



Discovering Sub-diseases with Model-based Machine Learning

Iain Buchan

Director, Farr Institute @ Health eResearch Centre

Director, Centre for Health Informatics, University of Manchester

23rd September 2015

Machine Learning for Personalized Medicine Summer School, Manchester

Endotype Discovery

Aim to identify **subgroups** (“endo-phenotypes” or “endotypes”) of disease risk or treatment outcome explained by a **distinctive** underlying **mechanism**

Foundation of **Stratified Medicine**, seeking better-targeted interventions



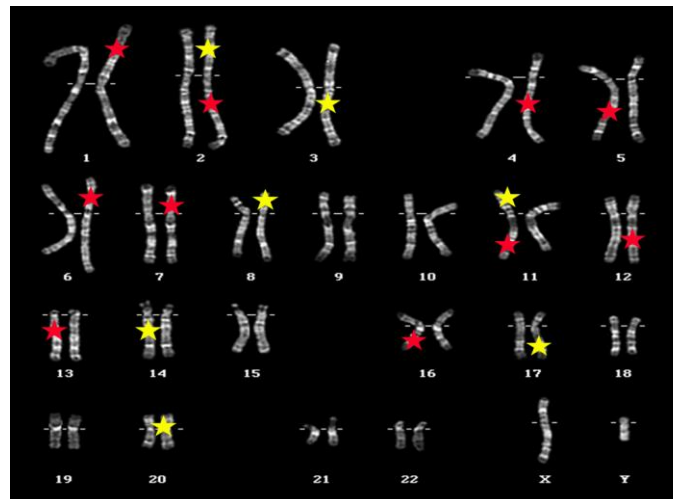
M.C. Escher
Order and Chaos, 1950

Asthma: Low GWAS Yield

- Legacy of non-replicated genetic epidemiology, typical of most common chronic disorders

★ Linkage in 1 study only

★ Linkage in >1 study



Asthma: Gene \cup Environment

- Important gene-environment interaction information may be averaged out in narrow or aggregated studies

CD14 Endotoxin Receptor

★ C allele associated

★ T allele associated

★ No association



Asthma: Heterogeneity

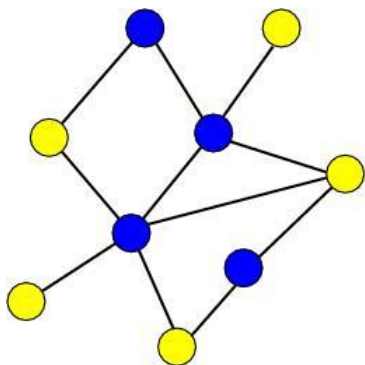
- Seems to be a collection of **several diseases**, each with distinctive pathophysiology, and environmental \cup genetic associates (“asthma endotypes”)
- Usually starts early in life, and may **progress**, **remit** or **relapse** over time
- Single cohort study **hypothesis-driven** epidemiology **lacks** temporal and environmental **complexity** needed to smoke out endotypes

Scaling-up

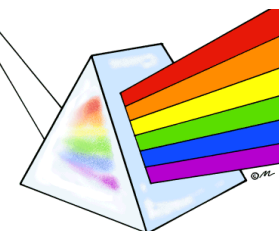
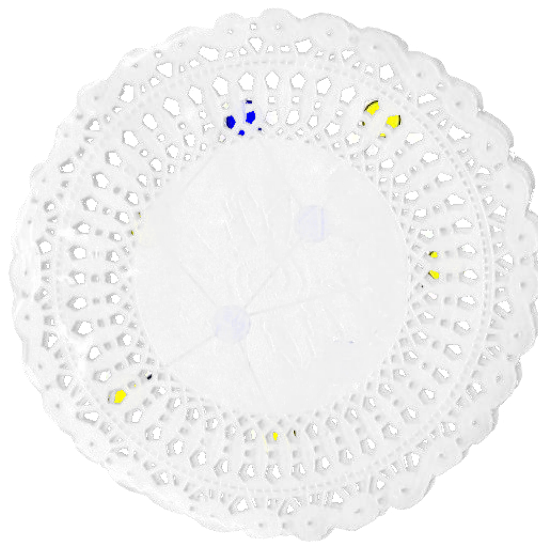
- **Multi-cohort**
 - Different windows on calendar **time**
 - Different windows on human **development**
 - Variety of **populations**
 - Variety of **environments**
- **Multi-disciplinary** and **multi-perspective**
 - Biostatistics (**deductive**) and machine-learning (**inductive**)
 - Tapestry of **reasoning** about mechanisms
- **Hypothesis forming** *and* **following**

