# All The World's A Sequencer: Global Hepatitis Outbreak Surveillance Technology (GHOST) Cloud Platform

## FY15-AMD-81 Project

S Sims[1], DS Campo[1], I Rytsareva[1], Y Zheng[2, 3], A Longmire[1], P Skums[1], Z Dimitrova[1], M Mirabito[4, 5], S Wang[4, 5], R Tracy[4, 5], T Sukalac[4, 5], C Lynberg[5], Y Khudyakov[1].

[1]Molecular Epidemiology and Bioinformatics Laboratory, Division of Viral hepatitis. [2]Eagle Medical Services, LLC, [3]AMD bioinformatics, NCEIZD, [4]Northrop Grumman. [5]NCHHSTP Informatics Office. [6]Research & Development, OD.

## Introduction

- The depth of sequencing afforded by recent Next-generation sequencing (NGS) technologies is capable of detecting large numbers of viral variants carried by infected individuals. This data offer novel prospects for understanding both intra- and inter-host pathogen evolution.
- As well as offering a wealth of epidemiologically relevant information, NGS provides an opportunity to implement a molecular surveillance of infectious diseases for accurate and comprehensive description of disease dynamics, detection of transmissions, monitoring of epidemic progress, and providing informed guidance for planning public health interventions.
- However, the greatest scientific strength of NGS systems is their greatest weakness in practice. Collection and processing of the volume of NGS data presents significant challenges and stress to legacy information technology (IT) systems. To be able to transfer, analyze, and store the torrent of data requires a fundamental rethink of the public health IT infrastructure status quo. We present GHOST, an IT system designed to collect, and analyze large scale data.
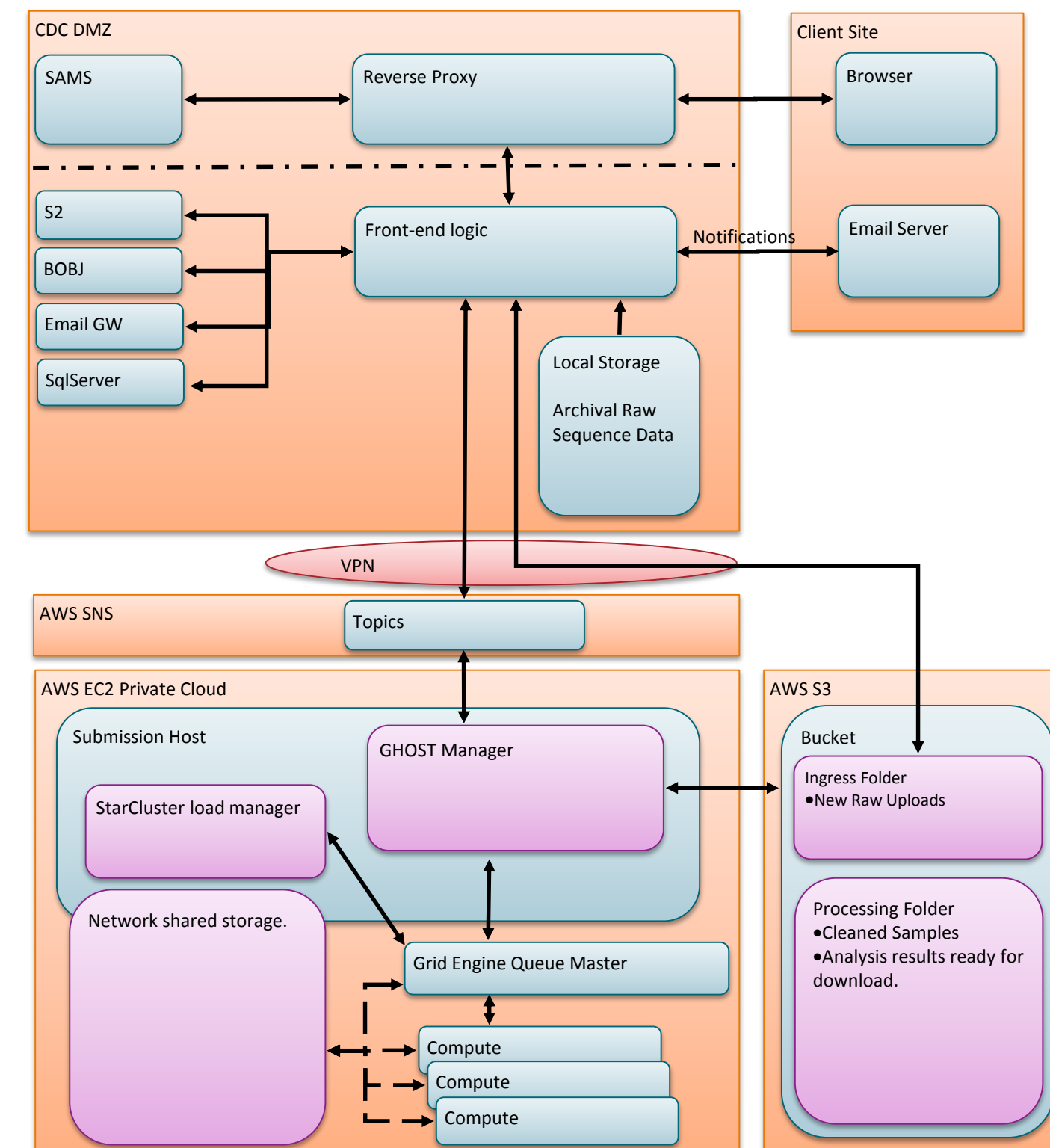
## Challenges

