

# Species and Quasispecies analysis using the public/FDA-HIVE platform

Raja Mazumder [mazumder@gwu.edu](mailto:mazumder@gwu.edu)


Assoc. Prof. Biochemistry and Molecular Medicine; Project Lead, public-HIVE;  
Co-director, The McCormick Genomic & Proteomic Center

GW



# High-performance Integrated Virtual Environment (HIVE)

- Main
- About HIVE
- Tutorial
- Publications
- Tools
- People
- Employment



## Welcome to HIVE

*High-performance Integrated Virtual Environment*

HIVE is a cloud-based environment optimized for the storage and analysis of extra-large data, like Next Generation Sequencing data, Mass Spectroscopy files, Confocal Microscopy Images and others.

HIVE uses a variety of advanced scientific and computational visualization graphics, to get the **MOST** from your HIVE experience you must use a supported browser. These include Internet Explore 8.0 or higher (Internet Explorer 9.0 is recommended), Google Chrome, Mozilla Firefox and Safari.

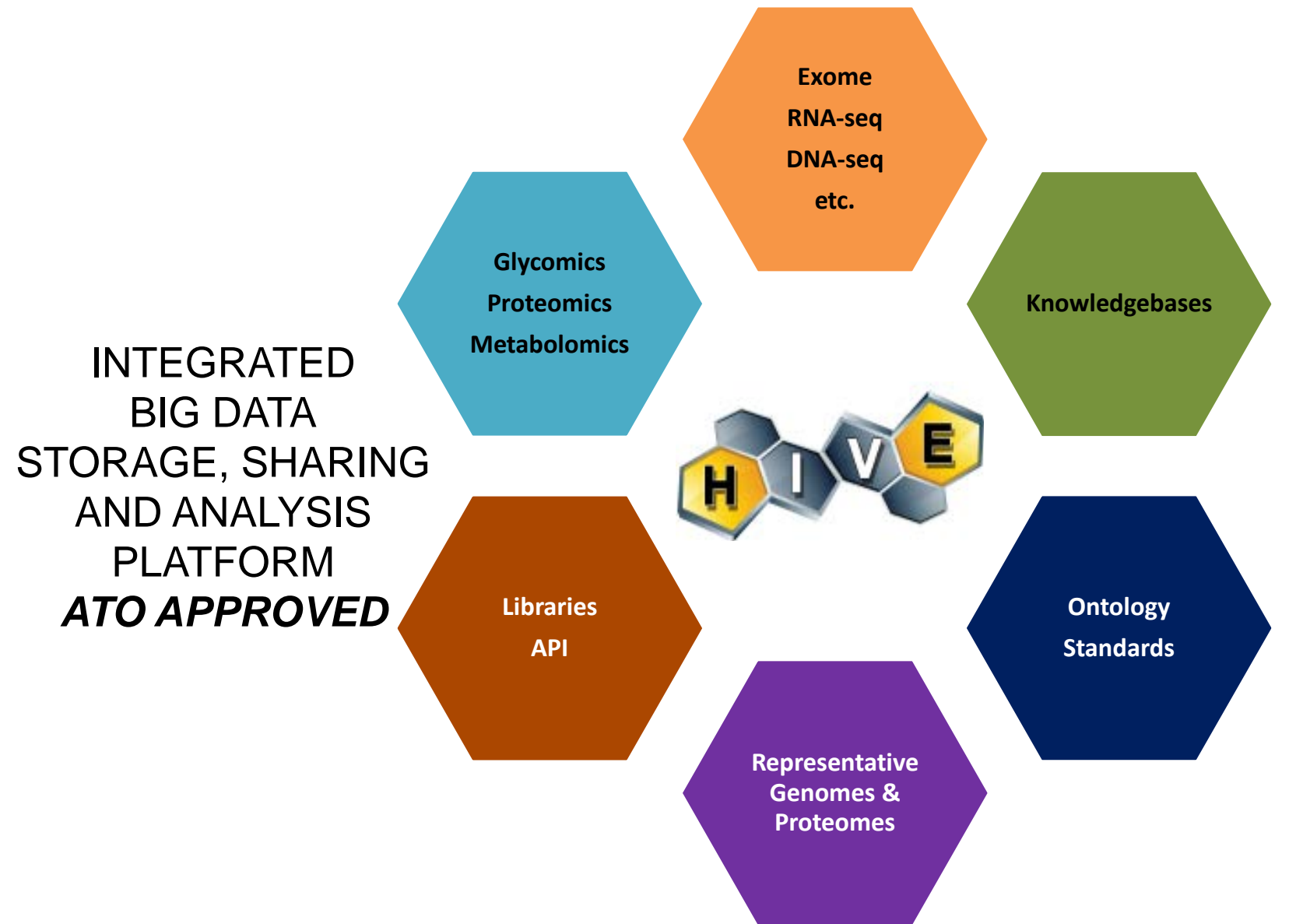
A few exemplary analytical outputs are displayed below for your enjoyment. But before you can take advantage of all that HIVE has to offer and create these objects for yourself, you'll need to [register](#).

