# ViPAR: a software platform for the Virtual Pooling and Analysis of Research data

**Associate Professor Kim Carter**
**Telethon Kids Institute**
**University of Western Australia**

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

"Any time scientists disagree, it's because we have insufficient data. Then we can agree on what kind of data to get; we get the data; and the data solves the problem.

Either I'm right, or you're right, or we're both wrong. And we move on ... "

Neil deGrasse Tyson

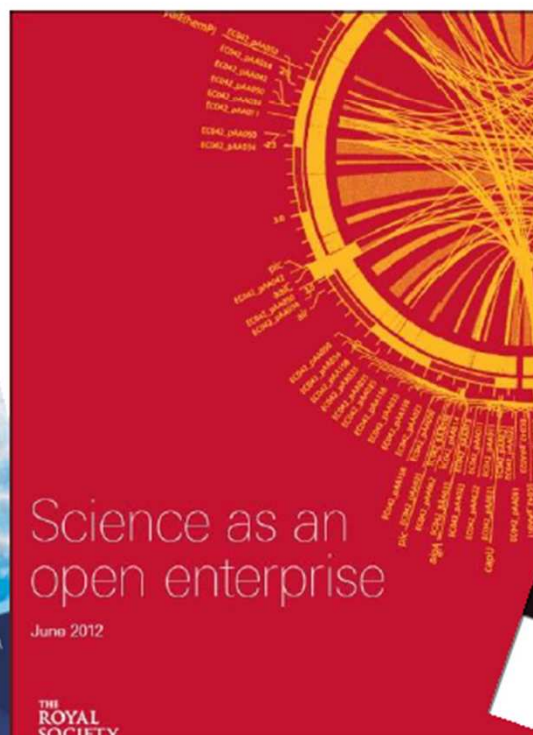# Telethon Kids Institute
# Perth, Western Australia

# Overview

- **Recent media on Data Sharing**
- **Technologies for Data Sharing**
- **How did I end up in this space ? (case study)**
- **Methods for Sharing and Analysing data together**
- **ViPAR**
- **Demo**

# Recent media on data sharing



SHARING CLINICAL RESEARCH DATA

Science as an open enterprise

June 2012

THE ROYAL SOCIETY

OECD Principles and Guidelines for Access to Research Data from Public Funding

OECD

nature

SPECIALS

Open Data White Paper

Unleashing the

#opendata
@uktransparency
@cabinetofficeuk

"Publicly funded research data are a public good, produced in the public interest, which should be made openly available with as few restrictions as possible in a timely and responsible manner that does not harm intellectual property."

RCUK Common Principles on Data Policy
http://www.rcuk.ac.uk/research/datapolicy/

INSTITUTE
Discover. Prevent. Cure.

# Media on data sharing (cont)

## THE LANCET

Comment

### Science as a public enterprise: the case for open data

Geoffrey Boulton[a], ✉ , Michael Rawlins[b], Patrick Vallance[c], Mark Walport[d]

[a] Grant Institute, Edinburgh University, Edinburgh EH9 3JW, UK

[b] National Institute for Health and Clinical Excellence, London, UK

[c] GlaxoSmithKline, London, UK

[d] Wellcome Trust, London, UK

## THE LANCET

Comment
### Sharing research data to improve public health

Mark Walport[a], ✉ and Paul Brest[b]

[a] Wellcome Trust, London NW1 2BE, UK

[b] Hewlett Foundation, Menlo Park, CA, USA

# Media on data sharing (cont)

"Increasing availability and promoting efficient data use to maximise public health benefits"

*Equitable:* balance needs of researchers, communities and funders

*Ethical:* protect individual privacy and dignity while recognising need to improve health using these data

*Efficient:* improve quality, value and contribution of research by building on existing, and reducing competition and duplication

# Unwanted sharing of data

arstechnica.com/security/2016/04/billion-dollar-bangladesh-hack-swift-software-hacked-no-firewalls-10-switches/

RISK ASSESSMENT —

## Billion dollar Bangladesh hack: SWIFT software hacked, no firewalls, $10 switches

The Bangladesh Bank's internal network security was sorely lacking.

