# Fighting the Tuberculosis Pandemic using Machine Learning

Kristin P. Bennett

Rensselaer Polytechnic Institute

TB-Insight Team
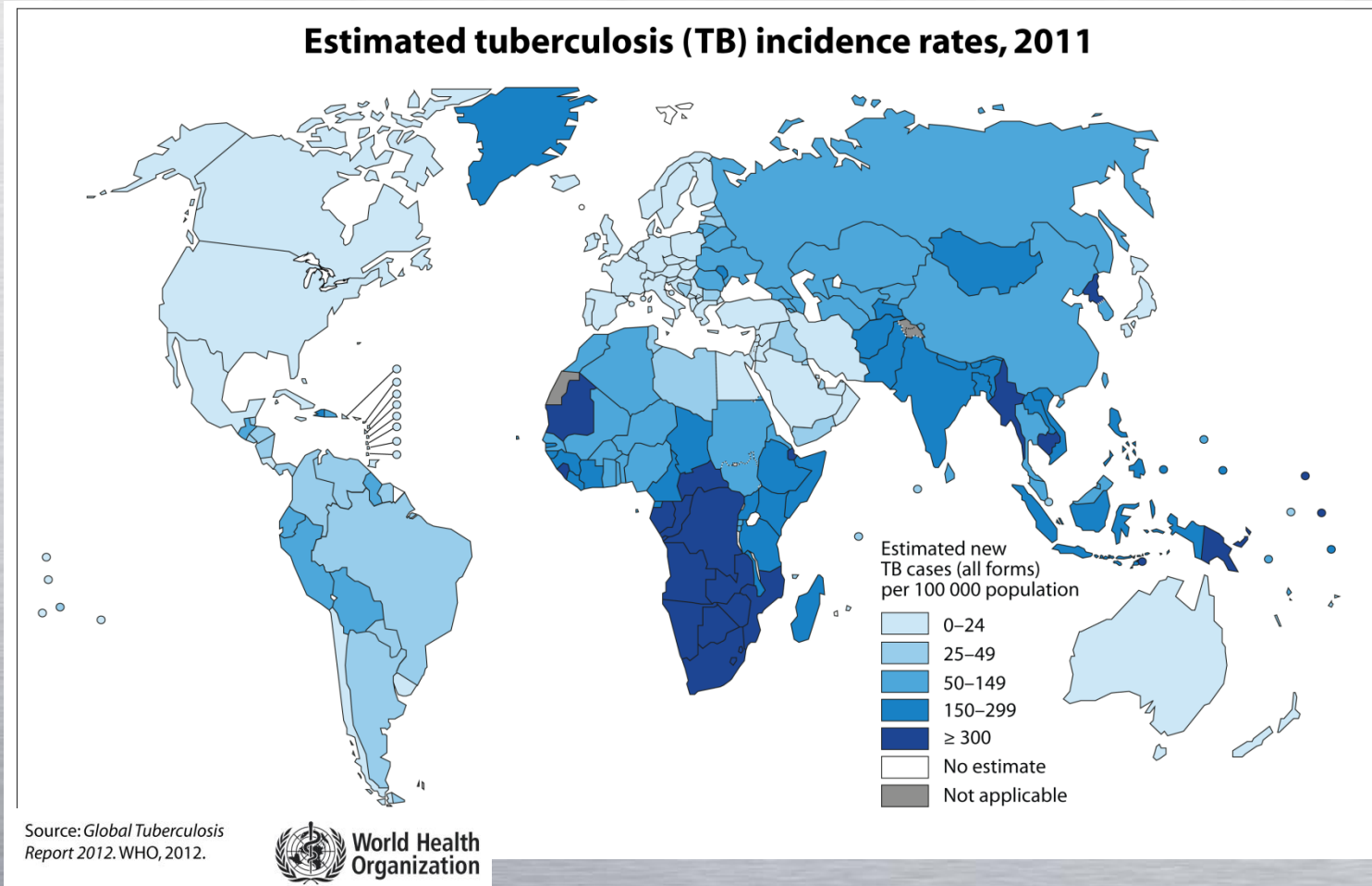
**TB** insight

Tuberculosis Tracking and Control

# 1/3 of World Latently Infected with TB



Estimated tuberculosis (TB) incidence rates, 2011

Estimated new TB cases (all forms) per 100 000 population
- 0–24
- 25–49
- 50–149
- 150–299
- ≥ 300
- No estimate
- Not applicable

Source: *Global Tuberculosis Report 2012.* WHO, 2012.