Try browsing the tabs above or click [here](#) for more information!

**2 portals: Public portal (GW) & a FDA only portal**

[mazumder@gwu.edu](mailto:mazumder@gwu.edu) | [vahan.simonyan@fda.hhs.gov](mailto:vahan.simonyan@fda.hhs.gov)

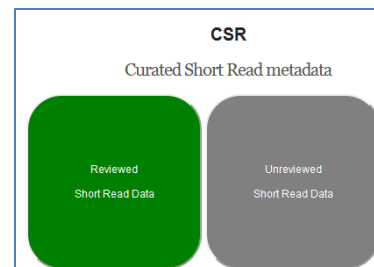
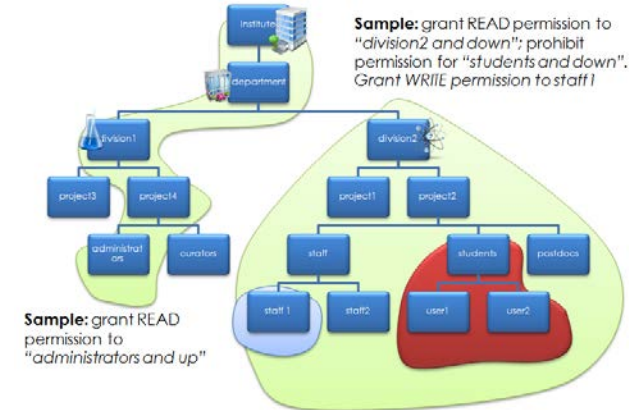
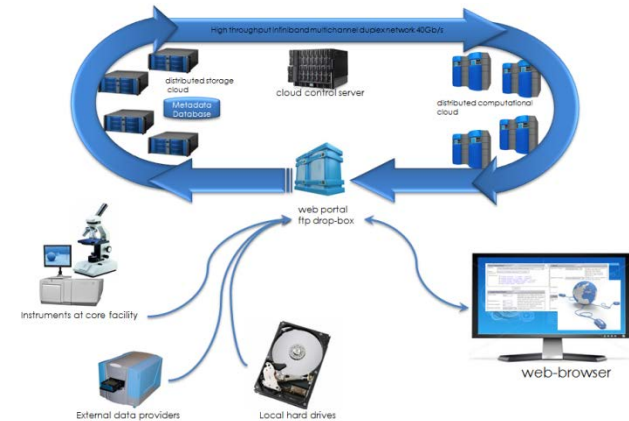
# HIVE tools + Community tools



References + Standards in collaboration with community

# HIVE key features and focus

- Distributed storage and computing
- Security features (FDA/industry/patient data in mind). Authorized to Operate (ATO) in Regulatory Environment
- Granular sharing capabilities
- Biocuration of content
- Compute speed

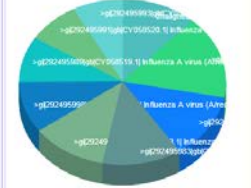


# HIVE: visualizations



**General Information**

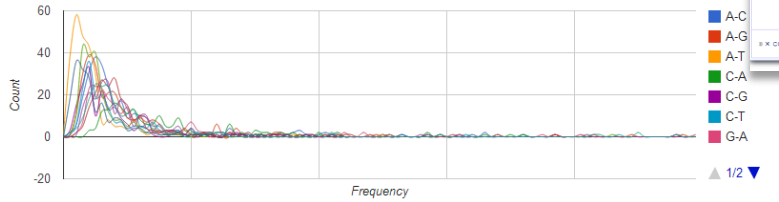
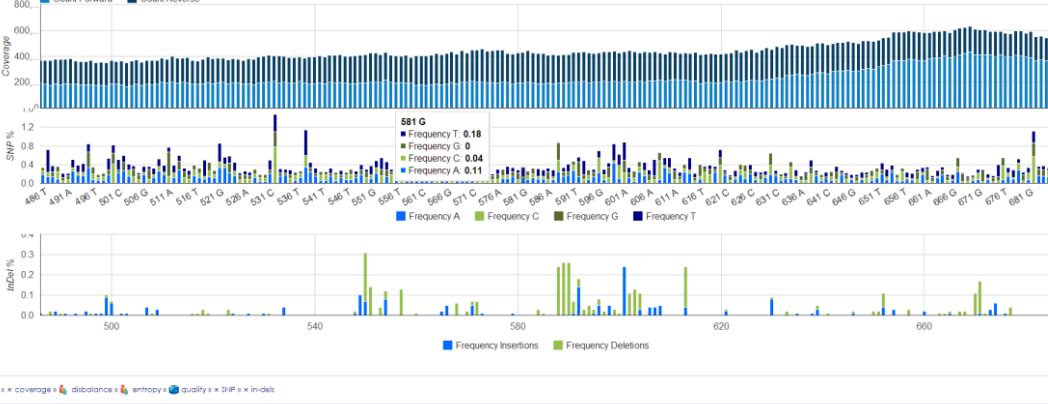
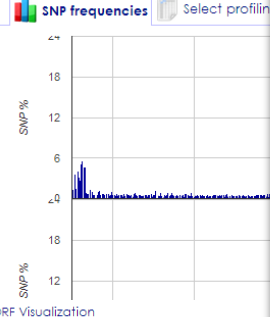
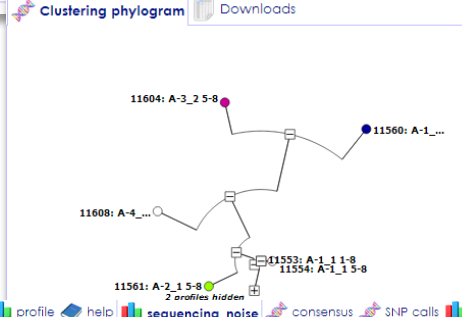
Reference Genome Length	934
Number of Reference Genomes	8
Mapped Regions	934
Total Contig Length	934
Mapped Coverage (% Reference)	100.00
Coverage on Contigs	381693
Total Number of Contigs	1
Unmapped Regions	0
Total Length of the Unmapped Regions	0.00
Unmapped Regions (% Reference)	0.00
Coverage on Gaps	0
Total Number of Gaps Found	0



