

Whole Genome, Physics-based Sequence Alignment for Pathogen Signature Design

Jason Gans¹, Wu-Chun Feng², Murray Wolinsky¹

¹Los Alamos National Laboratory

²Virginia Tech



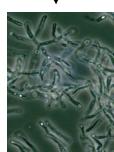
The need for pathogen detection



Specific vs broad spectrum signatures

Specific
Signature

GACTATACA ...



Broad Spectrum
Signature

ATGCCTAAT ...



Detecting multiple targets with a single experiment reduces cost



The World's Greatest Science Protecting America

Images courtesy of
<http://www.microbelibrary.org>



DNA-based detection assays

DNA signature development can exploit the growing number of sequenced genomes

- ~ 400 bacterial genomes available
- ~ 8000 viral genomes available

Mature technologies for high throughput detection

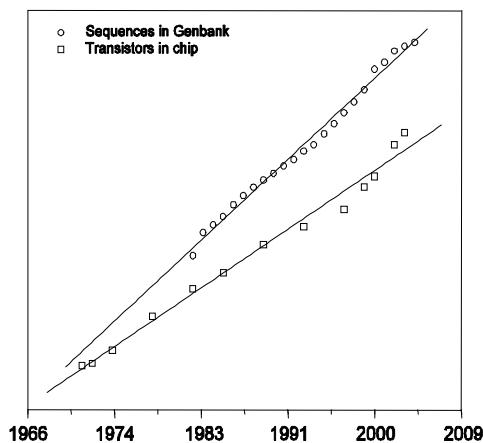
- Gene chip, PCR



The World's Greatest Science Protecting America



Parallel Computation (Moore's Law is not enough)



http://www.intel.com/museum/archives/history_docs/mooreslaw.htm
<http://www.ncbi.nlm.nih.gov/Genbank/genbankstats.html>



Signature Detection Pipeline

1. Partition DNA sequence data
2. Identify sequence fragments unique to target species (i.e. not found in background sequences)
3. Identify unique fragments that are found in multiple target species
4. Compute minimal set of unique fragments required to span target species set
5. Convert unique fragments into biological assays
6. Computationally and experimentally validate assays



The World's Greatest Science Protecting America



How to define sequence similarity/uniqueness?

The answer depends on how we ask the question

The assay format, e.g.
• PCR amplification
• DNA chip probe hybridization
• Single base extension
defines similarity.



How to define sequence similarity/uniqueness?

The answer depends on how we ask the question

The assay format, e.g.
• PCR amplification
• DNA chip probe hybridization
• Single base extension
defines similarity.

A sequence similarity metric that supports common assay formats
is DNA melting temperature, T_m

$T_m^{A,B} > T_m^H \Rightarrow$ Sequences A and B are equivalent

$T_m^{A,B} < T_m^L \Rightarrow$ Sequences A and B are dissimilar



How to define sequence similarity/uniqueness?

$$T_M = \frac{\Delta H^\circ}{\Delta S^\circ + R \ln\left(\frac{C_T}{N}\right)}$$

R = Gas constant
 C_T = Strand concentration
 N = sequence dependant constant

Nearest Neighbor Model

$$\Delta H^\circ = \Delta H_{init} + \Delta H_{sym} + \sum_{\substack{i=\{A,T,G,C\} \\ j=\{A,T,G,C\}}} n_{ij} \Delta H_{ij}$$

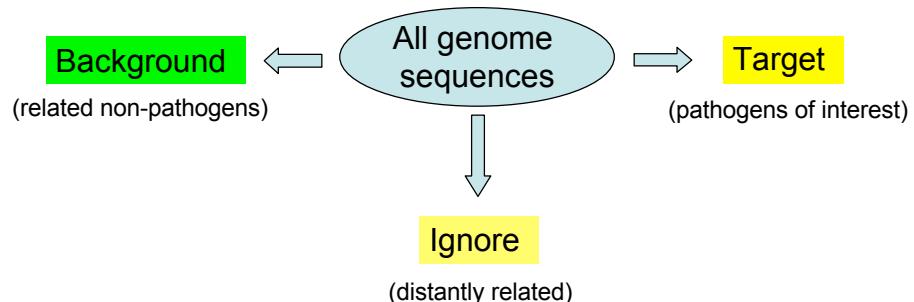
$$\Delta S^\circ = \Delta S_{init} + \Delta S_{sym} + \sum_{\substack{i=\{A,T,G,C\} \\ j=\{A,T,G,C\}}} n_{ij} \Delta S_{ij}$$