Health Data: No Mining Please

Problem Space



Observation Space



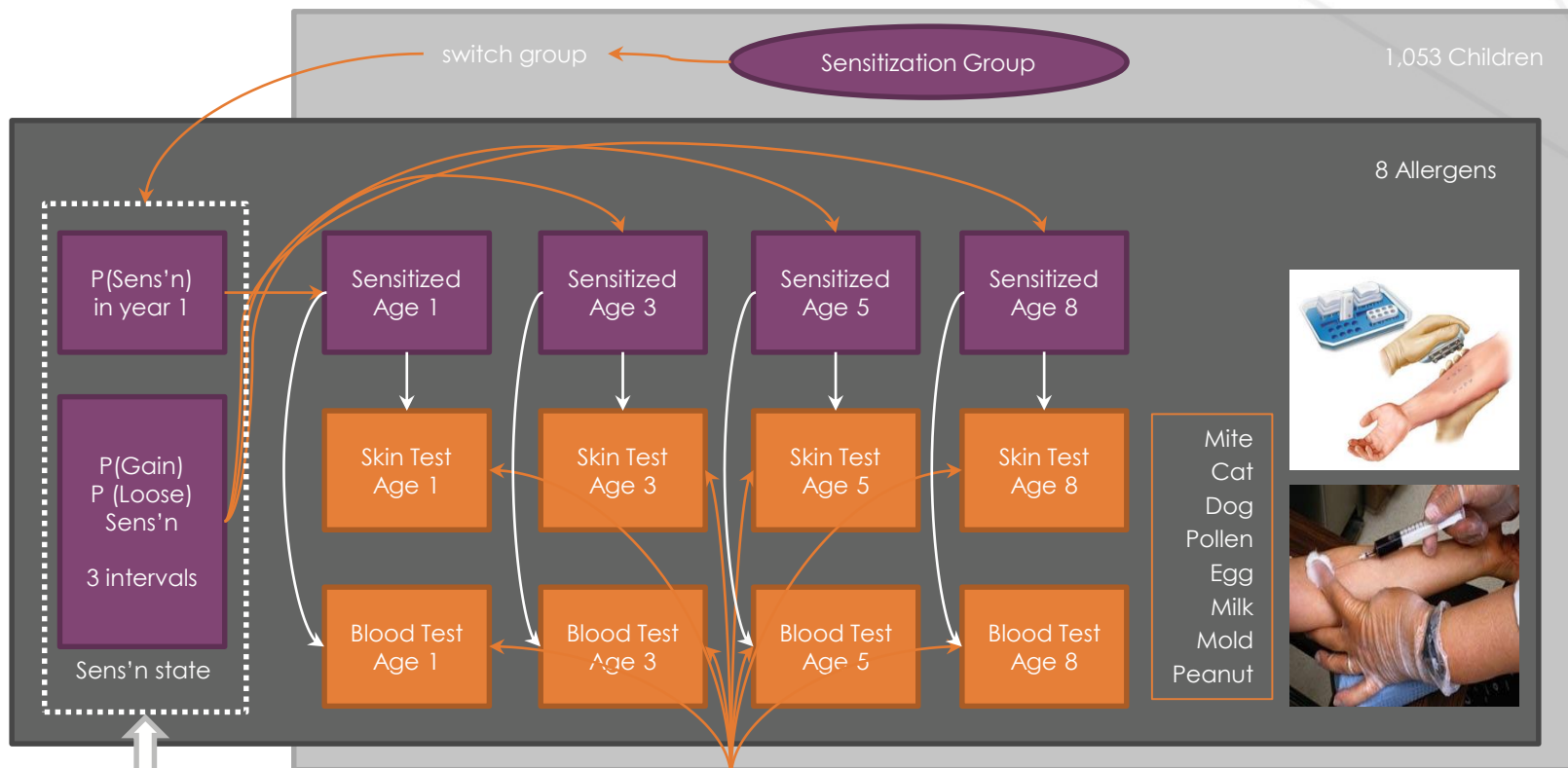
Data Space

Country	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	Total
Bangladesh	8.211	8.878	8.841	8.846	7.211	7.171	5.492	4.914	6.081	11.487	158.821
Canada	15.791	15.891	15.911	16.011	16.011	16.211	16.311	16.511	16.611	16.711	162.891
China	14.811	14.911	15.011	15.111	15.211	15.311	15.411	15.511	15.611	15.711	152.811
Colombia	14.111	14.211	14.311	14.411	14.511	14.611	14.711	14.811	14.911	15.011	147.811
Cuba	13.111	13.211	13.311	13.411	13.511	13.611	13.711	13.811	13.911	14.011	137.811
Dominican R	12.111	12.211	12.311	12.411	12.511	12.611	12.711	12.811	12.911	13.011	127.811
Ecuador	11.111	11.211	11.311	11.411	11.511	11.611	11.711	11.811	11.911	12.011	117.811
El Salvador	10.111	10.211	10.311	10.411	10.511	10.611	10.711	10.811	10.911	11.011	107.811
Germany	9.111	9.211	9.311	9.411	9.511	9.611	9.711	9.811	9.911	10.011	97.811
Guatemala	8.111	8.211	8.311	8.411	8.511	8.611	8.711	8.811	8.911	9.011	87.811
Hong Kong	7.111	7.211	7.311	7.411	7.511	7.611	7.711	7.811	7.911	8.011	77.811
India	6.111	6.211	6.311	6.411	6.511	6.611	6.711	6.811	6.911	7.011	67.811
Indonesia	5.111	5.211	5.311	5.411	5.511	5.611	5.711	5.811	5.911	6.011	57.811
Iran	4.111	4.211	4.311	4.411	4.511	4.611	4.711	4.811	4.911	5.011	47.811
Israel	3.111	3.211	3.311	3.411	3.511	3.611	3.711	3.811	3.911	4.011	37.811
Japan	2.111	2.211	2.311	2.411	2.511	2.611	2.711	2.811	2.911	3.011	27.811
Korea	1.111	1.211	1.311	1.411	1.511	1.611	1.711	1.811	1.911	2.011	17.811
Mexico	0.111	0.211	0.311	0.411	0.511	0.611	0.711	0.811	0.911	1.011	7.811
Nicaragua	0.011	0.011	0.011	0.011	0.011	0.011	0.011	0.011	0.011	0.011	0.811
Nigeria	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.011
Peru	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001
Philippines	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Poland	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Slovenia	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Thailand	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
USA	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Vietnam	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Yugoslavia	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Other	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Total	100.000	100.000	100.000	100.000	100.000	100.000	100.000	100.000	100.000	100.000	100.000

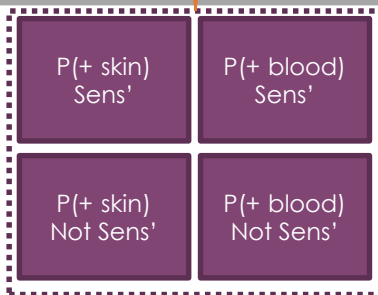
$$y = b_1x_1 + b_2x_2 + b_3x_3 + c$$

...Health is measured with **error** and **missingness**:
 Endo-phenotypes are resolved as if the researcher was looking through a prism and doyley at the problem

Hypothesising with Data



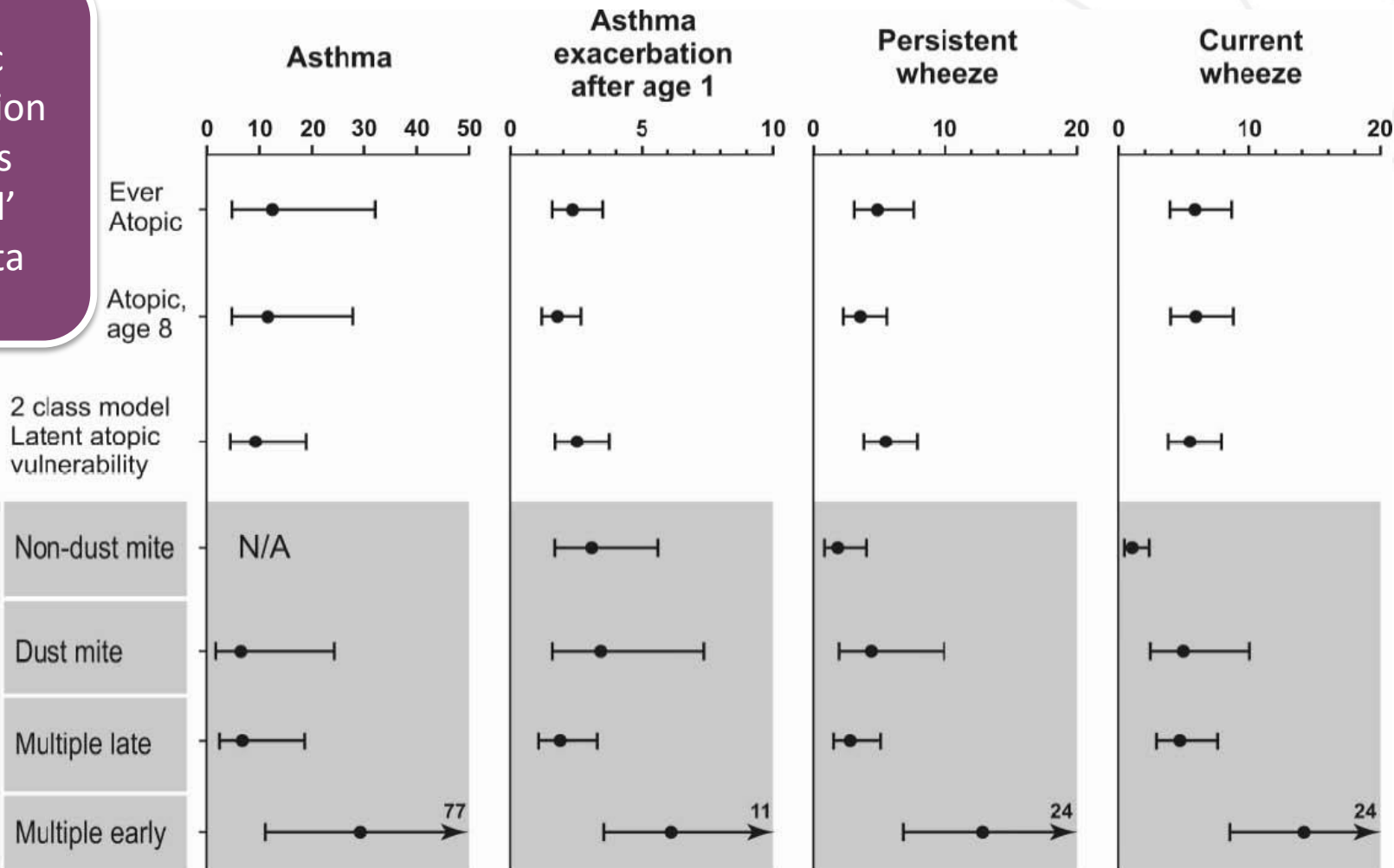
Machine-learning software & partial statistical models



