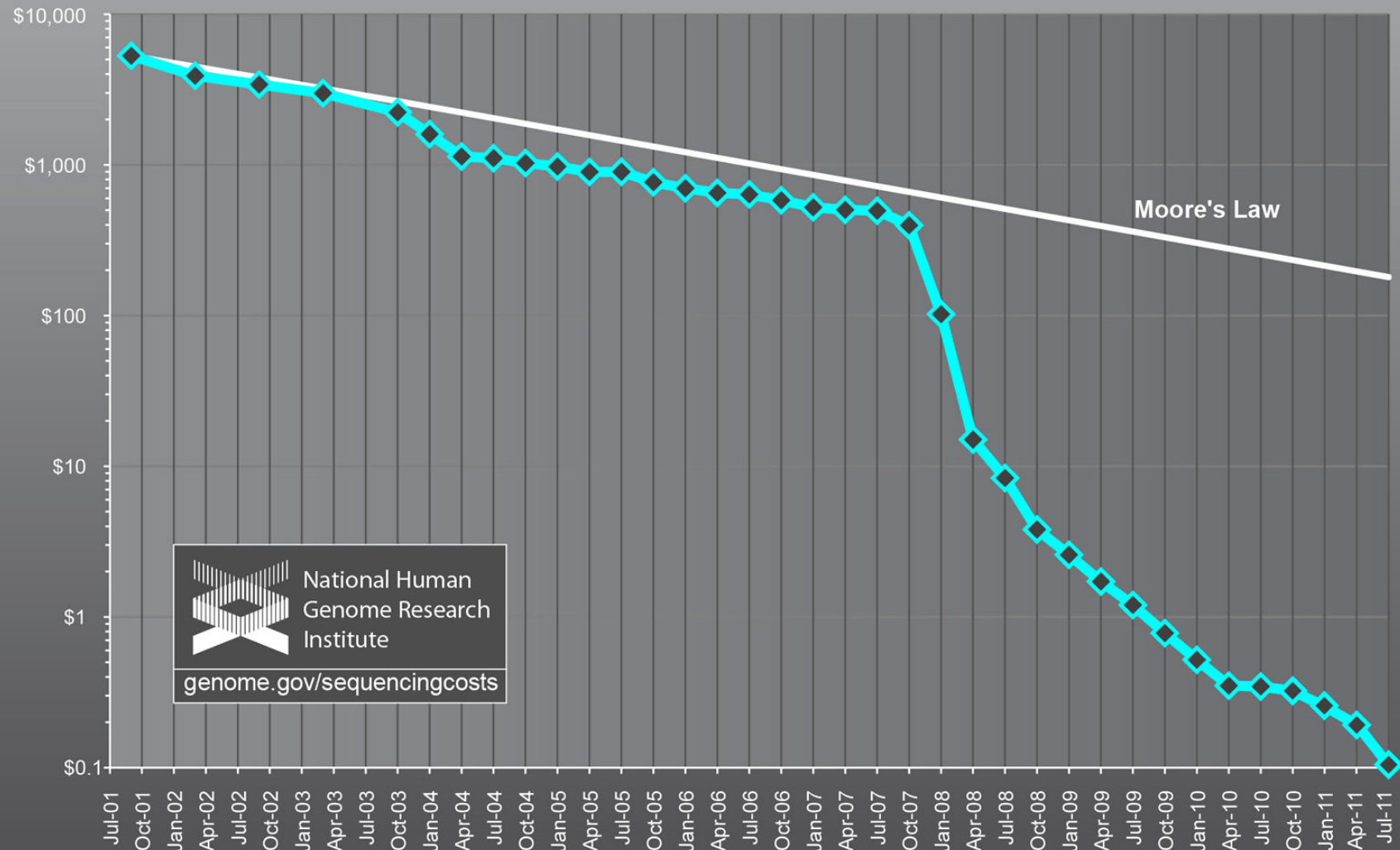# How will the control of infectious disease be improved by genomic analysis?