The World's Greatest Science Protecting America



Partition Target Sequences

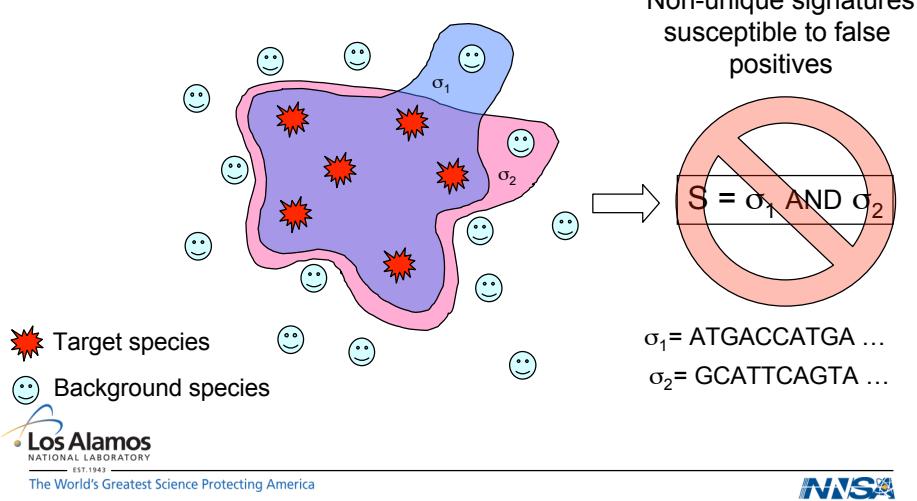


The World's Greatest Science Protecting America

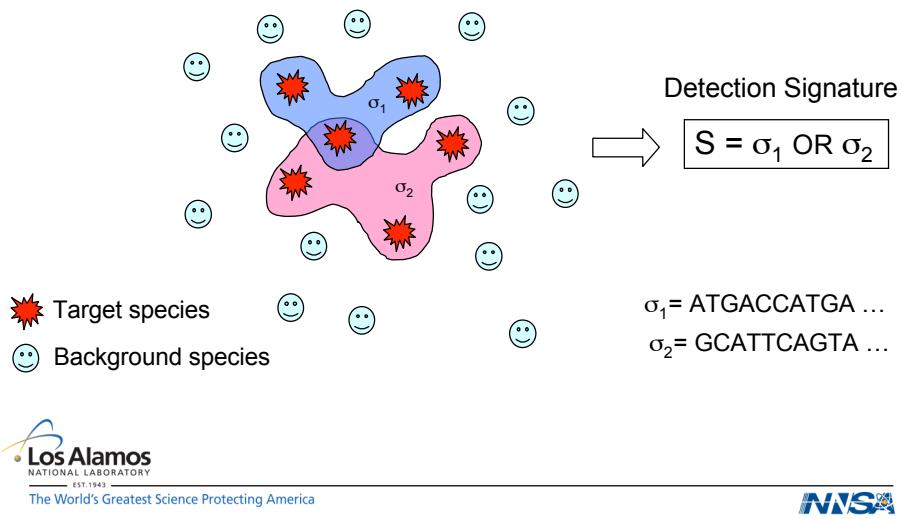


Isolate sequences unique to targets

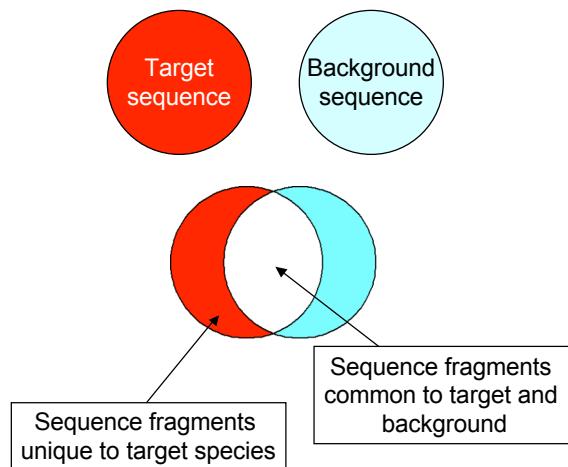
Why search for unique sequence?



Isolate sequences unique to targets



Isolate sequence fragments unique to target species



 Los Alamos
NATIONAL LABORATORY
EST. 1943
The World's Greatest Science Protecting America



Isolate sequences unique to targets: 1st generation

Using mpiBLAST to compute the set difference.

$$U = T - (T \cap B)$$

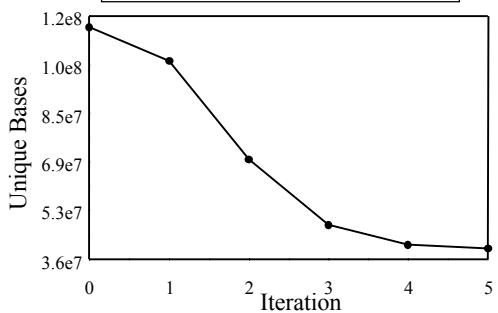
$$U \approx U_f$$

$$U_{i+1} = U_i - (U_i \tilde{\cap} B)$$

$$U_0 = T$$

$$U_f \approx U_3$$

- $\cap \sim$
1. Use mpiBLAST to identify regions of similarity between U_i and B .
 2. If a sequence fragment has a $T_m > 55^\circ\text{C}$ (over a ≤ 21 base window) then it is considered a match.