The Farr Institute of Health Informatics Research

New Asthma Risk Factor Found

Allergic sensitisation patterns 'learned' from data



Cross Cohort Team Research

Data & Harmonized
Metadata from Cohorts

MRC STELAR Consortium: www.asthmaelab.org

MAAS

SEATON

ASHFORD

ALSPAC

IOW

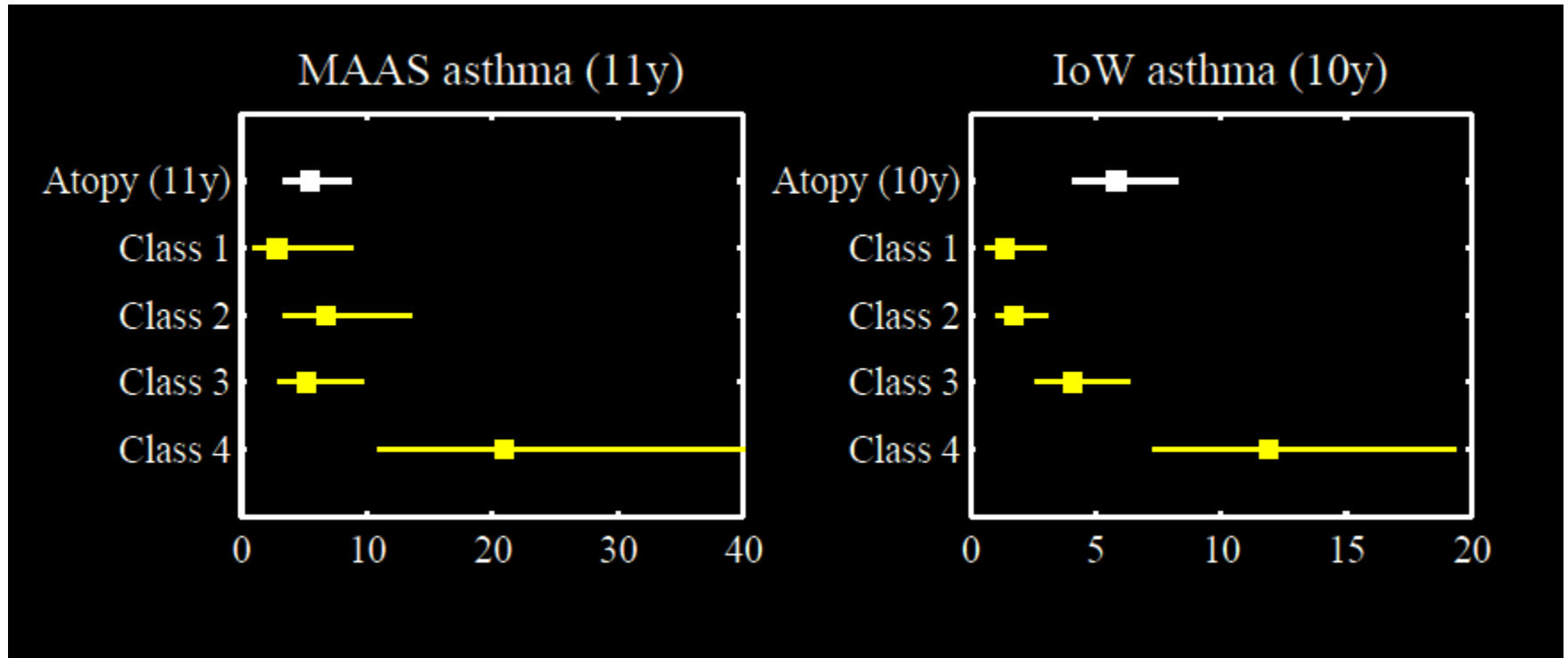
The screenshot displays the 'eLab Collaboration Site Dashboard' for the STELAR consortium. The page includes a navigation menu with options like 'My Dashboard', 'Sites', 'People', and 'Repository'. A welcome message for Adnan Custovic is at the top. Below this, there are sections for 'Site Members' (listing roles like Administrator, Collaborator, and Contributor), 'Site Activities' (showing recent document updates and deletions), 'Site Content' (listing documents like STELAR-5.docx and STELAR-1.docx), and 'Site Data Lists'. A sidebar on the left contains a 'Site Calendar' with dates and events.

Data Extracts

Modelling

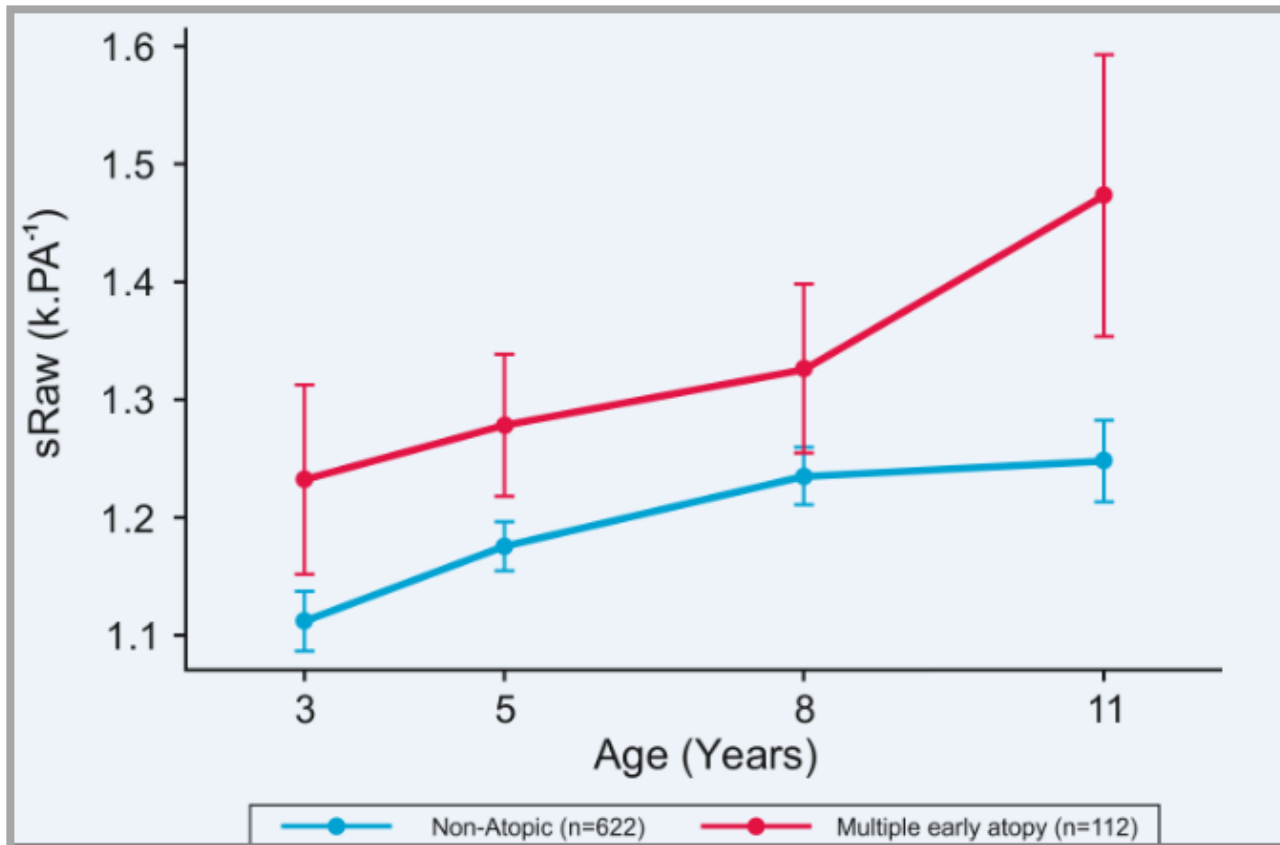
Networking:
Ideas, Activities,
Results, Meanings

'Learned' Atopy Classes Portable



Lazic et al, Allergy 2013; 68(6): 764-70

Risk Factor Development

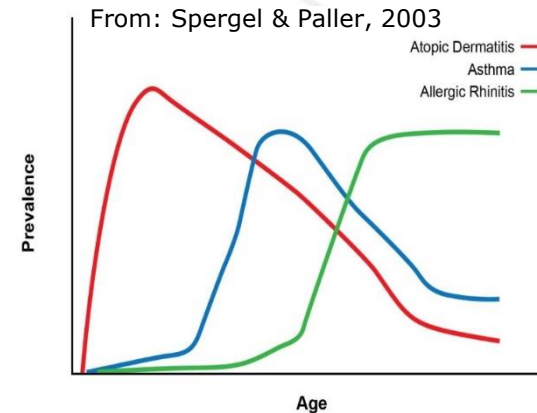


Risk group
'induced' from data
shows different
natural history of
airways resistance

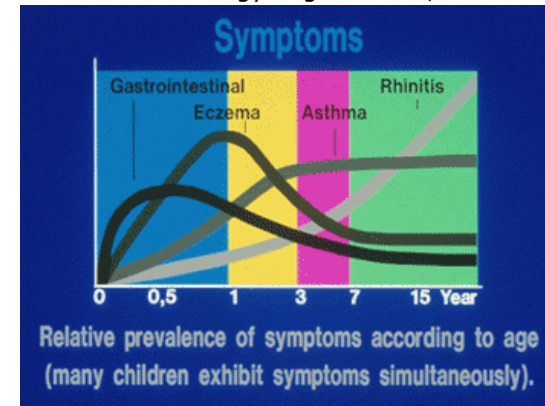
Belgrave et al, Am J Respir Crit Care Med 2014;189(9):1101-9

