

# Methods and tools for co-analysing biobanks

# **Paul Burton**

Professor of Genetic Epidemiology Departments of Health Sciences and Genetics University of Leicester







wellcometrust





SHIP Conference St Andrews, Sept 11th 2011









# DataSHIELD and DataSHaPER



Harmonized analysis of individual-

PROMOTING HARMONISATION OF PHOEBE EPIDEMIOLOGICAL BIOBANKS IN EUROPE



# Combined analysis: why bother?

THEORY AND METHODS

## Size matters: just how big is BIG?

# Quantifying realistic sample size requirements for human genome epidemiology

Paul R Burton,  $^{1,2,3_*,\dagger}$  Anna L Hansell,  $^{4,\dagger}$  Isabel Fortier,  $^{3,5}$  Teri A Manolio,  $^6$  Muin J Khoury,  $^{3,7}$  Julian Little  $^{3,8}$  and Paul Elliott  $^4$ 

- Sample size
- Depth of phenotyping
- Quality of measurement

**ALL critical** 

International Journal of Epidemiology 2009;38:263–273 doi:10.1093/ije/dyn147

# How big is BIG? ("typically" and "realistically")

- The direct effect of a gene
  - 2,000 cases minimum, 10,000 cases better
- Environmental and life-style factors
  - Highly context specific
- Gene-lifestyle and gene-gene "interactions"
  - Absolute minimum 10,000, usually need at least 20,000, a comprehensive platform needs at least 50,000
  - Scientifically fundamental

# The bottom line

- Effective access to well-measured data is vital
- Efficient valid analysis is crucial
- Effective <u>combined</u> analysis is also important
- Two fundamental challenges
  - Streamlined ethico-legally acceptable access to multiple data sets
  - Scientific harmonization
- New systems, methods and tools
  - DataSHIELD, DataSHaPER

# Combined analysis in practice

# Seven steps to combined analysis

- Define scientific question
- Identify information required
  - What is needed?
  - On how many people?
- Use catalogues to identify studies/cohorts/biobanks or biorepositories that can service the informational needs
- Develop "DataSchema" and "formal pairing rules"
- Create algorithms to generate required data in each study
- Generate the harmonized data sets
- Undertake an *appropriate* joint analysis (science and ELSI)

SHaPFF

Observatory

DataSHIELD

## http://www.p3gobservatory.org/



- Study catalogue
  - 162 studies
  - More than 12.5m participants
- Questionnaire catalogue
  - 76 questionnaires (cut and key-worded)
  - From 48 studies

## http://www.datashaper.org/ Dr Isabel Fortier – University of McGill

#### WELCOME TO DATASHAPER WEBSITE

The DataSHaPER (Data Schema and Harmonization Platform for Epidemiological Research) is both a scientific approach and a suite of practical tools. Its primary aims are to facilitate the prospective harmonization of emerging biobanks, provide a template for retrospective synthesis and support the development of questionnaires and information-collection devices, even when pooling of data with other biobanks is not foreseen.

#### • What is the DataSHaPER?

The development of these tools has been jointly funded by P<sup>s</sup>G, CPT **2**, PHOEBE **2**, and Generation Scotland **2**.











A DataSchema identifies and describes a thematic set of core variables that are of particular value in a specified scientific setting. It also contains associated support material including variable definitions, links to relevant Ontologies and Classifications, and access to reference questionnaires and operating procedures.



#### DATASCHEMA PLATFORM

A DataSchema contains a thematic set of core variables that are of significant relevance to the specific scientific area addressed by a specific study or group of studies.

The core variables in each DataSchema are grouped a four level nested hierarchy:

1. Module:

Assessment modes or type of element measured or collected. Each module subsumes one or more themes.

📲 2. Theme:

General area of interest. Each theme subsumes one or more domains.

🏐 3. Domain:

Risk factor or outcome of interest. Each domain subsumes a number of variables.

4. Variable:

Primary unit of interest for a statistical analysis.

Each DataSchema contains, where relevant, descriptive information to complements the list of variables itself, such as variable definitions, links to relevant standard classifications, reference questionnaires, operating procedures and so on.









