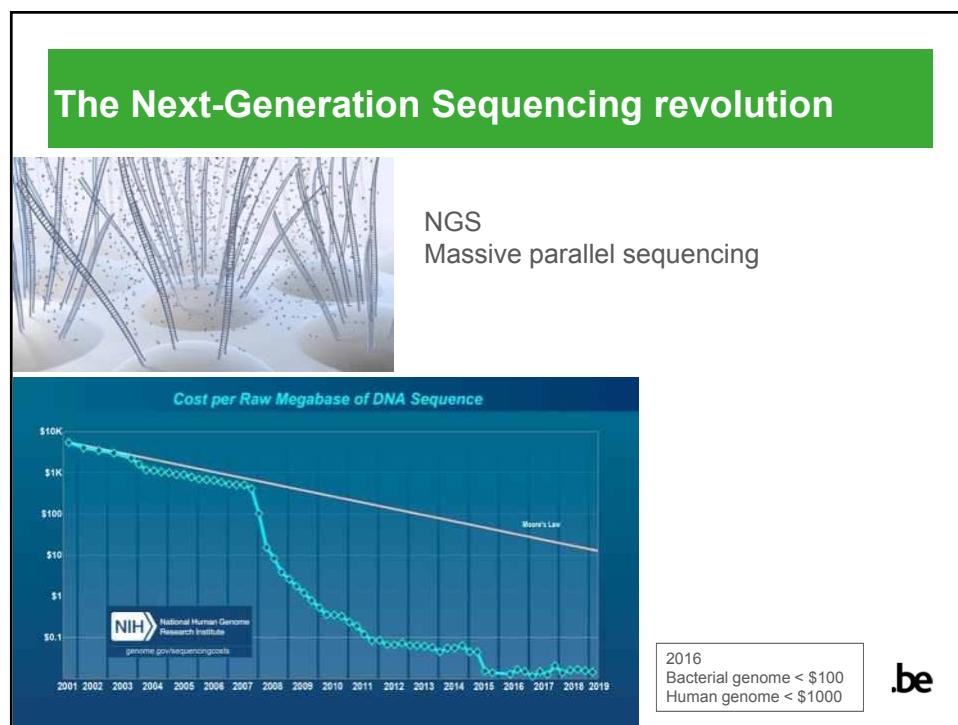
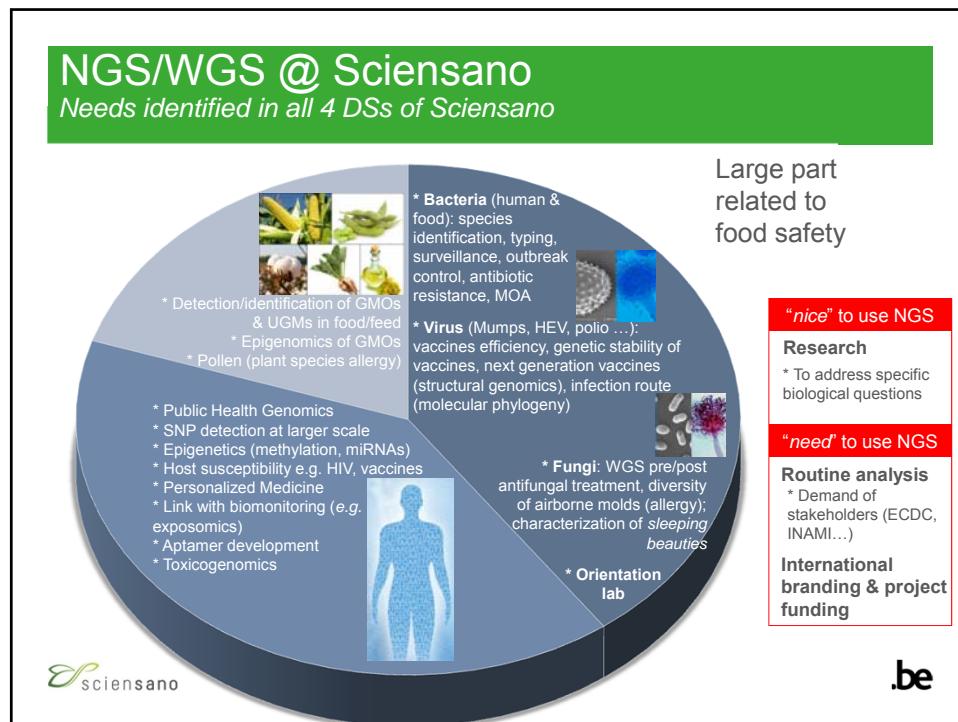
Christopher Quince^{1,7}, Alan W Walker^{2,7}, Jared T Simpson^{3,4}, Nicholas J Loman⁵ & Nicola Segata⁶



