Kevin Vanneste, PhD

Bioinformatics Platform
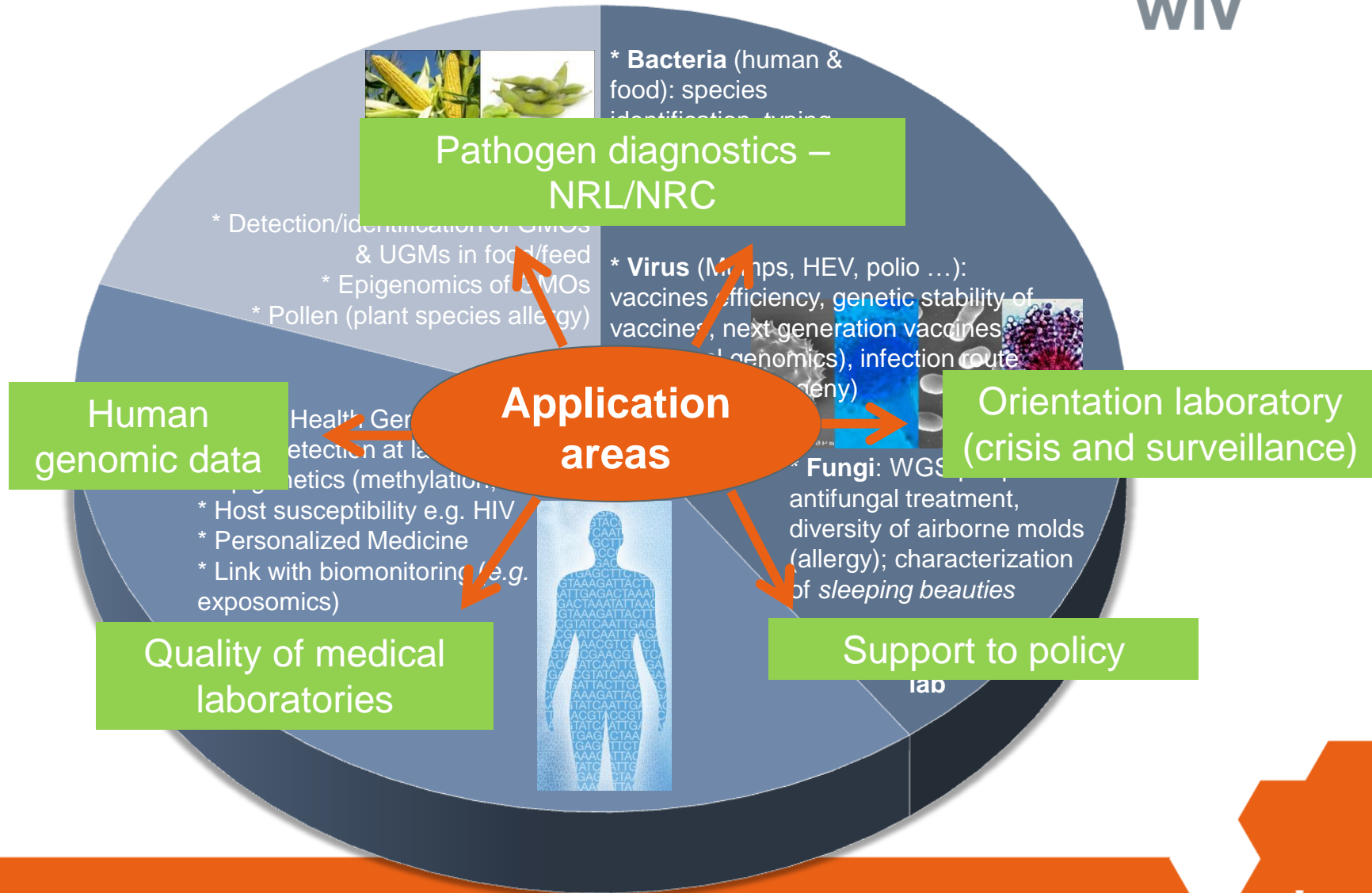
Platform Biotechnology and Molecular Biology

Department Expertise, Service Provision and Customer Relations

**Collaboration between the EURL-VTEC and the Platform for Biotechnology and Molecular Biology (WIV-ISP, Belgium) for the development of a bioinformatics pipeline for routine analysis of whole genome sequencing data for typing of STEC/*E. coli* using Galaxy**

12th Annual Workshop of the National Reference Laboratories for *E. coli* in the EU (12/10/2017-13/10/2017)

