



**St. Michael's**

Inspired Care.  
Inspiring Science.



UNIVERSITY OF TORONTO  
DALLA LANA SCHOOL OF PUBLIC HEALTH

# **'Big Data for Health' strategy at the Dalla Lana Faculty of Public Health**

**Prabhat Jha, on behalf of the big data for health working  
group of DLSPH  
(Chairs David Henry and Prabhat Jha)**

Prabhat.jha@utoronto.ca  
Twitter: Countthedeath



# 'Wide' and 'Deep' data

- **Deep**
  - **Genome wide analysis**
  - **Proteomics, metabolomics, microbiomics**
  - **Functional MRI**
  - **Geospatial-linked exposures**
- **Wide**
  - **Population data, inc mortality**
    - **Linked at the level of the individual**
    - **Administrative data**
    - **Electronic health records**
    - **Registries**
- **COMMON to Both: PLATFORMS to enable IMAGINATIVE linkages by diverse brains**

# MILLION DEATH STUDY IN INDIA

1. Visit 1 M homes (“true snapshot” of India) with a recent death & ask standard questions and get a narrative
2. Use non-medical surveyors (electronic entry + GPS)
3. Web-based double coding by 500 doctors (guidelines + adjudication and other strict quality control)
4. Study all diseases, work with census dept, keep costs

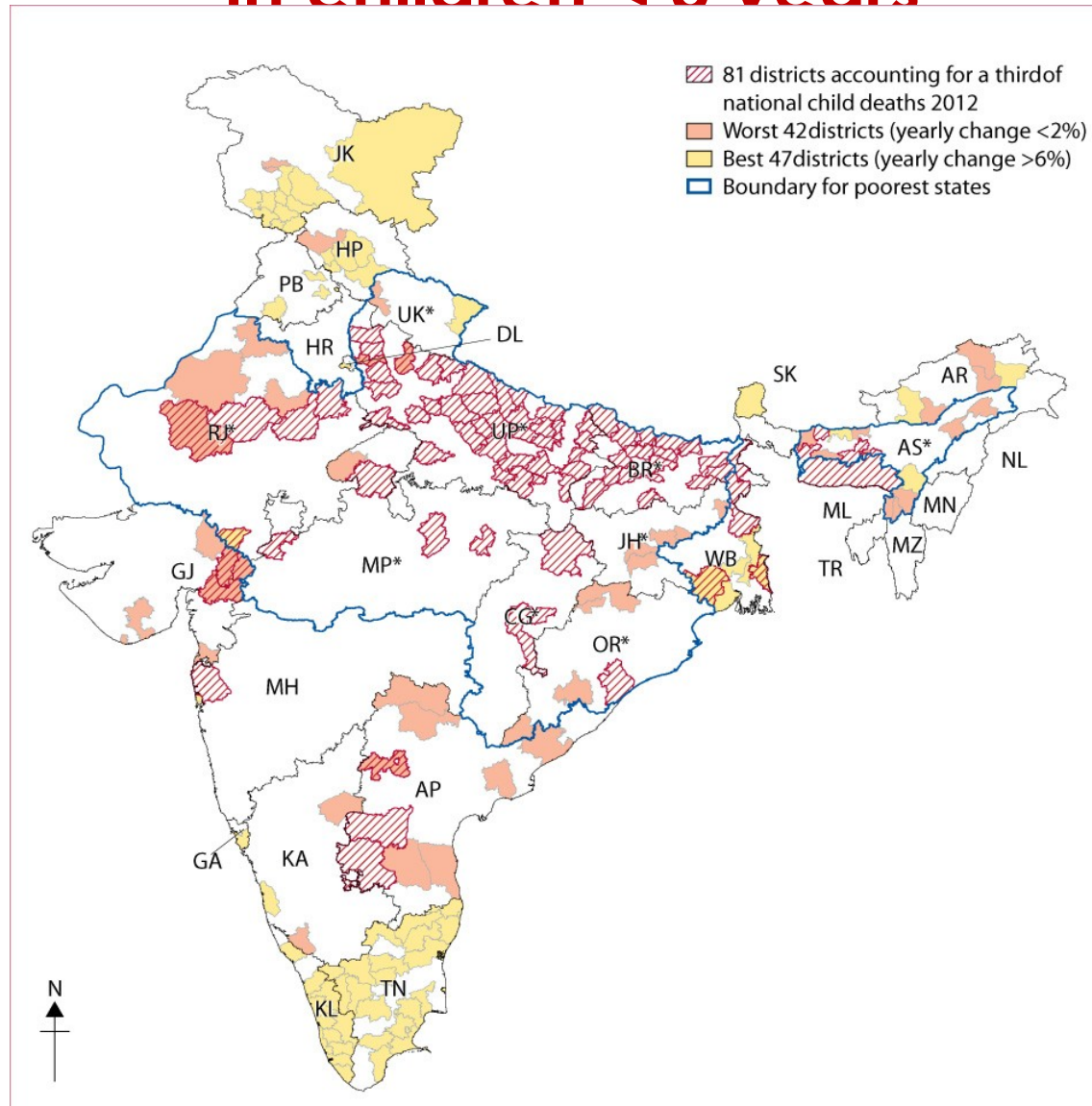


# **MILLION DEATH STUDY: selected results**

**(M=Millions, K=thousands)**

- **4-12M girls aborted before birth since 1980 (1/2 of these since 2000)**
- **1M smoking deaths (more than expected)** and 0.1M alcohol deaths
- **200K malaria deaths: WHO predicted only 15K**
- 100K HIV deaths: UNAIDS predicted 400K
- 60K pedestrian traffic deaths: Police estimate=9K
- 50K snakebite: WHO worldwide estimate=50K
- 33K cervical cancer: only 7K at Kashmir/Assam rate
- **Each common disease is rare somewhere in India, & hence is largely avoidable**