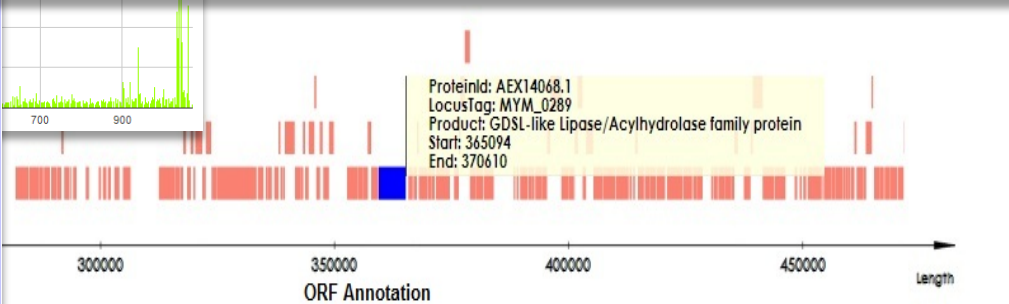
**Alignment**

```

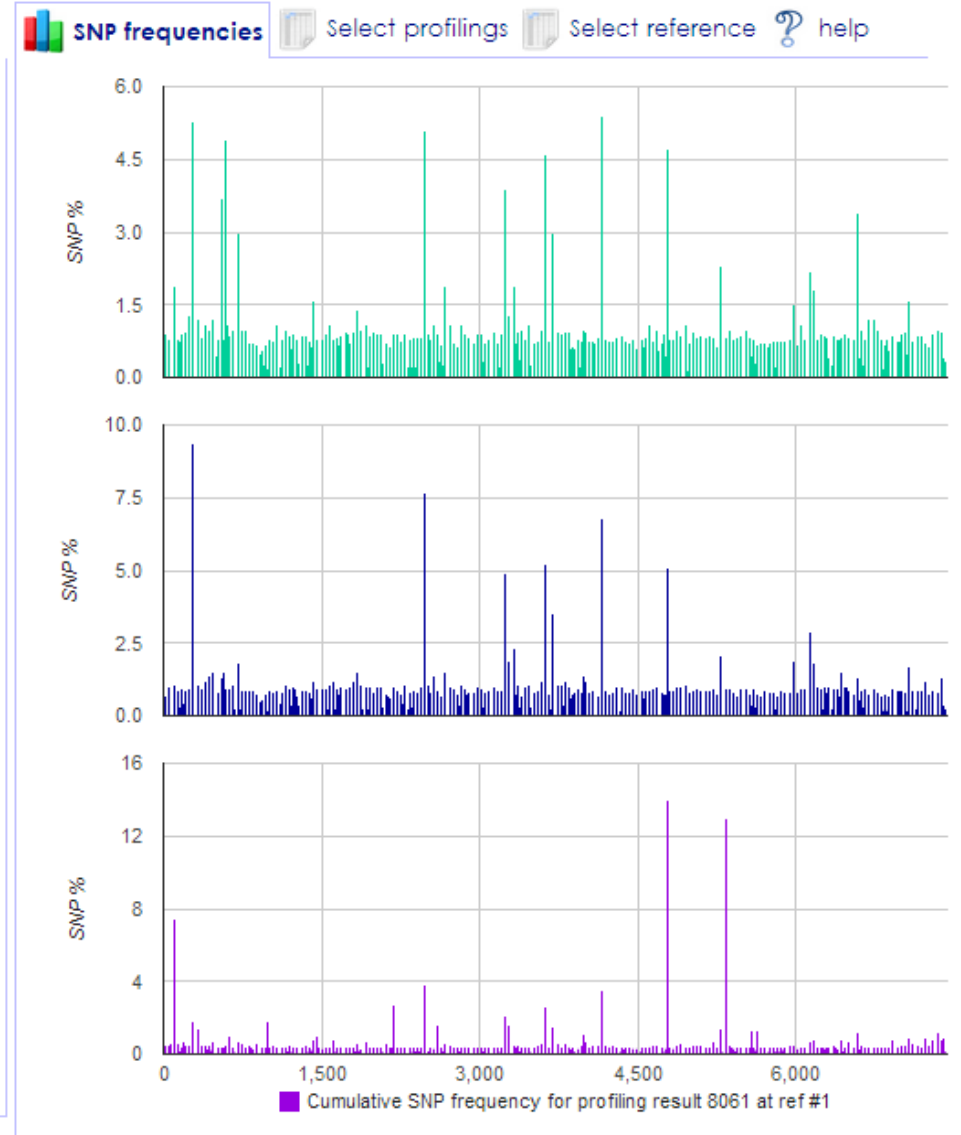
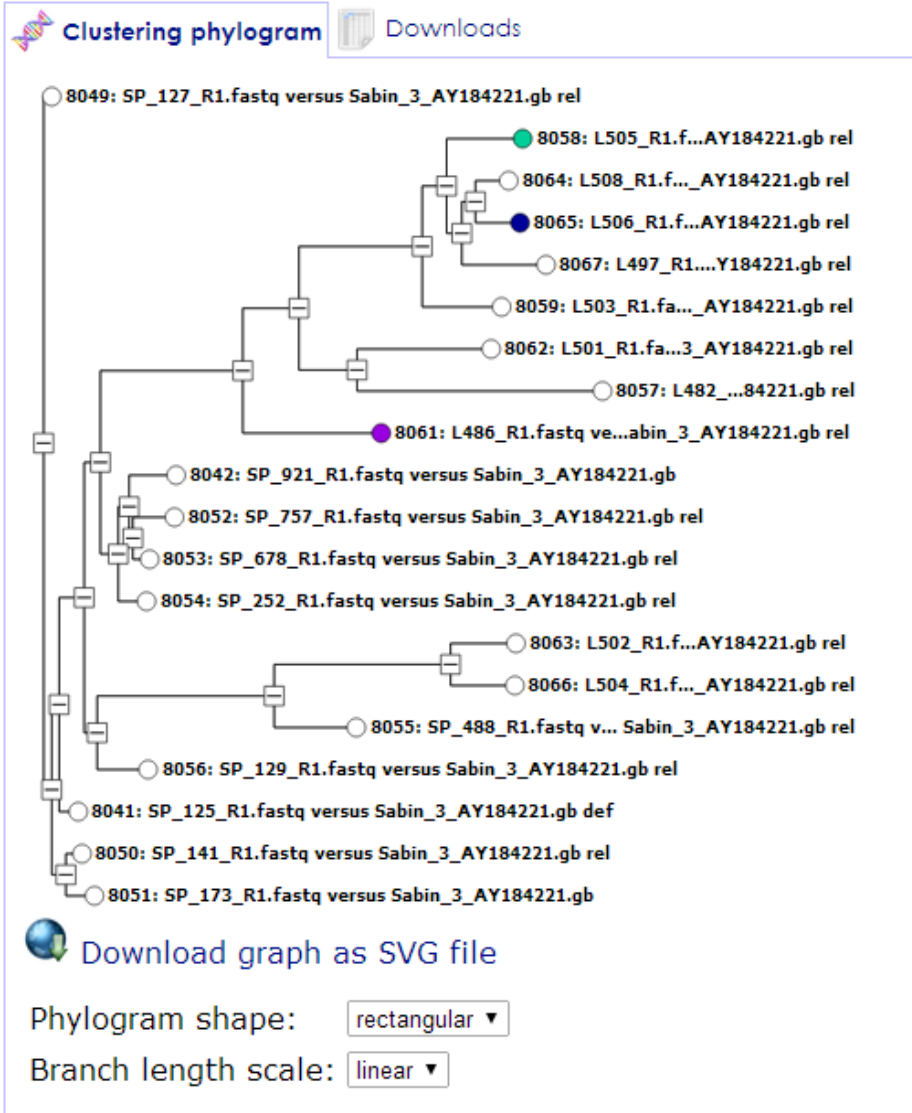
1 AGCGAAGCAGGTAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
19 AGCAAAAAGGTAAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
1 AGCGAAGCAGGTAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
16 AGCAAAAAGGTAAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
1 AGCGAAGCAGGTAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
10 AGCAAAAAGGTAAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
1 AGCGAAGCAGGTAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
14 AGCAAAAAGGTAAGATATTGAAAGATGAGTCTTCTAACCGAGTGGAAACCTAGCTCTCTATCATCCCGTCAGGCCCC
    
```