WGS will impact foodsafety worldwide

WGS will become standard methodology in food safety

- Foodborne disease surveillance
- Food inspection (testing) and monitoring
- Outbreak detection and investigation (higher resolution, more confident definition of clusters, linking sporadic cases...)
- Food source attribution studies
- Food technology developments
- + human & animal disease surveillance & outbreak response
- = thé One Health technology
- Several European countries already WGS implemented/ing in routine, supported by the competent authorities – Belgium?



Many applications

- Research tool
- Support proactive Public Health policy
 - ✓ Diagnosis current & emerging diseases
 - ✓ Outbreak control
 - ✓ Understand transmission patterns of microbial pathogens
 - ✓ Characterization of pathogens
 - ✓ Re-valorisation of collections (*sleeping beauties*)
 - ✓ Identify & prevent health risks food consumption & environment
 - ✓ GMO analysis
 - ✓ Vaccines
 - ✓ Public Health Genomics
 - ✓ ...
- Complement or replace current methods?

 .be

"Challenge in data analysis, not in data generation"

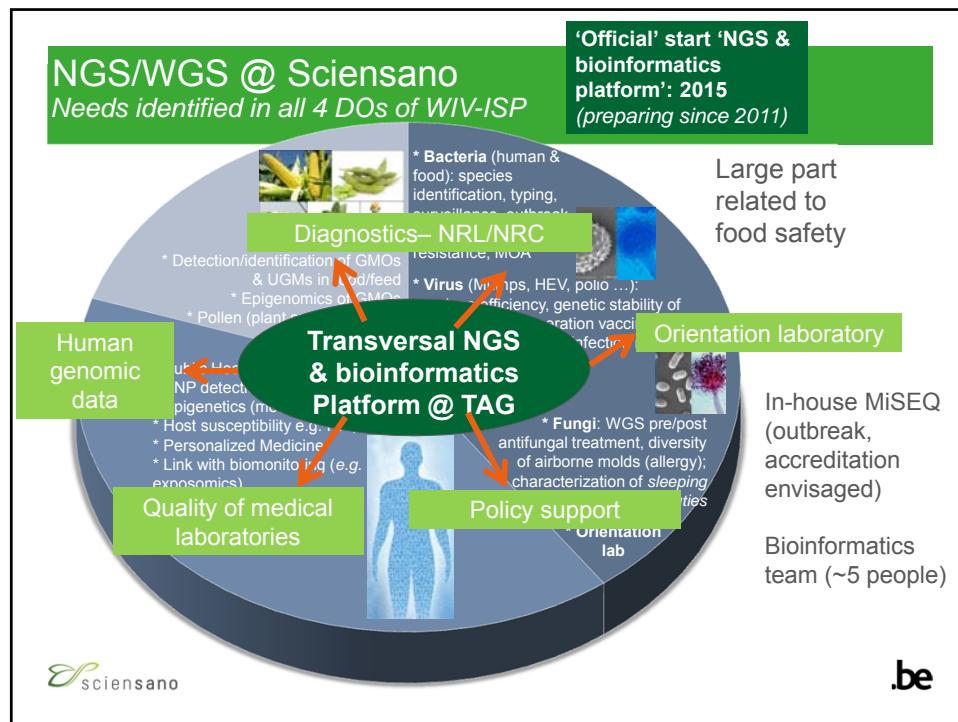
REFERENCE ALIGNMENT DE NOVO ASSEMBLY VARIANT DETECTION CHIP-SEQ RNA-SEQ

Need for NGS data analysis tools to rapidly analyse the massive amount of data to extract and interpret the required info correctly

« Without solving certain bioinformatics challenges, the NGS revolution will not be extensively available to public health professionals »

Picture taken from McPherson, Nature Methods 6, S2 - S5. 2009.

