Ana Rita Rebelo

*anrire@food.dtu.dk*

# Recommendations and proposals to develop harmonised protocols for CRE/CCRE surveillance and outbreak detection
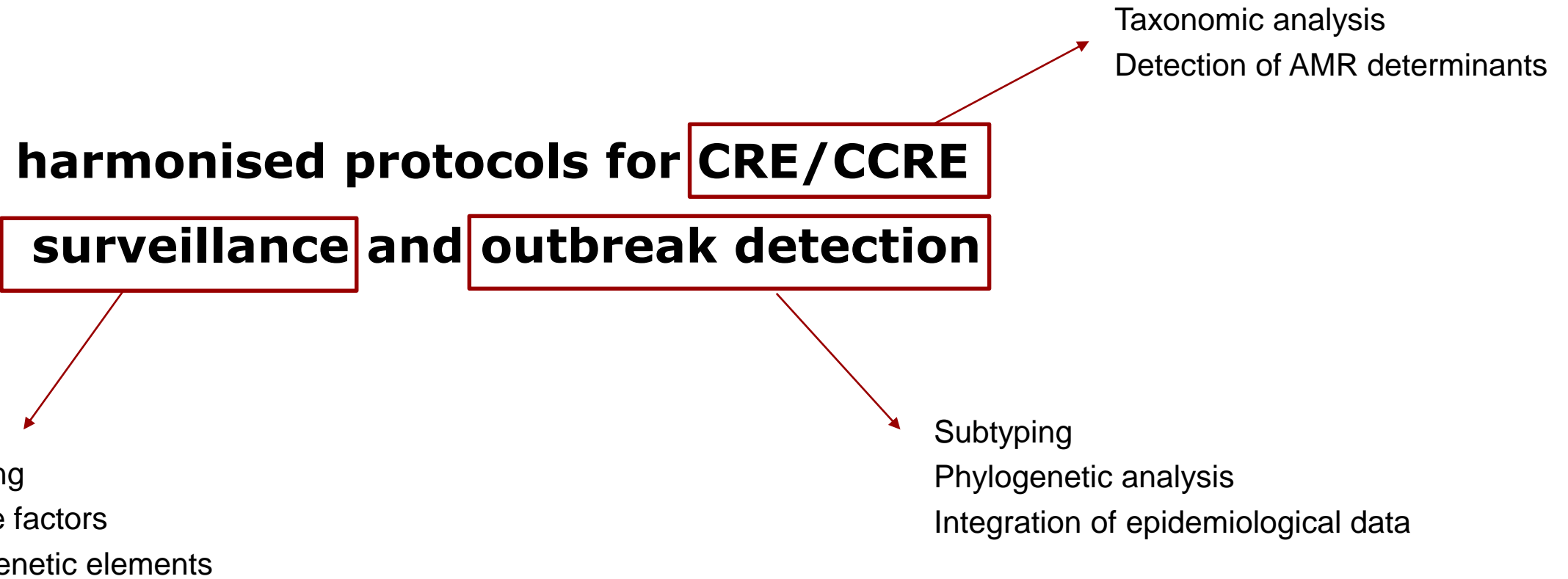
# Rationale – Why WGS?

*Advantages*

- Only one protocol
- Very large amount of data
- Higher discriminatory power
- Harmonised and automatic analysis
- Direct comparison
- Ease of storage
- Retroactive screening