Threshold	A-C	A-G	A-T	C-A	C-G	C-T	G-A	G-C	G-T	T-A	T-C	T-G
50 %	0.0007	0.0007	0.0003	0.0012	0.0006	0.0006	0.0008	0.0005	0.0008	0.0004	0.0007	0.0007
75 %	0.0012	0.0012	0.0005	0.0018	0.0008	0.0010	0.0013	0.0007	0.0011	0.0008	0.0012	0.0011
85 %	0.0017	0.0015	0.0009	0.0028	0.0011	0.0014	0.0019	0.0009	0.0015	0.0012	0.0014	0.0015
90 %	0.0023	0.0020	0.0011	0.0033	0.0013	0.0016	0.0021	0.0011	0.0020	0.0018	0.0018	0.0020
95 %	0.0037	0.0028	0.0017	0.0043	0.0027	0.0023	0.0038	0.0014	0.0033	0.0020	0.0022	0.0030
99 %	0.0066	0.0083	0.0041	0.0070	0.0039	0.0053	0.0068	0.0031	0.0077	0.0057	0.0042	0.0090



# NGS → Trees



HIVE SNP based hierarchal clustering and phylogenetic analysis

Sample Reads

12 Reference Strains

1 Upload to HIVE

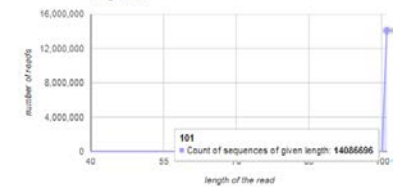


2 Automatic QC and Format Validation

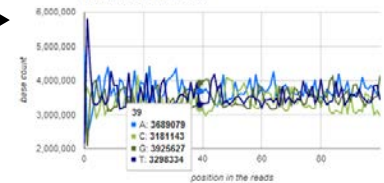
ACGT base count



Length count

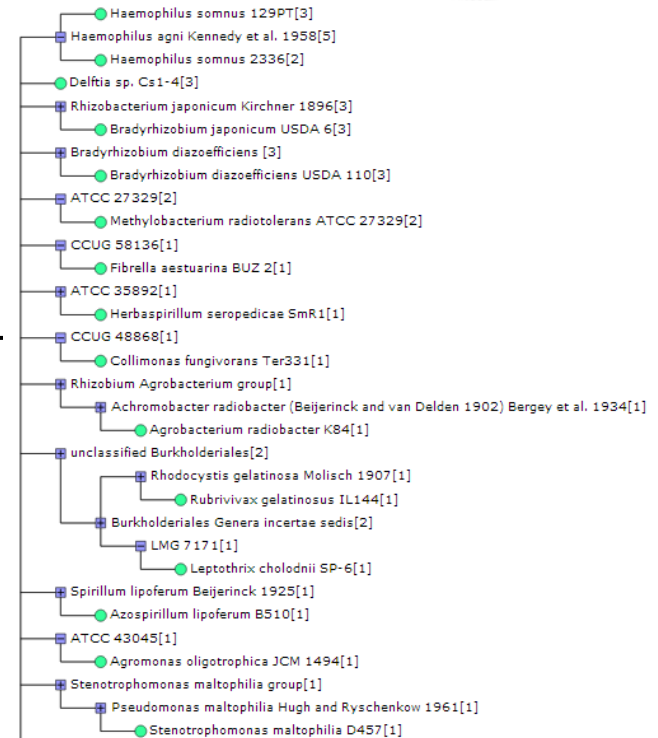


Lengthwise position count



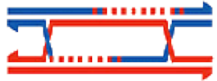
# HIVE: Pipeline

3 Analyze Contaminants using CensuScope

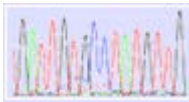


ITERATIVE LOOP: (Repeated within different parameters)

6a Recombinant mosaic discovery (Recombination Tool)



5 Align samples to mutual frame (HIVE-hexagon)



6b Generate SNV profile

6c Use SNV based phylogenetic tree to compare multiple samples (HIVE-phyloSNP)