 .be



MiSeq instrument for routine

Dr Ir senior Mol biologist (lead)
NGS expert technician
+ 2 other technicians trained for NGS

- NGS-instrument: in-house
 - Time- crisis (<- outsourcing: 6-8 weeks)
 - Technology: ~ stable
 - Strategic importance: central points for BE stakeholders
- Benchtop model: Illumina MiSeq
 - Short analysis time – acceptable amount of data
 - User-friendly
 - 'gold standard' for microbial NGS + cancer panels
- ISO17025 (ISO15189) accreditation (member of ISO working group (ISO/TC 34/SC 9/WG25): Whole genome sequencing for genomic typing and characterization)

Other technologies: R&D

- MinION
- PacBio (outsourcing)
- HiSeq (outsourcing)

sciensano .be

Bioinformatics Platform for massive data analysis

K. Vanneste
Bioinformatician
(lead)

Bioinformatician

Bioinformatician

Bioinformatician

Bioinformatician

Software engineer

NGS data analysis

Development and implementation of user-friendly bioinformatics pipelines & databases

g EBI Genomics Workbench

BioNumerics server

GNU/Linux

Galaxy

Commercial solutions

Microbial isolates
Mixed samples
Human samples

In-house developed solutions

Computational high-performance infrastructure

- Scientific storage
- In-house cluster
- Several (web)servers
- Highly controlled (versioning, logging of parameters/pipelines, stress tests...)

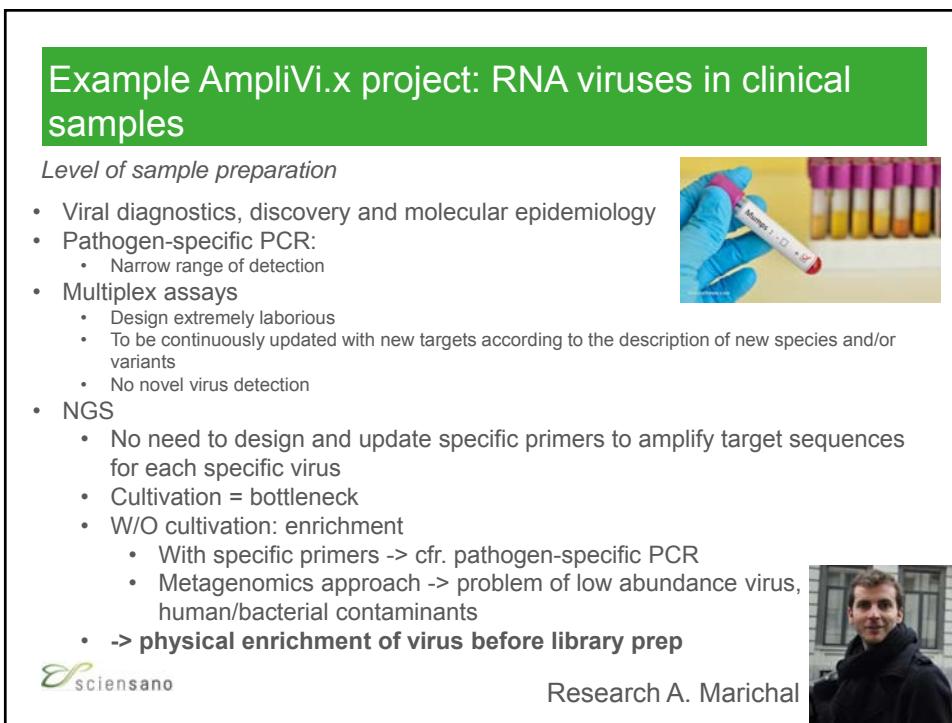
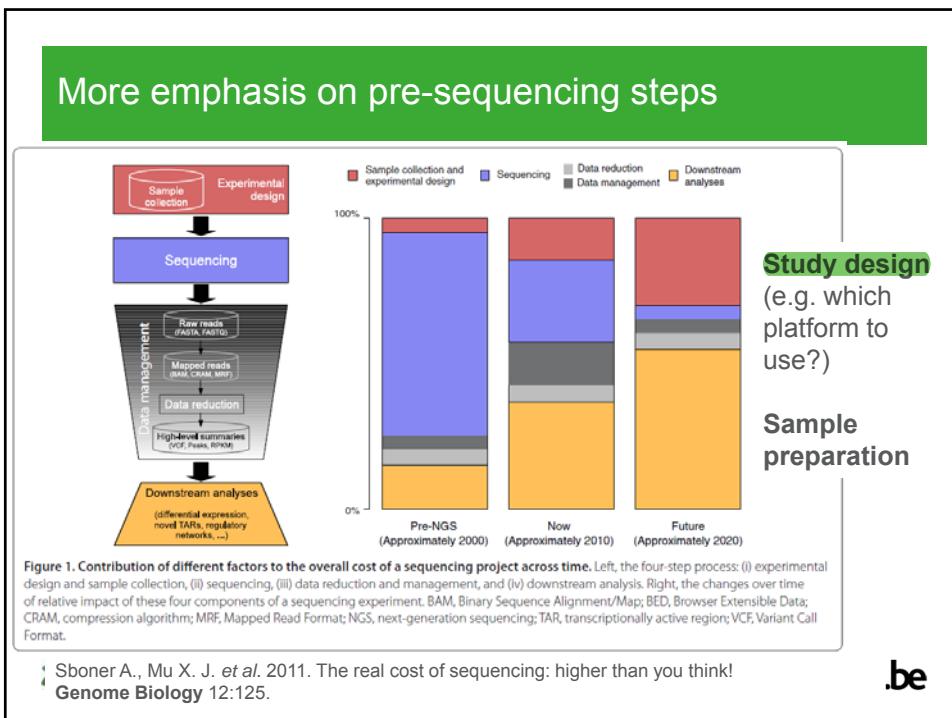
.be