- Build a system that is simple to use and requires minimal training.
- Large volumes of data to be stored.
- Work-load is in large bursts rather than spread out over time.
- Analysis step workload scales approximately with the square of the number of samples analyzed. i.e. analyzing 2x more samples requires 4x more work.
- Information Security considerations significantly limit the allowable system structure.
- High availability, fault tolerance, and disaster recovery.
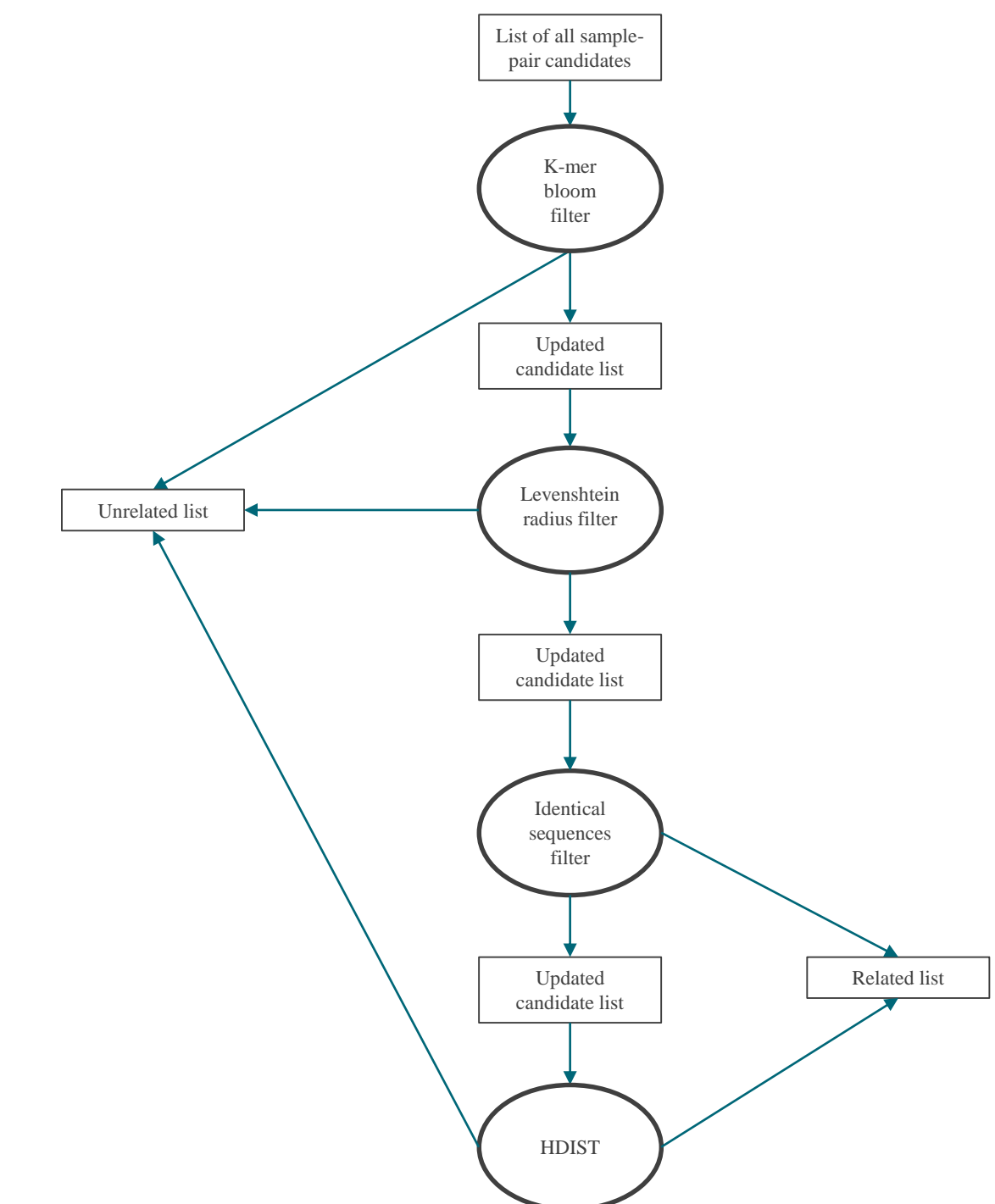- Pilot the new FedRAMP Amazon cloud with the help of ITSO.