# 81 districts are home to 37% of the national deaths in children < 5 years

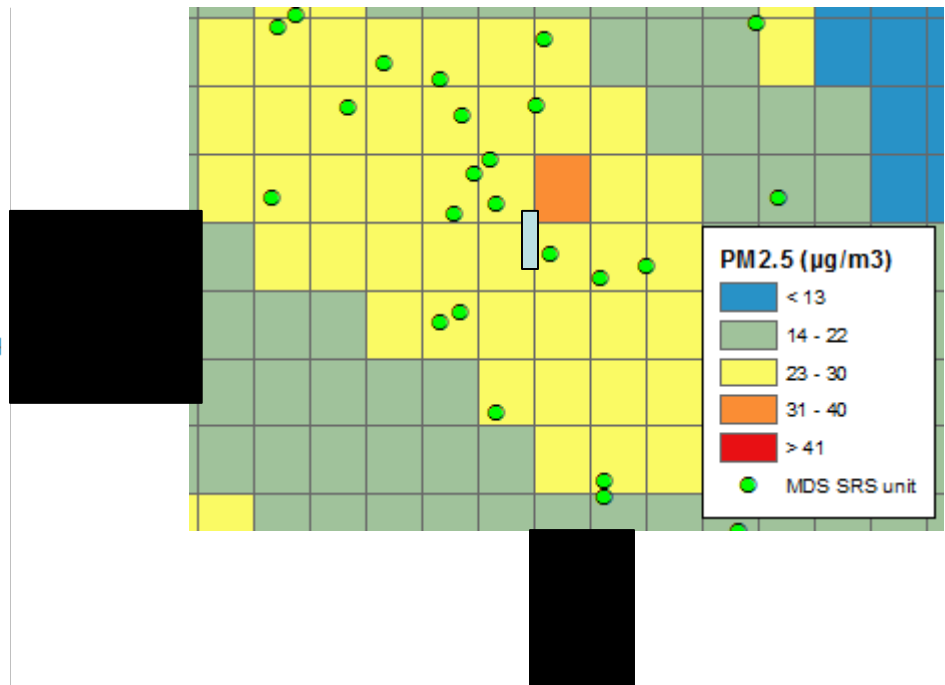
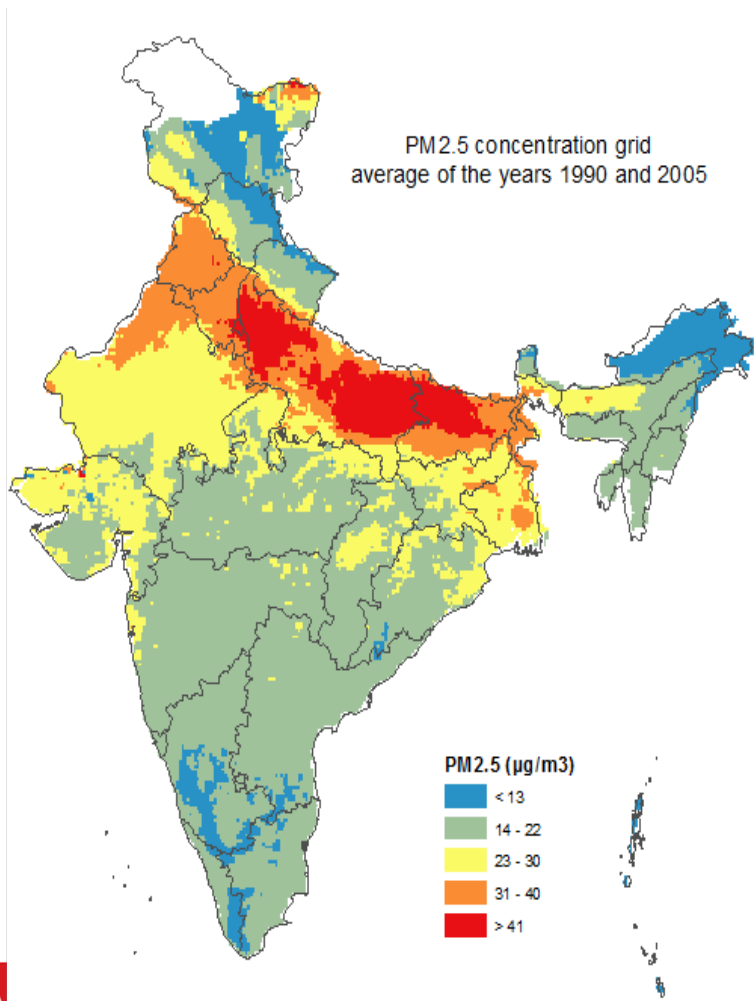


**68 of these  
81 districts  
are in  
poorer  
states**

**Figure 2:** Best and worst Indian districts according to change in under-5 mortality (2001–12), and districts with the highest third of under-5 mortality in 2012