Assumed Biology: Atopic March

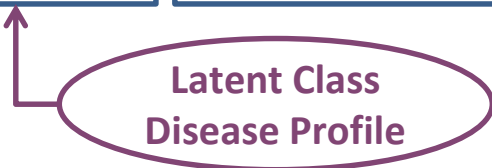
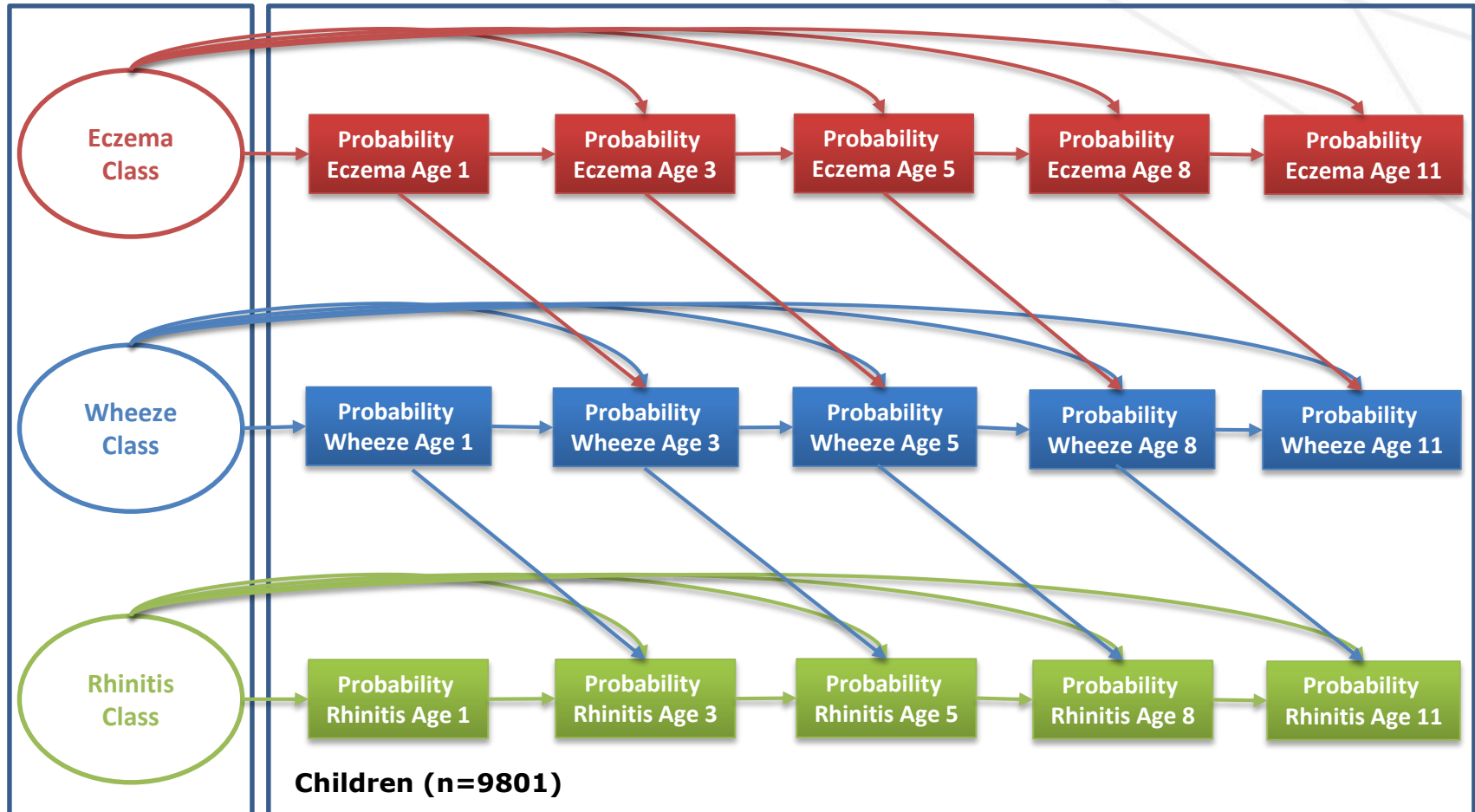
- Progression of **allergy**
Eczema → **Asthma** → **Rhinitis**
- Inferred from **population** summary →
- Assumed **causal** link between eczema – asthma & rhinitis
- Clinical response: **target** children with eczema to reduce progression to asthma



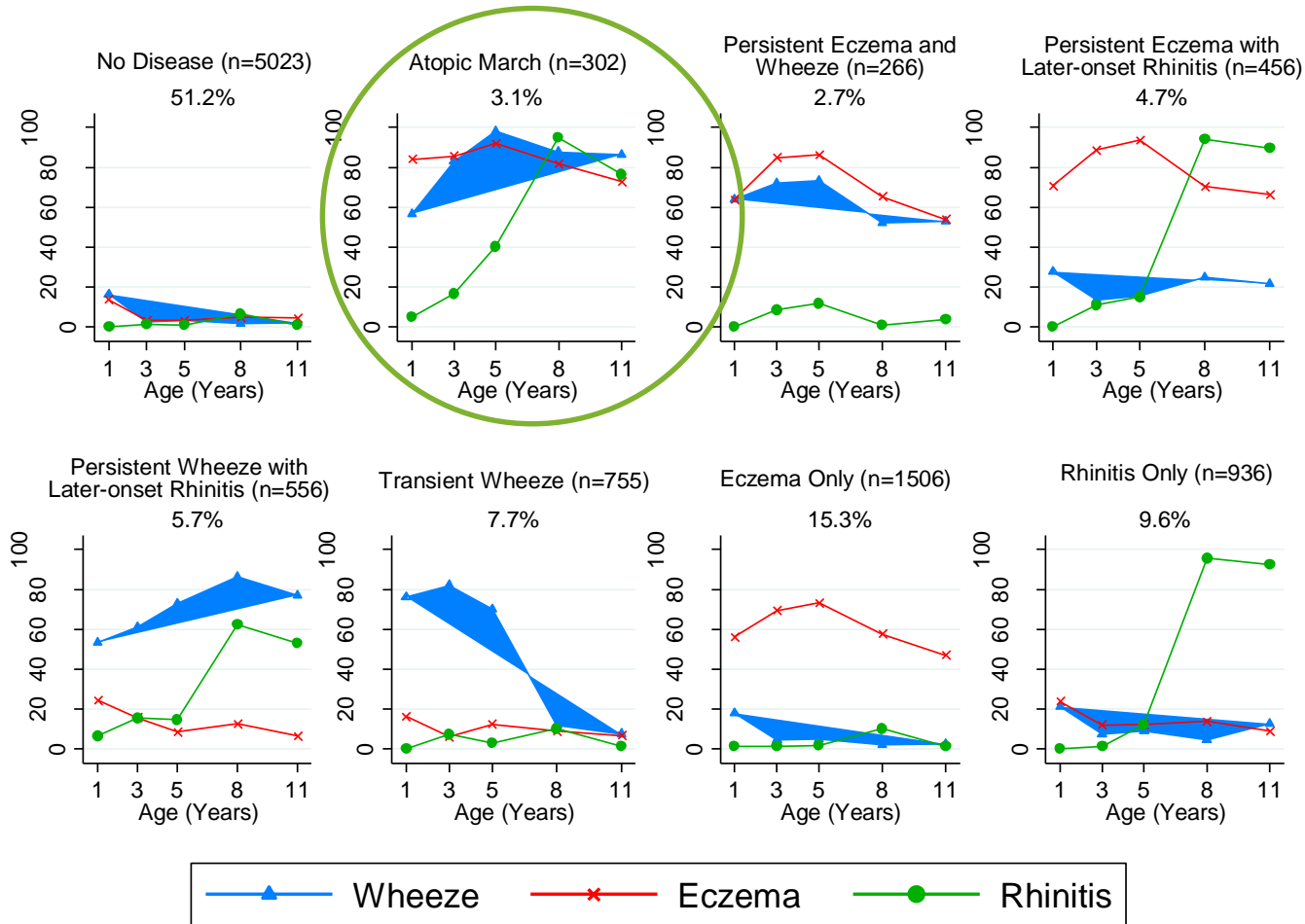
From: World Allergy Organization, 2014



Individual-level Longitudinal Analysis



Myth Bust by Learning from Data

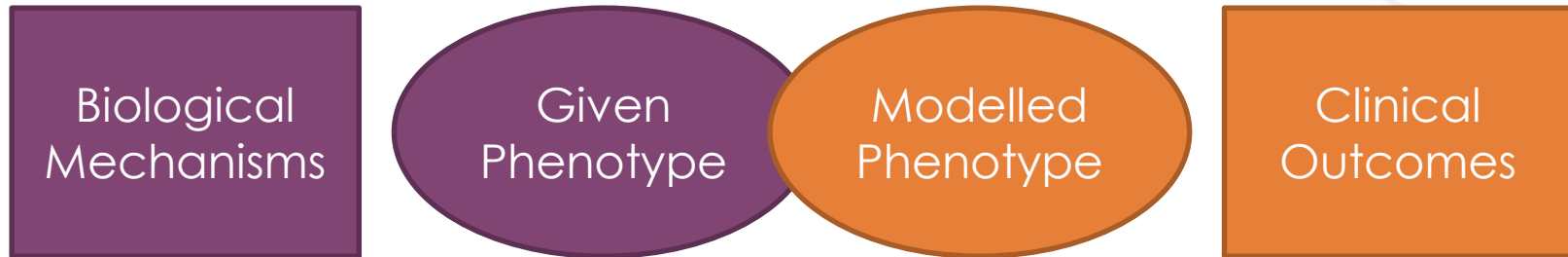


MRC STELAR consortium working at scale across MAAS and ALSPACS cohorts

From: Belgrave et al. Developmental Profiles of Eczema, Wheeze, and Rhinitis: Two Population-Based Birth Cohort Studies. PlosMedicine 2014