PETER BRIGHT - 4/26/2016, 6:15 AM

71

The Bangladesh central bank had no firewall and was using a second-hand $10 network when it was hacked earlier this year. Investigation by British defense contractor BAE Systems has also shown that the SWIFT software used to make payments was compromised, enabling the hackers to send money around the world without leaving any trace in Bangladesh.

In February, unknown hackers broke into the Bangladesh Bank and almost got away with just shy of $1 billion. In the event, their fraudulent transactions were cancelled after they managed to transfer $81 million when a typo raised concerns about one of the transactions. That money is still unrecovered, but BAE has published some of its findings.

The SWIFT organization is owned by 3,000 financial companies and operates a network for sending financial transactions between financial institutions. Institutions using the network must have existing banking relationships; SWIFT transactions do not actually send money but instead send payment orders that must then be settled by having the institutions involved moving money between accounts.

SWIFT's security stems from two major sources. Notionally, it's a private network, and most banks set up their accounts such that only certain transactions between particular parties are permitted. The network privacy means that it should be hard for someone outside a bank to attack the network, but if a hacker breaks into a bank—as was the case here—then that protection evaporates. The Bangladesh central bank has all the necessary SWIFT software and authorized access to the SWIFT network. Any hacker running code within the Bangladesh bank *also* has access to the software and network.

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Unwanted sharing of data

RISK ASSESSMENT —

## Billion dollar Bangladesh hack: SWIFT sof sw

The Ba

PETER BR

71

## Banner Health Identifies Cyber Attack

August 03, 2016

PHOENIX (August 3, 2016) - Banner Health announced today that it is mailing letters to approximately 3.7 million patients, health plan members and beneficiaries, food and beverage customers and physicians and healthcare providers related to a cyber attack. Banner Health immediately launched an investigation, hired a leading forensics firm, took steps to block the cyber attackers and contacted law enforcement.

On July 7, 2016, Banner Health discovered that cyber attackers may have gained unauthorized access to computer systems that process payment card data at food and beverage outlets at some Banner Health locations. The attackers targeted payment card data, including cardholder name, card number, expiration date and internal verification code, as the data was being routed through affected payment processing systems. Payment cards used at food and beverage outlets at certain Banner Health locations during the two-week period between June 23, 2016 and July 7, 2016 may have been affected. A list of the outlets that were affected can be found at www.BannerSupports.com. The investigation revealed that the attack did not affect payment card payments used to pay for medical services.

On July 13, 2016, Banner Health learned that the cyber attackers may have gained unauthorized access to patient information, health plan member and beneficiary information, as well as information about physician and healthcare providers. The patient and health plan information may have included names, birthdates, addresses, physicians' names, dates of service, claims information, and possibly health insurance information and social security numbers, if provided to Banner Health. The physician and provider information may have included names, addresses, dates of birth, social security numbers and other identifiers they may use. The investigation also revealed that the attack was

# Unwanted sharing of data



arstechnica.com/security/2016/04/billion-dollar

**RISK ASSESSMENT —**

## Billion dollar Bang
soft
swi

**Banner Hea**

The Ba

PETER BR

August 03, 2016

PHOENIX (August 3, 2016) –
patients, health plan memb
providers related to a cyber
firm, took steps to block the

On July 7, 2016, Banner Hea
systems that process payme
targeted payment card data
the data was being routed t
outlets at certain Banner He
have been affected. A list of
investigation revealed that t

On July 13, 2016, Banner He
information, health plan me
providers. The patient and h
names, dates of service, clai
provided to Banner Health. The physician and provider information may have included names, addresses, dates of
birth, social security numbers and other identifiers they may use. The investigation also revealed that the attack was

71

arstechnica.com/tech-policy/2016/08/australia-2016-census-personal-data-retention/

**CON CENSUS —**

## Australians threaten to take leave of their census

2016 Australian census stores names and addresses, prompting privacy, security outrage.

JENNIFER BAKER (UK) - 8/4/2016, 10:15 PM



Paramount

126

Next Tuesday is the day Australians must fill in—correctly—their census forms, or face a fine. However, many may be willing to take that risk as the Australian Bureau of Statistics (ABS) will rather extraordinarily be storing names and addresses in addition to the usual census results.