Rue Juliette Wytsmanstraat 14 | 1050 Brussels | Belgium
T +32 2 642 50 65 | F +32 2 642 52 92 | email: kevin.vanneste@wiv-isp.be | www.wiv-isp.be

.be

# Use of NGS & bioinformatics @WIV-ISP?



**Pathogen diagnostics – NRL/NRC**

**Application areas**

**Human genomic data**

**Orientation laboratory (crisis and surveillance)**

**Quality of medical laboratories**

**Support to policy**

* **Bacteria** (human & food): species identification, typing

* Detection/identification of GMOs & UGMs in food/feed
* Epigenomics of GMOs
* Pollen (plant species allergy)

* **Virus** (Mumps, HEV, polio …): vaccines efficiency, genetic stability of vaccines, next generation vaccines (genomics), infection route (deny)

* Health: Genetic detection at la... epigenetics (methylation)
* Host susceptibility e.g. HIV
* Personalized Medicine
* Link with biomonitoring (e.g. exposomics)

* **Fungi**: WGS... antifungal treatment, diversity of airborne molds (allergy); characterization of *sleeping beauties*

# The Next-Generation Sequencing revolution

**NGS as a well established research tool…**

- Universal method
- Single nucleotide resolution
- High-throughput, ultimate multiplex tool

**…but many challenges remain regarding data analysis and interpretation for routine applications in a public health setting!**

- Computational requirements
- Validation of standardized & optimized pipelines
- User-friendly access for non-experts
- Trade-off between quality and speed of analysis
- Traceability (databases, runs)

.be

# NGS & bioinformatics platform

**Mission: Utilizing NGS & bioinformatics**

- for the diagnosis, surveillance, control and characterisation of potentially harmful organisms
- to promote public health genomics by the effective integration into clinical use and public health policy

**Objectives: Develop and implement solutions and provide data acquisition and analysis tools to**

- complement the WIV-ISP laboratories services
- integrate the knowledge of genomics into public health policy

.be

ORGANISATION CHART WIV-ISP

# NGS & bioinformatics platform @WIV-ISP

**Mission**

PBB is a transversal scientific service using **molecular biology & bioinformatics** to conduct routine analysis as well as scientific and technologic research.

By building internal and external partnership, it generates new knowledge & customized tools to anticipate present and futures challenges affecting public health

.be

# Platform Biotechnology and Molecular Biology

**PBB** – N. Roosens

**GMOLAB**
N. Papazova

Scope: Molecular biology (GMO analysis of food & feed)

**BIOTECHLAB**
S. De Keersmaecker

Scope: Molecular biology (all fields in support to WIV-ISP)

**BIOIT Platform**
K. Vanneste

Scope: Bioinformatics (focus on analysis of NGS data in support to WIV-ISP)

**Sequencing Platform**

Sanger

NGS

NGS

Transversal

.be

# Activities of the bioinformatics platform



ispwIV

**BIOTECHLAB**   **NGS**

Novel applications

**Bioinformatics Platform**

**IT**

data generation

data analysis

Expertise

**Develop standards**
**Advise clients**
**Provide training**

Valid Implementation

Research Projects

**Centralized bioinformatics platform**
**User-friendly workflows**
**Obtain accreditation**

**Develop state-of-the-art pipelines**
**Collaborate (inter)nationally**
**Identify upcoming needs**

.be

# Develop and maintain centralized bioinformatics platform

## Bioinformatics platform



**Bioinformatician** *(lead)*

**Bioinformatician**

**Bioinformatician**

**Bioinformatician**

**Bioinformatician**

**Software engineer**

## Development and implementation of user-friendly bioinformatics tools, pipelines and databases



**Commercial solutions**

**In-house developed solutions**

**Microbial isolates**
**Mixed samples**
**Human samples**

.be

# The data analysis bottleneck



## Data storage?

~10s GB raw data (microbial)
~100s GB raw data (human)

# The data analysis bottleneck



**Data storage?**

**Data processing?**

~10s GB RAM
(microbial)

~100s GB RAM
(human)

.be

# The data analysis bottleneck

**Data storage?**
**Data processing?**
**Data analysis?**

SNPs?
Errors?