#### DataSchemas









#### Variable informations

Name : Occurrence of high blood pressure

Description : Occurrence of high blood pressure at any point during the life of the participant.

URI: http://www.datashaper.org/owl/2009/10/generic.owl#GENERIC\_326

≡ Categories		
Туре	Name	С
i⊒ Category	Never had high blood pressure	
i⊒ Category	Ever had high blood pressure	
🔁 Missing category	Prefer not to answer	
🔁 Missing category	Don't know	



The Harmonization Platform provides a template for the formal estimat ion of the potential to synthe size information between studies. At present, access to the Harmonization Platform is limited to collaborative context.



Harmonization Platform

#### HARMONIZATION PLATFORM

The process follows a rigorous approach including:

- 1. The development of rules providing a formal assessment of the potential for each individual study to generate each of the variables in a given DataSchema.
- 2. The application of these rules to determine and tabulate the ability of each study to generate each variable, thereby identifying the information that can be shared.

	Study A	Study B	Study C	Study D
Variable 1	IMPOSSIBLE	PARTIAL	COMPLETE	COMPLETE
Variable 2	COMPLETE	COMPLETE	COMPLETE	COMPLETE
Variable 3	PARTIAL	IMPOSSIBLE	IMPOSSIBLE	IMPOSSIBLE

- 3. The development and application of a processing algorithm enabling that study to generate the required variable in an appropriate form.
- 4. Access to the Harmonization Platform is limited to collaborative context. Please contact us to see how we can work together.



### Evaluating what can be harmonized

	KORA-gen	LifeGene	UK Biobank
Variable	Match	Match	Match
Occurrence of diabetes	Complete	Complete	Complete
Current treatment for high blood pressure	Complete	Impossible	Complete
Occurrence of diabetes in the family	Impossible	Impossible	Partial
Occurrence of high blood pressure in the family	Partial	Impossible	Partial
Current cigarette smoker	Complete	Partial	Complete
Current quantity of cigarettes smoked	Complete	Complete	Complete
Household income	Partial	Impossible	Partial
Country of birth-ISO-3166	Complete	Complete	Complete
Gender	Complete	Complete	Complete
Age at recruitment	Complete	Complete	Complete
Standing height	Complete	Complete	Complete
Weight	Complete	Complete	Complete
Waist circumference	Complete	Complete	Complete
Hip circumference	Complete	Complete	Complete
Body mass index	Complete	Complete	Complete
Diastolic blood pressure at rest	Complete	Complete	Complete
Systolic blood pressure at rest	Complete	Complete	Complete
Date of interview	Complete	Complete	Complete



Cccurrence of high blood pressure
🖃 Target
Participant
Disease
High blood pressure
Hypertension
Source of information
Participant
🖃 Period
Through all life
🖃 Format
Categorical
Collection mode
Questionnaire
🖃 Class
Occurrence
😹 Ontology references
URI
http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#Hypertension

## Quality, quantity and harmony: the DataSHaPER approach to integrating data across bioclinical studies

Isabel Fortier,<sup>1,2</sup>\* Paul R Burton,<sup>1,3</sup> Paula J Robson,<sup>4</sup> Vincent Ferretti,<sup>5</sup> Julian Little,<sup>1,6</sup> Francois L'Heureux,<sup>1</sup> Mylène Deschênes,<sup>1</sup> Bartha M Knoppers,<sup>1,7</sup> Dany Doiron,<sup>1</sup> Joost C Keers,<sup>8</sup> Pamela Linksted,<sup>9</sup> Jennifer R Harris,<sup>10</sup> Geneviève Lachance,<sup>1</sup> Catherine Boileau,<sup>11</sup> Nancy L Pedersen,<sup>12</sup> Carol M Hamilton,<sup>13</sup> Kristian Hveem,<sup>14</sup> Marilyn J Borugian,<sup>15,16</sup> Richard P Gallagher,<sup>15,16</sup> John McLaughlin,<sup>17</sup> Louise Parker,<sup>18</sup> John D Potter,<sup>19</sup> John Gallacher,<sup>20</sup> Rudolf Kaaks,<sup>21</sup> Bette Liu,<sup>22</sup> Tim Sprosen,<sup>23</sup> Anne Vilain,<sup>1</sup> Susan A Atkinson,<sup>3</sup> Andrea Rengifo,<sup>3</sup> Robin Morton,<sup>9</sup> Andres Metspalu,<sup>24</sup> H Erich Wichmann,<sup>25,26,27</sup> Mark Tremblay,<sup>28,29</sup> Rex L Chisholm,<sup>30</sup> Andrés Garcia-Montero,<sup>31</sup> Hans Hillege,<sup>32</sup> Jan-Eric Litton,<sup>33</sup> Lyle J Palmer,<sup>34</sup> Markus Perola,<sup>35,36</sup> Bruce HR Wolffenbuttel,<sup>37</sup> Leena Peltonen<sup>38</sup> and Thomas J Hudson<sup>5,39,40</sup>