Christophe Fraser

*MRC Centre for Outbreak Analysis*
*Dept of Infectious Disease Epidemiology*

# Desktop sequencing in every lab and in every hospital ward



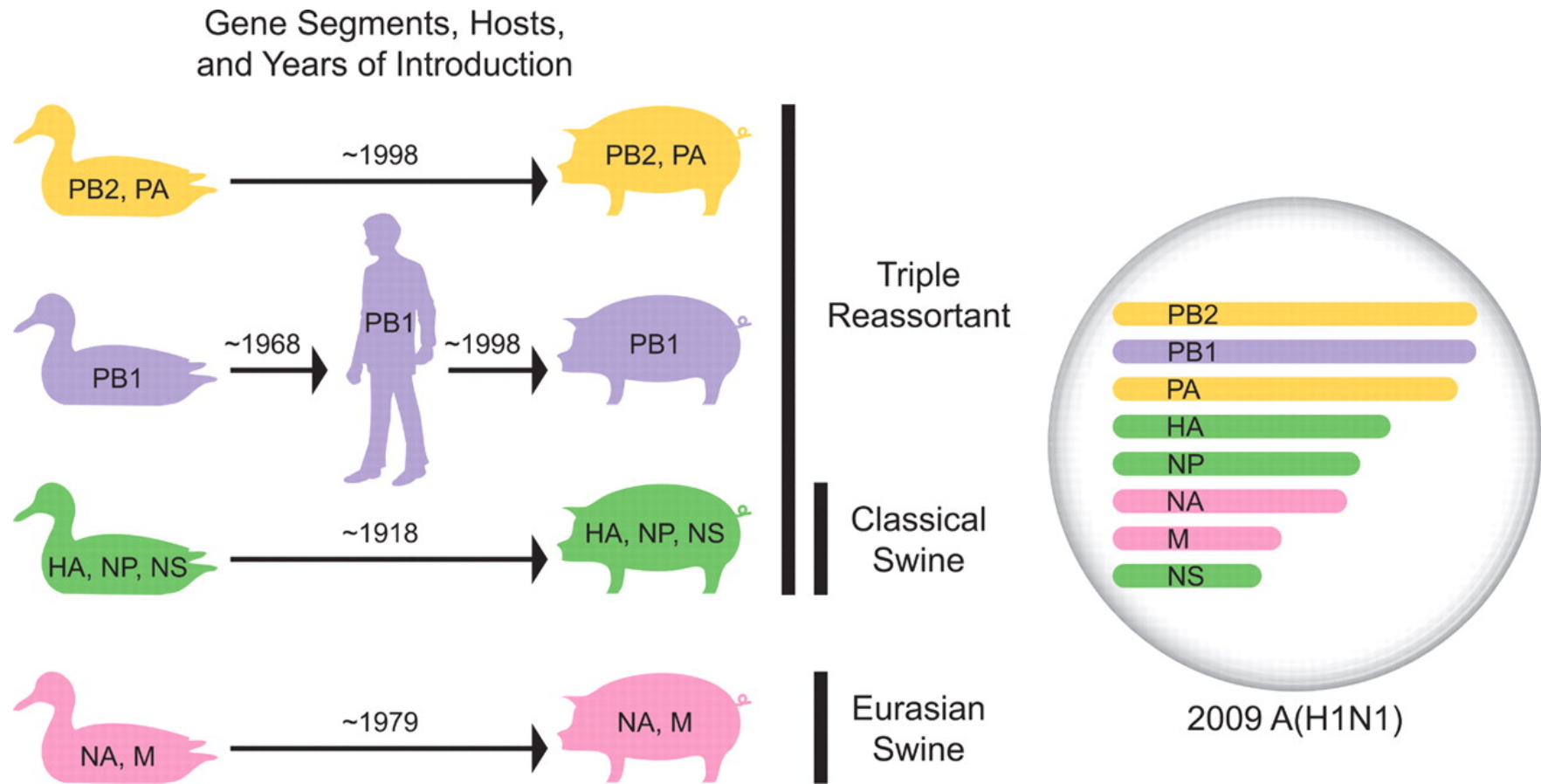## Sequencers are the new microscopes

# Why sequence?

- To track pathogens and outbreaks
- To identify factors important for disease
- To tailor treatment
- To understand evolution

# Two Examples:

- The 2009 H1N1 influenza pandemic

- The PMEN-1 lineage of *Streptococcus pneumoniae*

- The 2009 H1N1 influenza pandemic

- Tracking bacterial spread and evolution

# In April 2009, virological surveillance triggered a pandemic



Gene Segments, Hosts, and Years of Introduction

2009 A(H1N1)

Garten et al Science 2009, Smith et al Nature 2009

# Genetic data was shared in days

# Genomics revealed a complex evolutionary history



Still includes some essentially human-adapted core genes

*Gavin Smith et al, Nature 2009 & 2010*

## Number of full length HA sequences:

Currently over 2,600 full length HA sequences (for H1N1pdm virus)

# Earliest H1N1pdm trees (4 May)



Origin of outbreak strains
Estimated as 22/1/2009 (95% HPD
20/11/2009 to 13/3/2009)

We can already tell from tree structure that the epidemic wasn't growing very fast…

# Molecular epidemiology to detect local transmission (EpiInfo)

# PHYLODYNAMICS:



## BEAST (Bayesian Evolutionary Analysis Sampling Trees)

aaaagcaaca aaaatgaagg caatactagt agttctgcta
tgcagacaca ttatgtatag gttatcatgc gaacaattcat
actagaaaag aatgtaacag taacacactc tgttaaccttt
gaaactatgc aaactaagag gggtagcccc attgcatttg     +
ctggatcctg ggaaatccag agtgtgaatc actctccaca
tgtggaaaca tctagttcag acaatggaac gtgttaccca