## Solutions

- Build a system that is simple to use and requires minimal training
- Hybrid design with user management and presentation at CDC, together with high computational loads in the Amazon cloud.
- Simple web based interface with graph based visualization.
- Storage of raw data at CDC; cloud storage for computational data.
- Use the Amazon cloud's scalability and judiciously designed heuristic algorithms which reduce the workload.
- Standardized laboratory protocol with regional GHOST centers accepting samples.

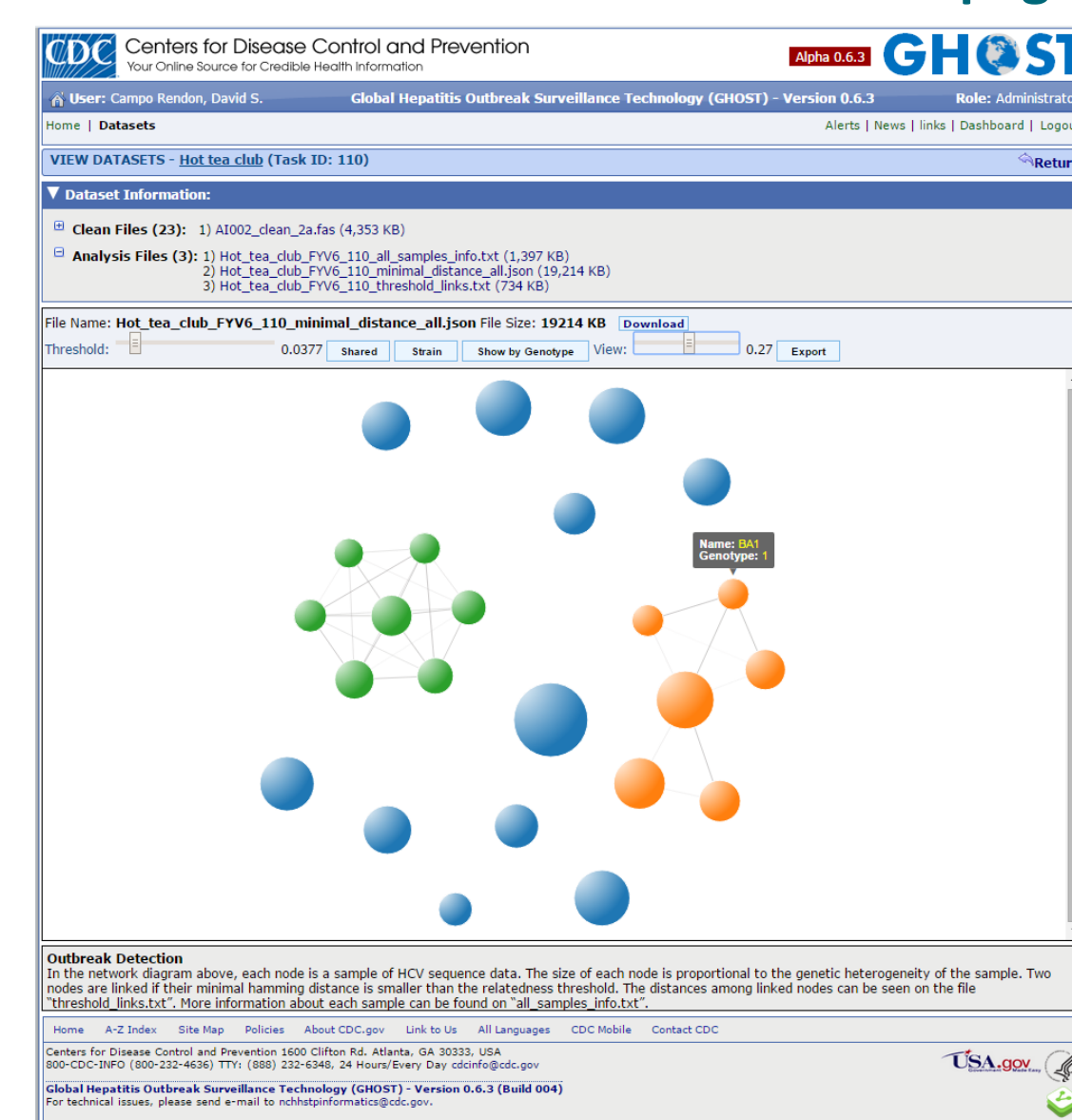### GHOST Cloud Architecture Diagram



## Filtering of candidate pairs

- The naïve, simplest form of our transmission link detection algorithm requires that the set of sequences for each endpoint of a putative transmission link be pooled and aligned. As the most time-consuming step of the pipeline, this represents a significant opportunity for algorithmic optimization.