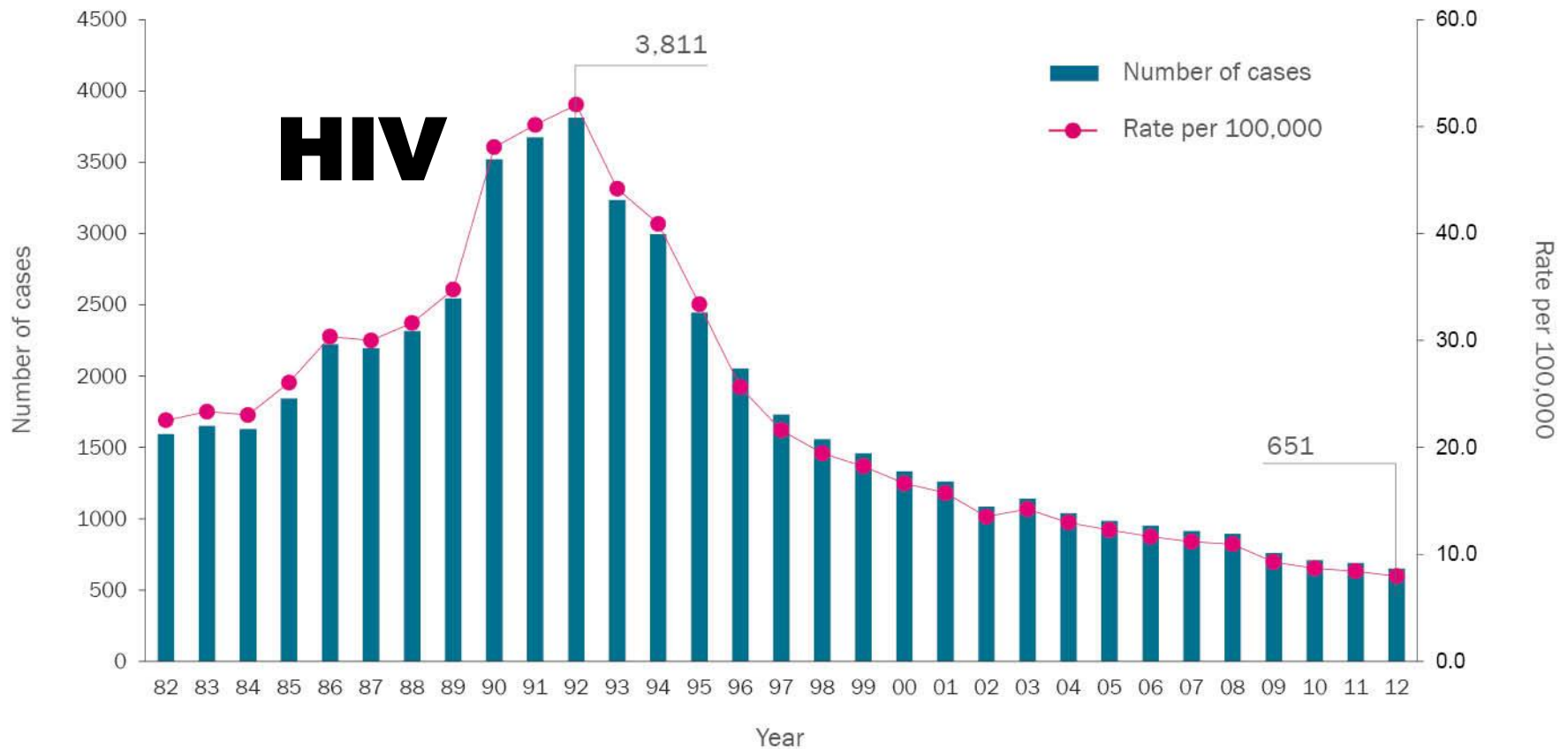World Health Organization

# 2.4 Million Deaths per Year

# Tuberculosis cases and rates, New York City, 1982-2012



1. Rates are based on official Census data.

Number of TB Cases in U.S.-born vs. Foreign-born Persons United States, 1993–2011*

*Updated as of June 25, 2012.

# Drug Resistance Threat



Countries that had notified at least one case of XDR-TB by the end of 2011

At least one case reported
No cases reported
Not applicable

Source: *Global Tuberculosis Report 2012*. WHO, 2012.

World Health Organization

- Susceptible
- Drug Resistant
- MDR-TB - Multi-Drug Resistant
- XDR-TB - Extremely-Drug Resistant
- TDR-TB?- Totally Drug Resistant
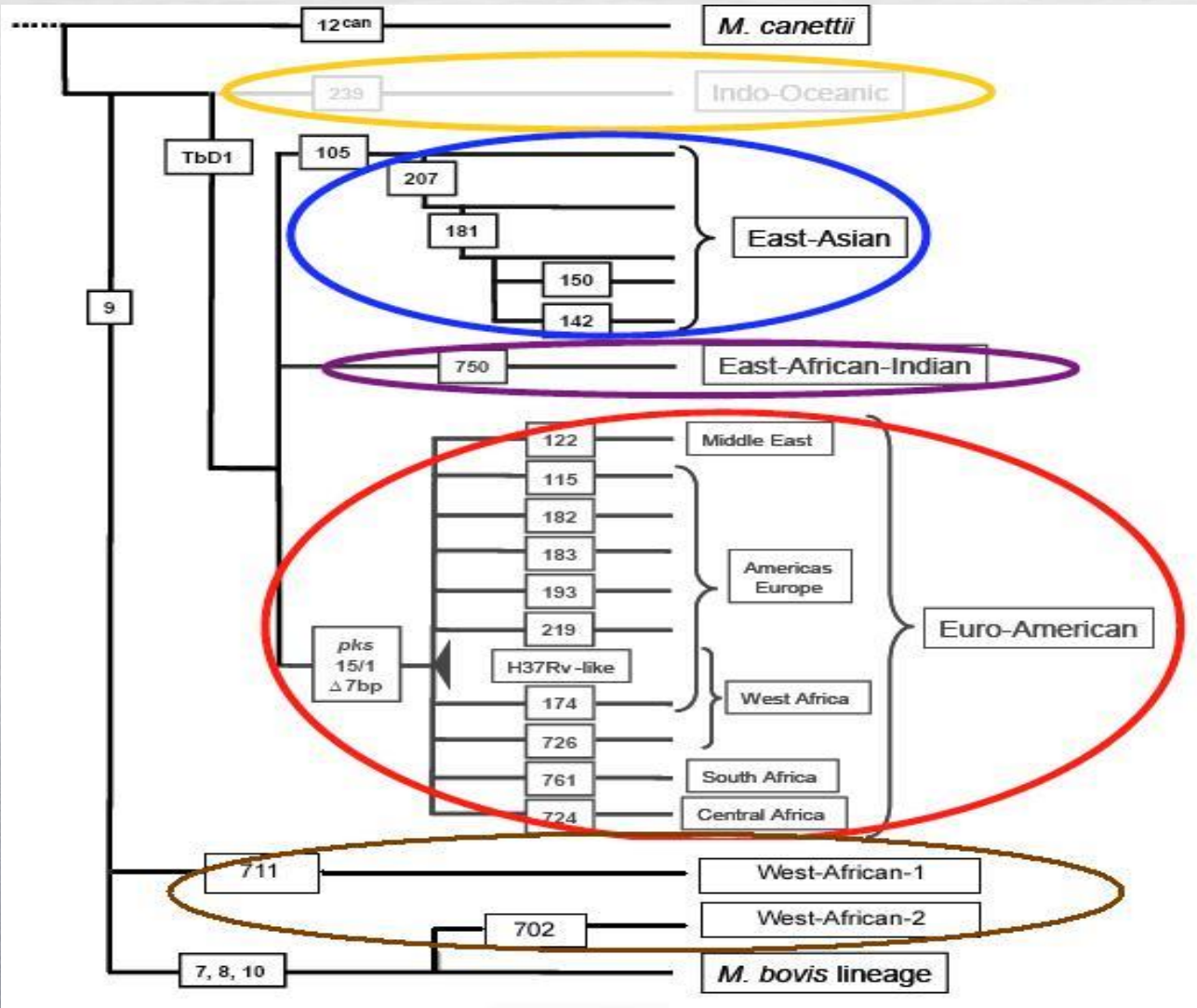
# Modern TB Control



TB Controller: Find source(s) of infection in order to identify people who need treatment and stop future transmission.

Tools:  **Contact Investigation**

**DNA Fingerprints of TB bacteria**

# Two or More DNA fingerprints gathered for every TB Patient in USA



- Spoligotype

1111111111111110000000000000000000000000000

- MIRU

  22531b153321

- RFLP

  BW90

# Major Phylogeographic Lineages of the MTBC



Determined by RD's

Predictable by
 Spoligotype

[Gagneux S et al. PNAS 2006]

# Spoligotype Genetic Diversity 37K Patients in US – 2004-2008

Spoligoforests labeled by CDC Expert Rules



LINEAGES

- East Asian (Beijing)
- East-African Indian (CAS)
- Euro-American
- Indo-Oceanic
- M. Caprae
- M. africanum
- M. bovis
- Unidentified