International Journal of Epidemiology 2010;1–11 doi:10.1093/ije/dyq139

## Is rigorous retrospective harmonization possible? Application of the DataSHaPER approach across 53 large studies

Isabel Fortier,<sup>1,2</sup>\* Dany Doiron,<sup>2</sup> Julian Little,<sup>1,3</sup> Vincent Ferretti,<sup>4</sup> François L'Heureux,<sup>2</sup> Ronald P Stolk,<sup>5</sup> Bartha M Knoppers,<sup>1,6</sup> Thomas J Hudson,<sup>4,7,8</sup> and Paul R Burton<sup>2,9,10</sup> on behalf of the International Harmonization Initiative<sup>†</sup>



International Journal of Epidemiology 2011;1–15 doi:10.1093/ije/dyr106 
 Table 1
 Pairing results (%) for selected variables presenting

 a high proportion of 'Complete/Partial Proximate' matches\*

Variable name	Complete/Partial proximate (%)
Occurrence of diabetes	89
Current use of alcohol	89
Level of physical activity	85
Occurrence of high blood pressure	83
Occurrence of menopause	83
Occurrence of cancer	81
Occurrence of stroke	81
Employment status	81
Menopause onset	75
Occurrence of asthma	74
Current quantity of cigarettes smoked	74
Occurrence of myocardial infarction	72
Living with partner	68
Type of cancer	66 Table
Standing height	66
Weight	66



Table 3 Univariate and multivariate models; variable characteristics and pairing results

\*For full-pairing results, please refer to supplementary ials (Supplementary data available at *IJE* online).

1200

		the percentage of variables presenting as:			
	Number of variables	Complete match	Partial Proximate match	Partial Tentative match	Impossible match
Univariate analysis					
Variable importance					
Essential (1)	38	60	2	15	23
Important (2)	45	38	5	13	44
Useful (3)	65	21	1	14	64
Targeted individual					
Participant (0)	101	46	3	12	39
Participant's family members (1)	47	14	1	19	66
Targeted period					
Current status (0)	40	50	4	11	36
All other periods (1)	108	31	2	15	52

Average across studies of



PROMOTING HARMONISATION OF **PHOEBE** EPIDEMIOLOGICAL BIOBANKS IN EUROPE DataSHIELD: full joint analysis without physically sharing the data

With: Philippe Laflamme Vincent Ferretti

GenomeCanada









wellcometrust

## Horizontally partitioned data



# A real and growing problem

- How can we undertake a joint analysis using multiple data sources if the data cannot physically be pooled?
  - Ethico-legal constraints
  - Physical size of the data objects
  - Intellectual property issues

#### cell meteroscience analysis an

PERSPECTIVE

## **On the Future of Genomic Data**

Scott D. Kahn



## DataSHIELD: resolving a conflict in contemporary bioscience—performing a pooled analysis of individual-level data without sharing the data

Michael Wolfson,<sup>1</sup> Susan E Wallace,<sup>2,3</sup> Nicholas Masca,<sup>4</sup> Geoff Rowe,<sup>1</sup> Nuala A Sheehan,<sup>4</sup> Vincent Ferretti,<sup>3,5</sup> Philippe LaFlamme,<sup>3,6</sup> Martin D Tobin,<sup>4</sup> John Macleod,<sup>7</sup> Julian Little,<sup>3,8</sup> Isabel Fortier,<sup>3,8,9</sup> Bartha M Knoppers<sup>2,3</sup> and Paul R Burton<sup>3,4,8,10</sup>\*