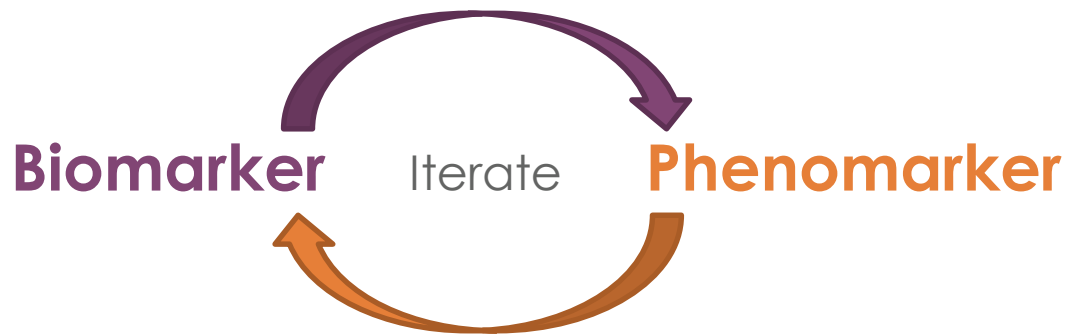
Joined-up Bio-Health Modelling



Bioinformatics

Health Informatics

Integrative Informatics needed



Portable Model
(not a fixed label/rule)

Real World Phenomarking

Total evidence-base predicts < 30% outcomes

?how to generate real world evidence on gliflozins

His GP...

Lifestyle factors:
diet, exercise

His nephrologist...

BP control

Primary Care

Renal Medicine

Diabetology



3 x evidence pipelines
1 x complex patient

Mr Jones...



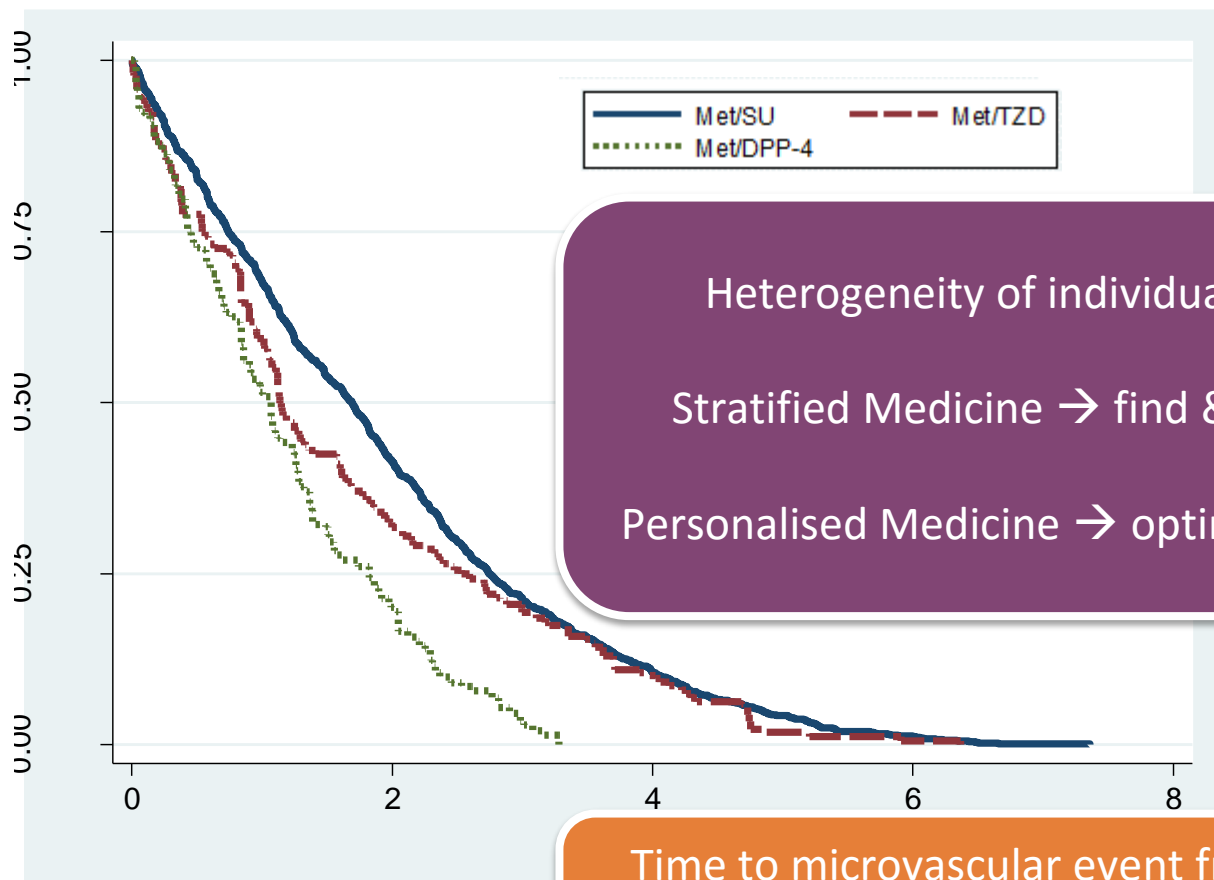
His diabetologist...

Glucose control

↑ Weight
→ ↑ BP

Valderas JM, Starfield B, Sibbald B, Salisbury C, Roland M. Defining comorbidity: implications for understanding health and health services. *Ann. Fam. Med.* 2009;7:357-363.

Dual Therapy for Diabetes



Heterogeneity of individual treatment response
Stratified Medicine → find & treat more subgroups
Personalised Medicine → optimise individual's response

Time to microvascular event from diagnosis of diabetes
Inverse probability weighted marginal structural model
Average causal effects of dual therapies

Data: Does Size Matter?

DATA



Vast volume,
velocity, variety...