# Geo-code this!

## Air pollution and mortality

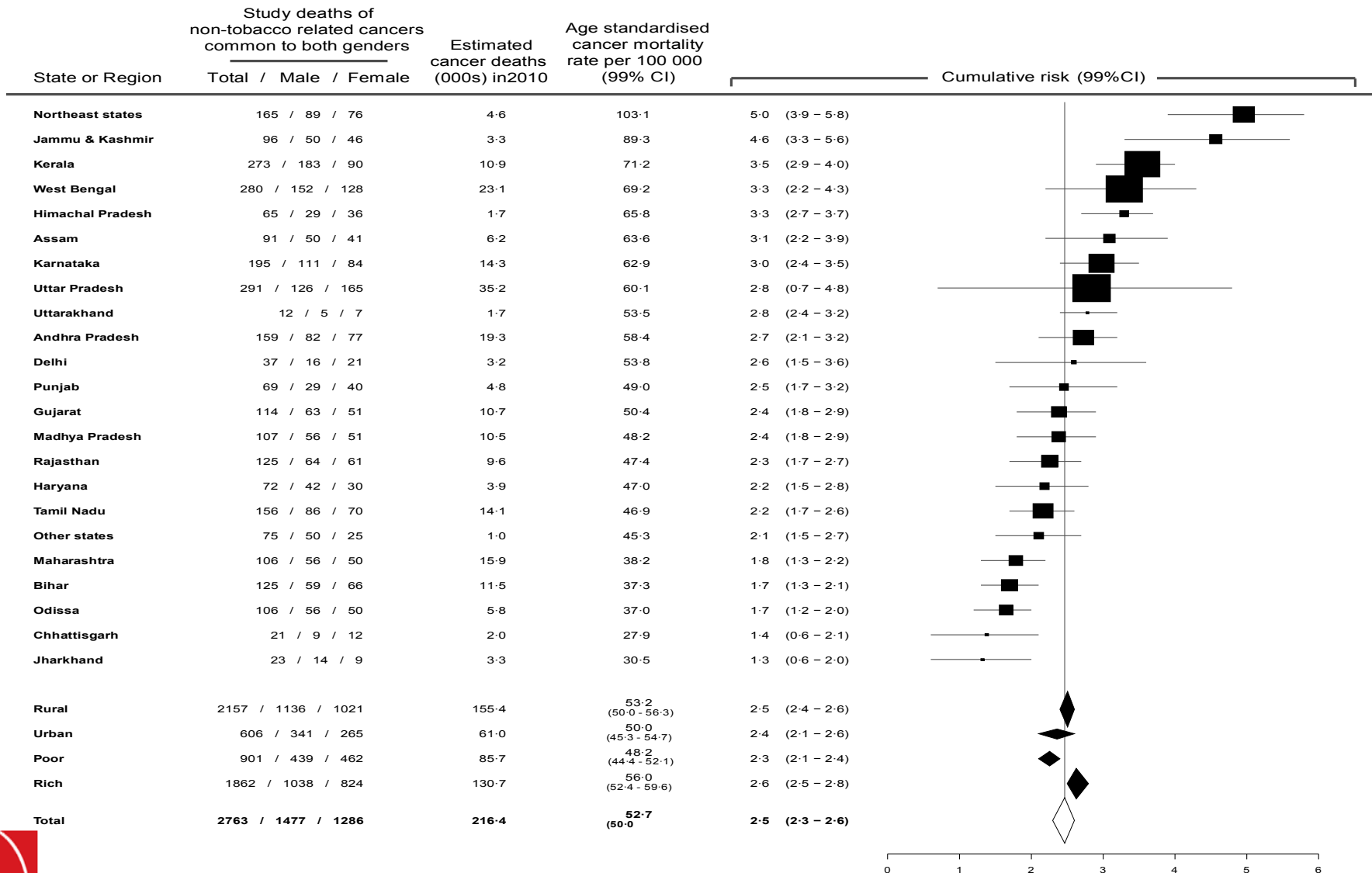


**Nationally representative estimates of ambient and indoor air pollution mortality (Half about WHO estimates!)**

Source: Yurie Maher, in press