# Bring data analysis into routine

**Tools**

Staple bioinformatics tools used and adopted by the scientific community

Tools can be combined to make pipelines

**'Push-on-the-button' pipelines**

Pipelines engineered by the BIOT platform

Case studies tackled according to priorities defined by direction committee

**E.g. *Neisseria meningitidis* pipeline**

**Computational requirements**

Tools and pipelines integrated directly into high-performance computational infrastructure WIV-ISP

**Validation of standardized & optimized pipelines**

Pipelines use validated parameters

**User-friendly access for non-experts**

Galaxy Workflow Management System for access to non-bioinformaticians

**Trade-off between quality and speed of analysis**

Different modes of analysis (e.g. surveillance versus outbreak)

**Traceability**

Automatically updated databases, logging of all parameters and runs

.be

# Offer a high-quality service platform

**Providing a high-quality service....**
Version control
Code review
Basic testing
Technical documentation (bioit Wiki)
User documentation (bioit Wiki)
Several in-take meetings with client to define needs
DTAP principle (**D**evelopment -> **T**esting -> **A**cceptance -> **P**roduction)

**Building up the quality system…**
2017: Benchmarking to consolidate internal quality system
2018: Obtain certification / accreditation

.be

# Combine tools to make pipelines

**Tweak parameters**



**Access using browser (simultaneous usage possible)**

**Standard tools**

**Custom tools**

**Workflows**

# Combine tools to make pipelines



**Tool output**

# Optimized 'push-on-the-button' pipelines

## User-friendly access

## Centralized computational infrastructure



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

**Automatically updated and traceable databases**

**Validated & optimized parameters**

.be

# Optimized 'push-on-the-button' pipelines



**Detailed output report**

ORGANISATION CHART WIV-ISP

**GENERAL DIRECTION**

Services of the Executive director
Support services
Managerial posts
Advisor posts

**COMMUNICABLE AND INFECTIOUS DISEASES**
- Food pathogens
- Bacterial diseases
- Viral diseases
- Mycology and aerobiology
- Immunology
- Platform of animal facility
- Laboratory of medical microbiology
- Platform BSL3 and Metrology

**FOOD, MEDICINES AND CONSUMER SAFETY**
- Medicines
- Chemical residues and contaminants
- Consumer safety
- Toxicology
- Health and environnement
- Platform for chromatography and mass spectrometry

**PUBLIC HEALTH AND SURVEILLANCE**
- Healthcare services research
- Surveys, lifestyle and chronic diseases
- Epidemiology of Infectious Diseases
- Healthcare associated infections & Antimicrobial resistance
- Cancer Centre
- EPI methods for Public health crisis
- Public health information system
- Special survey

**EXPERTISE, SERVICE PROVISION AND CUSTOMER RELATIONS**
- Biosafety and biotechnology
- Biological standardization
- Quality of medical laboratories
- Platform Biotechnology and molecular biology
- Coordination Cell for customer relations and orientation
- Dispatching Centre
- Healthdata.be
- Crisismanagement

# NRL STEC - Belgium

## Department communicable and infectious diseases

## Food pathogens

.be

# NRL STEC - Belgium



**Scientific Service Foodborne pathogens**
N. Botteldoorn

**Foodborne Outbreaks and toxines**
S. Denayer

- Foodborne Outbreaks
  - NRL Coagulase + Staphylococci
  - NRL VTEC (Food)
  - NRL FBO
- Botulism and other zoonotic clostridia
  - NRC C. botulinum & perfringens
  - NRL botulism

**Foodborne Bacteria and Viruses**
N. Botteldoorn

- Foodborne Bac...
  - NRL F... Micro...
  - NRL...
  - NR...
  - NR...
- Antimicr...
  - NRL A...
- Foodborne Vir...
  - NRC Norovirus
  - NRL Bivalve molluscs

**Foodborne parasites**
S. De Craeye

- ...genital
  - ...sis*)
  - ...nosis
- ...sortium

**Involved in research projects:**

- .Be Ready
- StEQIDEMIC.be
- EUROBIOTOX
- TOX-Detect

# NRL STEC - Belgium

E-mail: nrlvti-lnrtia@wiv-isp.be

# Intake BIOIT needs NRL STEC

- Requirement for a routine pipeline for the analysis of WGS data generated on VTEC and other *E. coli* isolates
- Pipeline should be 'push-on-the-button'
- All functionality should also be available as stand-alone tools
- **Pipeline should be streamlined with the functionality made available by the EURL-VTEC in their ARIES platform**

# Analysis BIOIT needs NRL STEC



**Galaxy WGS pipeline Analysis steps**

Input

Output

Ref. genomes

PE reads

agMLST db

Viruotyper db

Virulence Gene db

Serotyping db

Antibiotics Resistance Gene db

reads trimming & QC (FastQC, Trimmomatic)

Taxonomic Check (kraken)

reads mapping (BWA, bowtie2)

SNP calling (samtools, gatk)

SNP typing

agMLST (srst2, blastn)

Virulotyper (?) (srst2, blastn)

VirulenceFinder (srst2, blastn)

SerotypeFinder (srst2, blastn)

*de novo* assembly (Velvet, SPAdes, IDBA)

Antibiotics Resistance characterization (srst2, blastn)

HTML report

TSV report

**HReVAP tpying analysis steps**

PCR melting tempareture

HReVAP (?)