- As such, we have focused on classifying possible links between samples without resulting to the full distance measurement. We also take advantage of significant pre-computation in a common time/space tradeoff.
- The first filter uses a string metric, the edit distance, to define a sphere in sequence space that encompasses each set. By measuring the distance between the center of a pair of sets, and subtracting the radius of each hyper-sphere, sets that do not overlap are quickly pruned from the list.
- Secondly, for sequence sets to fall below the relatedness threshold they must also share a minimum length of k-mer. Any links that do not share any k-mers are pruned.
- We also check for sets that share an exact sequence, allowing us to add the pair to the set of links with a distance of zero without performing the full distance calculation.
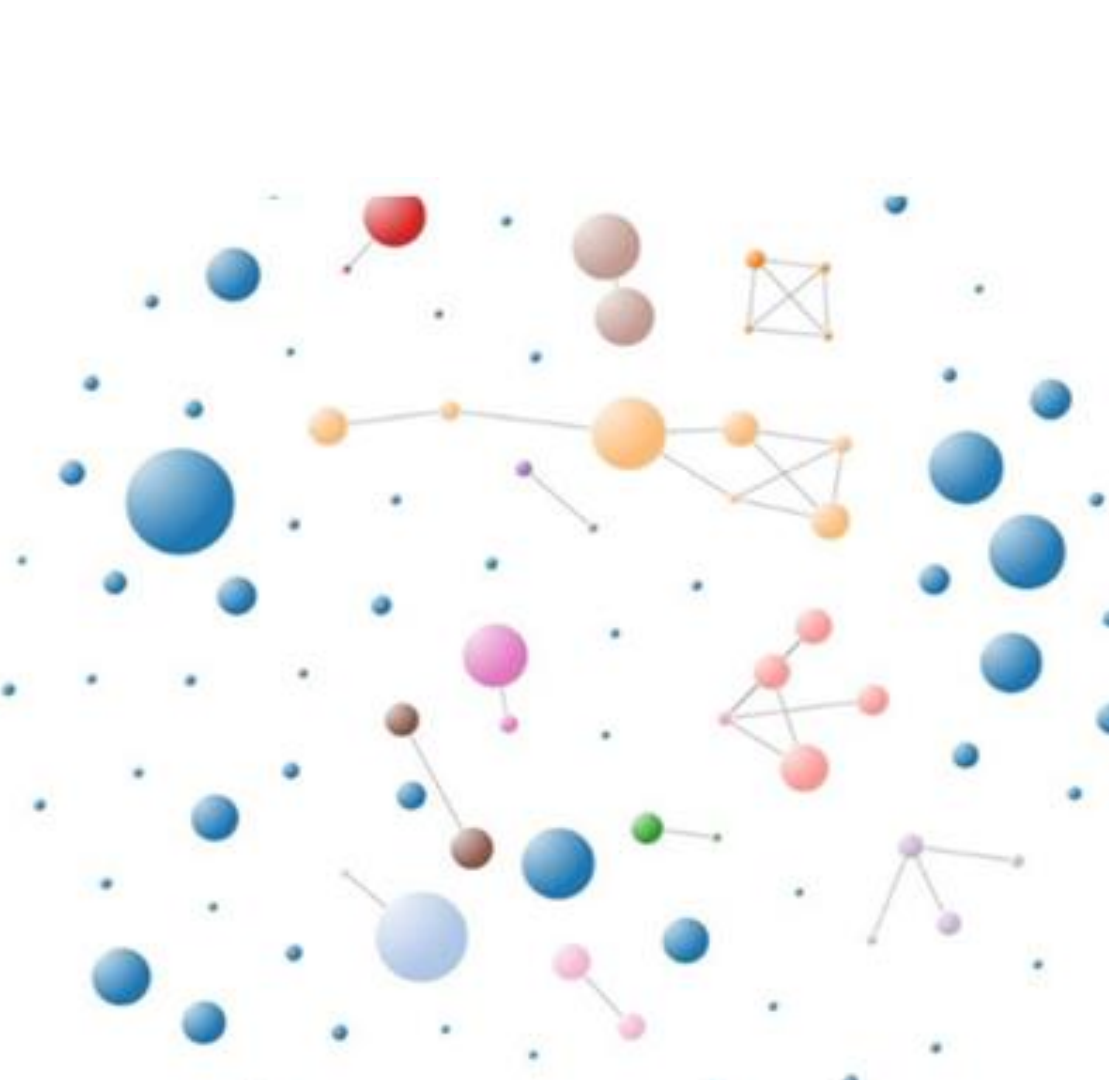


Our fast and efficient three-step filtering strategy removes 91.0% of all pairwise sequence comparisons, accurately establishing which pairs of HCV samples are below the relatedness threshold.

### Screenshot of GHOST's visualization page



- GHOST is available for beta testers at: https://webappx.cdc.gov/GHOST/

### Detail of the Indiana outbreak



- GHOST allowed timely and cost-effective analysis of HCV strains among PWID (n=240) in the recent HIV IN Outbreak.

### New visualization tools



- All the sequences of a linked pair are visualized with a k-step network.
- The sequences of patient 1 are shown in purple; the sequences of patient 2 are shown in red; the sequences shared by both patients are shown in yellow.
- This allows a qualitative assessment of the overlap between the sequences of two linked patients.
- k-step construction and visualization are done under 10 sec.