 Los Alamos
NATIONAL LABORATORY
EST. 1943
The World's Greatest Science Protecting America



Isolate sequences unique to targets: 1st generation

mpiBLAST

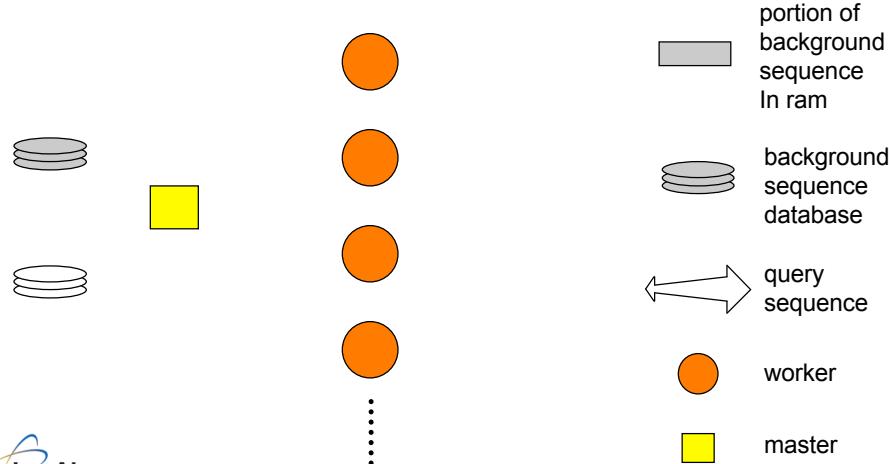
- Pro: Allowed signature design to use all available bacterial genomes
- Con: Required iteration and separate computation of set difference

2nd generation implementation

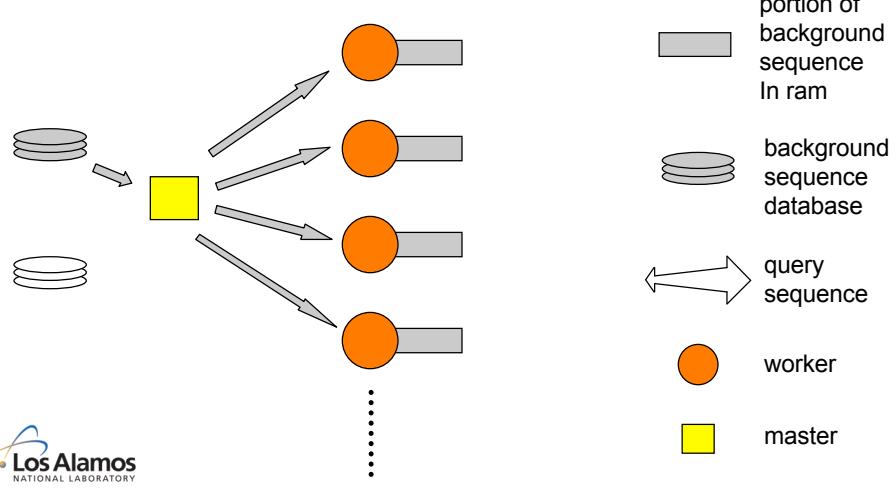
- Use high level structure of mpiBLAST (i.e. master/worker, database segmentation, query segmentation, MPI/C++)
- Directly compute set difference: *in silico* subtractive hybridization



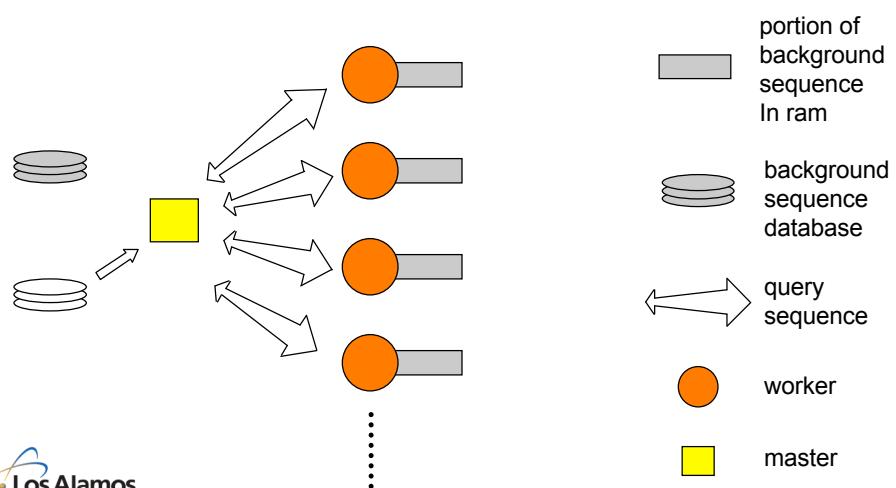
In silico subtractive hybridization



In silico subtractive hybridization



In silico subtractive hybridization



Target fragment alignment

In silico subtractive hybridization produces fragments unique to every target species

species A → {ACTGGATCGATC, GGCTGGATTCTAGG, TAGGCTTAGGCTTA, ATTGGGCCAGATAG, ...}
species B → {GCTTCTAGACAAAC, ATGGCGATTAGCCA, GCTAAGCCTAGCTA, TTTGACTAGATCAC, ...}
species C → {GGCTAGGCCCA, GGCTCTAGGATATAAC, CGCTCTAGGCTTATAT, CGGATTCTGGCTTAG, ...}