# Unwanted sharing of data



arstechnica.com/tech-policy/2016/08/australia-2016-census-personal-data-retention/

## Sony Makes it Official: PlayStation Network Hacked

By Keir Thomas, PCWorld    Apr 23, 2011 7:35 AM

When Sony's PlayStation Network was taken offline three days ago, all eyes fell on the Anonymous group, who've taken a dislike to Sony over its treatment of hardware hacker George Hotz. The network allows online play between PlayStation 3 consoles and boasts 70 million users, so this is no small inconvenience.

## Electronic health files secure despite NBN hacking, minister insists

SEAN PARNELL AND BEN PACKHAM   The Australian   July 28, 2011 12:00AM

the data was being routed t

## Gordon Brown's shock that his family medical records were hacked

Rebekah Brooks, then editor of The Sun, contacted the Browns, informing them that she had obtained medical details about their four-year-old son Fraser

By The Independent Reporting Team, Cahal Milmo, Martin Hickman, Oliver Wright and Ian Burrell

Tuesday, 12 July 2011          ⌄ SHARE  |  🖶 PRINT  |  ✉ EMAIL  |  A A A TEXT SIZE

birth, social security numbers and other identifiers they may use. The investigation also revealed that the attack was

ELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Funder expectations (UK)

Legend: ● Full Coverage  ◐ Partial Coverage  ○ No Coverage

| Research Funders | Policy Coverage | | Policy Stipulations | | | | | Support Provided | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Published outputs | Data | Time limits | Data plan | Access/ sharing | Long-term curation | Monitoring | Guidance | Repository | Data centre | Costs |
| AHRC | ● | ● | ● | ● | ● | ◐ | ○ | ● | ○ | ◐ | ◐ |
| BBSRC | ● | ● | ● | ● | ● | ● | ● | ● | ● | ◐ | ● |
| CRUK | ● | ● | ● | ● | ● | ● | ● | ◐ | ● | ○ | ○ |
| EPSRC | ● | ● | ● | ◐ | ● | ● | ● | ◐ | ○ | ○ | ● |
| ESRC | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ◐ |
| MRC | ● | ● | ● | ● | ● | ● | ○ | ◐ | ● | ○ | ◐ |
| NERC | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ◐ |
| STFC | ● | ● | ● | ● | ● | ● | ● | ◐ | ● | ◐ | ◐ |
| Wellcome Trust | ● | ● | ● | ● | ● | ● | ● | ● | ● | ◐ | ● |

http://www.dcc.ac.uk/resources/policy-and-legal/overview-funders-data-policies

# Journal expectations

**An increasing number of journals expect data to be publicly available at time of publication**



PLOS | ONE

Publish | About | Browse | Search 🔍

advanced search

Acceptable Data-Sharing Methods

Unacceptable Data Access Restrictions

Explanatory Notes and Guidance

Recommended Repositories

FAQs for Data Policy

## Data Availability

**The following policy applies to all of PLOS journals, unless otherwise noted.**

PLOS journals require authors to make all data underlying the findings described in their manuscript fully available without restriction, with rare exception.

When submitting a manuscript online, authors must provide a *Data Availability Statement* describing compliance with PLOS's policy. If the article is accepted for publication, the data availability statement will be published as part of the final article.

Refusal to share data and related metadata and methods in accordance with this policy will be grounds for rejection. PLOS journal editors encourage researchers to contact them if they encounter difficulties in obtaining data from articles published in PLOS journals. If restrictions on access to data come to light after publication, we reserve the right to post a correction, to contact the authors' institutions and funders, or in extreme cases to retract the publication.