*Classification Models*

**TB-Lineage**
- Rule-based model: 2012.
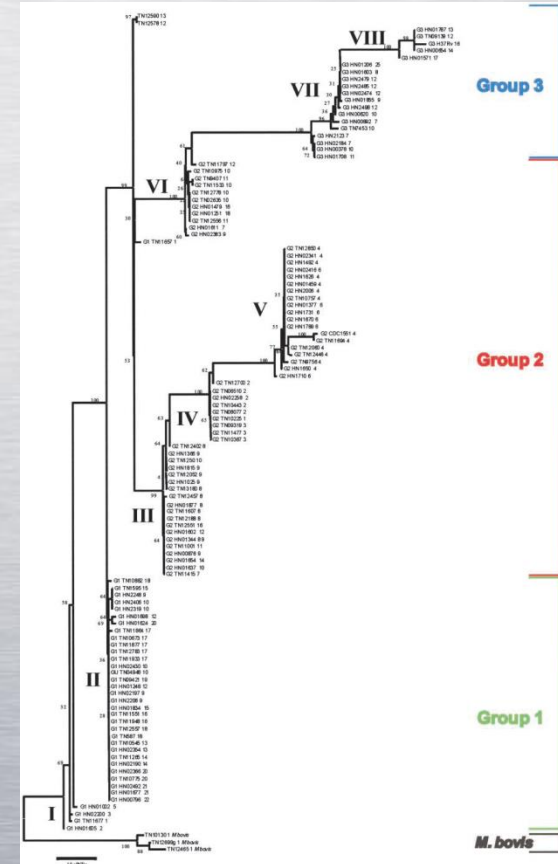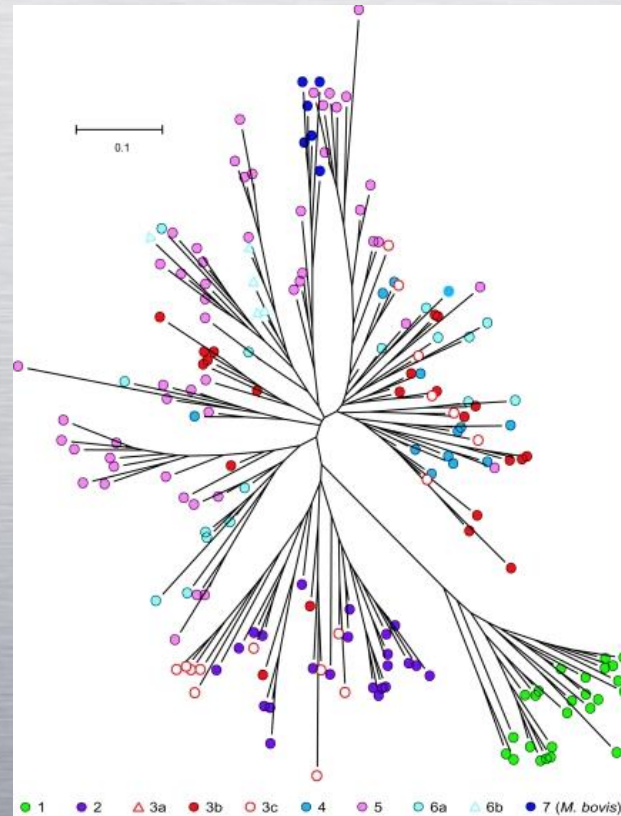- Bayesian Network: 2010.

# Sublineages-varying opinions

62 sublineages based on spoligotype signatures
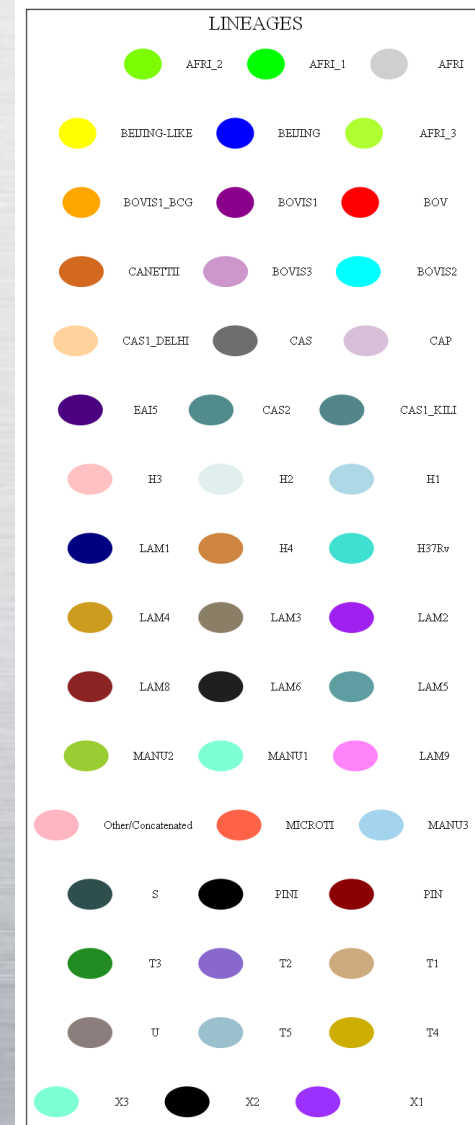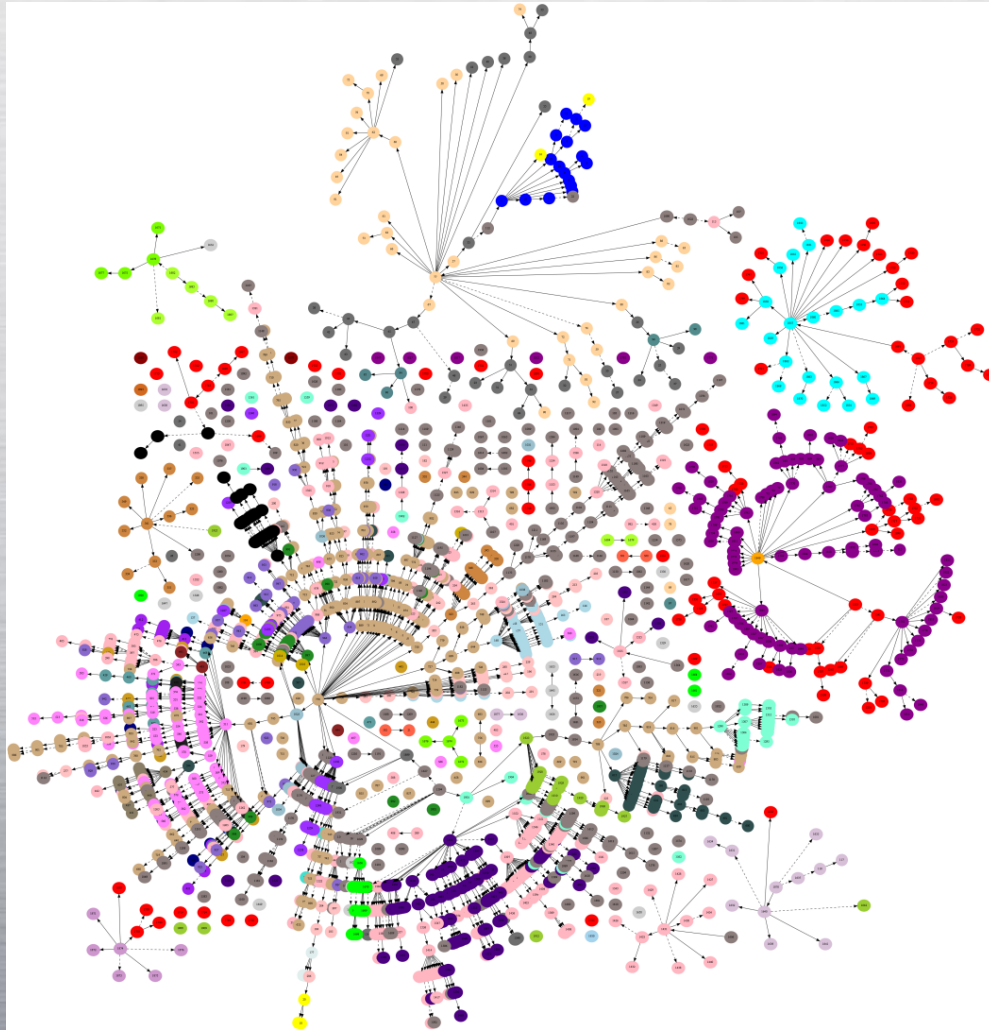[Brudey et al BMC Microbiol. 2006]