TSUNAMI

METHODS & OUTPUTS



Supra-linear growth
in papers & tools

BLIZZARD

EXPERTISE



Similar number of
analysts

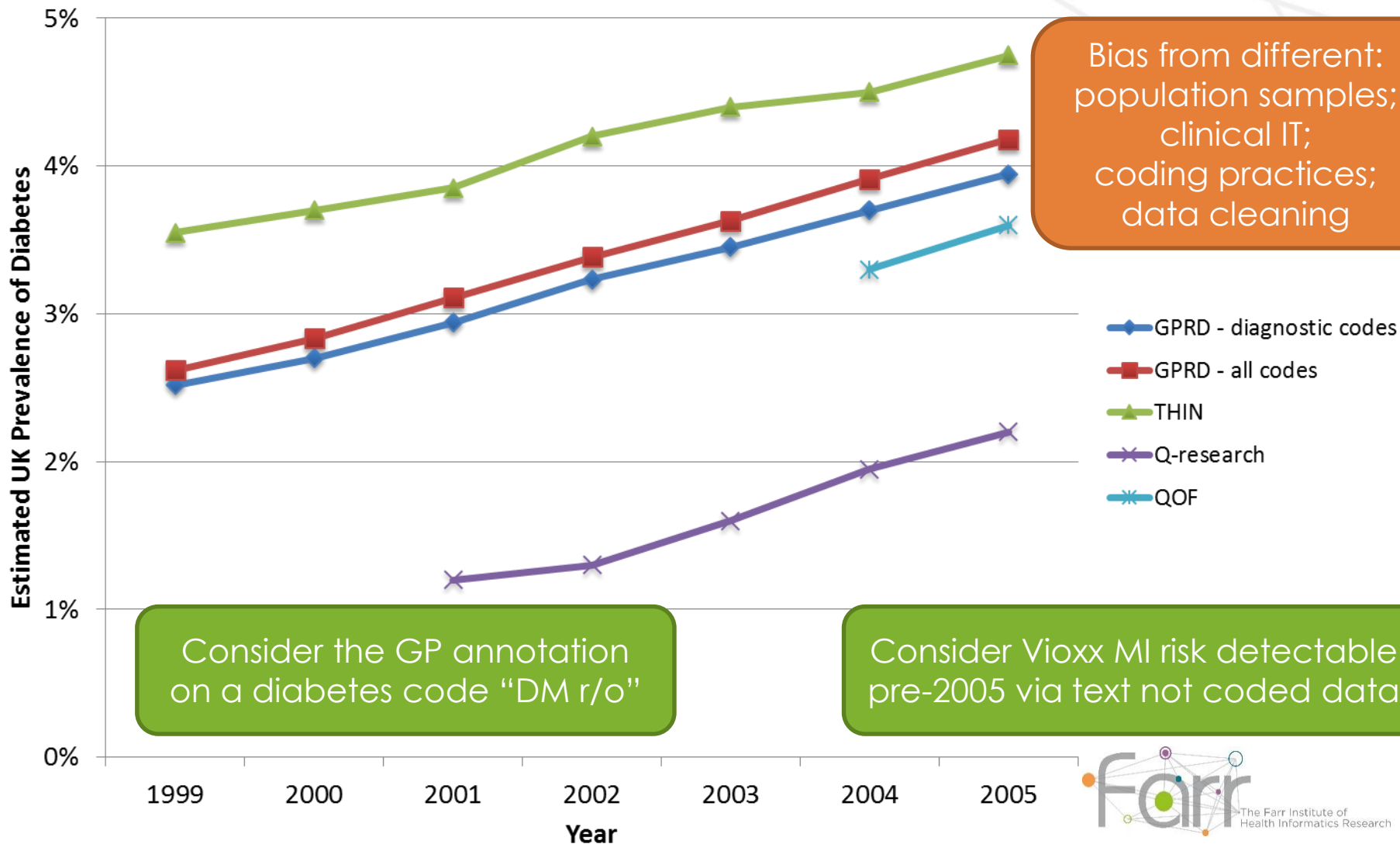
DROUGHT

More data * small-scale research = more small-scale research
>n with >heterogeneity can *reduce* 'power'

Ioannidis JPA. Why most published research findings are false.
PLoS Med. 2005 Aug;2 (8):e124.

Overhage JM, Ryan PB, Schuemie MJ, Stang PE. Desideratum for
evidence based epidemiology. Drug Saf. 2013 Oct;36 Suppl 1:S5-14.

Beyond In-licensed Data

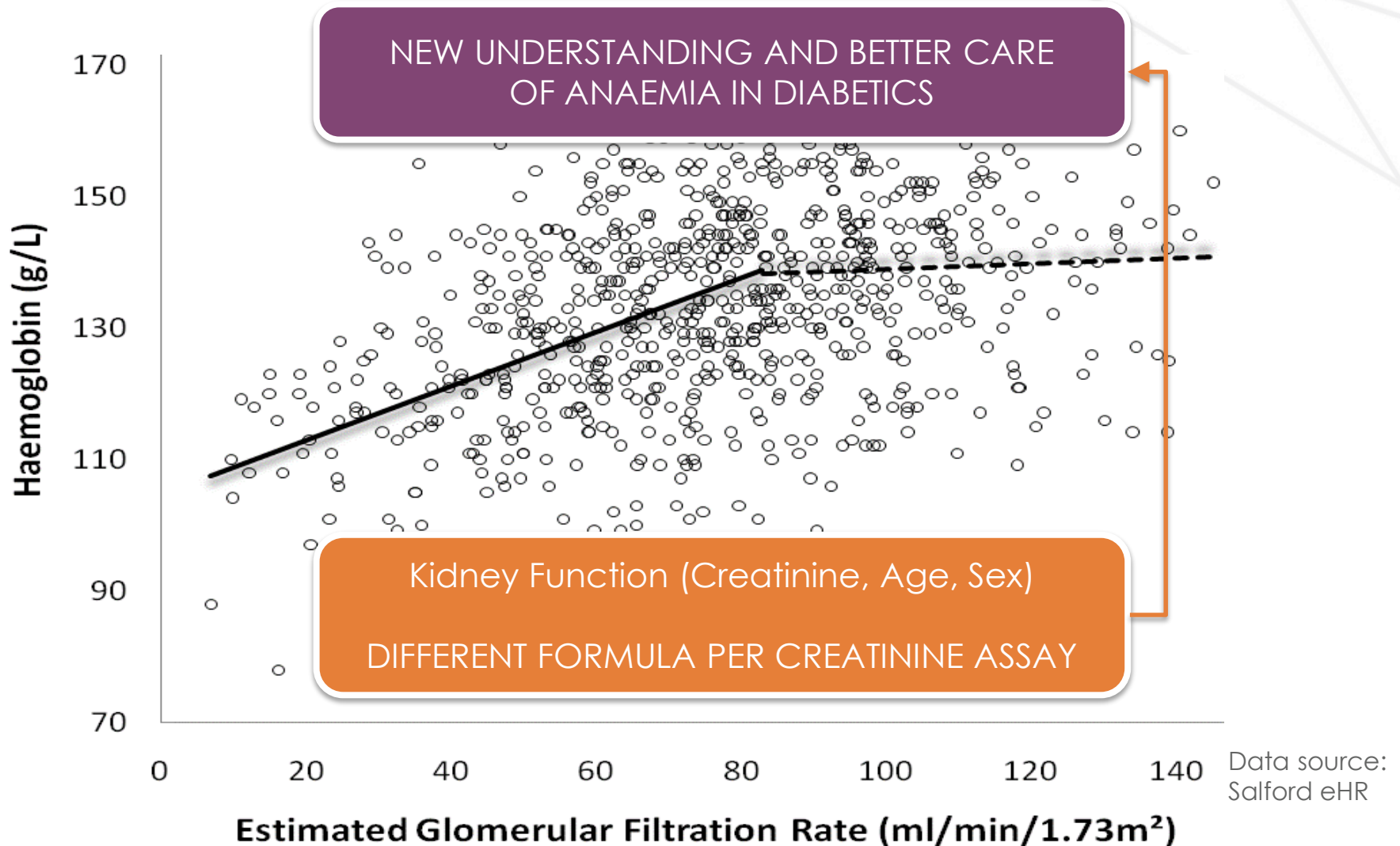


Bias from different:
population samples;
clinical IT;
coding practices;
data cleaning

Consider the GP annotation
on a diabetes code "DM r/o"

Consider Vioxx MI risk detectable
pre-2005 via text not coded data

Key Local Metadata



Farr @ Health e-Research Centre



Missed Opportunities Detector

Identify patients with target disease

Example
CKD

Exclude if quality standard inappropriate

Example
Terminal illness

Exclude if quality standard achieved

Example
BP target

Identify how care could be improved

Example
Rx optimization

Salford Clinical Commissioning Group

Salford Integrated Record  234k popln.

Web Interface

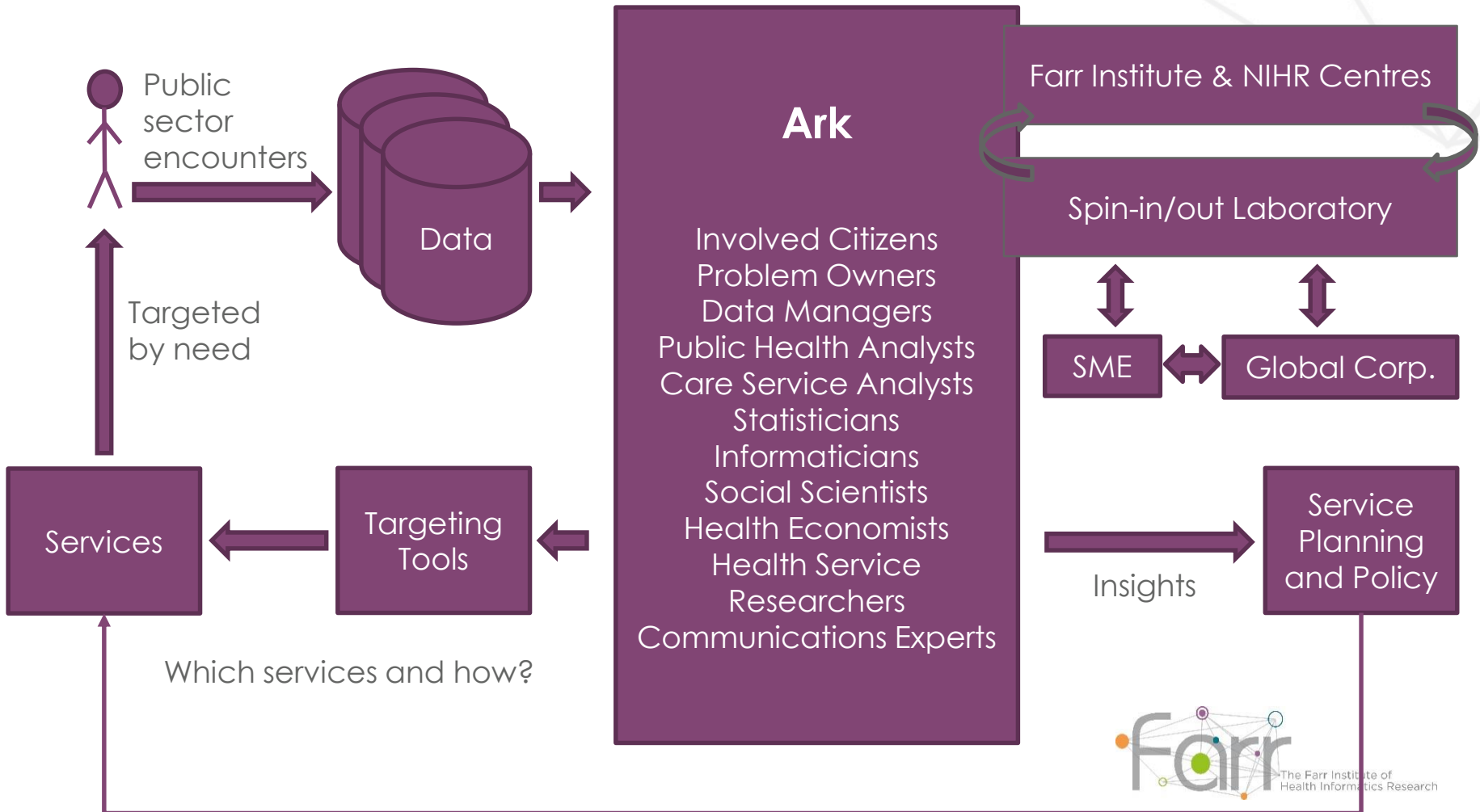
Practice-level Audit + ? Patient-level Decisions