International Journal of Epidemiology 2010;39:1372–1382







# Two approaches to data synthesis

- Study level meta-analysis (SLMA)
  - Estimate a summary statistic from each study. Calculate an appropriately weighted mean and standard error across all studies
  - "Conventional meta-analysis"
- Individual level meta-analysis (ILMA)
  - Pool all of the individual level data from each of the studies into one large data set and then analyse that data set as if it was one single study (with parameters for heterogeneity)
  - "Direct pooling"

# Study level meta-analysis

- If summary statistics available, quick, cheap, and convenient (e.g. from published literature)
  - RCTs, Health Care Assessment etc
- If the required summary statistics are straightforward to define and easy to calculate, also very convenient
  - Genome Wide Association Studies (GWASs)

# Study level meta-analysis

- Quick, easy and it works
- But SERIOUS lack of flexibility for example:
  - One million SNPs on a GWA chip are successfully analysed
  - But, then you want to study interaction of all apparently associated SNPs with age and sex
  - Impossible unless these analytic results provided up-front
- Contemporary bioscience is getting more complex
- Exploratory analysis needs flexibility

ILMA (direct data pooling) therefore preferable

# **ELSI restrictions on sharing**

## Exemplar wording

Wallace S, Lazor S, Knoppers BM. Chapter in Kaye J and Stranger M. Principles and Practice in Biobank Governance. Ashgate, Farnham 2009

- Use of data restricted to researchers participating in the original study
- Use of data restricted to researchers in one country
- The need to obtain ethico-legal and scientific permission to access the data
  - Often a protracted and time consuming process
  - Often needs multiple clearances

# **Competing public goods**

- Analytic flexibility greatly favours ILMA
- ELSI can prohibit or discourage ILMA
  - $\rightarrow$  Most current GWASs based on SLMA
  - BUT: this situation is not sustainable as things become more complex, unpredictable and exploratory

# **DataSHIELD: a novel solution**



Use <u>parallel processing</u> to fit equivalent models on each study simultaneously

Parallel analyses linked together by transmitting entirely <u>non-identifying</u> <u>summary statistics</u>

In many settings this produces mathematically <u>identical results</u> to fitting a single model to all the data held in one pooled data set

# **DataSHIELD: a novel solution**



#### Analysis commands (1)

b.vector<-c(0,0,0,0)

glm(cc~1+sex+snp+bmi, family=binomial, start=b.vector, maxit=1)


#### **Summary Statistics (1)**

#### Score vector Study 5

[36, 487.2951, 487.2951, 149]

#### Information Matrix Study 5

500	70.56657	70.56657	297
70.56657	7646.29164	7646.29164	65.39412
70.56657	7646.29164	7646.29164	65.39412
297	65.39412	65.39412	382



#### **Summary Statistics (1)**

#### Score vector Study 5

[36, 487.2951, 487.2951, 149]

#### Information Matrix Study 5

500	70.56657	70.56657	297
70.56657	7646.29164	7646.29164	65.39412
70.56657	7646.29164	7646.29164	65.39412
297	65.39412	65.39412	382



Analysis commands (2)

b.vector<c(-0.322, 0.0223, 0.0391, 0.535)

glm(cc~1+sex+snp+bmi, family=binomial, start=b.vector, maxit=1)



#### Summary Statistics (2)

#### **Score vectors**

#### **Information Matrices**



#### Summary Statistics (2)

#### **Score vectors**

#### **Information Matrices**

### and so on .....

11110

11-12



#### **Summary Statistics (4)**



#### **Updated parameters (4)**



#### **Updated parameters (4)**

#### **Final parameter estimates**

Coefficient	Estimate	Std Error
Intercept	-0.3296	0.02838
BMI	0.02300	0.00621
BMI.456	0.04126	0.01140
SNP	0.5517	0.03295