⋮

Need to associate every unique target fragment with one or more target genomes

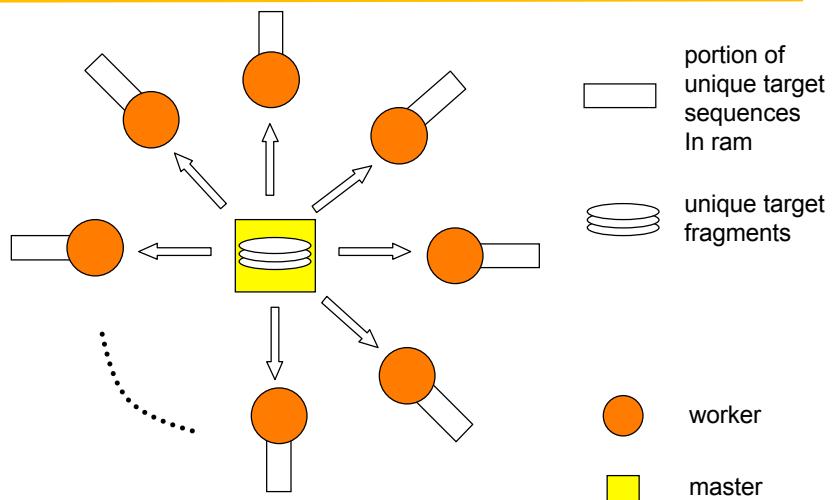
CGTAGGCTTAGGATAT	AGGCTATAGCGGA	CAATAGGCTTATAGGATTA	...
↓	↓	↓	
{species A, species B, species C}	{species B, species D, species E, species F}	{species A, species F}	



The World's Greatest Science Protecting America



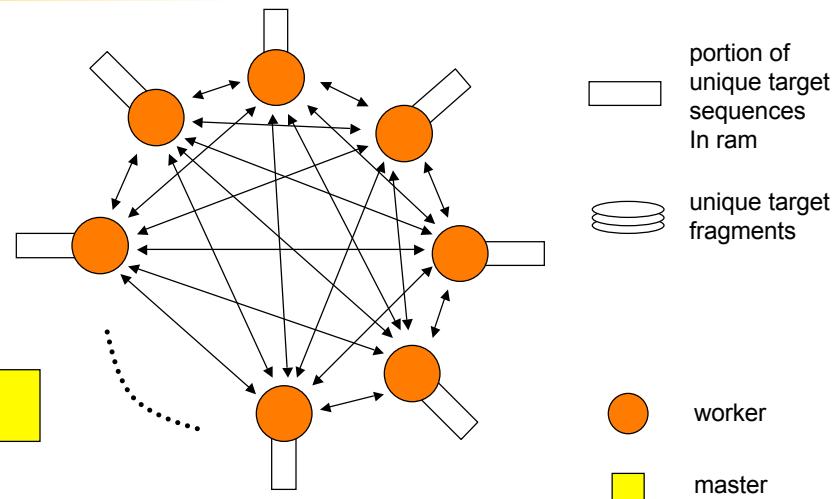
Target fragment alignment



The World's Greatest Science Protecting America



Target fragment alignment

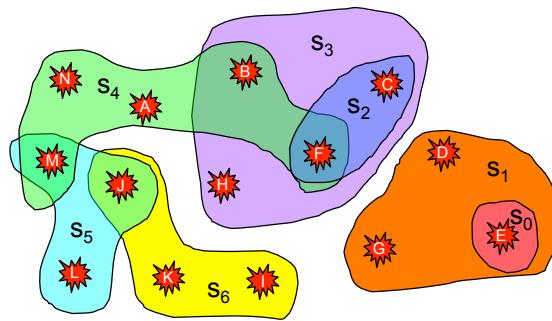


Los Alamos
NATIONAL LABORATORY

EST. 1943 The World's Greatest Science Protecting America



Find the minimal set of signatures that will detect all target species



S_i is a set of species that contain a common sequence fragment

Los Alamos
NATIONAL LABORATORY

EST. 1943 The World's Greatest Science Protecting America



Find the minimal set of sequence fragments that will detect all target species

Given:

$$\begin{aligned} T &= \{A, B, C, D, E, F, G, \dots\} \\ S_0 &= \{E\} \\ S_1 &= \{D, E, G\} \\ S_2 &= \{B, C, F, H\} \\ &\vdots \end{aligned}$$

Our task is to find $G = \{S_\alpha, S_\beta, S_\delta, \dots\}$

Where $T = \bigcup_i G_i$
and
 $|G|$ is a global minimum

The set coverage problem is NP hard!