53 GP Practices + 1 Hospital

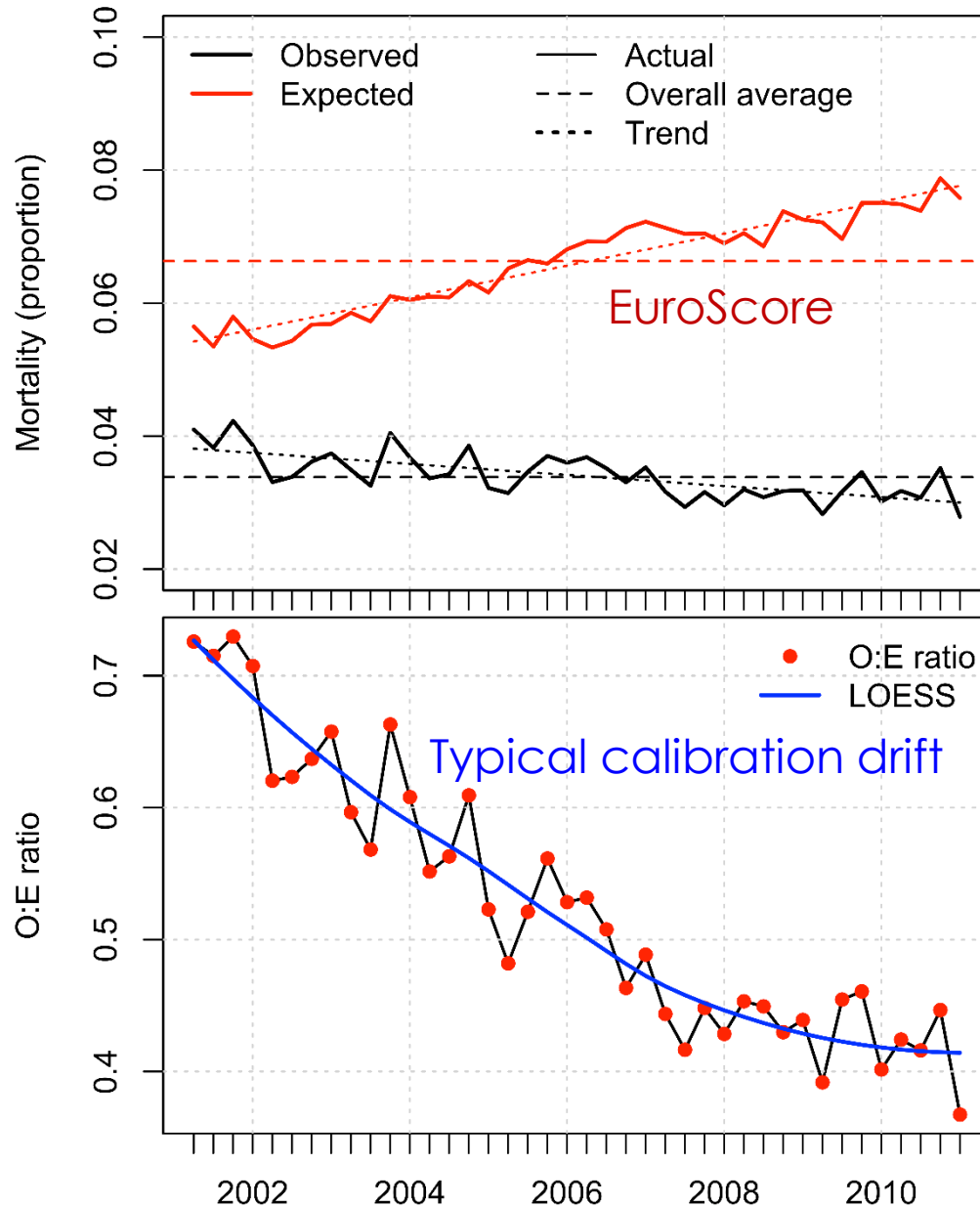
EHR EHR EHR EHR

Care Professionals

Connected Health Cities



Clinical Outcome Prediction

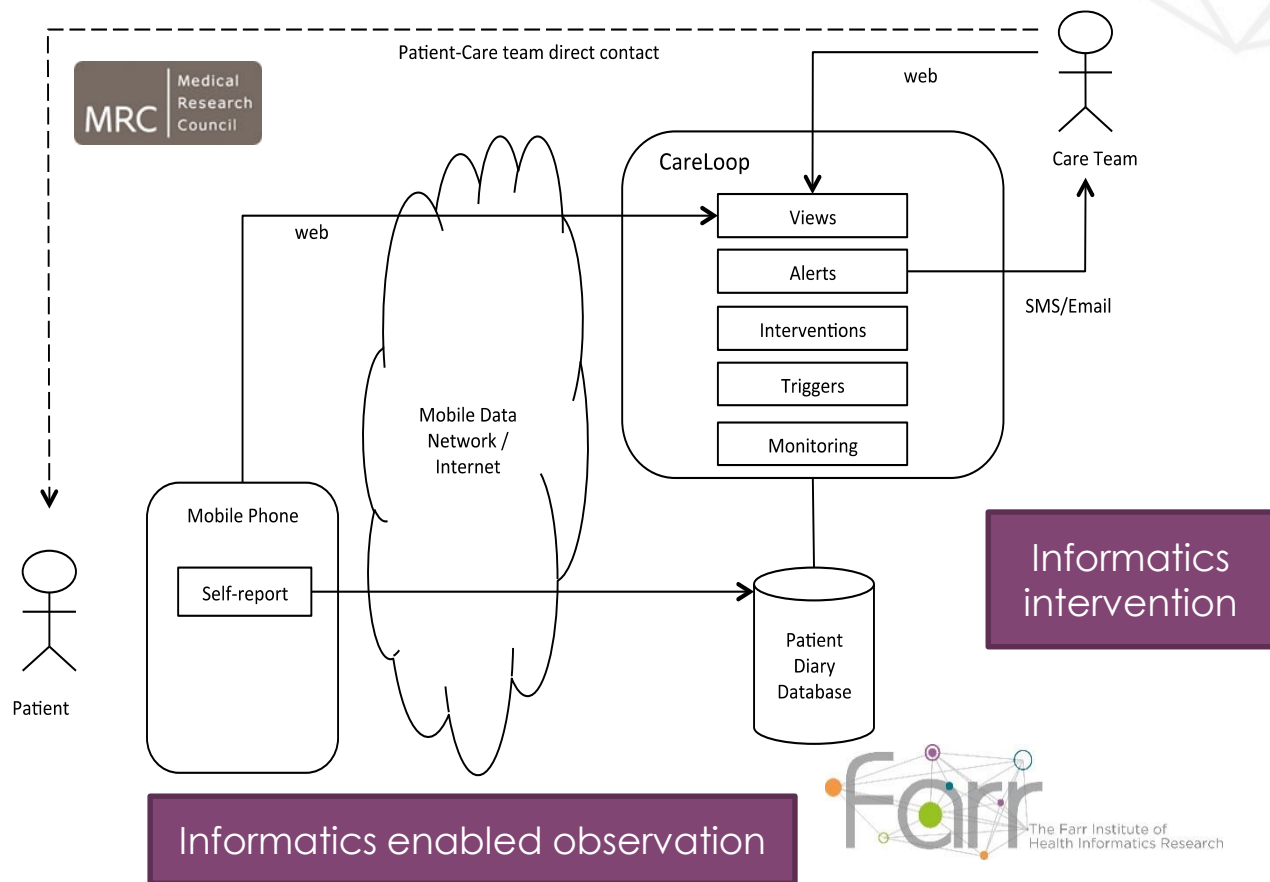
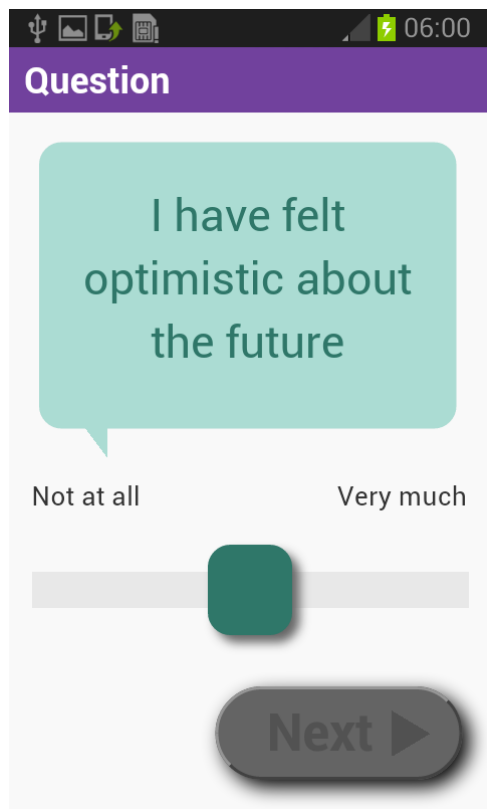