Methods acceptable to PLOS journals with respect to data sharing are listed below, accompanied by guidance for authors as to what must be indicated in their data availability statement and how to follow best practices in reporting. If authors did not collect data themselves but used another source, this source must be credited as appropriate. Authors who have questions or difficulties with the policy, or readers who have difficulty accessing data, are encouraged to contact the relevant journal office or data@plos.org.

http://journals.plos.org/plosone/s/data-availability

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Data sharing journals



Current model of open data sharing.

Data paper

Journal → Reviewer → Error reporting

$$ $$

Citations!

Others papers

Closing the circle between data generators and users.

Gorgolewski, Milham, and Margulies, 2013 Front. Neurosci.

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Large Consortia

- To elucidate subtle genetic & environmental effect on common diseases ➡ larger sample sizes req'd

- Power comes from the pooled sample size

   (ie an individual level meta-analysis

- Multi-site, multi-national means potential ethical, legal and privacy barriers to sharing research data

# Disincentives to sharing?

- Policies aren't always compelling

- Labor not always recognised / rewarded

- Authorship: Nominally only First / Last really counts

- Tenure: Metrics Reward Individual not teams

- HIPAA and similar: Scary and significant individual penalties for data loss

- Deidentification: Is it truly possible?

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Different Perspectives and Drivers

## Technological

## Economic



## Legal

## Personal

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# How do you actually share data?

# Data sharing by physical media



Usually require local collaborator (proximity) for transfer (manual) or has potentially significant time delays (eg postage)

# Technologies for sharing data (by electronic means)

Removes the barrier of geography to allow shared "storage" of data. Doesn't necessarily make it any easier to do anything once you have the data

# Examples of projects facilitating sharing of research data (Aus)





**ARCS** – provides data storage, data transfer, collaboration and conferencing facilities

**Biogrid** – platform for integrating and analysing clinical, imaging and biospecimen data across jurisdictions

**NeCTAR** – cloud services for workflows, tools and servers

# Examples of successful data sharing meta-resources

# Examples of successful data sharing projects (repositories)

# Examples of successful data sharing projects (repositories)



Data sharing 1.0

# Data warehouses vs Databases

## Database vs Data Warehouse

| | Database | Data Warehouse |
|---|---|---|
| Purpose | Data retrieval, updating and management | Data analysis and decision making |
| Application | OLTP (Online Transaction Processing) | Reporting and OLAP (Online Analytical Processing) |
| Format | Normalized | Denormalized |
| Time Frame | Current/Real-Time | Historical |

Note that you can use a database as a data warehouse. It depends on, for which purpose you have design the database

www.edureka.co/data-warehousing-and-bi

# Data warehouses vs Databases

"Bioinformatics issues for conducting sophisticated

multiple-pass genome sequence analysis"

Kim W. Carter

Murdoch University, 2004.

# How/why did I get interested in data sharing technologies?



FIG. 5.  THE SIDE STEP BY WHICH AN OPPONENT IS ELUDED.

# The International Consortium for Autism Registry Epidemiology (iCARE)



iCARE annual meeting, Sweden 2011

# Vision for iCARE and autism research

Motivated and willing collaborators wanted to join forces to create a resource that allows for analyses that:

- cannot be performed with a single existing system

- enhance the analytic potential of a single data system

- allow direct comparison of findings across different data systems with eg geographic, population, or data collection variation.

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Overarching aims of iCARE

1. **Funded (Autism Speaks 2009-2013) to setup a <u>multi-national virtual data set</u> for the study of pre- and perinatal risk factors for autism.**

   - Establish a minimal database requirement.
   - Establish and conduct data preparation protocols across registries.
   - Enable and test centralized data access protocols to multiple registries.
   - Establish a collaborative framework and guidelines for a working relationship among sites/investigators.

2. **Demonstrate the utility of the resource**

TELETHON
**KIDS**
INSTITUTE
Discover. Prevent. Cure.