10 groups based on 212 single nucleotide polymorphism (SNP) markers
[Filliol J. Bacteriol 2006]

8 groups based on 230 sSNPs
[ Gutacker Genetics 2002]

# What's the story with sublineages?



Spoligoforest – drawn using GraphViz Twopi
Data labels –SpolDB4 [Brudey et al 2006]

# SPOTCLUST: Hidden-Parent Bayesian Network for Spoligotypes (Vitol et al, 2006)

Standard resource for **tuberculosis sublineage identification** used in over 96 publications.

Strains of *Mycobacterium tuberculosis* from Western Maharashtra, India, Exhibit a High Degree of Diversity and Strain-Specific Associations with Drug Resistance, Cavitary Disease, and Treatment Failure

Characterization of multiple and extensively drug resistant *Mycobacterium tuberculosis* isolates with different ofloxacin-resistance levels

Distinct clinical and epidemiological features of tuberculosis in New York City caused by the RD$^{Rio}$ *Mycobacterium tuberculosis* sublineage

*Mycobacterium bovis* infection in livestock workers in Ibadan, Nigeria: evidence of occupational exposure

High prevalence of subclinical tuberculosis in HIV-1-infected persons without advanced immunodeficiency: implications for TB screening

Whole cell & culture filtrate proteins from prevalent genotypes of *Mycobacterium tuberculosis* provoke better antibody & T cell response than laboratory strain H$_{37}$Rv

# Spoligotyping

- Direct repeats (DR) separated by variable spacers



- Contiguous on chromosome, order well conserved

- Forty three spacers used

- Presence of a spacer is detected: 1- present (▪), 0 - absent (□)

| Strain | Binary description of spoligotypes |
|---|---|
| *M. tuberculosis* Beijing |  |
| *M. bovis* |  |
| |  |
| Indo-Oceanic |  |

# Rule-Based Method: (Sub)Lineage Visual Rules

*M. africanum*

*Haarlem 2*

- Determined by human experts.
- Ill-defined.
- Incomplete.
- Frequently ambiguous.
- No precedence.
- May or may not correspond to actual evolutionary groups.

# First Try – Naïve Bayes
## *M. tuberculosis* Haarlem2 Family

- Prototype = probabilities



S - spacer present

N - spacer absent

- Bernoulli Mixture Model



- Biology is wrong!!!

# Unsupervised Hidden Parent Multivariate Bernoulli Mixture Model

Model child spacer $S$, given unobserved parent spacer $H$

- With very high probability child matches parent

- Children are much more likely to lose spacer than gain

  - $P(S=1|H=1) = 0.99$, $P(S=0|H=1) = 0.01$

  - $P(S=0|H=0) = 1\text{-}1e\text{-}7$ $P(S=1|H=0) = 1e\text{-}7$

$C$

$H_i$

$S_i$

43

## SPOTCLUST (2006)

- Unsupervised except

*34 SPOLDB2 Visual Rules* used to initialize clusters

- Trained using 535 spoligotypes

- Number of sublineages (36) picked by MCCV



| Family | Total (n) | Stability | Description |
|---|---|---|---|
| EAI3 | 112 | 0.96 | |
| LAM3 | 138 | 0.95 | |
| Haarlem1 | 236 | 0.94 | |
| Beijing | 985 | 0.92 | |
| X2 | 364 | 0.88 | |
| CAS | 283 | 0.87 | |
| LAM4 | 146 | 0.84 | |
| T4 | 67 | 0.83 | |
| X3 | 469 | 0.81 | |
| EAI5 | 171 | 0.80 | |
| M. bovis BCG | 109 | 0.78 | |
| Family34 | 60 | 0.76 | |
| Family33 | 119 | 0.75 | |
| EAI2 | 153 | 0.73 | |
| M. africanum | 60 | 0.71 | |
| Family36 | 46 | 0.68 | |
| T3 | 56 | 0.67 | |
| LAM9 | 534 | 0.67 | |
| LAM8 | 58 | 0.63 | |
| Family35 | 31 | 0.59 | |
| Haarlem2 | 74 | 0.58 | |
| T1 | 1084 | 0.58 | |
| LAM10 | 73 | 0.57 | |
| Haarlem3 | 603 | 0.50 | |
| H37Rv | 122 | 0.49 | |
| T2 | 57 | 0.45 | |
| X1 | 395 | 0.41 | |
| LAM7 | 55 | 0.40 | |
| EAI1 | 22 | 0.40 | |
| EAI4 | 70 | 0.34 | |
| S | 134 | 0.27 | |
| LAM1 | 142 | 0.24 | |
| LAM2 | 94 | 0.16 | |
| LAM5 | 43 | 0.15 | |
| M. microti | 3 | 0.08 | |
| LAM6 | 2 | 0.02 | |

# New Challenges

- More data: 119,684 isolates from  US CDC, NYDOH, NY State DOH, and Institut Pasteur de Guadeloupe, MIRUVNTRPlus
- More types of DNA fingerprints – Spoligotypes   and MIRU
- More proposed sublineages (70?)
- Putative labels from multiple experts
- Missing Data

# Who's Right?

| Ctop | Cmid | Csub |
|------|------|------|
| Indo-Oceanic | Bangladesh | EAI6-BGD1 EAI7-BGD2 |
| | India | EAI3-IND |
| | Manila | EAI2-Manila |
| | Mexico | EAI-Mexico |
| | Nonthaburi | EAI2-nonthaburi |
| | Vietnam | EAI4-VNM |
| | Unknown Mid-level | EAI1-SOM EAI2 EAI8-MDG |
| *Mycobacterium africanum* | West African 1 | AFRI_2 AFRI_3 |
| | West African 2 | AFRI_1 |
| *Mycobacterium bovis* | *Mycobacterium bovis* | BOV_1 BOV_2 BOV_3 |
| *Mycobacterium canettii* | *Mycobacterium canettii* | Canettii |
| *Mycobacterium caprae* | *Mycobacterium caprae* | Caprae |
| *Mycobacterium microti* | *Mycobacterium microti* | Microti |
| *Mycobacterium mungi* | *Mycobacterium mungi* | M. mungi |
| *Mycobacterium pinnipedii* | *Mycobacterium pinnipedii* | Pini1 Pini2 |

# Semi-supervised Hierarchical Lineage Model

A: 12 Major Lineages

B: 22 Mid-level Lineages

C: 70 + 9 Sub-lineages

- Estimated 92% Cross-validated Accuracy

# Major Lineage Results



- Balanced Classification Rate about 98%
- No changes in major lineages
- MANU Modern?