sciensano

R&D versus routine applications in a public health setting, many challenges

- Sample preparation

.be



**Example MycoMOLAIR project:
Indoor air quality**

Level of sample & library preparation

Indoor time: 75-90%
Mold exposure: 18-50%

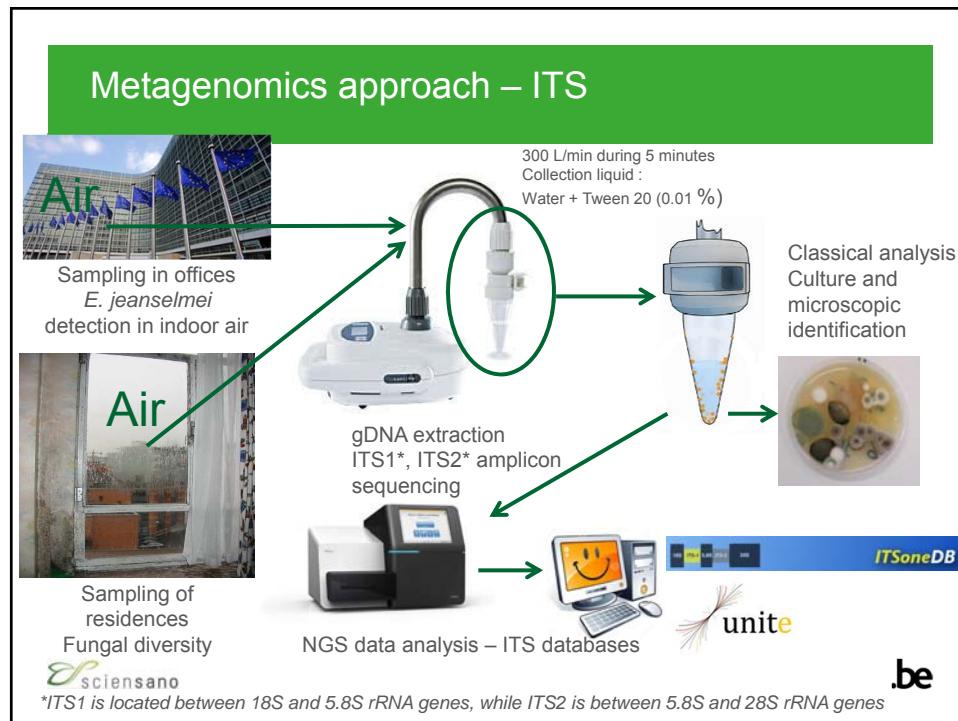
- > Severe public health problems (e.g. asthma) – causal relationship?
- > consequences