## **Conventional analysis**

#### **Coefficients:**

	Estimate	Std. Error	
(Intercept)	-0.32956	0.02838	
BMI	0.02300	0.00621	
BMI.456	0.04126	0.01140	
SNP	0.55173	0.03295	

### DataSHIELD analysis

### Does it work?

Parameter	Coefficient	Standard Error
<b>b</b> <sub>intercept</sub>	-0.3296	0.02838
<b>b</b> <sub>BMI</sub>	0.02300	0.00621
<b>b</b> <sub>BMI.456</sub>	0.04126	0.01140
<b>b</b> <sub>SNP</sub>	0.5517	0.03295



# Current Implementation





# Vertically partitioned data





# **Safeguards**

- Ethical approval for all contributing studies
- Confidentiality agreements for all
- Record *all* information passed to and from each DC *at the DC*
- No new models until summary statistics fully understood
  - At present GLMs OK, random effects no!!
- Harmonization essential

# Safeguards

- For the future a software wrapper
  - Monitor, scrutinise and interpret all incoming and outgoing information flows and block and identify/record any request, or series of requests, that might – either by accident or by design – generate potentially disclosive information

### For the moment

• Suggest more than one statistician (from more than one centre) closely follows the analysis

# Current state of play

- Main funding, BioSHARE-eu and ALSPAC
- DataSHIELD implemented in R
  - Refining at present
- R based implementation inserted into OPAL
  - Part of OBiBa suite
- First workshop held in Lyon, April 2011
- Analytic validity confirmed on real IPRI data
- Three pilot studies moving forward

# Thank you for listening

# **Combined analysis** in principle





**GEN2PHEN** 

- Workable access
  - Integrative IT platforms
  - Working within ELSI guidelines and edicts particularly across jurisdictions
- Scientific compatibility
  - Prospective v Retrospective
  - Stringent v Flexible









Numb	er of S	tudies	5									162			
With	Sumn	hary in	nform	ation	(to be	e valio	dated)	)				70			
With	Comp	lete in	nform	ation	(valid	ated I	oy stu	dy inv	vestig	ators)		92	-		
AU 0.0	٨	B	С	D	E	F	G	н	I	J	к	L	М	N	(



### Catalogue current content

Number of studies	48
Number of questionnaires	76

//www.p3gobservatory.org/questionnaireblock/viewBlock.htm?questionnaireId=6&blockId=77 - Windows Internet Explorer

http://www.p3gobservatory.org/questionnaireblock/viewBlock.htm?questionnaireId=6&blockId=77

1120

#### UK Biobank: a large-scale prospective epidemiological resource [UKB] / UK Biobank Touchscreen Questionnaire / Block 77

L3	Has a doctor ever told you that you have diabetes?	Select one from
		- YE Yes
		- NO No
		- UN Do not know
		- DA Prefer not to answer
L3A	Did you only have diabetes during pregnancy?	Select one from
		- YE Yes
		- NO No
		- NA Not applicable
		- UN Do not know
		- DA Prefer not to answer
L3B	What was your age when the diabetes was first diagnosed?	Enter number
		OR
		UN Do not know
		OR
		DA Prefer not to answer
T 3 C	Did you start inculin within one year of your diagnosis of diabetes?	Select one from

# **Combined analysis** in principle





**GEN2PHEN** 

- Workable access
  - Integrative IT platforms
  - Working within ELSI guidelines and edicts particularly across jurisdictions
- Scientific compatibility
  - Prospective v Retrospective
  - Stringent v Flexible









# Context specific!!

Identify subjects with a doctor diagnosis of diabetes for an aetiological study



### How long is 'LONG'?



Disease





### DataSHaPER website

#### \_ @ 🗙 ContaSHaPER - Windows Internet Explorer pa 😽 🗙 🛛 Google $\Theta \Theta$ http://www.datashaper.org/ 2-🔵 paul.g... + 🔌 + 🚓 -Google 🚼 Search 🔹 🖓 More ≫ 📥 🝷 🔂 Page 👻 🙆 Tools 🙆 DataSHaPER Data ØG **SHaPER** Observatory DataShape Secretariat Home Contact Us **Terms And Conditions**

#### WELCOME TO DATASHAPER WEBSITE

The DataSHaPER (Data Schema and Harmonization Platform for Epidemiological Research) is both a scientific approach and a suite of practical tools. Its primary aims are to facilitate the prospective harmonization of emerging biobanks, provide a template for retrospective synthesis and support the development of questionnaires and information-collection devices, even when pooling of data with other biobanks is not foreseen.