Production line of clinical prediction models is broken

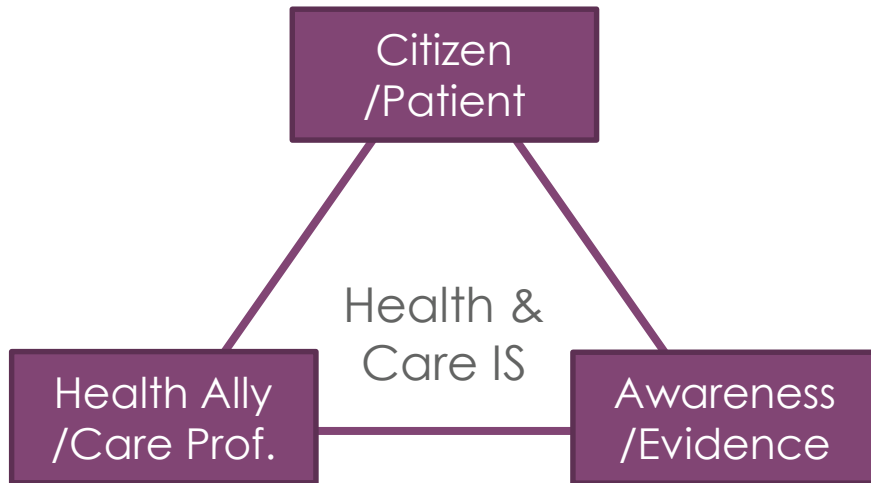
EU Directive 2007/47:
The law now sees algorithms as medical devices

Health System Extension: Mobile

Aim: To Reduce Relapse in Schizophrenia via Smartphone
Drug + behaviour (information * psychological endotype) = outcome



Digital Health Triangle



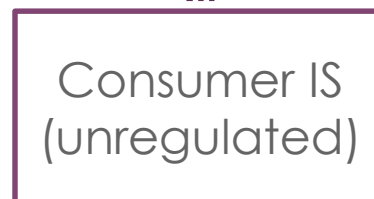
Digital Health Economy (20**?)

Ubiquitous technologies
Unifying models
Usable interfaces (avatar etc.)
Stratification ↔ personalisation
Actionable micro-evidence
Service ↔ research

Records
Guidelines

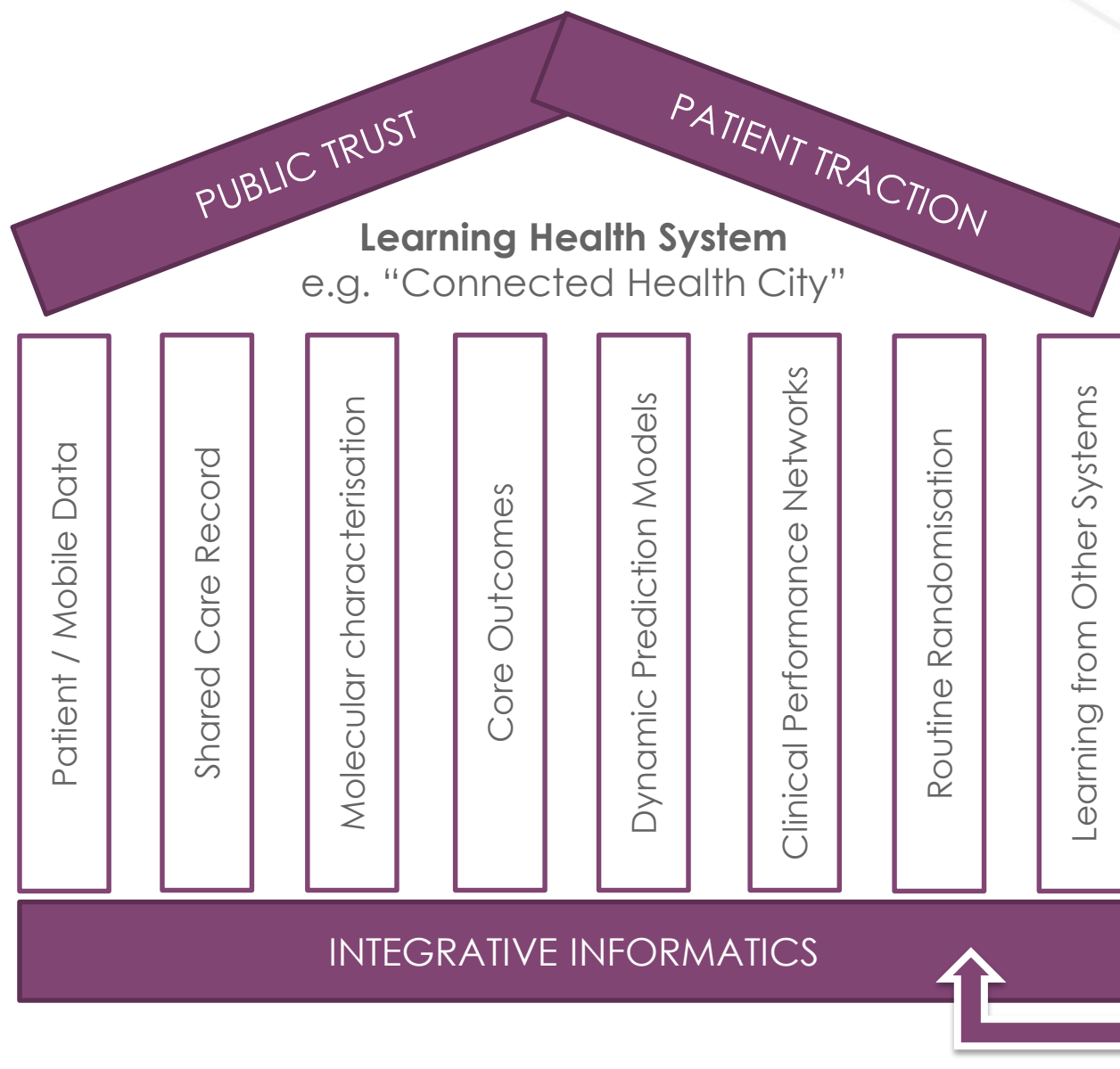


Two Worlds of Healthcare (2015)



Apps & sensors
Healthcare information
Social networks

Depth, Trust and Scale



From care + data
→ research
→ translation



To a federation of research in care:
pulled & shaped by local communities,
with academic and industry partners



Sub-disease Research Tips

1. Target

- Plausible diagnostic **aggregation**
- Unexplained **variation** in clinical outcomes

2. Data

- Multiple populations/settings (**heterogeneity & replication**)
- Useful **temporal** structure

3. Analytics

- Multi-perspective ML **pre-model framing** – don't rush in!
- Heuristic **phenomarker-biomarker** resolution