*Why now?*

- Increase in sequencing accuracy
- Decrease in cost
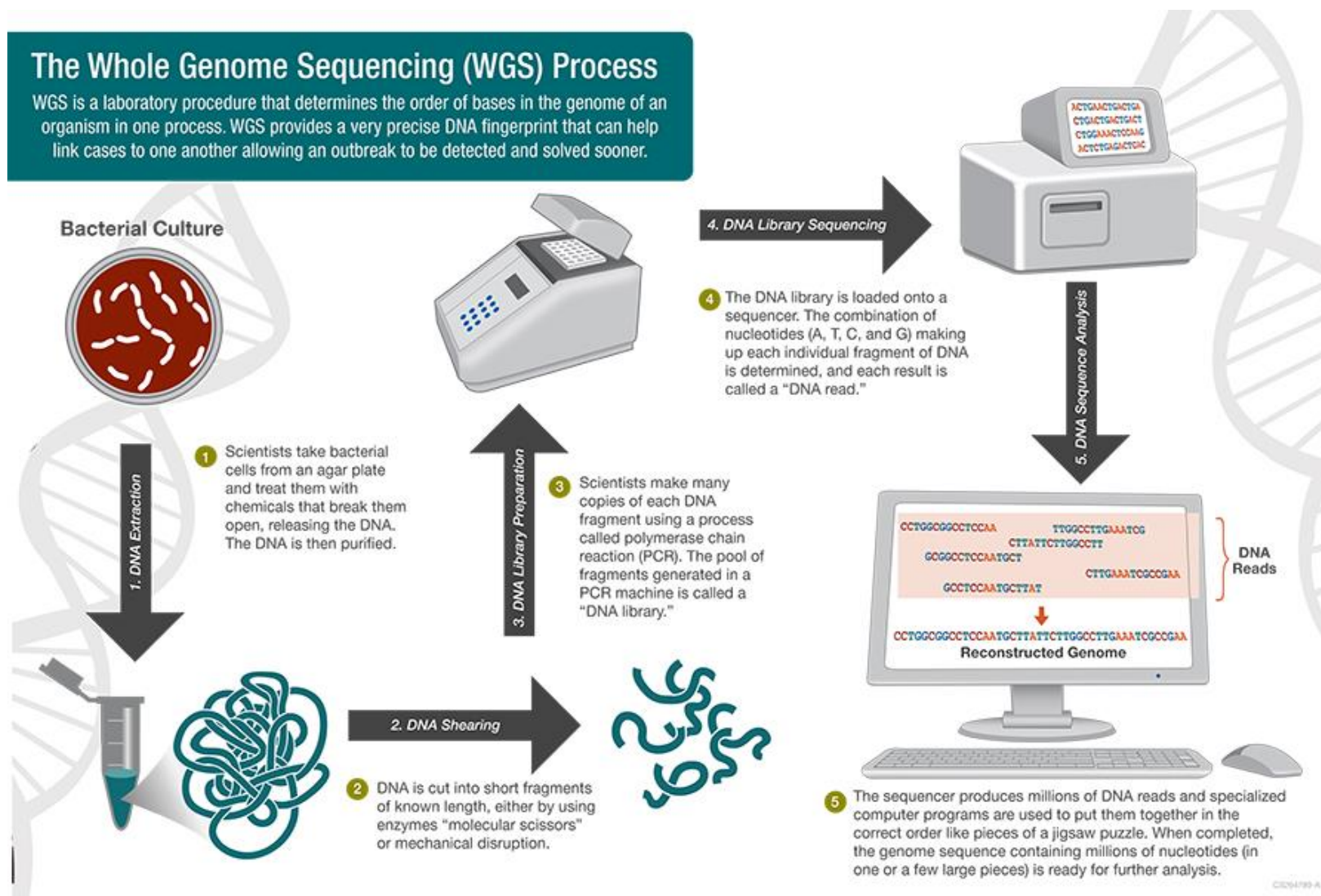- Coordinated efforts throughout Europe