8 Reinitiate loop with different parameters

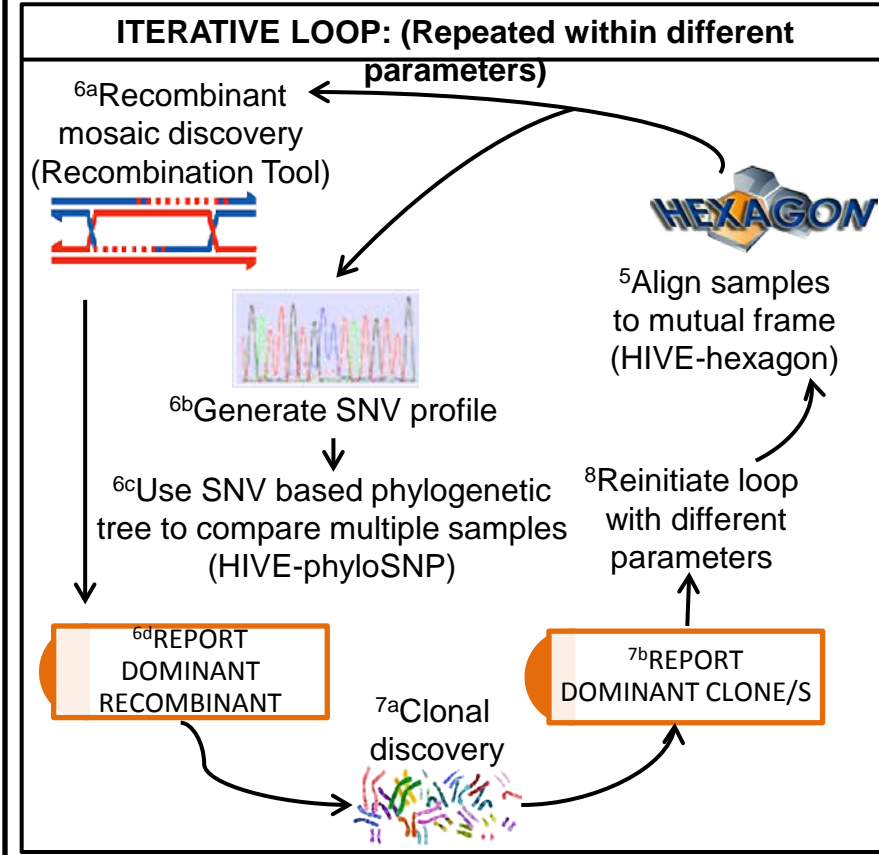
6d REPORT DOMINANT RECOMBINANT

7b REPORT DOMINANT CLONE/S

7a Clonal discovery

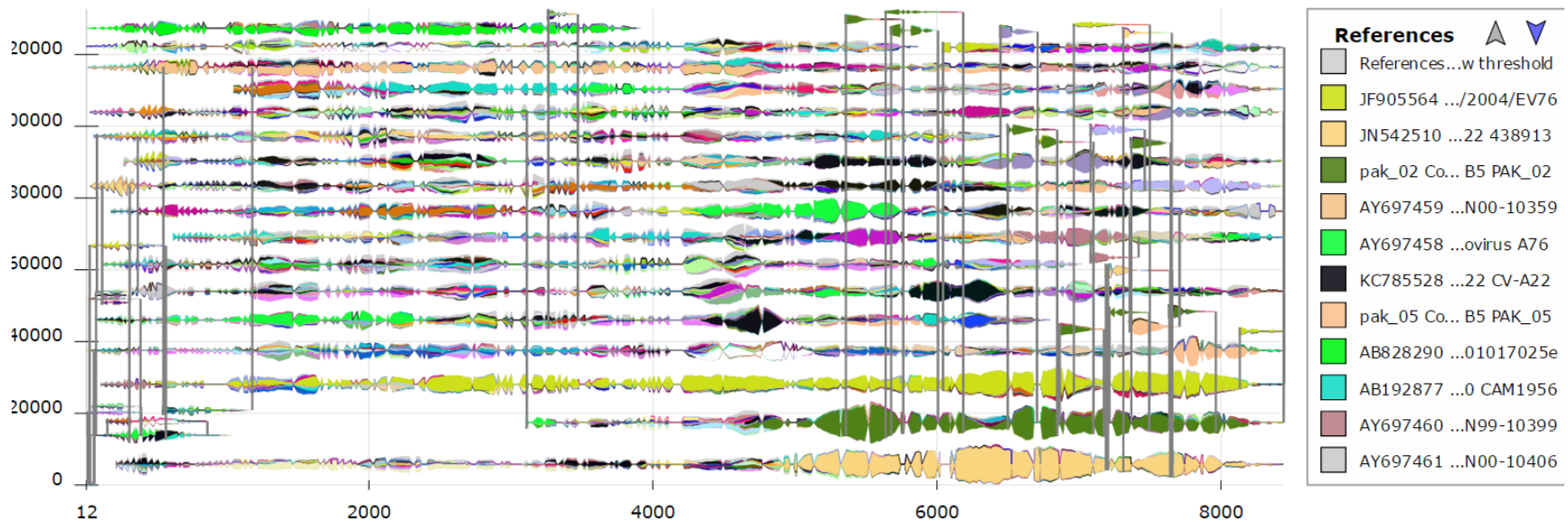


4 Remove host/contaminant



# Environmental sample

- Stool sample aligned against database of 600 enteroviruses.
- Paired end reads of 300bp.
- Only 0.1% of reads aligned to the reference set
- Haplotype reconstruction performed with 1% mutation cutoff
- ORF validation revealed 12 out of 15 biologically meaningful global haplotypes

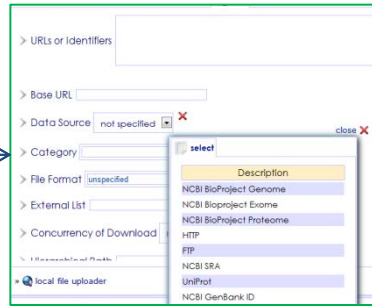
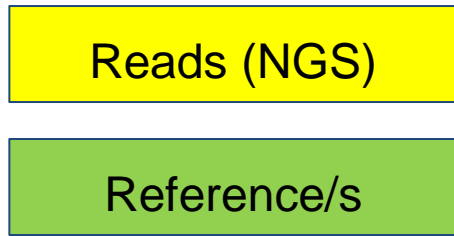