Find the minimal set of sequence fragments that will detect all target taxa

Approximate solution

Solve for G using Metropolis Monte-Carlo driven simulated annealing

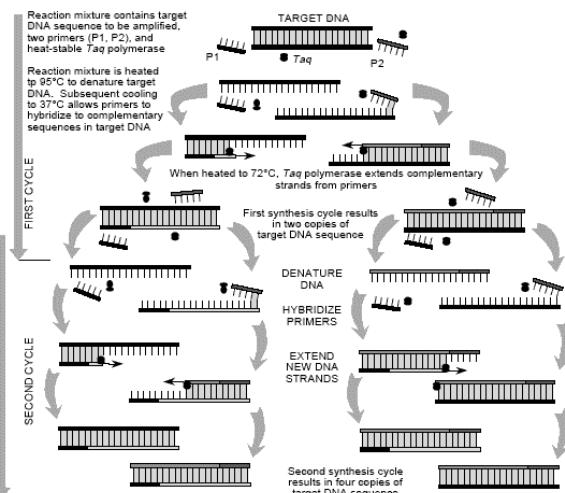
Trial move generation: $G_i \rightarrow G'_{i+1}$
Add random ($[0, n]$) number of sets
and
Delete random ($[0, n]$) number of sets

Acceptance criteria: $G_{i+1} =$
 G'_{i+1}
 $|G'_{i+1}| \leq |G_i|$
or
 $R \leq \exp(-|G'_{i+1}| - |G_i|)/T$
 $R \in [0, 1]$



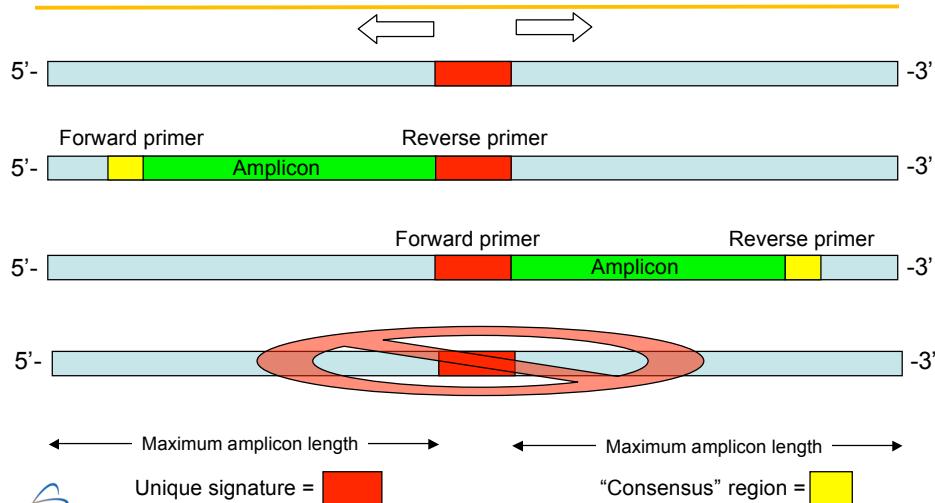
Filter signatures with experimental constraints: PCR Primers