• What is the DataSHaPER?

A

The development of these tools has been jointly funded by PSG, CPT , PHOEBE , and Generation Scotland



Doc...

🙆 C:\...

🐼 Cal...

🖂 RE:...



A DataSchema identifies and describes a thematic set of core variables that are of particular value in a specified scientific setting.

It also contains associated support material including variable definitions, links to relevant Ontologies and Classifications, and access to reference questionnaires and operating procedures.



🔞 🐼 👎 📓 👩 🏉 🭕

🚰 start

The Harmonization Platform provides a template for the formal estimation of the potential to synthetized information between studies. At present, access to the Harmonization Platform is limited to collaborative context.

€ P³G...

🖾 Doc...

🤮 Ado...

🔏 htt...

💽 Pau...

🖭 Pre...

#### What's New ?

- The Renal DataSchema is now online. This DataSchema addresses Chronic Kidney Disease (CKD) and was developed in collaboration with GENECURE<sup>GP</sup>;
- A new search tool that lets you track "elements" across different DataSchemas is now available:
- When applicable, we have linked some DataSchemas elements (themes/domains/variables) to PhenX<sup>dP</sup> elements;
- The new "What is the Data SHaPER" page contains technical documentation with helpful descriptions of the Data SHaPER.

Visit the P<sup>s</sup>G Observatory<sup>GP</sup> website to see recent additions and improvements.

🔏 Dat...

21:18

Ø

2 P<sup>3</sup>G Observatory - LifeGene

¥

(PG »	
<mark>Study</mark>	
Catalogue	

ارژه http://www.p3gobservatory.org/	catalogue.htm?studyId=21		🛱 🗸 🔀 🖌 Google	ABP
÷				```
			Legend 🖌 Yes 🗶 No	? Unknown
General design				
Study design	Cohort			
Type of participants	Families			
Target or final number of participants	500000			
Target or final number of families	No information available			
Target or final number of DNA	No information available			
Participant selection / Characteristics	of the population			
Selection criteria				
🗸 Age	Maximun: 55 years			
Country of residence	Sweden			
Recruitment procedures Initially, enrollment of subjects will subjects through contacts via mas	be performed by use of national d s e-mails, advertisements, for enri	atabases such as Swedis chment of the cohort with	sh Twin Registry, followed by enro a under-studied ethnic populations	ollment of a, etc.
Data Sources				
		Cross-sectional	Repeated/continuous	
Questionnaires to participants/	respondents	Х	1	
<ul> <li>Direct physical measures</li> </ul>		х	1	
Biological samples		1	х	
Medical paper				
<ul> <li>Electronic databases</li> </ul>				
X Genealogical records				
Collection procedures				

Participants will be requested to fill out web-based questionnaires and answer questions by use of cell-phones etc. Data on exposures and phenotypes will be collected repeatedly over time. Linkage to comprehensive registers with health care information will be performed yearly. Blood sampling for DNA extraction and analyses of key biomarkers will be performed according to standardized procedures. Data collection will be for research purpose only. All data will be obtained in coded format. Subject identification data registered in conjunction to biological samples will be stored in such a way that unauthorized persons will have no access and be able to link data and/or samples to the donors. Information transmitted by electronic means will be encrypted. Subject data transferred to a third country (e.g. outside EU) must be de-identified. If data cannot be de-identified before transfer, data security issues must be agreed upon in writing in advance.