https://www.cdc.gov/pulsenet/pathogens/protocol-images.html#wgs

# WGS-BASED ANALYSIS OF BACTERIA - REQUIREMENTS

- Expertise on DNA extraction methods

- Expertise on library preparation methods

Not too technically demanding

Ideally a dedicated room

- Access to sequencing platform

- Access and expertise on bioinformatics tools

Main challenges: cost, implementation

- Data management infrastructure

Main challenges: cost, compatibility

| Tool | Reference database | Description and output |
|---|---|---|
| **Tools for taxonomic analysis and typing** | | |
| **KmerFinder [150,151]** | KmerFinder | Provides hits of the query genome against whole reference genomes, the respective % of identity and % of coverage |
| **SILVA [152]** | SILVA | Collection of 16S rRNA genes, also possible to perform phylogenetic analysis and obtain phylogenetic trees |
| **MLST [153]** | PubMLST | MLST schemes, provides the sequence type |
| **rMLST [126]** | rMLST | rMLST schemes, provides the predicted species and respective allelic support metric |
| **SerotypeFinder [154]** | SerotypeFinder | Serotype, specific for E. coli |
| **SeqSero [** | | |
| **PneumoC** | | |
| **Tools for** | | |
| **cgMLST [** | | be used exclusively for typing but also clustering |
| **CSIPhylog** | | through FastTree |
| **Evergreen** | | |
| **Tools for detection of antimicrobial resistance determinants** | | |
| **ResFinder [129]** | ResFinder, PointFinder | Provides hits against reference ARGs and PMs and the respective % of identity and % of coverage, position in genome and predicted phenotype |
| **KmerResistance [150,151]** | KmerResistance | Provides hits of the query genome against reference genomes, as well as the detected ARGs and respective % of identity and % of coverage |
| **CARD/RGI [158]** | CARD | Provides hits against reference ARGs and respective % of identity and % of coverage. Other options are possible and the service is highly focused on ontology and standardization |
| **AMRFinder [159]; AMRFinderPlus [160]** | NCBI RefSeq | Provides hits against reference ARGs and respective % of identity and % of coverage |
| **ARIBA [161]** | CARD, ResFinder, NCBI Bacterial AMR Reference Gene DB, ARG-ANNOT, MEGARes, PubMLST, others defined by user | Provides hits against reference ARGs and respective % of identity and % of coverage |
| **Tools for detection of virulence factors** | | |
| **VirulenceFinder [162]** | VirulenceFinder | |
| **Victors [131]** | Victors | |
| **Tools for detection and analysis of mobile genetic elements** | | |
| **PlasmidFinder [163]** | PlasmidFinder | |
| **Platon [164]** | Platon | Provides hits against reference plasmids and respective % of identity and % of coverage, as well as relevant genes |
| **pMLST [153]** | PubMLST | Plasmid typing schemes |
| **MobileElementFinder [135]** | MobileElementFinder | Provides type and reference sequences of MGEs, respective % of identity and % of coverage, as well as associated ARGs and VFs |
| **Pipelines for extensive analyses** | | |
| **NCBI Pathogen Detection** | NCBI DBs | Detects ARGs and VFs, provides SNP-based phylogenetic analysis |
| **Pathogenwatch [165]** | Pathogenwatch, tools' DBs | Performs taxonomic analysis, determines MLST and cgMLST and provides cgMLST-based phylogenetic clustering |
| **BIGSdb [166]** | PubMLST BIGSdb | Performs annotation and taxonomic analysis, detects ARGs and plasmids, determines MLST, rMLST and cgMLST, provides phylogenetic and spatio-phylogenetic analysis |
| **PATRIC [167]** | PATRIC, but also includes others such as CARD, NDARO and VFDV | Performs assemblies, quality control, annotation and taxonomic analysis, detects ARGs and performs phenotype prediction, detects VFs and MGEs, provides phylogenetic analysis, variation analysis and genome alignments |

This is just a subset….

….and just from tools for analysis .

**Accounting for future priority pathogens and priority AMR profiles for surveillance.**

# RECOMMENDATION 1

**MAIN OBJECTIVE:** the harmonised WGS approach is streamlined for surveillance and analyses of CRE/CCRE bacteria

**BUT:** should already account for the future integration of two additional epidemic-prone healthcare-associated antimicrobial-resistant bacterial pathogens of public health priority

**BECAUSE:** current European guidelines and Commission Implementing Decisions might at any time be reviewed if new epidemiological situations become established, changing the scope of organisms under surveillance

**AND:** we should avoid working in silos

**EXAMPLES:**

o selection of ARG databases that contain ARGs and point mutations for other microorganisms besides CRE/CCRE

o defining QC parameters at are not exclusive to *Enterobacterales*

*establishing the acceptable genome size deviation as a percentage of the expected genome size, and not as a numerical deviation from the expected 5 million base-pairs.*

**Choose well defined AMR genotypes for validation of the WGS approach.**

# RECOMMENDATION 2

**MAIN OBJECTIVE:** Predict AMR phenotypes (which are often regulated by a combination of several genotypic determinants)

**BUT:** many of the genetic mechanisms are currently not well defined and prove difficult to be detected through WGS approaches

**SOLUTION:** Defining a subset of well-studied ARGs and point mutations allows for comparison of sampling methods, laboratory protocols and bioinformatics workflows

**EXAMPLE:** Adequately dectecting *mcr*-genes in colistin-resistant *Enterobacterales* could correspond to the benchmark indicating an adequate approach, even knowing that other mechanisms of resistance (e.g. PMs) exist

**Establish the control parameters to be used.**


**Establish the thresholds for the control parameters.**

Many different:

- DNA extraction kits
- Sequencing platforms
- Bioinformatics approaches
- Bioinformatics tools

Is harmonization feasible?