# Cancer (non tobacco/non infection):

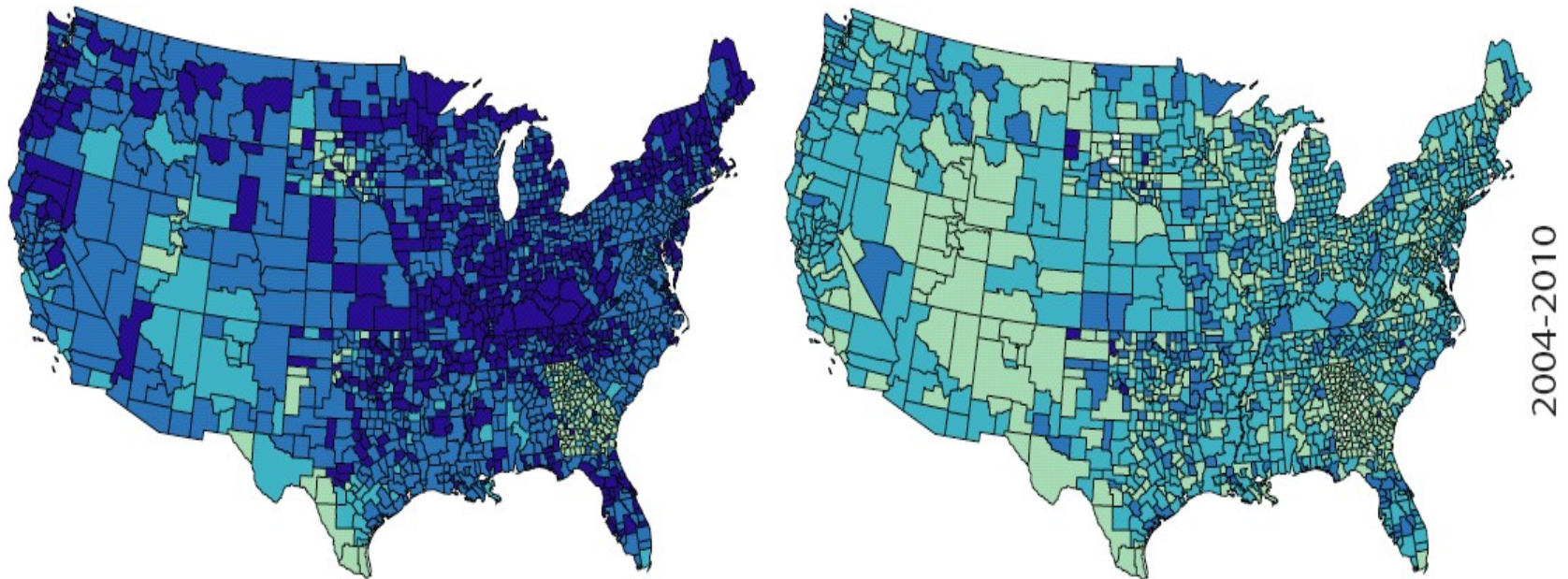
**BOTH GENDERS aged 30-69 years**



Source: Dikshit et al, Lancet 2012



# Tobacco deaths (% of total), WOMEN, US, 