DNA Amplification Using PCR



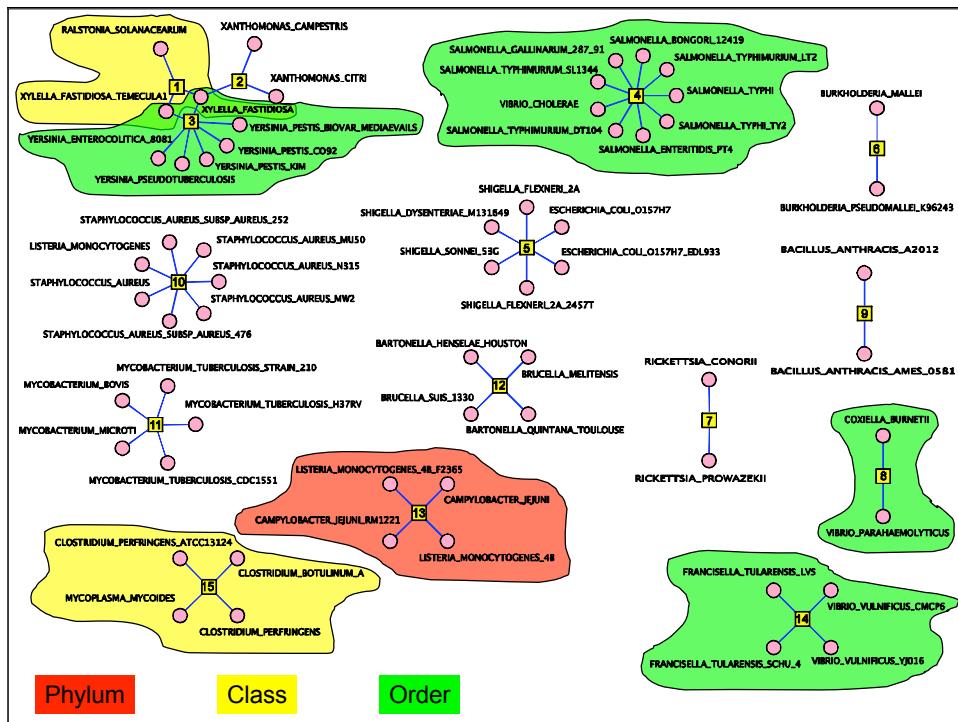
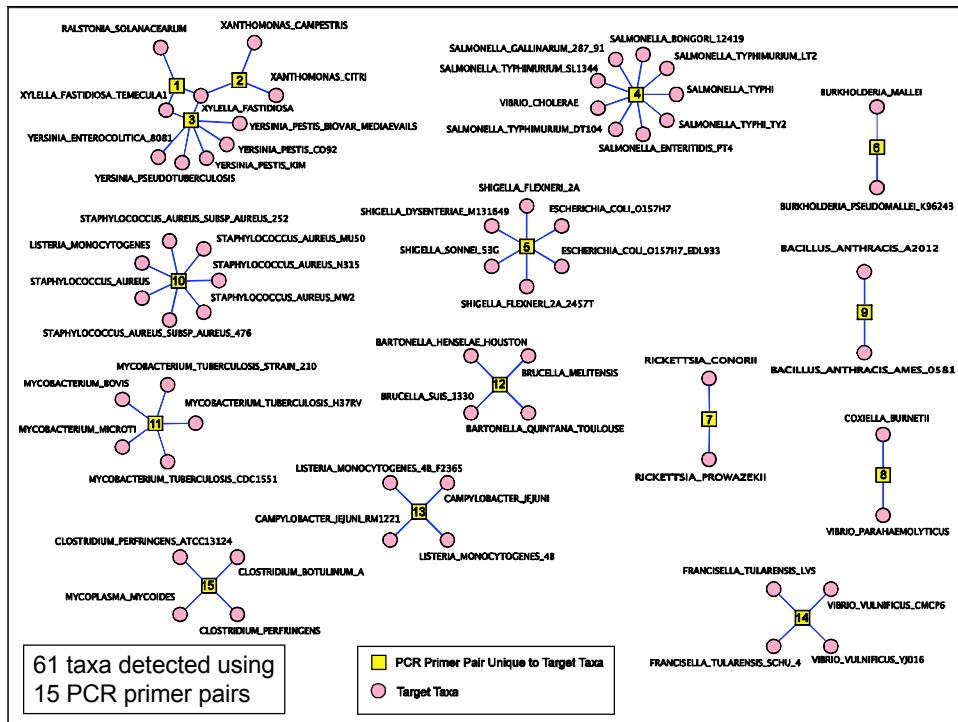
Source: DNA Science, From DOE Human Genome pamphlet, *Primer of Molecular Genetics*

