

# Interview

# Folkhälsomyndigheten

Marcela Pereira

*Implementation Note: Following assignment guidance for limited computational resources, pipeline architecture and fractioned database for assembly (2% sized to allow local run) of results demonstrate production workflow for target cloud environment.*

*Pipeline git link: [https://github.com/marcelladane/NF\\_nextflow\\_pipeline.git](https://github.com/marcelladane/NF_nextflow_pipeline.git)*

# Pipeline design rationale

## Workflow technologies selection

- Nextflow
  - Native SLURM integration
  - Good AWS support (including Batch and S3)
  - Built-in containerisation support
  - Better for HPC environments than Snakemake
  - Active slack community for help troubleshooting

### ABRicate database:

- **Database subsetting strategies** (random sampling portions of fasta file).
  - Decreased size allowed for local run.
  - Also advantageous when running new tools during the test phase
    - Decreased computational demands.
    - Allows for fast validation (CI/CD processes).











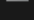

## Pipeline tools

- Quality control:
  - FastQC + MultiQC
    - Standard for both Illumina and Nanopore
    - Lightweight and fast execution
- Assembly:
  - Illumina: SPAdes (Quite robust, good for bacterial genomes)
  - Nanopore: Flye (less computationally intensive and faster than the other option (Canu), also has simple parameter tuning and is very good for small genomes)
- AMR Annotation:
  - ABRicate: good for mutation detection and good python compatibility.
  - Note: in case of metagenomes - ARM++ instead (Rationale: Quantitative, handles mixed communities; gives abundance metrics)
- Additional reporting:
  - Quast + Bandage: Assembly QC + figures

# Pipeline development






- Pipeline troubleshooting

How many troubleshooting scripts can that take? A lot ...

 abricate_debug_script	2025-06-09 21:56	SH Source File
 abricate_dependency_debug	2025-06-09 22:06	SH Source File
 create_small_nanopore	2025-06-09 17:15	SH Source File
 database_content_debug	2025-06-09 22:18	SH Source File
 extract_real_sequences	2025-06-09 22:15	SH Source File
 fixed_sequence_extraction	2025-06-09 22:22	SH Source File
 klebsiella_database	2025-06-09 22:01	SH Source File
 main.nf.backup	2025-06-09 17:27	BACKUP-fil
 main.nf.backup_abricate	2025-06-10 09:34	BACKUP_ABRICAT...
 morning_diagnostic	2025-06-10 09:25	SH Source File
 quast_diagnostic	2025-06-10 11:26	SH Source File
 setup_fohm_databases	2025-06-09 15:06	SH Source File

Best practice on CI/CD and reproducible development:

(Do small updates for each step developed with a well explained comment – easier to track back and return to previous versions)

<b>update readme</b>  marcelladane pushed 1 commit to <b>main</b> • 2936ca6...5a1a230 • yesterday
<b>update readme</b>  marcelladane pushed 1 commit to <b>main</b> • 2633efa...2936ca6 • yesterday
<b>Update QC pipeline with demonstration outputs</b>  marcelladane pushed 1 commit to <b>main</b> • ae74101...2633efa • yesterday
<b>Add FastQC quality control and MultiQC reporting:</b>  marcelladane pushed 1 commit to <b>main</b> • d2d76d5 • yesterday
<b>Remove pipeline_n</b>  marcelladane pushed 1 commit to <b>main</b> • d2d76d5 • yesterday
<b>feat: implement da</b>  marcelladane created <b>main</b> • d2d76d5 • yesterday

Add FastQC quality control and MultiQC reporting;