#### Follow up procedures

No information available

#### Baseline principal variables of interest

#### Health information

- Diseases history ICD-10 All conditions
- Familial disease history
- Early life
- Women's health
- Quality of life

#### Physical / Biochemical measures

- Body structure measures
- Body function measures

#### Sociodemographic Characteristics

Birth location

- Biochemical measures
- Parents birth location

Language restricting the scope of data sharing (Wallace et al, 2009)

- Use of data restricted to researchers participating in the original study
  - "All research data are confidential ... they will only be used in medical research and [will] remain in the sole use of the participating researchers."
- Use of data restricted to researchers in one country
  - "Blood and DNA samples may...be distributed to laboratories...around [country] for further research."
  - "Research using the anonymous samples will be done by [researchers] ... throughout [country]."

# The need to obtain both scientific and ethical approval

- "The [Project] gives approved researchers access to data and samples.... All researchers will only have access to coded data or samples, in order to protect your privacy. They also have to obtain prior scientific and ethical approval, as described above, and their research must fit the purpose of the resource/biobank."
- "The [Project] expects to receive requests and, if approved, provide access to data and samples to overseas researchers and international collaborators. These researchers must follow the same procedures as all other researchers. All access is subject to the strictest scientific and ethical scrutiny..."

# ELSI restrictions on third-party sharing

- Substantive reasons:
  - Relevant values held throughout many societies
  - Disclosure of identity
  - Sensitive information
  - Intellectual property
- Exemplar wording
  - Wallace S, Lazor S, Knoppers BM. Chapter in Kaye J and Stranger M. Principles and Practice in Biobank Governance. Ashgate, Farnham 2009



#### The Healthy Obese project (HOP): A Core Project in BioSHaRE-eu

 Create HOP DataSCHEMA for variables to be shared across three major studies in BioSHaRE-eu







Y13D	Has/did your father ever suffer from?	Select from	If DA go to
	(You can select more than one answer)	- HE Heart disease	Y16
		- ST Stroke	Otherwise
		- BP High blood pressure	Go to Y13E
		- DB Diabetes	
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	
Y13D	Has/did your adopted father ever suffer	Select from	If DA go to
AD	from? (You can select more than one	- HE Heart disease	Y16AD
	answer)	- ST Stroke	Otherwise
		- BP High blood pressure	Go to
		- DB Diabetes	Y13EAD
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	

Y16D	Has/Did your mother ever suffer from?	Select from	If DA go to
	(You can select more than one answer)	- HE Heart disease	Y17
		- ST Stroke	Otherwise
		- BP High blood pressure	Go to Y16E
		- DB Diabetes	
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	
Y16D	Has/Did your adopted mother ever suffer	Select from	If DA go to
AD	from? (You can select more than one	- HE Heart disease	Y17AD
	answer)	- ST Stroke	Otherwise
		- BP High blood pressure	Go to
		- DB Diabetes	Y16EAD
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	



Y19	Have any of your brothers or sisters	Select from	If DA go to
	suffered from any of the following	- HE Heart disease	next section
	illnesses? (You can select more than one	- ST Stroke	Otherwise
	answer)	- BP High blood pressure	Go to Y20
		- DB Diabetes	
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	
Y19A	Have any of your adopted brothers or	Select from	If DA go to
D	sisters suffered from any of the following	- HE Heart disease	next section
	illnesses? (You can select more than one	- ST Stroke	Otherwise
	answer)	- BP High blood pressure	Go to
		- DB Diabetes	Y20AD
		- EM Chronic	
		bronchitis/emphysema	
		- AZ Alzheimers	
		disease/dementia	
		- NN None of the above	
		- UN Do not know	
		- DA Prefer not to answer	

### Strategic investment

- Effective joint data analysis involving large-scale biobanks will be an essential element of tomorrow's bioscience
- Generic tools and methods exist *now* that make it much easier to construct and then use harmonized data sets
- Strategic national and international funding is required to ensure that such tools continue to evolve, and are implemented in effective multipurpose user IT environments that are cheap and easy to use
- Some encouragement from funders may be needed to ensure they are actually used.

### Strategic investment

- The generic integrative work I have described is undertaken and promoted by large international consortia such as P<sup>3</sup>G, PHOEBE, BBMRI and BioSHaRE-eu
- Investment in such platform projects should therefore be international - all nations should contribute to a truly collaborative endeavour
- This investment should be equitably shared to ensure long-term sustainability