2004-2010 by county (0.4M deaths/year)

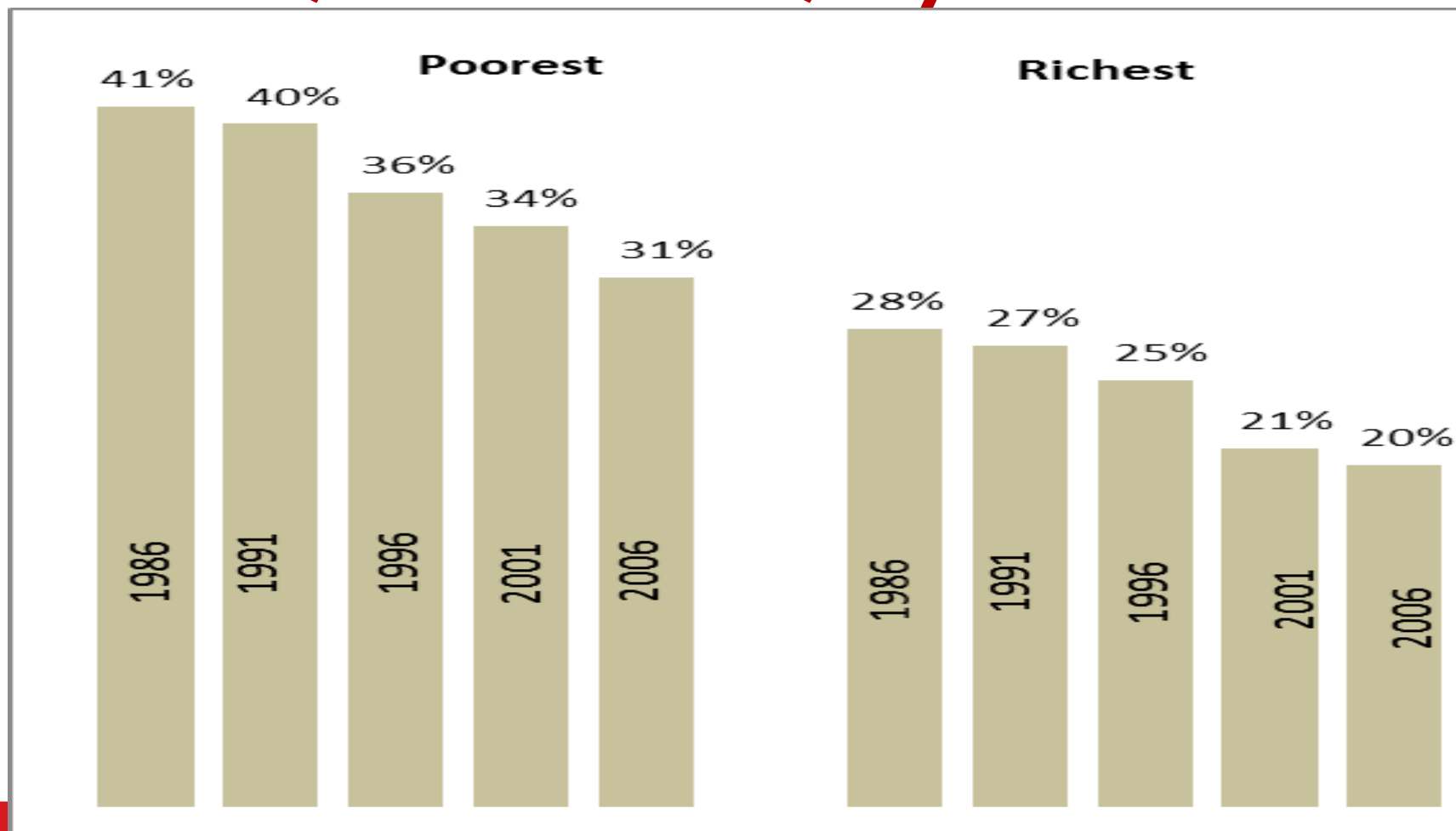


**Evaluate: smoking, obesity, health insurance (Obamacare) on mortality changes**

**Extend to 1.7 M deaths in Ontario and 4 M deaths in