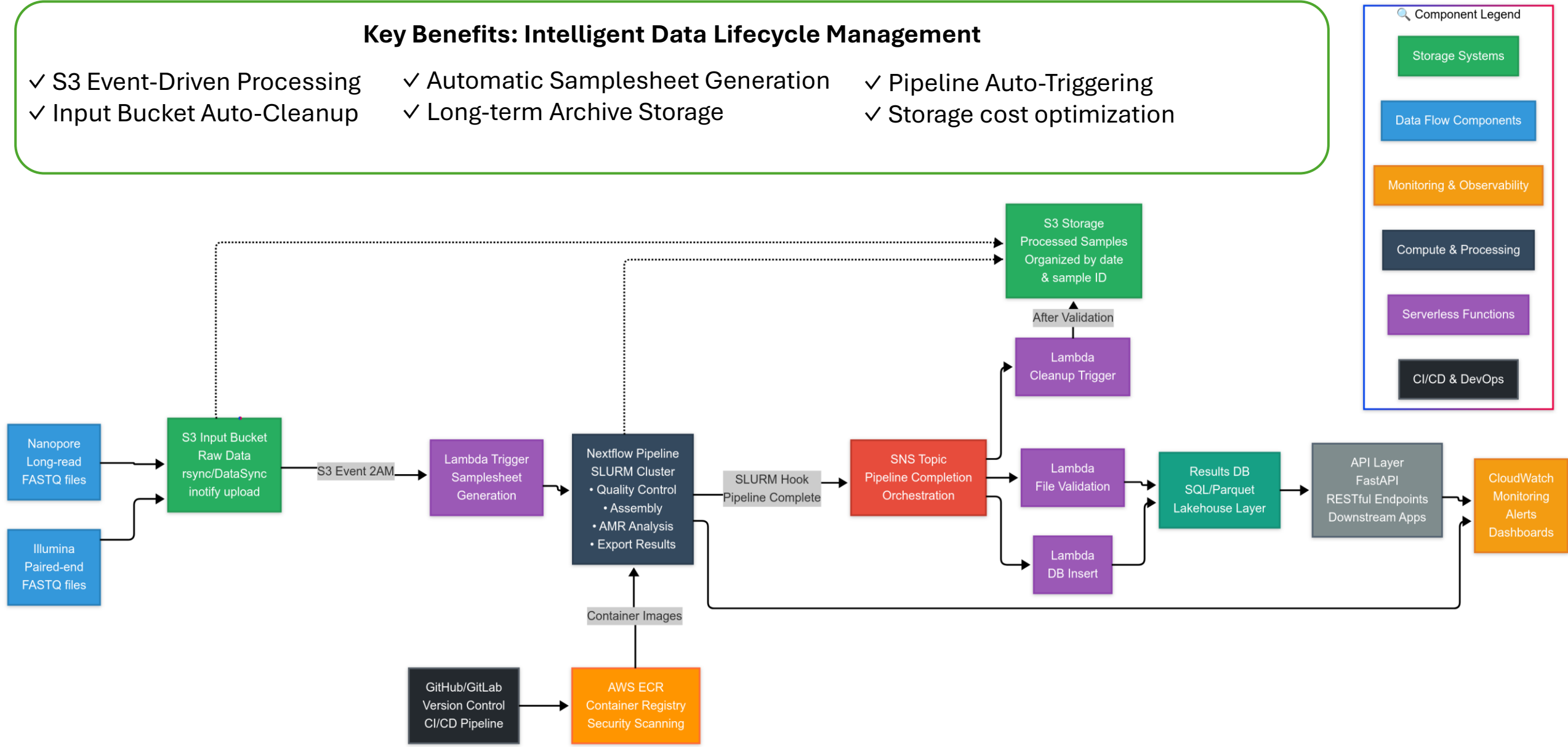
- Implement FastQC analysis for both Illumina and Nanopore reads
- Add MultiQC module for aggregated quality control reports
- Include error handling for FastQC configuration issues
- Add multiple deployment profiles (local, production, cloud)
- Update resource management for different environments

# Data flow automation

- Technologies:** Nextflow workflows + SLURM + Docker containers + Lakehouse architecture
- Data flow:** Illumina (paired-end) + Nanopore (single-end) → Unified analysis pipeline

## Key Benefits: Intelligent Data Lifecycle Management

- ✓ S3 Event-Driven Processing
- ✓ Automatic Samplesheet Generation
- ✓ Pipeline Auto-Triggering
- ✓ Input Bucket Auto-Cleanup
- ✓ Long-term Archive Storage
- ✓ Storage cost optimization



# Production Deployment - Serverless + HPC Integration

- Lambda Function Examples:

## Samplesheet Generation Trigger

```
1 # lambda_exemple.py
2 # This AWS Lambda function processes S3 events to handle FASTQ file uploads,
3 # extracts metadata, creates a samplesheet, and triggers a Nextflow pipeline.
4
5 import boto3
6 import json
7
8 def lambda_handler(event, context):
9     # S3 event triggers when new FASTQ uploaded
10    bucket = event['Records'][0]['s3']['bucket']['name']
11    key = event['Records'][0]['s3']['object']['key']
12
13    # Extract metadata and create samplesheet
14    sample_id = extract_sample_id(key)
15    platform = detect_platform(key) # illumina/nanopore
16
17    # Generate samplesheet.csv
18    create_samplesheet(sample_id, platform, bucket, key)
19
20    # Trigger Nextflow pipeline
21    trigger_slurm_job(sample_id)
22
```

## Database update validation

```
23 # Lambda function to validate and cleanup after processing
24 def validate_and_cleanup(event, context):
25     sample_id = event['sample_id']
26
27     # Multi-step validation
28     amr_present = check_amr_results(sample_id)
29     qc_complete = verify_qc_metrics(sample_id)
30     assembly_valid = validate_assembly_stats(sample_id)
31
32     if all([amr_present, qc_complete, assembly_valid]):
33         # Move to archive and cleanup
34         move_to_archive(sample_id)
35         cleanup_input_bucket(sample_id)
36         send_completion_notification(sample_id)