# Sub-level Results



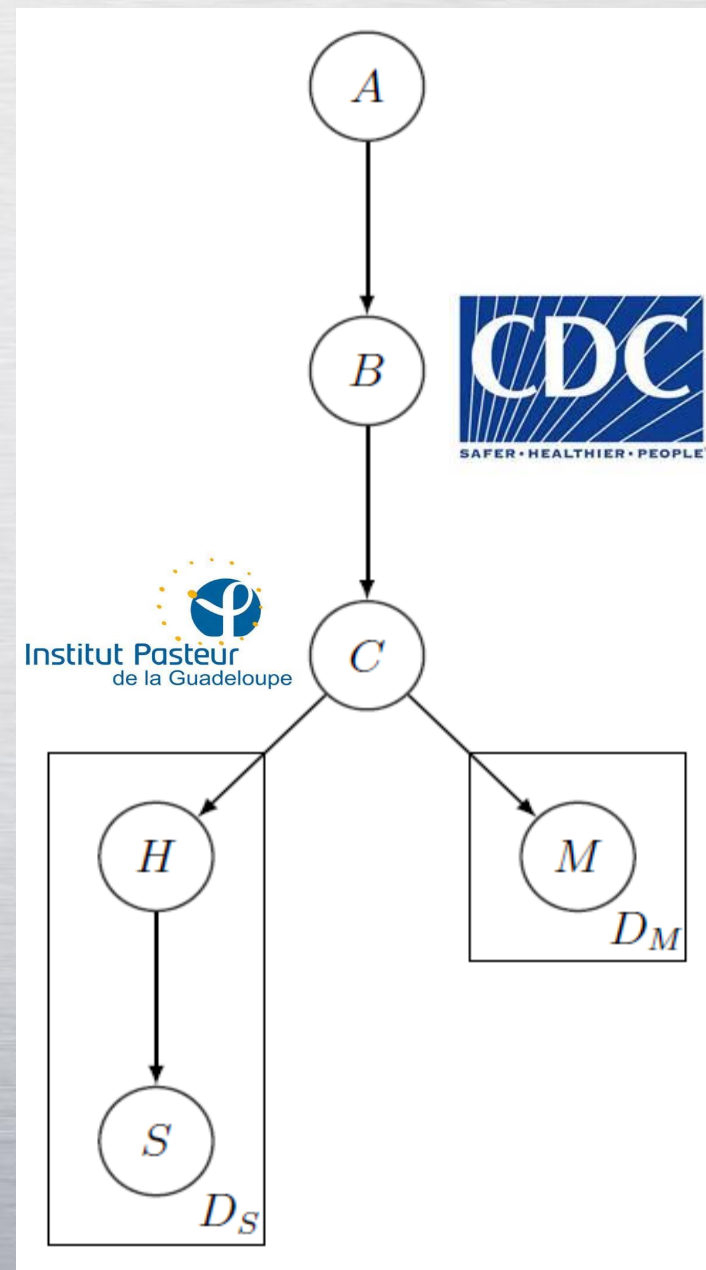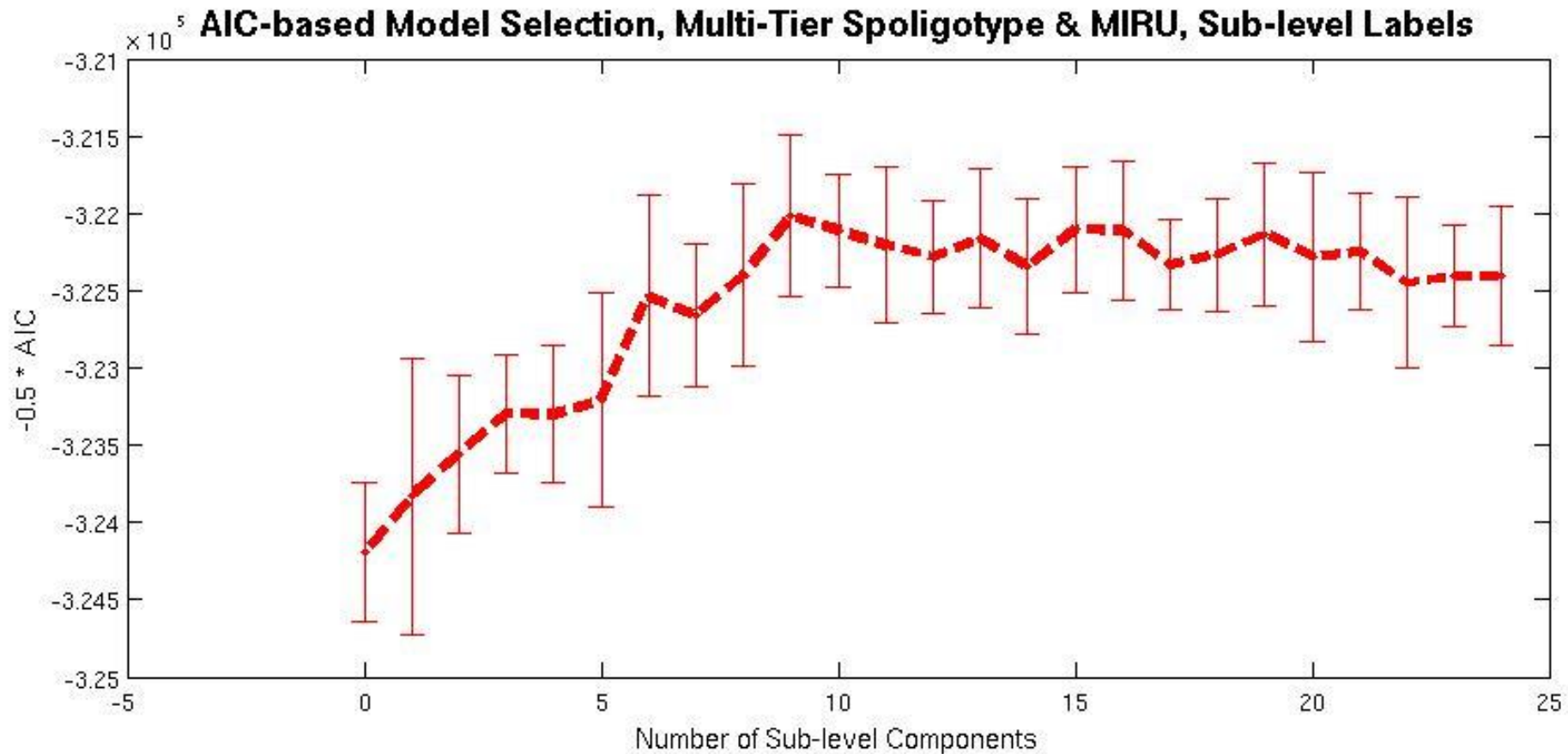| Top Pred. Label | Prob. | Mid Pred. Label | Prob. | Sub Pred. Label | Prob. | Size | Spoligotype Probabilities | MIRU2 | MIRU4 | MIRU10 | MIRU16 | MIRU20 | MIRU23 | MIRU24 | MIRU26 | MIRU27 | MIRU31 | MIRU39 | MIRU40 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M. africanum | 0.034 | West African 2 | 0.423 | AFRI_1 | 1 | 1719 | | | | | | | | | | | | | |
| M. africanum | 0.034 | West African 1 | 0.577 | AFRI_2 | 0.663 | 1491 | | | | | | | | | | | | | |
| M. africanum | 0.034 | West African 1 | 0.577 | AFRI_3 | 0.043 | 58 | | | | | | | | | | | | | |
| M. africanum | 0.034 | West African 1 | 0.577 | OtherSub6 | 0.283 | 722 | | | | | | | | | | | | | |
| M. bovis | 0.065 | M. bovis | 1 | BOV_1 | 0.501 | 4035 | | | | | | | | | | | | | |
| M. bovis | 0.065 | M. bovis | 1 | BOV_2 | 0.456 | 3248 | | | | | | | | | | | | | |
| M. bovis | 0.065 | M. bovis | 1 | BOV_3 | 0.038 | 375 | | | | | | | | | | | | | |
| M. canettii | 0.002 | M. canettii | 0.999 | Canettii | 0.998 | 212 | | | | | | | | | | | | | |
| M. caprae | 0.011 | M. caprae | 1 | Caprae | 1 | 1373 | | | | | | | | | | | | | |
| M. microti | 0.002 | M. microti | 1 | Microti | 0.998 | 232 | | | | | | | | | | | | | |
| M. mungi | 0.004 | M. mungi | 1 | M. mungi | 0.999 | 500 | | | | | | | | | | | | | |
| M. pinnipedii | 0.006 | M. pinnipedii | 1 | Pini1 | 0.832 | 631 | | | | | | | | | | | | | |
| M. pinnipedii | 0.006 | M. pinnipedii | 1 | Pini2 | 0.165 | 152 | | | | | | | | | | | | | |
| East Asian (Beijing) | 0.108 | East Asian (Beijing) | 1 | Beijing | 1 | 12922 | | | | | | | | | | | | | |
| East-African Indian | 0.044 | East-African Indian | 1 | CAS1-Delhi | 0.842 | 4359 | | | | | | | | | | | | | |
| East-African Indian | 0.044 | East-African Indian | 1 | CAS2 | 0.047 | 286 | | | | | | | | | | | | | |
| East-African Indian | 0.044 | East-African Indian | 1 | CAS1-Kili | 0.104 | 552 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | X1 | 0.362 | 4145 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | X2 | 0.371 | 4282 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | X3 | 0.202 | 2574 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | LAM8 | 0.015 | 318 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | OtherSub5 | 0.01 | 164 | | | | | | | | | | | | | |
| Euro-American | 0.617 | X | 0.161 | OtherSub7 | 0.033 | 534 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | H1 | 0.297 | 4899 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | H3 | 0.522 | 7096 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | Ural-1 | 0.062 | 1006 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | Ural-2 | 0.043 | 526 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | H2 | 0.056 | 1086 | | | | | | | | | | | | | |
| Euro-American | 0.617 | Haarlem | 0.214 | OtherSub8 | 0.018 | 653 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM9 | 0.36 | 7215 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM6 | 0.049 | 1048 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM4 | 0.043 | 313 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM1 | 0.092 | 1674 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM5 | 0.04 | 330 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM3 | 0.165 | 3041 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | T5-RUS1 | 0.018 | 1050 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM7-TUR | 0.058 | 895 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM11-ZWE | 0.063 | 1058 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM2 | 0.084 | 1585 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | LAM12-Madrid1 | 0.004 | 97 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | EAI2 | 0.01 | 388 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | OtherSub1 | 0.003 | 102 | | | | | | | | | | | | | |
| Euro-American | 0.617 | LAM | 0.259 | OtherSub9 | 0.003 | 166 | | | | | | | | | | | | | |
| Euro-American | 0.617 | EuroAm-African | 0.096 | LAM10-CAM | 0.321 | 1923 | | | | | | | | | | | | | |
| Euro-American | 0.617 | EuroAm-African | 0.096 | S | 0.493 | 3620 | | | | | | | | | | | | | |
| Euro-American | 0.617 | EuroAm-African | 0.096 | T2-uganda | 0.184 | 1121 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T1 | 0.691 | 14862 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T4 | 0.011 | 46 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T3 | 0.052 | 827 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T2 | 0.109 | 2973 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T-tuscany | 0.004 | 44 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T3-OSA | 0.004 | 158 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T1-RUS2 | 0.008 | 281 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T5 | 0.025 | 179 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T5-Madrid2 | 0.015 | 224 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T4-CEU1 | 0.017 | 372 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | T3-ETH | 0.021 | 454 | | | | | | | | | | | | | |
| Euro-American | 0.617 | T | 0.271 | H37Rv | 0.02 | 521 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Manila | 0.461 | EAI1-SOM | 0.203 | 1302 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Manila | 0.461 | EAI2-Manila | 0.755 | 4100 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Manila | 0.461 | OtherSub3 | 0.01 | 68 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Manila | 0.461 | OtherSub4 | 0.01 | 113 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Nonthaburi | 0.018 | EAI2-nonthaburi | 0.993 | 215 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | India | 0.146 | EAI3-IND | 0.736 | 1271 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | India | 0.146 | EAI8-MDG | 0.258 | 766 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Bangladesh | 0.205 | EAI6-BGD1 | 0.365 | 1036 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Bangladesh | 0.205 | EAI7-BGD2 | 0.634 | 1325 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Mexico | 0.02 | EAI-Mexico | 0.998 | 287 | | | | | | | | | | | | | |
| Indo-Oceanic | 0.107 | Vietnam | 0.151 | EAI4-VNM | 0.997 | 2215 | | | | | | | | | | | | | |
| Manu | 0.001 | Manu | 0.999 | OtherSub2 | 0.989 | 272 | | | | | | | | | | | | | |