**Well defined set of QC parameters**

- For the raw data

   *E.g. nr. and length of raw reads, depth of coverage*

- For the assembled genomes

   *E.g. N50, nr. of contigs, genome size*

- For the performance of the tools

   *E.g. accurately detect PMs and ARGs in sets of benchmarking data*

# RECOMMENDATION 3 AND RECOMMENDATION 4

We propose two sets of QC parameters that will allow the users to maintain flexibility in their choice of WGS platform and bioinformatics tools:

o   Sequencing QC parameters

o   Data management QC parameters

**Sequencing QC parameters:**

any sequencing platform, protocol and bioinformatics tool

$\downarrow$

raw data and assembled genomes with quality equal or above limited thresholds for defined control parameters

Current general consensus:

o   "number and length of raw reads"

o   "depth of coverage"

o   "number of contigs in the assembled genome"

o   "N50 and deviation from expected genome size"

# RECOMMENDATION 3 AND RECOMMENDATION 4

**Sequencing QC parameters:**

o the depth of coverage (both of the raw data and also of the assembled genome) should be at least of 30 times (30X)

*15X has proven sufficient to an adequate detection of ARGs and point mutations in* E. coli

o the number of contigs in the assembled genome must be lower than 1,000, and ideally lower than 500

*short-read technologies*

o the size of the assembled genome is dependent on the target species. To account for genome plasticity and mobile genetic elements we suggest that a maximum threshold of 10% of variation in the number of base-pairs (BPs) should be adopted.

*CRE/CCRE: genome size 5 million BPs = a variation of plus/minus 0.5 million BPs would be acceptable*

**Data management QC parameters:**

Ensure that all workflows respect the same data management directions and there is:

-traceability of data and methods

-compatibility of data types and formats between different bioinformatics approaches

-comparability of results between settings

**Data management QC parameters:**

Parameters:

 ?

ISO standard will soon become available and help guide us

Examples from ISO draft:

o   defining the <u>minimum metadata requirements</u> and respective <u>adequate descriptors</u>

o   describing the proper <u>registry methods</u> for the DNA extraction and sequencing protocols, WGS platforms and bioinformatics analyses.

**Data management QC parameters:**

--secret--

# BENCHMARKING DATASETS

## European Commission's Joint Research Centre

**2018: "The challenges of designing a benchmark strategy for bioinformatics pipelines in the identification of antimicrobial resistance determinants using next generation sequencing technologies"**

**2021: "A roadmap for the generation of benchmarking resources for antimicrobial resistance detection using next generation sequencing"**

o Benchmarking approaches to validate sequencing and bioinformatics workflows

o Ensure that different pipelines can be used while at the same time adhering to the same minimum standards of performance

# BENCHMARKING DATASETS

## HOW?

o creating platform-specific validation datasets

o using simulated data that complies with specific certifications (not biases towards the creating platform)

o accepting the *fastq* format as the standardized input for analysing the performance of bioinformatics tools

o **accuracy should be dependent on the ability of the workflow to identify the correct AMR determinants that are introduced in the dataset (and not on agreement with phenotypic results)**

o these AMR determinants should include species and mechanisms which are selected based on international priority lists

Define a set of bioinformatics tools and databases as potential candidates to be included in the harmonised approach.

# BIOINFORMATICS TOOLS - OVERVIEW

**Purpose**

Quality control
Assembly
Taxonomic analysis
Phylogeny
Serotyping/Subtyping
Detection of AMR determinants
Detection of other determinants