# How did I get involved with iCARE?

Setup a **multi-national virtual data set** for the study of pre- and perinatal risk factors for autism.

The group had funding and ethics approval for virtual pooling but the original IT group had to withdraw

# iCARE consortium characteristics

| Site | Population Size | Birth Years | Births/Year | Coverage | Health Care Provision |
|------|-----------------|-------------|-------------|----------|-----------------------|
| Denmark | 5.5 million | 1980-2007 | 62,000 | National | Public |
| Finland | 5.4 million | 1987-2008 | 60,000 | National | Public |
| Israel | 7.6 million | 1987-2006 | 86,000 | National | Public |
| Norway | 4.8 million | 1980-2005 | 55,000 | National | Public |
| Sweden | 9.4 million | 1980-2008 | 107,000 | National | Public |
| Western Australia | 1.9 million | 1983-1999 | 24,000 | State | Public and private |

Population-based registries ➡ Large samples, unbiased, prospective, with ability to link with other population databases

# Data dictionary and harmonisation

**Critical to the success of the project (and any data sharing project) is making data comparable across sites**

**Challenge**

Diverse data availability and formats, over time and by site

**Solution: Harmonisation**

generic term for procedures that create comparability between data derived from different sources

# Data dictionary and harmonisation

| Minimum Variable List | Data Availability Surveys | Data Dictionary | Certification | Archives |
|---|---|---|---|---|

Data dictionary V1 (Sep 2010) – 48 variables

Data dictionary V2 (Mar 2011) – 52 variables

Data dictionary V3 (Mar 2013) – 58 variables

# Recap: So why do we want to share research data?

**Benefits**

- Increased power, leading to subgroup analyses and interactions
- Ability to compare outcomes and validate models across sites

**Pitfalls**

- Cost of harmonisation – time and $
- Difficult to validate/error check when anonymous
- Potentially more complex analysis
- Ethico-legal issues with privacy & consent
- Requires strong collaborations

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Are there electronic methods that facilitate data sharing and analysis together?

# Methods for Sharing and Analysing Research Data

## 1. Traditional Meta Analysis

Combining existing results of analyses on similar outcomes and predictors using similar methodology

- needs published data on similar outcomes and predictors using similar methodology
- can suffer from "resolution"
As individual level data may not be used

Meta-Analysis Results



Mean and 95 Pct CI (Study_0 is Pooled Results)

# Methods for Sharing and Analysing Research Data (cont)

## 2. Manual Data Pooling

Analysis of harmonised pooled data from multiple sites, sent to a single analysis site

- (manually) unified methodology
- needs harmonised data
- consent and ethics for sharing
- significant effort for a single analysis centre
- requires strong collaboration

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Methods for Sharing and Analysing Research Data (cont)

## 3. Automated Data Pooling (federation)
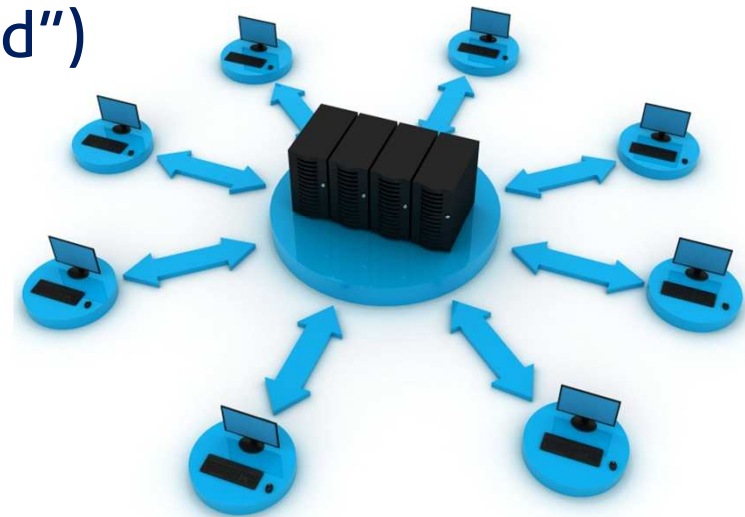Analysis of harmonised pooled data automatically from multiple sites (pooling in the "cloud")