HReVAP type

**SNP phynolgeny analysis steps**

PE reads

SNP phylogeny (BioNumerics)

sample SNPs

simple SNP phylogeny

SNP tree

# Data processing

- Quality control:
    - FastQC for generation of raw read reports
    - Trimmomatic for trimming of raw reads
    - FastQC for generation of trimmed read reports
- Contamination check
    - Kraken for identification/confirmation of species based on *kmer* counting, using the entire NCBI RefSeq Microbial database (updated automatically)
- *De novo assembly*
    - SPAdes or Velvet(Optimiser) for assembly of trimmed reads
    - QUAST for quality control of assembly
- 'Advanced' quality control
    - A series of custom checks to ensure adequate quality for functional interpretation (%cgMLST genes found, median coverage, N-content…)
    - Provide 3 outcomes: 'pass' (all checks passed), 'warning' (questionable quality but OK for interpretation), and 'fail' (sample should be re-sequenced)

.be

# Sequence typing

- MLST, cgMLST, wgMLST, agMLST
  - SRST2 for direct read mapping or BLAST+ for checking assembly
  - MLST (Pasteur/Warwick), cg/wgMLST (Enterobase), agMLST (ARIES) - updated automatically
- Serotyping
  - SRST2 for direct read mapping or BLAST+ for checking assembly
  - SerotypeFinder (DTU/CGE) - updated automatically
- Virulence typing
  - SRST2 for direct read mapping or BLAST+ for checking assembly
  - VirulenceFinder (DTU/CGE) - updated automatically
- Antibiotics resistance
  - SRST2 for direct read mapping or BLAST+ for checking assembly
  - ResFinder (DTU/CGE), CARD (McMaster University), ARG-annot (University of Marseille) - updated automatically
- Plasmid typing
  - SRST2 for direct read mapping or BLAST+ for checking assembly
  - PlasmidFinder(DTU/CGE) - updated automatically

.be

# SNP typing/phylogeny

- SNP typing
  - SnapperDB for SNP typing
  - Italian/UK/Belgian database(s) for typing (?)
- SNP phylogeny
  - PHEnix pipeline, CFSAN pipeline, in-house implementation using Samtools for determining phylogeny based on SNPs
  - Output that consists out of a newick tree file and basic visualization of the resulting phylogenetic tree
  - Due to the nature of the these tools, they cannot be integrated in the 'push-on-the-button' pipeline

.be

# Others

- HReVAP (ARIES)
- Community requirements?

**VTEC Pipeline 0.2 a pipeline for the characterization of VTEC isolates (Galaxy Version 0.2)**   ▾ Options

## Input

**Sample name**

[                                                                                    ]

If no sample name is entered, the system will try to detect one based on the input read files. [WARNING] Sample name can NOT be changed afterwards.

**Forward reads**

[ 📄 ] [ 📑 ] [ 📁 ]   [ 5: 3902_S43_L001_R1_001.fastq                              ▾ ]

**Reverse reads**

[ 📄 ] [ 📑 ] [ 📁 ]   [ 5: 3902_S43_L001_R1_001.fastq                              ▾ ]

**Assembler**

[ VelvetOptimiser                                                                    ▾ ]

**Type of analysis**

[ Fast: allele detection based on Blastn alignment (DNA) and Blastx alignment (Peptide)   ▾ ]

**Library kit**

[ Nextera                                                                            ▾ ]

## Resistance Characterization

**ResFinder**

[ Yes ] [ No ]

**ARG-ANNOT**

[ Yes ] [ No ]

**CARD**

[ Yes ] [ No ]

## Virulence Characterization

**VirulenceFinder**

[ Yes ] [ No ]

## Serotype Determination

**SerotypeFinder**

[ Yes ] [ No ]

## Plasmid Replicon Detection

**PlasmidFinder - Enterobacteriaceae**

[ Yes ] [ No ]

## Sequence Typing

**Classic MLST - Pasteur**

[ Yes ] [ No ]

**Classic MLST - Warwick**

[ Yes ] [ No ]

**cgMLST (From Enterobase)**

[ Yes ] [ No ]

✔ Execute

# Collaboration ISS and WIV-ISP



Data bases **+** Source code

Pipeline

Belgian community     EU community

.be

**Thank you for your attention!**
**Questions?**