Mexico**

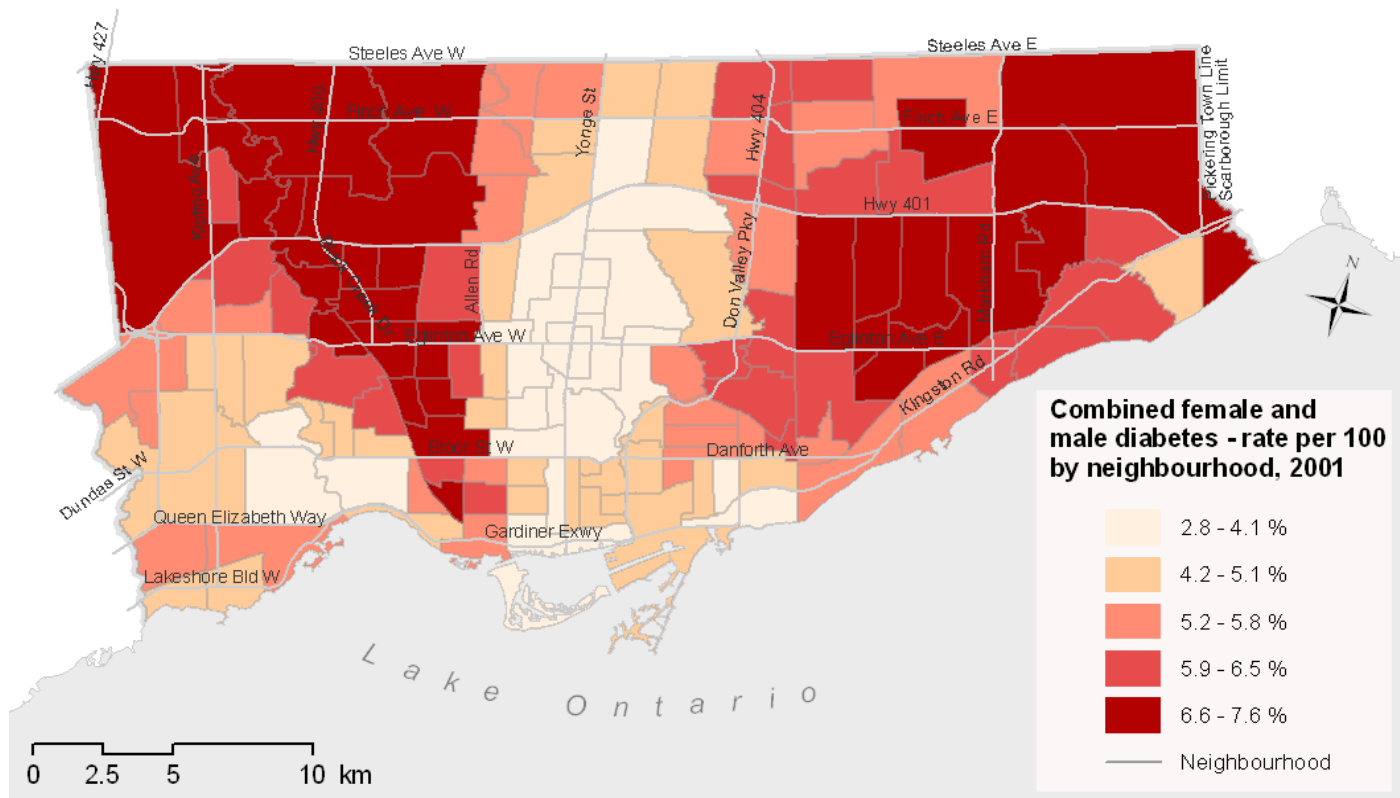


# Tobacco deaths (% of total), MEN, Canada, 1986-2006, by income

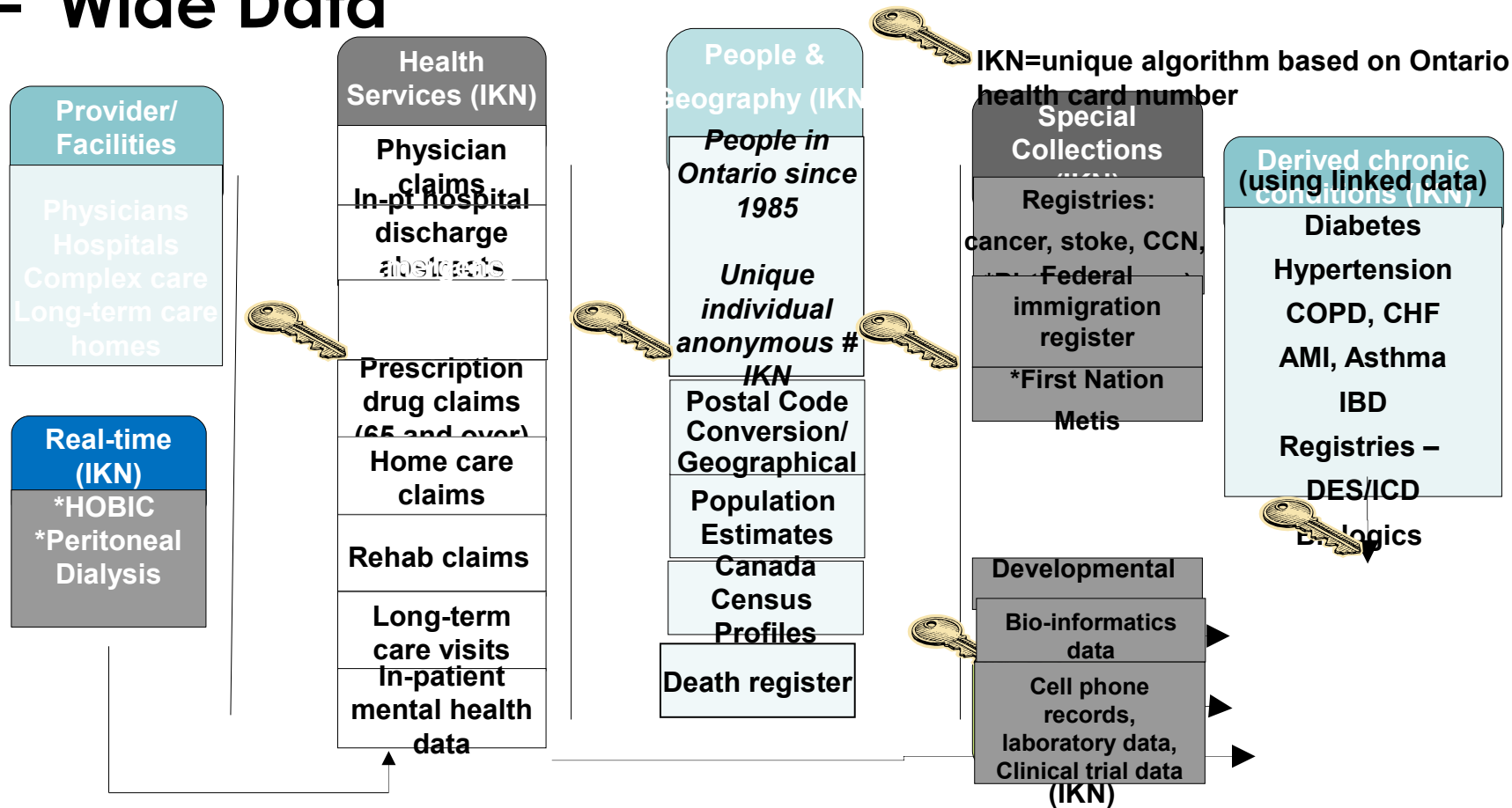


Source: Jha et al, CMAJ submitted