- unified methodology
- needs harmonised data
- consent and ethics for sharing
- strong collaboration
- requires informatics infrastructure to enable the federation process
- Requires some infrastructure at all data-contributing sites

# Methods for Sharing and Analysing Research Data (cont)

## 4. Automated decentralised analysis – eg DataSHIELD

Analysis of harmonised data automatically from multiple sites by pooling statistics
(ie no data transfer)



- unified methodology
- needs harmonised data
- may need consent and ethics for sharing
- potentially limited types of analysis
- Requires significant informatics infrastructure at all data-contributing sites

# ViPAR: Virtual infrastructure for pooling and analysis of research data

# What is Federation?



?

# WHAT IS DATABASE FEDERATION?



**Tools that transparently integrates multiple autonomous databases into a single virtual entity**

# Federation within RDMS

## Site 1

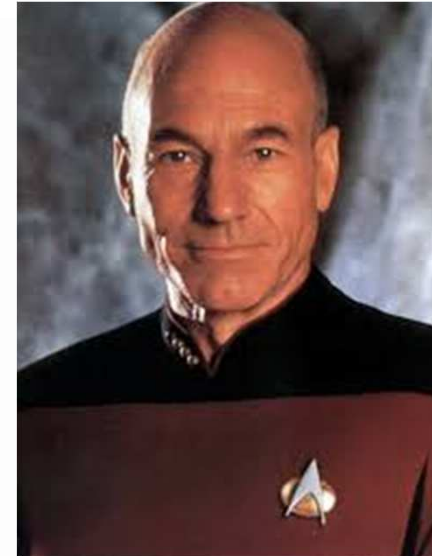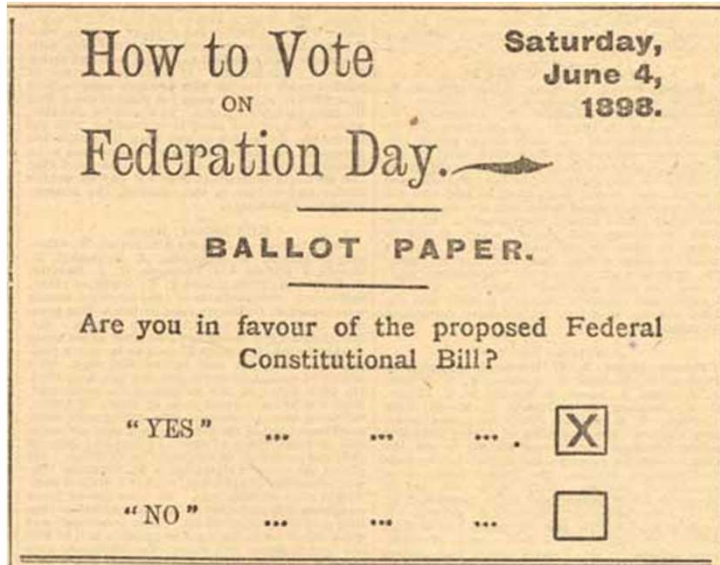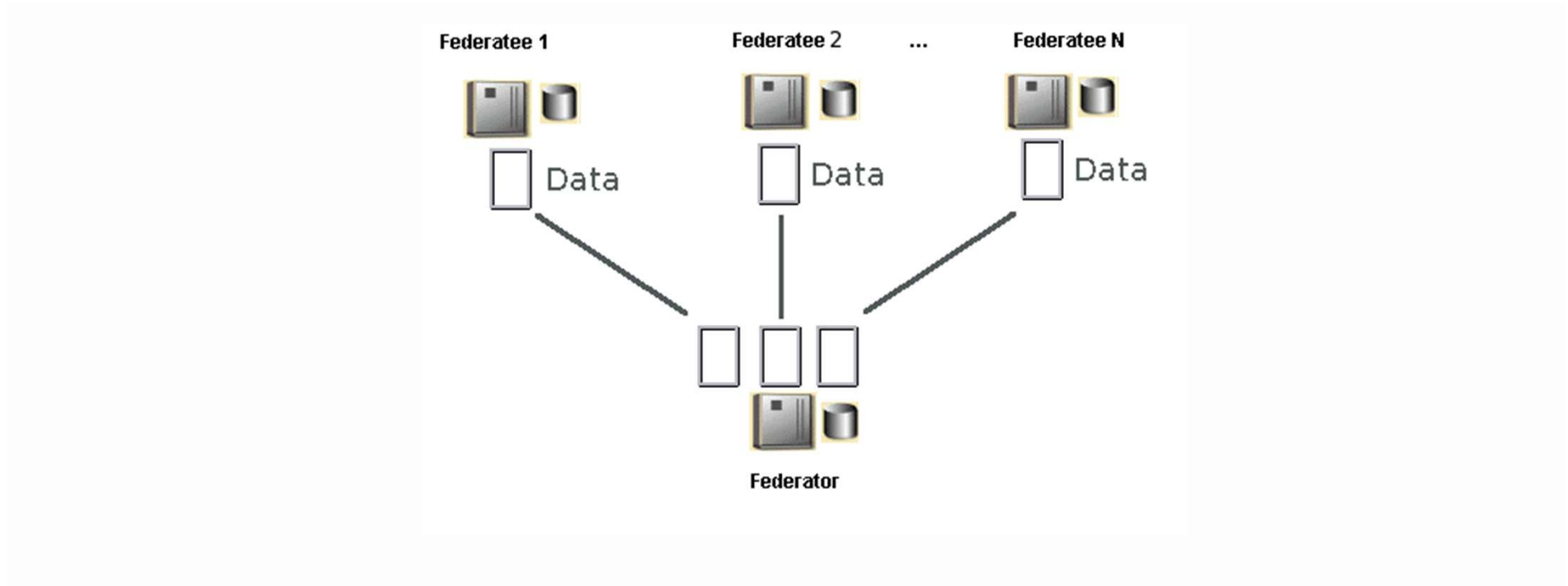| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