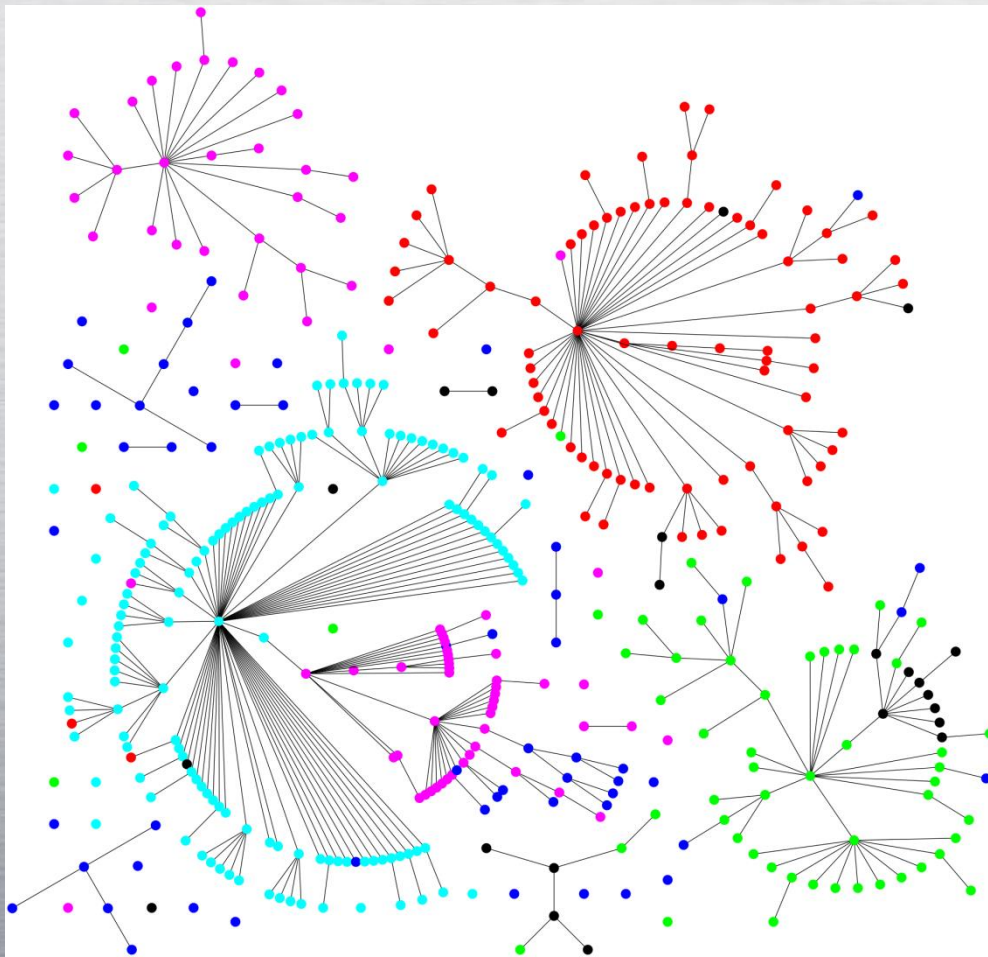# Model Adds Sublineages but not Mid or Major Lineages