# Diabetes is more common in neighbourhoods with poor walkability

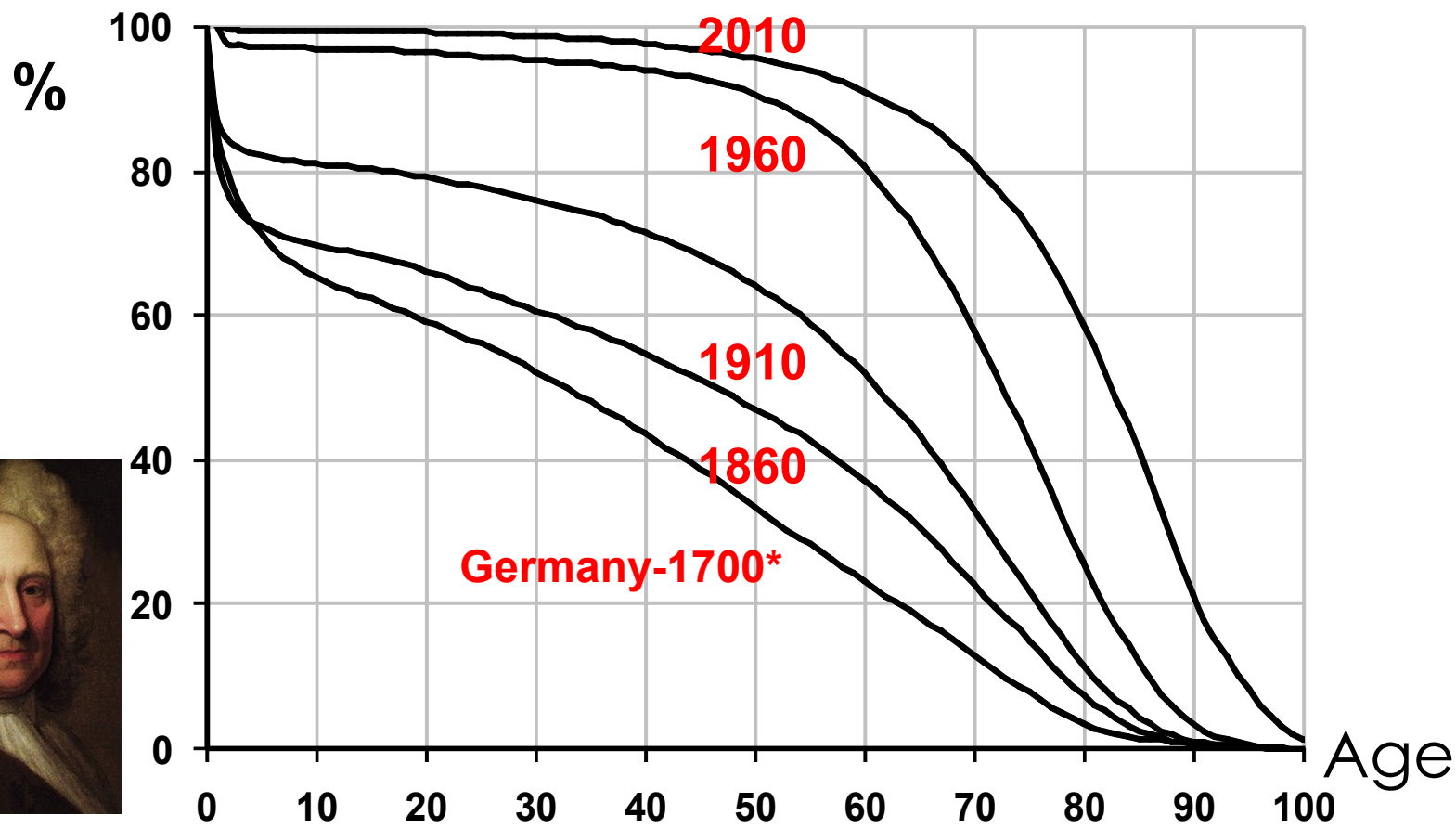


# – ‘Wide Data’



**Project Data-Set**

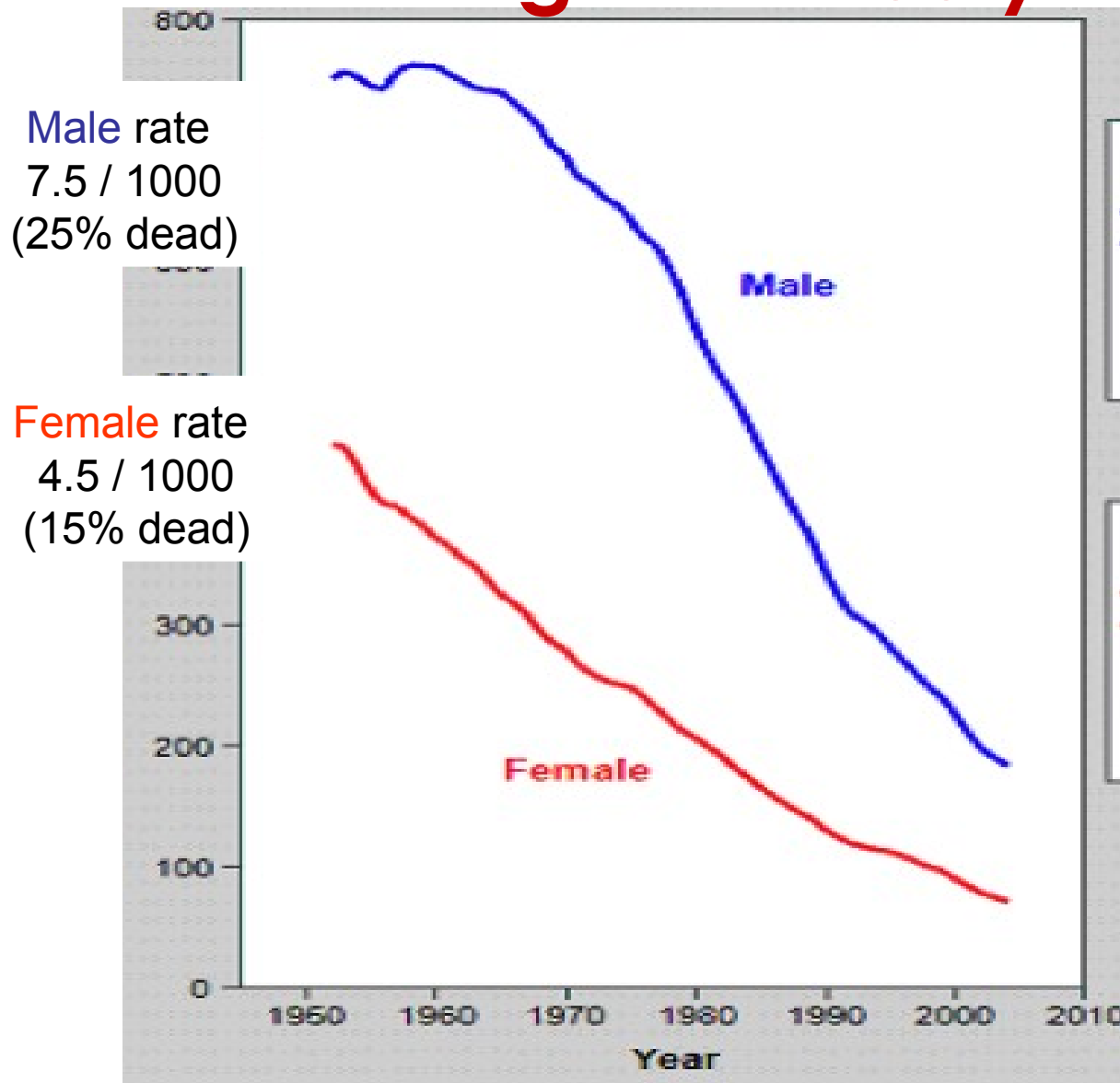
# Males, England & Wales, % survival at period rates