**Accessibility**

Web-based
Command line
        Local vs. server

**Maintenance**

Benchmarked
Curated

**Cost**

Open access
Subscription

**Data**

Raw data as input
Assemblies as input
Integration of metadata

# EXAMPLE: BIOINFORMATICS TOOLS FOR PREDICTION OF AMR

| Tool | Target species | Reference database | Output | Comments |
|---|---|---|---|---|
| **SRST2** | All | CARD, PubMLST, or others defined by user | Reference sequences and respective % of coverage, depth | Also taxonomy, phylogeny, VFs, plasmids, other, depending on provided databases |
| **ARIBA** | All | CARD, ResFinder, ARG-ANNOT, MEGARes, NCBI Bacterial AMR Reference Gene Database, PubMLST, or others defined by user | Reference sequences and respective % of identity, % of coverage | Also phylogeny, VFs and plasmids, depending on provided databases (such as plasmidfinder, VFDB, VirulenceFinder) |
| **KmerResistance** | All | Own | Reference genomes, ARGs and respective % of identity, % of coverage | Also taxonomy |
| **ResFinder** | All | Own | ARGs and respective % of identity, % of coverage, position in genome, predicted phenotype | NA |
| **PointFinder** | Limited | Own | Mutated gene, protein translation, predicted phenotype | Included in ResFinder but can be used locally by itself. Currently under development for Klebsiella spp. |
| **RGI** | All | CARD | Reference sequences and respective % of identity, % of coverage, other options | Integrated in the Galaxy server; allows proteome analysis |
| **AMRFinder; AMRFinderPlus** | All (Limited PMs) | NCBI RefSeq | Reference sequences and respective % of identity, % of coverage | Included in NCBI Pathogen Detection |
| **SSTAR** | All | Own (created by merging ResFinder and ARG-ANNOT) | Reference sequences and respective % of coverage, depth | Can be used with other reference databases |
| **ABRicate** | All | CARD, ResFinder, ARG-ANNOT, MEGARES, NCBI AMRFinderPlus, or others defined by user | Reference sequences and respective % of identity, % of coverage | Also VFs and plasmids, depending on provided databases Also VFs and plasmids, depending on provided databases (such as plasmidfinder and VFDB) |
| **CARD** | All | Own | ARGs and point mutations, respective prevalence and predicted phenotype | Highly focused on ontology and standardization. VFs and mobile genetic elements currently being added |

# EXAMPLE: BIOINFORMATICS PIPELINES FOR PREDICTION OF AMR

| Tool | Target species | Reference database | Output |
|------|---------------|--------------------|--------|
| Pathogenwatch | Limited | Own, tools' databases | Taxonomy, MLST, cgMLST and clustering. Other functionalities for Klebsiella spp. derived from Kleborate |
| Enterobase | Limited | Tools' databases | Genome assembly and annotation, serotyping, MLST, cgMLST, rMLST, phylogenetic analysis |
| BIGSdb | All | PubMLST BIGSdb | Annotation, taxonomy, ARGs, plasmids, MLST, rMLST, cgMLST, phylogenetic and spatio-phylogenetic analysis, comparative genomics |
| NCBI Pathogen Detection | Limited | Own | ARGs, VFs, SNP-based phylogenetic analysis |
| PATRIC | All | Own, but also includes others such as CARD, NDARO and VFDV | Assemblies, QC, annotation, taxonomy, ARGs, phenotype prediction, VFs, mobile elements, phylogenetic analysis, variation analysis, genome alignments, comparative genomics, other options |
| Ridom SeqSphere+ € | All | Own, tools' databases; Includes NCBI AMRFinder and VFDB | Assemblies, QC, taxonomy, ARGs, VFs, MLST, cgMLST, phylogenetic analysis |
| Bionumerics € | All (wgMLST schemes available for limited species) | Own, tools' databases, others provided by user | Assemblies, QC, annotation, taxonomy, ARGs, plasmids, MLST, rMLST, cgMLST, wgMLST, phylogenetic analysis, comparative genomics, other options for E. coli (ARGs, PMs, VFs, plasmids, serotypes) |

o should be **open-access**, in order to respect probable budget limitations of certain users

o should be available as **online interfaces** (or be part of other interfaces) to avoid the need for expensive computing resources and specific professionals

o should also be **downloadable for local** usage for users with the resources and interest in doing so

o should be **benchmarked**, and should be transparently and continuously **curated**

o The owners should be available for **collaboration** with the European Commission, ECDC, and this provider to facilitate implementation and testing, and to potentially coordinate events such as suspension of updates during External Quality Assessments, and to provide scientific and technical support to the users

**Consider the results of other tasks incurred in the context of this project.**

o WGS approaches currently used by the NRLs

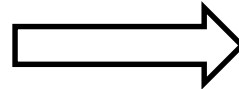o minimum conflict as possible with the ones currently used

anrire@food.dtu.dk

# Thank you!