# Ten New Putative Sublineages

- Discovers new sublineages and assigns them mid-level lineage
- New lineage characterized by ''long'' deletions frequently covering "typical" lineage deletions
- Covers many previously unlabeled isolates.

| X | | 0.161 | X1 | 0.362 | 4145 | |
|---|---|-------|-----|-------|------|---|
| X | | 0.161 | X2 | 0.371 | 4282 | |
| X | | 0.161 | X3 | 0.202 | 2574 | |
| X | | 0.161 | LAM8 | 0.015 | 318 | |
| X | | 0.161 | OtherSub5 | 0.01 | 164 | |
| X | | 0.161 | OtherSub7 | 0.033 | 534 | |

# Visualizing X Lineage Spoligoforest



Graphvis Two-PI

LAM8
OtherSub5
OtherSub7
X1
X2
X3

# Multi-Objective Embedding Methods

Good visualizations
minimize *stress* and *edge crossings*

# Edge Crossing Constraints as Classification

No Crossing =

Correct Classification

Crossing =

Incorrect  Crossing

# Multi-Objective Graph Embedding

Add classifier for each potential edge/node crossing parameterized by *U*

embedding error on *X* (MDS stress)

+ misclassification error based on *U* (SVM max margin)

$$\min_{X,u} Stress(X) + \sum_{i=1}^{m} \rho_i [\| (-A_i(X)u_i + 1)_+ \|_1 + \| (B_i(X)u_i + 1)_+ \|_1^1]$$

Optimize by an alternating algorithm on *X, U* using scalable classification and embedding algorithms

# Comparison: LAM sublineages



a) *Stress Majorization with Edge-crosssing penalties*

b) *Stress Majorization*

| | |
|---|---|
| ● | LAM1 |
| ● | LAM2 |
| ● | LAM3 |
| ● | LAM5 |
| ● | LAM6 |
| ● | LAM8 |
| ● | LAM9 |

c) *Graphviz Twopi*

d) *Laplacian Eigenmap*

# Visualizing X Lineage adding MIRU and Spoligotype distances



Legend:
- LAM8
- OtherSub5
- OtherSub7
- X1
- X2
- X3

Shabbeer et al, 2012 Multiobjective Embedding, MDS +SVM

# Spoligotype Genetic Diversity within 4.3K TB Patients in NYC – 2001-2007



CDC FAMILIES

East Asian

East-African Indian

Euro-American

Indo-Oceanic

Mycobacterium africanum

Mycobacterium bovis

# Host-Pathogen Graphs

Patients = Circles
DNA Fingerprint = Box

Boxes nested to indicate multiple DNA fingerprints
- Other box = spoligotypes
- Inner box = RFLP.

Color by Patient property = Region of Birth

Split by lineage



32

S5624

HH8

HH10

M97

M11

U

U2

001

AREAS

Americas

Central Africa

East Africa

East Asia

Europe

Indian Subcontinent

North Africa / Middle East

Oceania

South Africa

Southeast Asia

United States

West Africa

33

AREAS

- Americas
- Central Africa
- East Africa
- East Asia
- Europe
- Indian Subcontinent
- North Africa / Middle East
- Oceania
- South Africa
- Southeast Asia
- United States
- West Africa

S5370

AH

AH44

34

Euro American

# Disease as Stock Market

Companies=Bacteria

CDC

UNITED STATES FEDERAL TRADE COMMISSION BUILDING

Buyers = Patients

# Stock Market Tree Map



http://robslink.com/SAS/democd9/sp500.htm

Split by Lineage – East African Indian



| Patient | ID | Biomarker 1 | Biomarker 2 | TB continent |
|---------|------|-------------|-------------|---------------------|
| Patient 1 | 105 | S00669 | MY8 | Indian Subcontinent |
| Patient 2 | 2443 | S00210 | GD139 | Indian Subcontinent |
| Patient 3 | 2452 | S00210 | MY44 | Indian Subcontinent |
| Patient 4 | 2487 | S00247 | NO12 | East Africa |

# NYC - East Asian

Euro American

# NYC - Euro-American

# Spoligotype – S00241
# BW90 with Patient Age



Ongoing transmission in the US that became resistant to INH antibiotic

# NYC - M. bovis

# M. Bovis with Age of Patient/Country



Extremely Young: Outbreak in US-born children of Mexican Parents likely due to unpasteurized cheese.

# NYC *M. bovis* (2001-2007)



- Extra pulmonary *M. bovis* strikes
  - Mexican Immigrants
  - US-born children of Mexican Immigrants
- Hypothesized caused: Unpasturized cheese

# Indo-Oceanic with Time in US



Large clusters with few US patients and no found epi-links.