Moisture – energy saving houses?
Accurate identification needed
Problems with classical (culture & microscopy based) analysis
Dead/uncultivable fraction: allergenic?-> fungal diversity?
Need for molecular methods

PhD Xavier Libert

.be

sciensano





 **microorganisms**



Article

Exploiting the Advantages of Molecular Tools for the Monitoring of Fungal Indoor Air Contamination: First Detection of *Exophiala jeanselmei* in Indoor Air of Air-Conditioned Offices

Xavier Libert^{1,2}, Camille Chasseur³, Ann Packeu³, Fabrice Bureau², Nancy H. Roosens¹ and Sigrid C. J. De Keersmaecker^{1,*} 

¹ Transversal activities in Applied Genomics, Sciensano, J. Wytsmanstraat 14, 1050 Brussels, Belgium; libert.xavier85@gmail.com (X.L.); nancy.roosens@sciensano.be (N.H.R.)

² Cellular and Molecular Immunology, Groupe Interdisciplinaire de Génoprotéomique Appliquée (GIGA), Université de Liège (ULg), Avenue de l'Hôpital, 1 (B34), 4000 Sart-Tilman, Belgium; fabrice.bureau@ulg.ac.be

³ Mycology and Aerobiology, Sciensano, J. Wytsmanstraat 14, 1050 Brussels, Belgium; Camille.Chasseur@sciensano.be (C.C.); ann.packeu@sciensano.be (A.P.)

* Correspondence: Sigrid.dekeersmaecker@sciensano.be; Tel.: +32-2-642-5257

 .be



R&D versus routine applications in a public health setting, many challenges

- Sample preparation
- Standardized, harmonized generation of NGS data for different applications
- Computational requirements
- Validation of standardized & optimized pipelines, for different applications in public health (isolates, clinical samples, metagenomics...)
- User-friendly access for non-experts (application in NRL/NRC)
- Trade-off between quality and speed of analysis (crisis...)
- Traceability (databases, runs)
- Integration with other high-throughput technologies

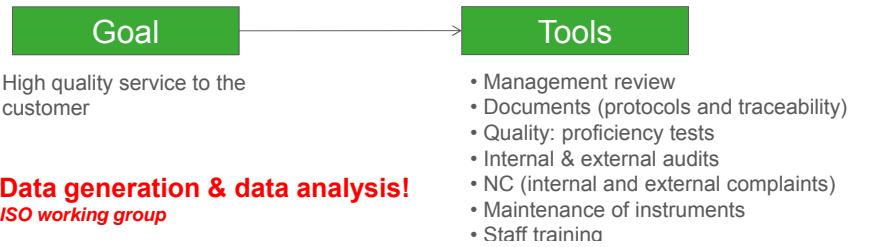
 .be

Accreditation: ISO 17025

- General requirements for the competence of testing and calibration laboratories

General considerations

- management system ensuring adequate qualification and training of staff;
- metrology system ensuring periodical calibration of equipment;
- method verification demonstrating that laboratory is meeting the performance characteristics



.be

Bioinformatics & ISO

- No ISO standard yet
- Best practices:
 - Version control (codebase + databases) (GIT)
 - Code review
 - DTAP approach (Development, Testing, Acceptance, Production)
 - Technical documentation (Wiki)
 - User documentation (Wiki + SOP)
 - Traceability: automatically updated databases + logging of all parameters (including version of databases used) and runs
 - All relevant information to repeat an analysis with a pipeline is stored:
 - Input files (checksums), Parameters, Output files (checksums), Database versions, Username, Date, ...
 - Extensive validation of pipelines
 - Repeatability
 - Reproducibility
 - Part of complete workflow within pathogen profiling