## Site 2

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

## RDMS federated table

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

TELETHON KIDS INSTITUTE
Discover. Prevent. Cure.

# Federation within RDMS (cont)

## RDMS federated table

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

Select firstname,age from table

| First Name | Age |
|---|---|
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |

# Federation within RDMS (cont)

## RDMS federated table

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

Select firstname,age from table

| First Name | Age |
|---|---|
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |

# Federation within RDMS (cont)

**RDMS federated table**

Select firstname,age
from table

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

| First Name | Age |
|---|---|
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |
| Mickey | 73 |
| Bat | 54 |
| Wonder | 39 |
| Donald | 65 |
| Bugs | 58 |
| Wiley | 61 |
| Cat | 32 |
| Tweety | 28 |

Intermediate copy of entire database tables
made on the fly at the central site

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Federation within RDMS (cont)

## RDMS federated table

| First Name | Last Name | Address | City | Age |
|---|---|---|---|---|
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |
| Mickey | Mouse | 123 Fantasy Way | Anaheim | 73 |
| Bat | Man | 321 Cavern Ave | Gotham | 54 |
| Wonder | Woman | 987 Truth Way | Paradise | 39 |
| Donald | Duck | 555 Quack Street | Mallard | 65 |
| Bugs | Bunny | 567 Carrot Street | Rascal | 58 |
| Wiley | Coyote | 999 Acme Way | Canyon | 61 |
| Cat | Woman | 234 Purrfect Street | Hairball | 32 |
| Tweety | Bird | 543 | Itotltaw | 28 |