# *Indo-Oceanic Anomaly*



Hypothesized cause:  IO strains have longer latency phenotype

# Survival curves for all lineages

Proportion of cases not yet activated

# Surveillance Data can reveal novel phenotypes and genotypes



Estimated new TB cases (all forms) per 100 000 population

0–24

## Whole Genome Sequencing versus Traditional Genotyping for Investigation of a *Mycobacterium tuberculosis* Outbreak: A Longitudinal Molecular Epidemiological Study

Andreas Roetzer, Roland Diel, Thomas A. Kohl, Christian Rückert, Ulrich Nübel, Jochen Blom, Thierry Wirth, Sebastian Jaenicke, Sieglinde Schuback, Sabine Rüsch-Gerdes, Philip Supply, Jörn Kalinowski, Stefan Niemann

# Challenges of Disease Control using Molecular Epidemiology

- **Coupling human intelligence and analytics to help public health care workers control disease**
- **Informing local decisions with global data**
- **Allocating scarce control resources effectively by predicting disease dynamics**
- **Incorporating rapidly evolving data**
  - Contact Investigations
  - New biomarkers for pathogen/host
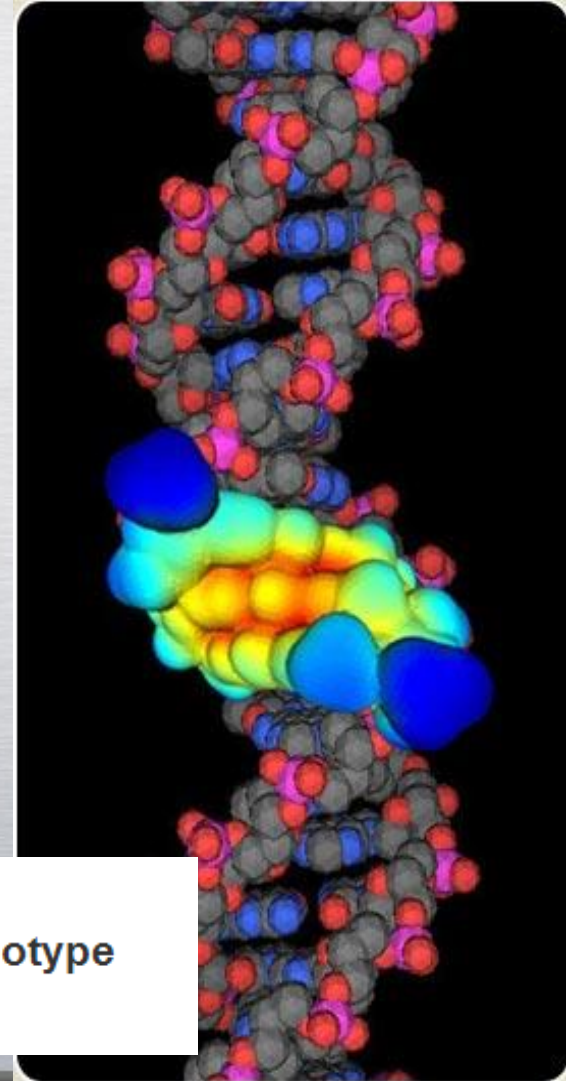  - Electronic Medical Records
  - Social media
- **Getting the biology right**
- **Preserving Privacy**

# Personalized Medicine based on Host and Pathogen DNA



- Discovery of Host/Pathogen Coadaptation

- Control and treatment efforts guided by host and pathogen DNA

- Better models for drug development, etc

Ethnicity and mycobacterial lineage as determinants of tuberculosis disease phenotype
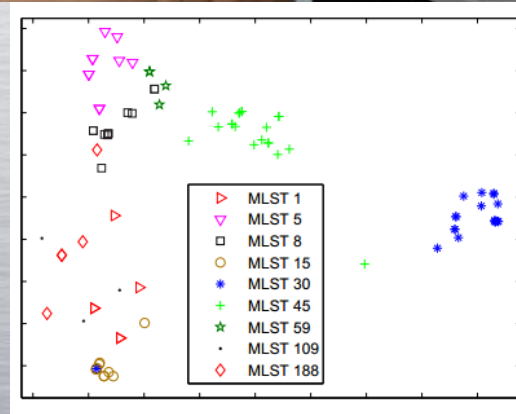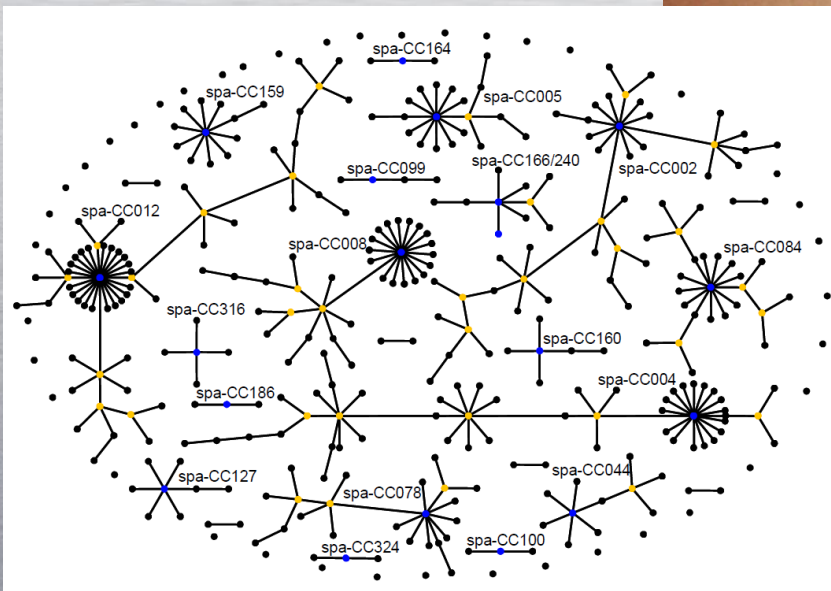
Thorax, 2013

# Other Diseases/Pathogens

**MRSA**

Mellmann *et al. BMC Microbiology* 2007, CDC, Agius et al, *IEEE/ACM Trans. on Comp. Bio.*, 2007.

# TB-Insight Project Team Supported by NIH R01-LMN009731

RPI:

- Professors: K. Bennett and B. Yener
- Graduate Research Assistants:

  **Amina Shabbeer**, **James Blondin**, **Inna Vitol,** Janani Ranganathan, Srivatsan Raghavan, Cagri Ozcaglar, Chris Gatti

- Postdoc:  **Minoo Aminian**
- Undergraduate Researchers:  Eric Dubois, Kane Hadley, Michael

## During this talk approximately 998 people developed active TB 276 people died of TB

Lauren Cowan, Jeff Driscoll, *US Centers for Disease Control*
Natalia Kurepina,  Barry Kreiswirth, *Public Health Research Institute*
**Nalin Rastogi,**  Phillip Supply  *Institut Pasteur*
Vincent Escuyer,  New York State Department of Health
Shama Ahuja,  Bianca Perri, Jeanne Sullivan
*New York City Department of Health*

**TB** insight
Tuberculosis Tracking and Control