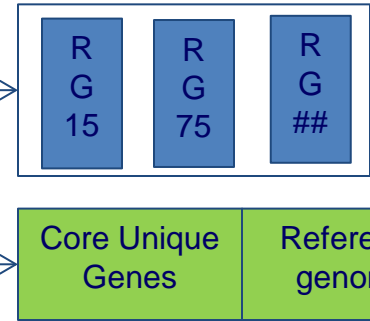


# Detection pipeline

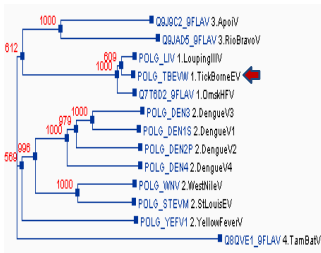
Evolving complex system  
Needs to be highly flexible



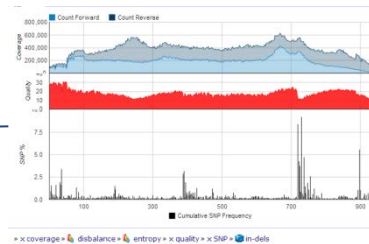
HIVE dmDownloader  
NCBI eutils/Webservices



Reference/Representative  
genome filters



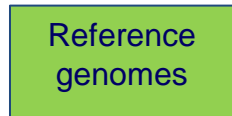
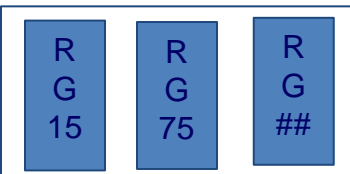
Phyetic placement



SNV profiling

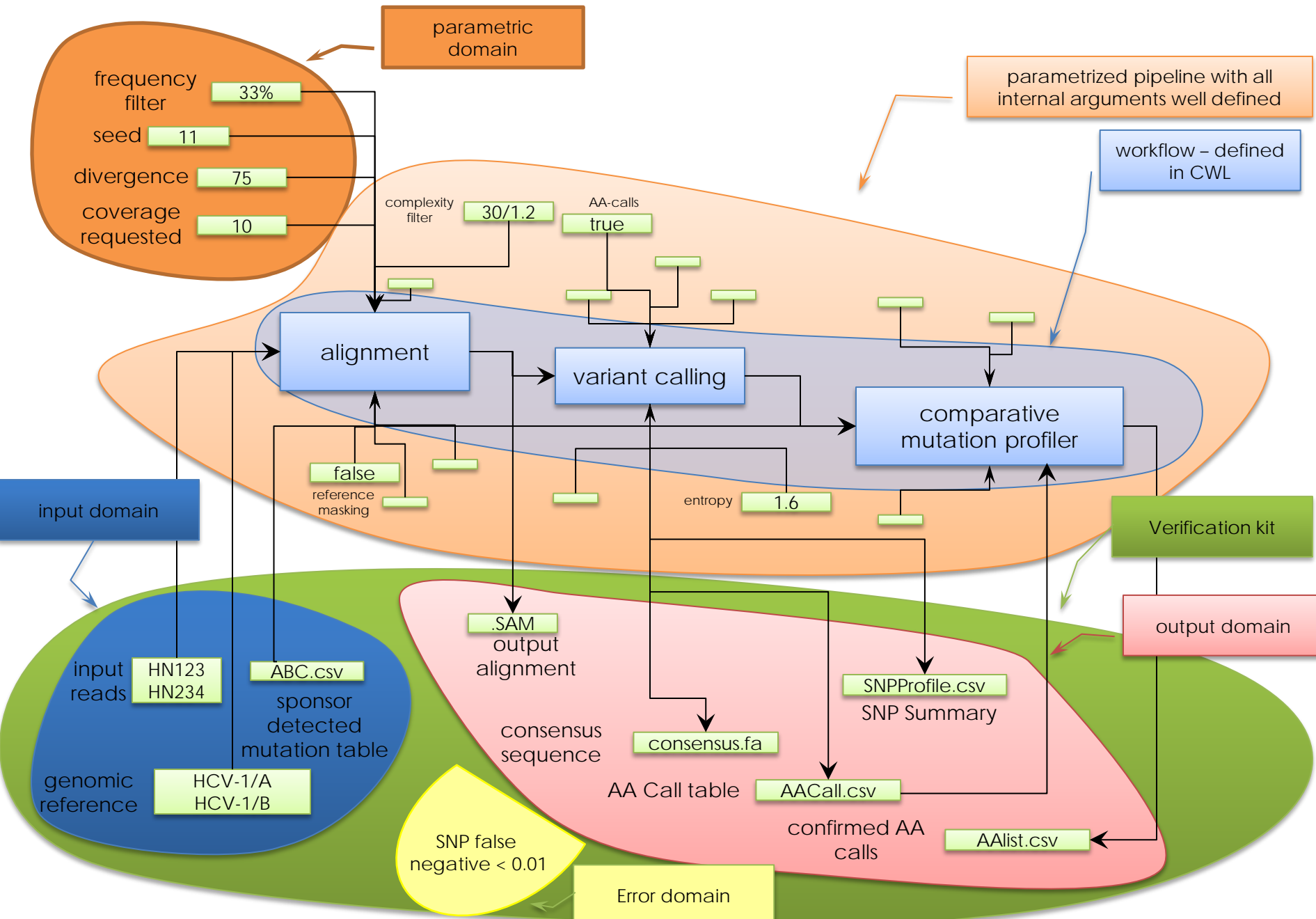
#	Sequence	Repeats	Direction	Start	Alignment
1	hg 127763381 ref NC_005364.2 ...	(*)	1	1	TTTGTAAATATACTATGTTATAGAGATTTTCGAC
965192	@g 127763381 ref NC_005364.2 ...	(*)	1	1	TTTGTAAATATACTATGTTATAGAGATTTTCGAC
1	hg 127763381 ref NC_005364.2 ...	(*)	5	1	TATGTAATCTATGTTATAGAGATTTTCGACGTT
831702	@g 127763381 ref NC_005364.2 ...	(*)	1	1	TATGTAATCTATGTTATAGAGATTTTCGACGTT
1	hg 127763381 ref NC_005364.2 ...	(*)	10	1	TACTATGTTATAGAGATTTTCGACGTTTCGAC
359220	@g 127763381 ref NC_005364.2 ...	(*)	1	1	TACTATGTTATAGAGATTTTCGACGTTTCGAC
1	hg 127763381 ref NC_005364.2 ...	(*)	13	1	CTATGTTATAGAGATTTTCGACGTTTCGACAT