Bayesian
Skyline
Plot

BSP

# The coalescent with variable population size

# The coalescent

Consider 2 infected people randomly chosen from $N_t$ total infected people

# The coalescent

Probability they shared a common ancestor time T ago is $\dfrac{1}{N_T}\displaystyle\prod_{x=1}^{(T-1)}\left(1-\dfrac{1}{N_x}\right)$



Time is measured in generations of infection

The coalescent can thus be used to reconstruct the number infected from a phylogeny…



O coalescent events

…provided a suitable model is used to relate mutations to time (*molecular clock model*)

# Earliest H1N1pdm trees (4 May)



Can already tell that the epidemic wasn't growing very fast
Based on <u>assuming</u> exponential growth: strong assumptions → robust estimates

# 9 May: Updated trees

Epidemiological model:



Population genetic model:



Orange – first iteration
Blue – updated – more effort to obtain 'random' unlinked sample

*Fraser et al, Science 2009*

Back-calculation

Outbreak investigation

*Effective reproduction number*

Phylodynamics

Time-series analysis

- The 2009 H1N1 influenza pandemic

- Tracking bacterial spread and evolution

# Streptococcus pneumoniae



- Gram +ve, commonly carried.
- Near ubiquitous in children.
- Causes otitis, pneumonia, invasive disease, meningitis, ….
- Causes 10% of all paediatric mortality.
- Diverse patterns of virulence and resistance.
- Antigenically diverse (92 serotypes).
- 7- and 13-valent vaccines now available.
- Naturally competent – recombinogenic.

# Vaccine-caused Serotype Replacement

- Serotype replacement was complete in US in <10 years (Hanage et al, Epidemics 2009).
- Disease levels decreased approx 2/3 in US, but very little decline in the UK (HPA UK).
- Huge ecological perturbation, with unknown effect on antibiotic resistance and virulence factors
- What are implications for global roll-out?

# Probing the core genome with Multi-Locus Sequence Typing (MLST)

*Now cheaper to sequence all 2,150,000 base pairs (whole genome with next generation sequencing) than 3,500 base pairs (MLST with capillary sequences).*

# A population genomic analysis

- Focus on PMEN-1 lineage:
  - Earliest recognised multi-drug resistant lineage of *Streptococcus pneumoniae* (penicillin, chloramphenicol, tetracycline, occasionally: fluoroquinolones & rifampicin, …)
  - Predominantly serotype 23F/ST81
  - Caused 40% of invasive disease in USA in 1990s
  - Member of the highly mosaic cluster (*based on BAPS/MLST analysis*)

- Full genomes from 241 isolates

# Rapid Pneumococcal Evolution in Response to Clinical Interventions



Global spread:

Western Europe
Eastern Europe
North America
South Africa
Southeast Asia
Central and South America

Vaccine escape:

Spanish origin (1970):

A

19F †

19F

15B

time

*Croucher et al, Science 2011*

# Multiple acquisition & loss of antibiotic resistance

- Whole lineage is resistant to penicillin, chloramphenicol & tetracycline.
- Fluoroquinolone resistance mutations acquired & lost, seemingly random.
- Macrolide resistance cassettes acquired repeatedly through horizontal gene transfer

**Key:**
- — Western Europe
- — Eastern Europe
- — North America
- — South Africa
- — South-East Asia
- — Central and South America

Recombination density

Number of events encompassing a base: >10, 9, 8, 7, 6, 5, 4, 3, 2, 1

pspA    cps locus    ICESpn23FST81    Prophage MM1    psrP    pspC

0.0030

90% of polymorphisms acquired by recombination, covering 75% of genome
(identified by SNP density)

Recombination not uniform along genome, but marked by hotspots

Some suggestion of lineages having experienced hyper-recombination?

**Key:**
- Western Europe
- Eastern Europe
- North America
- South Africa
- South-East Asia
- Central and South America

Recombination density
Number of events encompassing a base

0.0030

*pspA*    *cps* locus    ICESpn23FST81    Prophage MM1    *psrP*    *pspC*

Hyper-recombination accounts for over 50% of polymorphisms – role still unclear

# Phylogeny reveals a rich adaptive history



*Blue: hyper-recombination events*
*Yellow: serotype switches*
*Purple: Acquisition of macrolide resistance cassettes*
*Red: Fluoroquinolone resistance mutations*
*White: Abrogation of competence.*

With whole genome data, we can appreciate that genetic, genomic and selective events all occur concurrently, not in isolation. This will require new theory.

# Summary

- PMEN-1 lineage defined by acquisition of an accessory multidrug-resistant gene (ICE), and rapid spread since 1970

- Further changes in lineage driven by rapid switching in accessory genome (inc. antibiotic resistance and vaccine escape)

- New understanding of mechanisms of recombination

# Conclusions

- Sequencing is cheap, and will soon be standard.

- Most infections will be 'typed'.

- If nothing else, cheapest & quickest way of determining antibiotic resistance profile.

- Interpretation & analysis will remain challenging.

- Data storage & sharing, and linkage to meta-data will be extremely challenging, but necessary.