

Advanced Molecular Detection Southeast Region Bioinformatics

Where to Start Guide Sam Marcellus, MPH 07/22/2024



Updates

- New Staff Introductions
 - Nikhil Reddy
 - Arnold Rodríguez-Hilario
- On-going Florida computer issues to be resolved soon

IT Department Communication



Communicating Needs



Understanding Capabilites



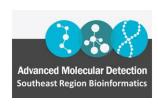
Defining Tasks



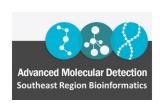
Bioinformatics Resources

Communicating with IT: The Basics

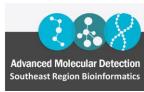
- Buy-in from leadership
 - Lab Leadership
 - Upper-Level Agency Leadership
- Identify allies in IT
 - Who is familiar with Linux?
 - Does anyone have experience with Cloud or computing clusters?
 - Who has worked well with you in the past?
- Are there any existing contracts to be aware of?
- Clearly (and repeatedly) state your needs (in writing)



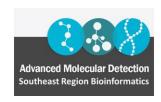
- Could range from bullet points in an email to a formal memo
- Next-Generation Sequencing has become an essential tool
 - Define NGS
 - List the uses of NGS in your lab
 - Surveillance
 - Clinical
 - Outbreak Investigation
 - Name all the groups or pathogens using NGS
 - Note that uses will only expand in the future
 - Discuss the increase in staff on the lab side to handle sequencing



- NGS and Bioinformatics have unique computational needs
 - Clouds "attached" to sequencers (BaseSpace)
 - Large amounts of data
 - Need larger than usual computational power
- Sample Tracking
 - Is current system working for you?
 - What's your ideal set up?
 - Can you offer what another state is using as an example?
- Current computing set up and why it's insufficient
 - HiPerGator?
 - It's down two weeks per year
 - Share resources with the region



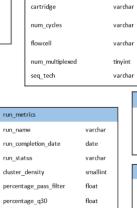
- List your sequencers and their capacities
 - GB of data per run
 - Samples per run
- List dry lab positions and their responsibilities
- Briefly describe epis and what they do with the sequencing data
- Any useful diagrams can be added as appendices
- Short-, medium-, and long-term goals for your team and its infrastructure





Samples_on_runs	
run_name	varchar
bphl_accession	char
unique_id	uniqueidentifier

run_requeues	
run_name	varchar
requeue_reason	text
corrective_action	text
requeue_date	date
requeue_tech	varchar



varchar

run_name

run_date

run_type

instrument_type

instrument_id

library_prep

percent_phix

indexes

bphl instrument name

enrichment_protocol

library_conc_loaded

,		cluster_id	,
rcha	r	strain_id	,
	repeat_s	samples	
	bphl_accession		char
	repeat_reason		text
	repeat_run		varchar
	repeat_r	runs	
	run_nam	ne_original	varchar

varchar

text

run_name_repeat

repeat_reason

smartgene_hiv unique_id

bphl_accession

analysis_pipeline

analysis_date

genotype

ghost_hep

unique_id

bphl_accession

analysis_pipeline

analysis_date

genotype

run_name

analyst

analyst

varchar

date

varchar

varchar

varchar

varchar

varchar

varchar

varchar

float

		bacterial_wgs	
uniqueidentifier		unique_id	uniqueidentif
char		bphl_accession	char
varchar		run_name	varchar
varchar		analyst	varchar
varchar		analysis_pipeline	varchar
date		analysis_date	date
varchar		organism_id_wgs	varchar
varchar		assembly_qc	varchar
		raw_reads	int
		dean_reads	int
uniqueidentifier		species_id_mash	varchar
char		nn_accession_mash	varchar
varchar		mash_distance	float
varchar		species_id_kraken2	varchar
varchar		kraken2_percent	float
date		mlst_scheme	varchar
varchar		st	smallint
varchar		mean_read_length	int
varchar		mean_read_quality	tinyint
varchar		est_genome_depth	int
		num_contigs	smallint
		longest_contig	int
		N50	int
		L50	smallint
ar		assembly_length	int
_		gc	float
		annotated_cds	mediumint
har		amr_report	(file)
har		kpc	varchar
		vim	varhcar
		ndm	varchar
		оха	varchar

imp

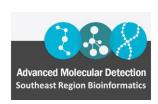
varchar

covid_wgs	
unique_id	uniqueidentifier
bphl_accession	char
run_name	varchar
analyst	varchar
analysis_pipeline	varchar
analysis_date	date
reference_seq	varchar
start_pos	int
end_pos	int
raw_reads	int
dean_reads	int
mapped_reads	int
percent_map_clean_reads	float
cov_bases_mapped	int
percent_genome_cov_map	float
mean_depth	int
mean_base_qual	tinyint
mean_map_qual	tinyint
assembly_length	int
num_n	int
percent_ref_genome_cov	float
vadr_flag	varchar
qc_flag	varchar
pangolin version	varchar
lineage	varchar



Understanding Capabilities

- Current policies in place
 - IT Policies
 - Data Retention Policies
 - Privacy Policies
 - If it's not written down, it's not a real policy
- List of available software and current partnerships
 - Any cloud agreements with other agencies?
 - IT staff familiar with Linux working in a different agency
- Overarching Data Modernization Plans



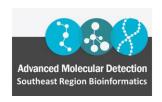
Defining Tasks

Laboratory Tasks

- Day-to-day management of resources
- Data analysis
- Implementation and use of bioinformatic pipelines

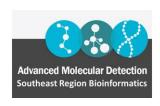
IT Tasks

- Install of new resources on state computers
- Data encryption
- Data protection
- Database construction
- Technical Support



Bioinformatics Resources

- Trainings helpful for new staff or IT staff interested in bioinformatics
- StaPH-B (State Public Health Bioinformaticians)
 - https://staphb.org/training.html
- CDC Genomic Epidemiology Toolkit
 - https://www.cdc.gov/advanced-molecular-detection/php/training/index.html
- Data Camp
 - https://app.datacamp.com/
 - Data Scientist in R Track
 - Data Analyst in SQL Track



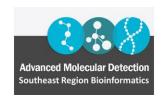
Bioinformatics Resources

PulseNet

- https://www.cdc.gov/pulsenet/hcp/about/next-gen-wgs.html
- An Overview of PulseNet Databases
 - https://doi.org/10.1089/fpd.2019.2637

Bioinformatics in Public Health

• Libuit KG, Doughty EL, Otieno JR, Ambrosio F, Kapsak CJ, Smith EA, Wright SM, Scribner MR, Petit Iii RA, Mendes CI, Huergo M, Legacki G, Loreth C, Park DJ, Sevinsky JR. Accelerating bioinformatics implementation in public health. Microb Genom. 2023 Jul;9(7):mgen001051. doi: 10.1099/mgen.0.001051. PMID: 37428142; PMCID: PMC10438813.





Advanced Molecular Detection Southeast Region Bioinformatics

Questions?

bphl-sebioinformatics@flhealth.gov

Sam Marcellus, MPH
Bioinformatician
Samantha.marcellus@flhealth.gov

Nikhil Reddy, MS
Bioinformatician
Nikhil.reddy@flhealth.gov

Molly Mitchell, PhD

Bioinformatics Supervisor

Molly Mitchell @ flhealth.gov