HIVE Hexagon mapping



Cannot be done in isolation  
Requires collaborative effort

Communicating computational analysis to FDA



```

{id": "obj.      ",
"type": "antiviral_resistance_detection",
"name": "HCV1a [taxID:31646] ledipasvir [PubChem:67505836] resistance SNP [SO:0000694] detection",
"version": "1.1",
"digital_signature": "905d7fce3f3ac64c8ea86f058ca71658",
"verification_status": "unreviewed",
"publication_status": "draft",
"usability_domain": ["Identify baseline single nucleotide polymorphisms (SNPs[SO:0000694]),
insertions[SO:0000667], and deletions[SO:0000045] that correlate with reduced ledipasvir[PubChem:67505836]
antiviral drug efficacy in Hepatitis C virus subtype ", "Identify treatment emergent amino acid
substitutions[SO:0000048] that correlate with antiviral drug treatment failure", "Determine whether the
treatment emergent amino acid substitutions[SO:0000048] identified correlate with treatment failure
involving other drugs against the same virus"],
"authors": [{"name": "Eric Donaldson"}, {"orcid": "0000-      -      -      "}],
"description_domain": {
  "xref": [{"SO":      ", "SO":      ", "PubChem":      ", "taxID":      " ", "PMID":      ", "PMID":
  "}],
  "keywords": ["antiviral resistance", "SNP"],
  "pipeline_steps": {"HIVE_hexagon": {...}, "HIVE_heptagon": {...}},

```

**Metadata**  
 (not needed for computation)

parametrized pipeline with all internal arguments well defined

```

"execution_domain": {
  "platform": "hive",
  "pipeline_version": "0.1",
  "driver": "shell",
  "url": "https://hive.biochemistry.gwu.edu/workflows/",

```

workflow - defined in HIVE

```

  "env_parameters": {"hive_vir": "1"},
  "script": "hive://workflows/antiviral_resistance.py",
  "prerequisites": [
    {"name": "HIVE_hexagon", "version": "1.3"}, {"name": "HIVE_heptagon", "version": "1.3"}]

```

```

"parametric_domain": {
  "hexagon_minimum_coverage": "0.15",
  "hexagon_seed": "14",
  "hexagon_minimum_match_len": "66",
  "heptagon_freq_cutoff": "0.10",
  "heptagon_divergence_threshold_percent": "30"},

```

parametric domain

output domain

Verification Kit

input domain

```

"io_domain": {
  "reference_uri": ["http://www.ncbi.nlm.nih.gov/nuccore/CP000139.1"],
  "input_uri_list": ["hive://nuc-read/      ", "hive://nuc-read/      "],
  "output_uri_list": ["hive://data/514769/dnaAccessionBased.csv", "hive://data/514801/SNP.csv"]},
"error_domain": ["false negative discovery < .0 ", "false positive discovery < .01"]}

```

# Biocompute workshop speakers and panelists



# Acknowledgements

**HIVE TEAM MEMBERS (GW+FDA), COLLABORATORS, USERS**

Funding sources for our projects

NIH, FDA, GW, NSF, Pharma

Contact

[mazumder@gwu.edu](mailto:mazumder@gwu.edu)

Special thanks

Konstantinos Karagiannis, Vahan Simonyan, Alin Voskanian,  
Konstantin Chumakov, Viswanath Ragupathy, Naila Gulzar, Krista Smith



HIVE jamboree  
2016