More information: Lambert et al. 2017. Journal of AOAC International Vol. 100, No.3 DOI: 10.5740/jaoacint.16-0269

:

Validation of bioinformatics pipelines

LAMBERT ET AL.: JOURNAL OF AOAC INTERNATIONAL VOL. 100, NO. 3, 2017 721

FOOD BIOLOGICAL CONTAMINANTS

Baseline Practices for the Application of Genomic Data Supporting Regulatory Food Safety

DOMINIC LAMBERT
Canadian Food Inspection Agency, Ottawa Laboratory (Carling), Ottawa, ON, Canada

ARTHUR PIGHILING
U.S. Food and Drug Administration, Center for Food Safety and Applied Nutrition, College Park, MD

EMMA GREGORY
Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC, Canada

GARY VAN DOMSELAAR
Public Health Agency of Canada, National Microbiology Laboratory, Winnipeg, MB, Canada

PETER EVANS
U.S. Department of Agriculture, Food Safety and Inspection Service, Washington, DC

SHARON BERKELIET
Canadian Food Inspection Agency, Ottawa Laboratory (Carling), Ottawa, ON, Canada

DUNCAN CRAG
Food Standards Australia New Zealand, Canberra, Australia

P. SCOTT CHAPOK
Commonwealth Scientific and Industrial Research Organisation, Melbourne, Australia

ROBERT STONES
Food and Environment Research Agency, Sand Hutton, York, United Kingdom

FIONA BRINKMAN
Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC, Canada

ALEXANDER ANGERS-LOVSTAD and JOACHIM KELLYA
Eurofins Commission, Joint Research Center, Ispra, Italy

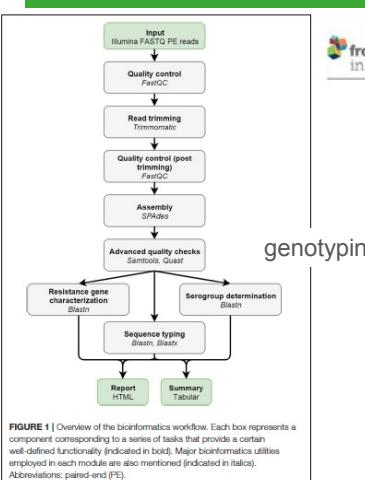
WILDA TONG
U.S. Food and Drug Administration, National Center for Toxicological Research, Little Rock, AR

BERTON BLAIS[†]
Canadian Food Inspection Agency, Ottawa Laboratory (Carling), Bldg 22, 960 Carling Ave., Central Experimental Farm, Ottawa, ON, Canada K1A 0Y9





Neisseria case study



genotyping

Validation of a Bioinformatics Workflow for Routine Analysis of Whole-Genome Sequencing Data and Related Challenges for Pathogen Typing in a European National Reference Center: *Neisseria meningitidis* as a Proof-of-Concept

OPEN ACCESS
Published: 06 March 2019
doi: 10.3389/fmicb.2019.00362



PhD Bert Bogaerts

Edited by:
Ruth Ann Lewis,
University of Arizona,
United States;
Reviewed by:
Marta de Jong,
Belgium

Bert Bogaerts¹, Raf Winand¹, Qiang Fu¹, Julien Van Braekel¹, Pieter-Jan Coeynsens¹, Wesley Mattheus², Sophie Bertrand², Sigrid C. J. De Keersmaecker¹, Nancy H. C. Roosens¹ and Kevin Vanneste^{1*}