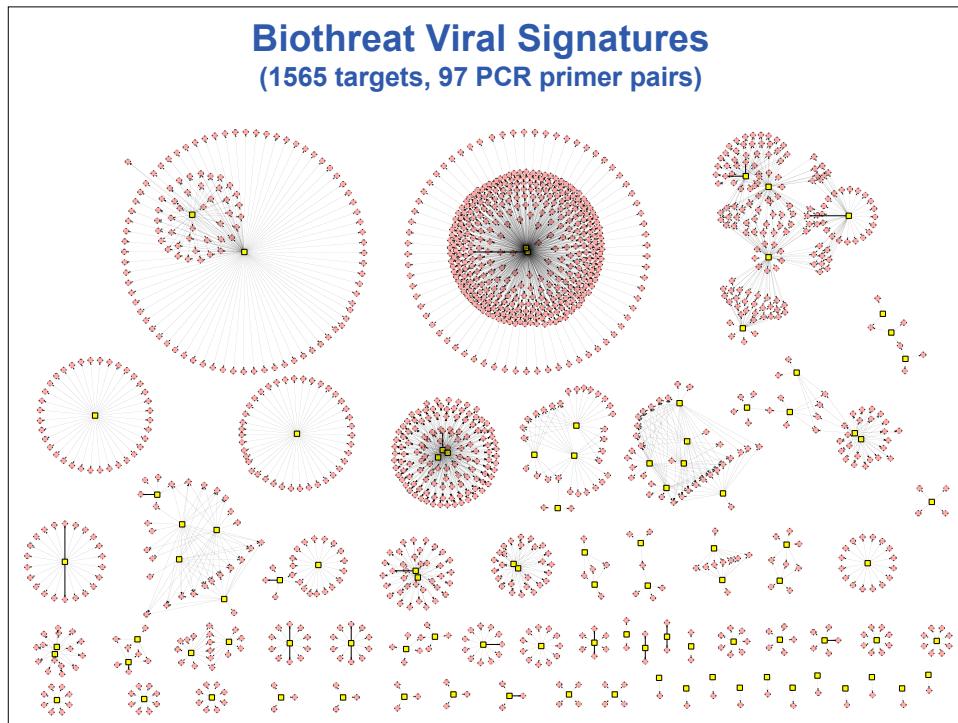
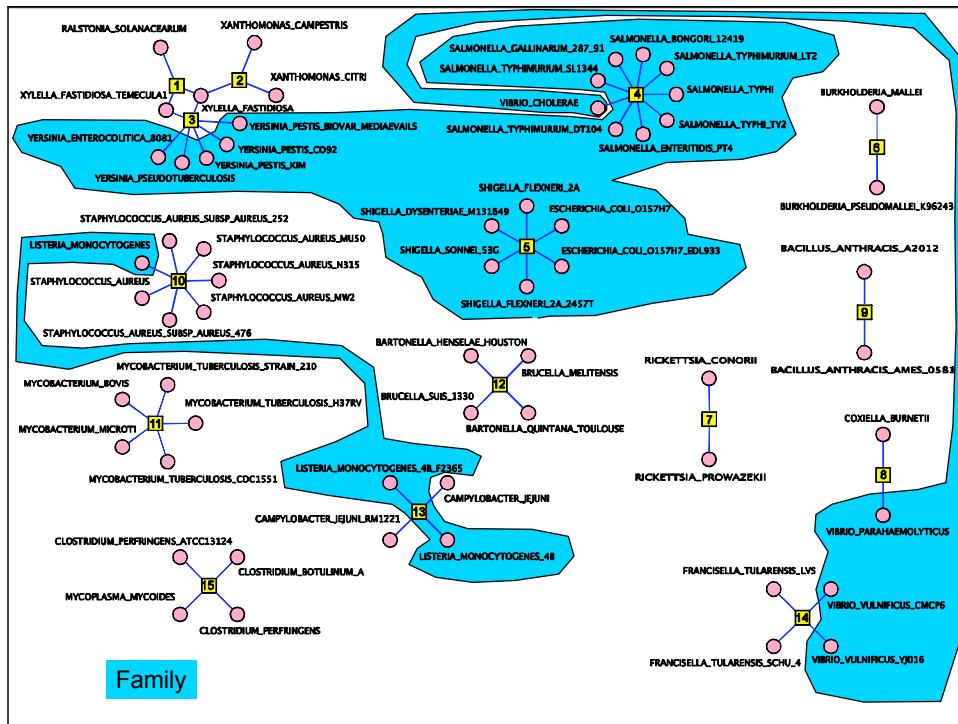
Filter signatures with experimental constraints: PCR Primers