```

- SLURM integration example script:

```
1 bash#!/bin/bash
2
3 #SBATCH --job-name=fohm-amr-pipeline
4 #SBATCH --cpus-per-task=16
5 #SBATCH --mem=64G
6 #SBATCH --time=06:00:00
7 #SBATCH --mail-type=FAIL
8 #SBATCH --mail-user=bioinformatics-team@folkhalsomyndigheten.se
9 #SBATCH --output=/logs/pipeline_%j.out
10 #SBATCH --error=/logs/pipeline_%j.err
11
12 # Load modules
13 module load nextflow/23.04.0
14
15 # Set up error handling
16 set -e
17 trap 'send_failure_alert $SLURM_JOB_ID' ERR
18
19 # Run pipeline with auto-scaling
20 nextflow run main.nf \
21     --input s3://fohm-input/samplesheet.csv \
22     -profile production \
23     -with-tower \
24     -resume
25
26 # Success notification (optional)
27 echo "Pipeline completed successfully for job $SLURM_JOB_ID" | \
28     mail -s "FOHM AMR Pipeline Success" ops-team@folkhalsomyndigheten.se

```

# API Layer - Data Access & Integration

- Example code for API implementation

```
39 # FastAPI Application for AMR Data Access
40 from fastapi import FastAPI, Query
41 from fastapi.responses import StreamingResponse
42 import pandas as pd
43 import pyarrow.parquet as pq
44 import io
45
46 app = FastAPI(title="FOHM AMR API", version="1.0.0")
47
48 @app.get("/samples/")
49 async def get_samples(
50     date_from: Optional[str] = None,
51     date_to: Optional[str] = None,
52     bacteria: Optional[str] = None,
53     format: str = Query("json", enum=["json", "csv", "parquet", "excel"])
54 ):
55     # Query lakehouse Parquet tables
56     df = query_parquet_tables(date_from, date_to, bacteria)
57
58     if format == "csv":
59         return StreamingResponse(
60             io.StringIO(df.to_csv(index=False)),
61             media_type="text/csv",
62             headers={"Content-Disposition": "attachment; filename=amr_results.csv"}
63         )
64     elif format == "parquet":
```

```
1 # Key API Endpoints:
2 bash
3 # 1. Date-based Queries with Format Selection:
4 GET /api/v1/samples?date_from=2025-01-01&date_to=2025-01-31&format=csv
5 # Returns: CSV file download with all January 2025 samples
6
7 GET /api/v1/samples?date_from=2025-01-01&format=parquet
8 # Returns: Parquet file for data science workflows
9
10 # 2. Tool-specific Results:
11 GET /api/v1/analysis/spades?min_n50=100000&format=excel
12 # Returns: Excel file with SPAdes assemblies N50 > 100kb
13
14 # 3. Bacteria-specific AMR:
15 GET /api/v1/amr?bacteria=escherichia_coli&gene=blaCTX-M&format=json
16 # Returns: Excel file specific for AMR analysis
```










```
76 # Data Query Backend:
77 def query_parquet_tables(date_from, date_to, bacteria):
78     # Direct Parquet querying with PyArrow
79     table = pq.read_table('s3://fohm-lakehouse/amr_results.parquet')
80     df = table.to_pandas()
81
82     # Apply filters
83     if date_from:
84         df = df[df['date_processed'] >= date_from]
85     if bacteria:
86         df = df[df['bacteria_species'] == bacteria]
87
88     return df
```

# AMR Analysis Results & Platform Comparison

- **Illumina Results (Demonstrated)**
  - **Assembly Quality:** N50: 261,925 bp, 50 contigs, 5.23 Mb total
  - **AMR Detection Framework:** ABRicate with CARD database (mock implementation)
  - **Expected Detection:**  $\beta$ -lactamase, aminoglycoside, tetracycline resistance genes
  - **Advantages:** >99.9% accuracy, cost-effective, excellent for surveillance
- **Nanopore Potential (Production Environment)**
  - **Unique Capabilities:** Complete plasmid assemblies, resistance gene clusters
  - **Structural Insights:** Integron analysis, insertion sequences, chromosomal integration
  - **Clinical Value:** Superior for outbreak investigation and novel resistance mechanisms

Technology	Known Genes	Novel Variants	Structural Variants	Use case
Illumina	95-98%	70-80%	20-30%	Routine surveillance
Nanopore	90-95%	85-95%	80-90%	Outbreak investigation
Hybrid	98-99%	90-95%	85-95%	Comprehensive profiling

# Production Readiness & Next Steps

- **Immediate Implementation (Weeks 1-2):**
-  **Illumina Pipeline Deployment:** Production-ready with existing infrastructure
-  **SLURM Integration:** Automated job submission and resource management
-  **Basic API:** Core endpoints for surveillance team access
- **Short-term Optimization (Months 1-2):**
-  **Nanopore Integration:** Cloud resources for Flye assembly
-  **Real ABRicate Database:** Full CARD database implementation
-  **Complete Containerization:** All tools in production containers
- **Medium-term Enhancement (Months 3-6):**
-  **ML Integration:** Resistance prediction models
-  **Dashboard Development:** Real-time surveillance monitoring
-  **Clinical Integration:** LIMS system connectivity