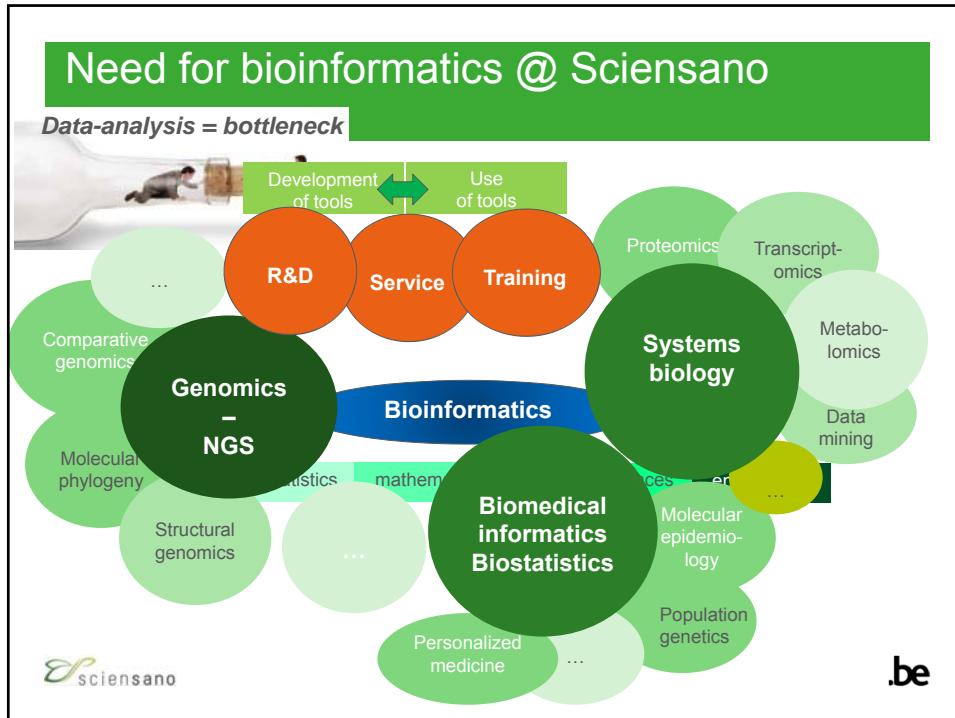
¹Transversal Activities in Applied Genomics, Sciensano, Brussels, Belgium; ²Bacterial Diseases, Sciensano, Brussels, Belgium

FIGURE 1 | Overview of the bioinformatics workflow. Each box represents a component corresponding to a series of tasks that provide a certain well-defined functionality (indicated in bold). Major bioinformatics utilities employed in each module are also mentioned (indicated in italics). Abbreviations: paired-end (PE).

Push-button pipeline = user-friendly, for NRL/NRC Report







Acknowledgements

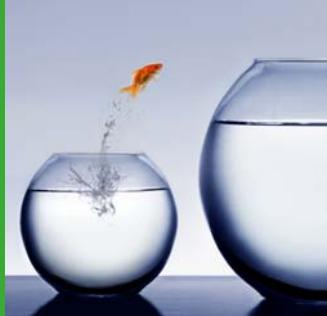
- Sciensano
 - Transversal activities in Applied Genomics: Nancy Roosens, Kevin Vanneste, Véronique Wuyts, Assia Saltykova, Stefan Hoffman, Maud Delvoye, Els Vandermassen, Dirk Van Geel, Loïc Lefèvre, Bas Berbers, Stéphanie Nouws, Florence Buytaers, Qiang Fu, Bert Bogaerts, Raf Winand, Thomas Delcourt, Julien Van Braekel, Axel Marichal, Xavier Libert
 - Food Pathogens: Katrijne Dierick, Nadine Botteldoorn, Sarah Denayer, Bavo Verhaegen, Cristina Garcia-Graells
 - Bacterial Diseases: Sophie Bertrand, Wesley Mattheus, Vanessa Mathys, Pieter-Jan Ceyssens
 - QSE: Anne-Marie Vanherle, Patricia Cliquet
 - Ugent: Kathleen Marchal
 - ISO working group TC 34/ SC 9/ WG 25 Whole-genome sequencing for typing and genomic characterization

 **sciensano** healthy all life long

Thank you for your attention!
Questions?

Contact

Sigrid.dekeersmaecker@sciensano.be



Sciensano • Rue Juliette Wytsmanstraat 14 • 1050 Brussels • Belgium
T +32 2 642 51 11 • T Press +32 2 642 54 20 • info@sciensano.be • www.sciensano.be **.be**