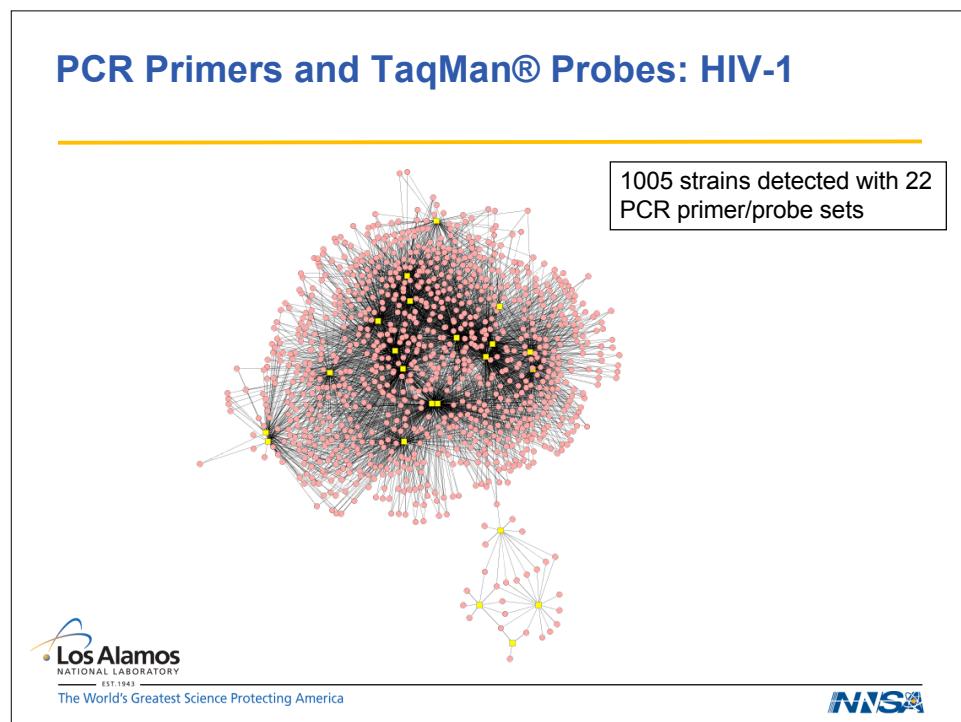
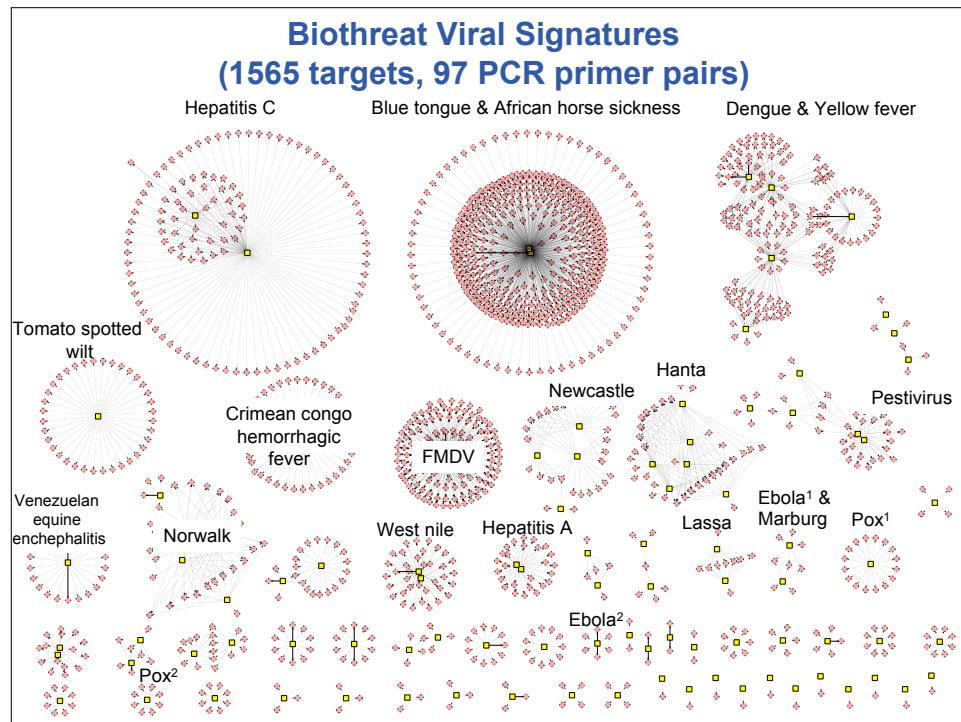


The World's Greatest Science Protecting America

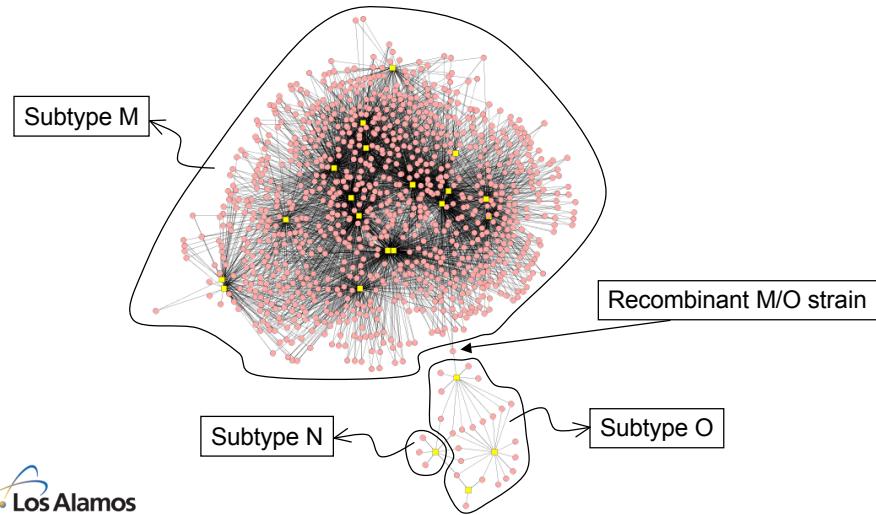








PCR Primers and TaqMan® Probes: HIV-1



Acknowledgements

Informatics Team

Byron Delaney
Norman Doggett
Rob Leach
Jian Song
Chris Stubben
Murray Wolinsky*

Experimental Validation

John Dunbar*
Lance Green
Scott White

mpiBLAST

Lucas Carey
Aaron Darling
Wu-Chun Feng

Funding

DHS
TSWG



The World's Greatest Science Protecting America