Source: Gary Whitlock, CTSU from Registrar-General reports and Human Mortality Database

\* Males and females combined- from Edmond Halley, 1693 for Breslaw, Germany

# Canada: vascular death rates in middle age over 50 years

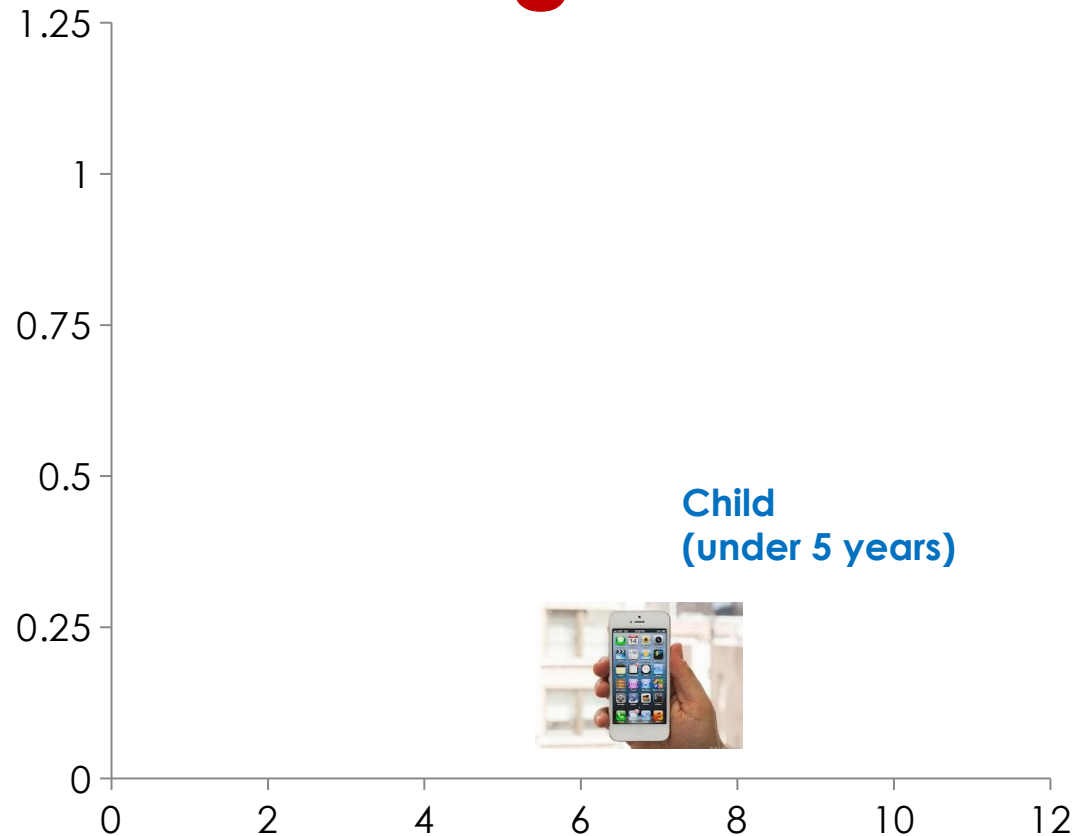


Vascular death  
at ages 35-69,  
Canada 2005:  
7% Male,  
3% Female

# Marginal costs for maximal child survival are falling



GDP per capita (\$2005, PPP, per day)



“Critical” incomes is real \$ needed to achieve  $\frac{1}{2}$  of maximal survival (in that year)  
from 1970 to 2007

# Marginal costs for maximal adult survival are rising



“Critical” incomes is real \$ needed to achieve  $\frac{1}{2}$  of maximal survival (in that year)

from 1970 to 2007; note higher adult costs due in part to HIV and tobacco

# Linkage of Wide and Deep data

- With large cohorts (100,000s) the administrative data are an efficient way of determining outcomes
- Identifying phenotypes in the administrative data to enable true population-based approaches
- Linking EMR/Admin data to bio-banked samples has great potential
- Overlaying environmental exposures or nutritional maps on geocoded disease distributions
- Great opportunity to study the 'Exposome'



# Analytical approaches

- **Complex algorithms to analyse genomic/ proteomic/ metabolomic data**
- **Data Mining: extract information from a data set and transform it into an understandable structure for further use**
- **Complex algorithms and models: image analysis, facial recognition, weather prediction**
- **Machine learning**
- **Measuring associations between exposures and outcomes**
  - **Propensity scores**
  - **Inverse probability of treatment weighting**
  - **Instrumental variable analysis**

# Discussions at the U of Toronto

## ·Big Data Committees in Faculty of Medicine and in the DLSPH

- Medicine – concentration on deeper analytical capability: genomics proteomics....
- DLSPH – discussions around linkage of broad datasets: surveys (OHS, CCHS etc), public health data, social services data, education data etc
- DLSPH – Focus on exposomics
- Also the potential for establishment of a Health Observatory in Ontario using linked electronic health and administrative data
- Support for the data platform from CIHR/SPOR

# A few big themes of this meeting

- Substantial capacity in genomics/biological correlates
- Unique representative population (Ontario, n=13M)
- Converging vision- ICES to population, genomics to large datasets
- Need for computing/analytic other PLATFORMS (think LHC plus computing grid)
- GLOBAL relevance- esp with chronic disease/mental health
- Need to train young scientists across disciplines