Federation engines in MySQL and Postgres appear to be poorly implemented:
* suitable for LAN only
* suitable for smallish datasets
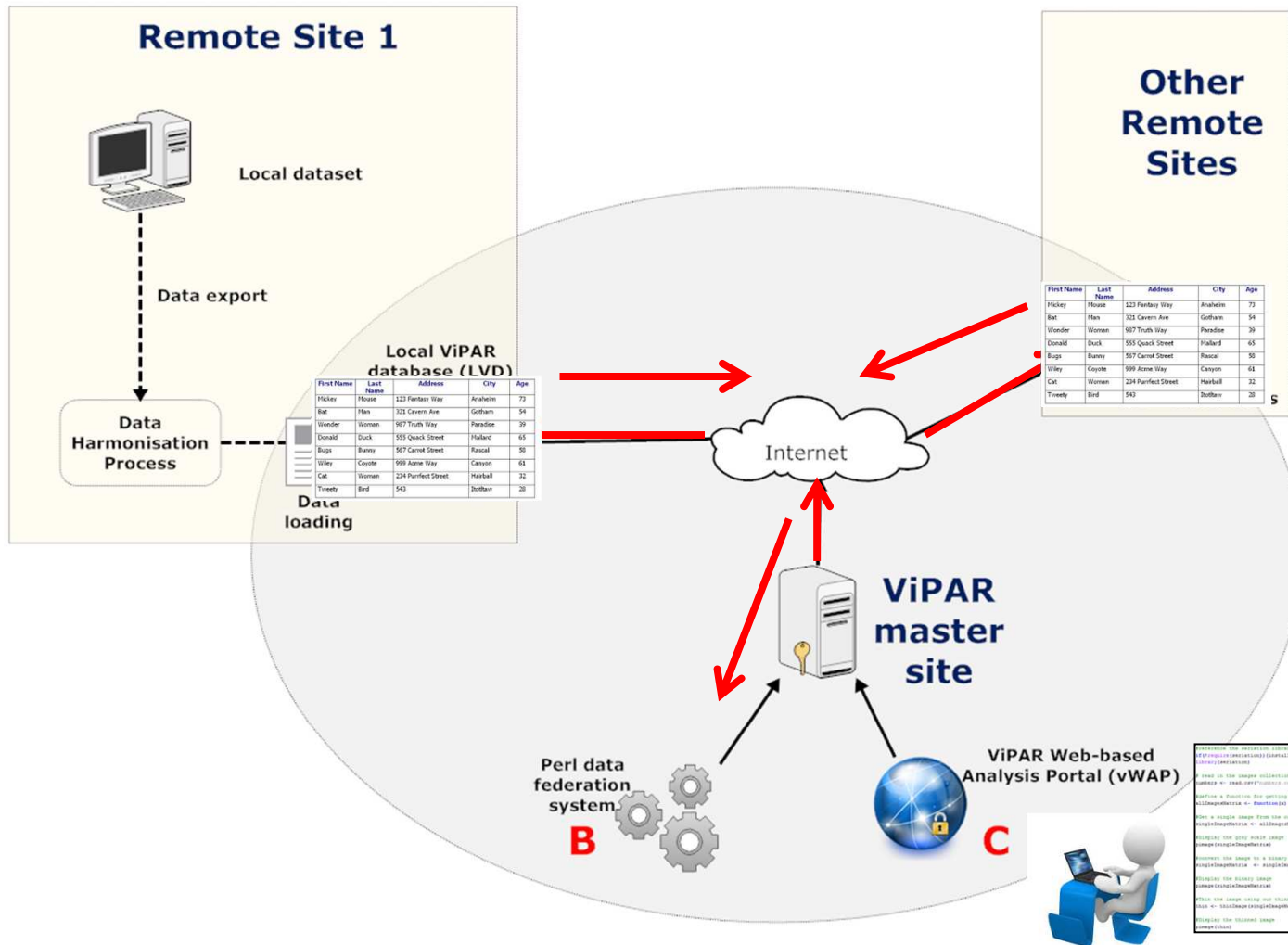
Problem: needs to work securely over Internet with large datasets

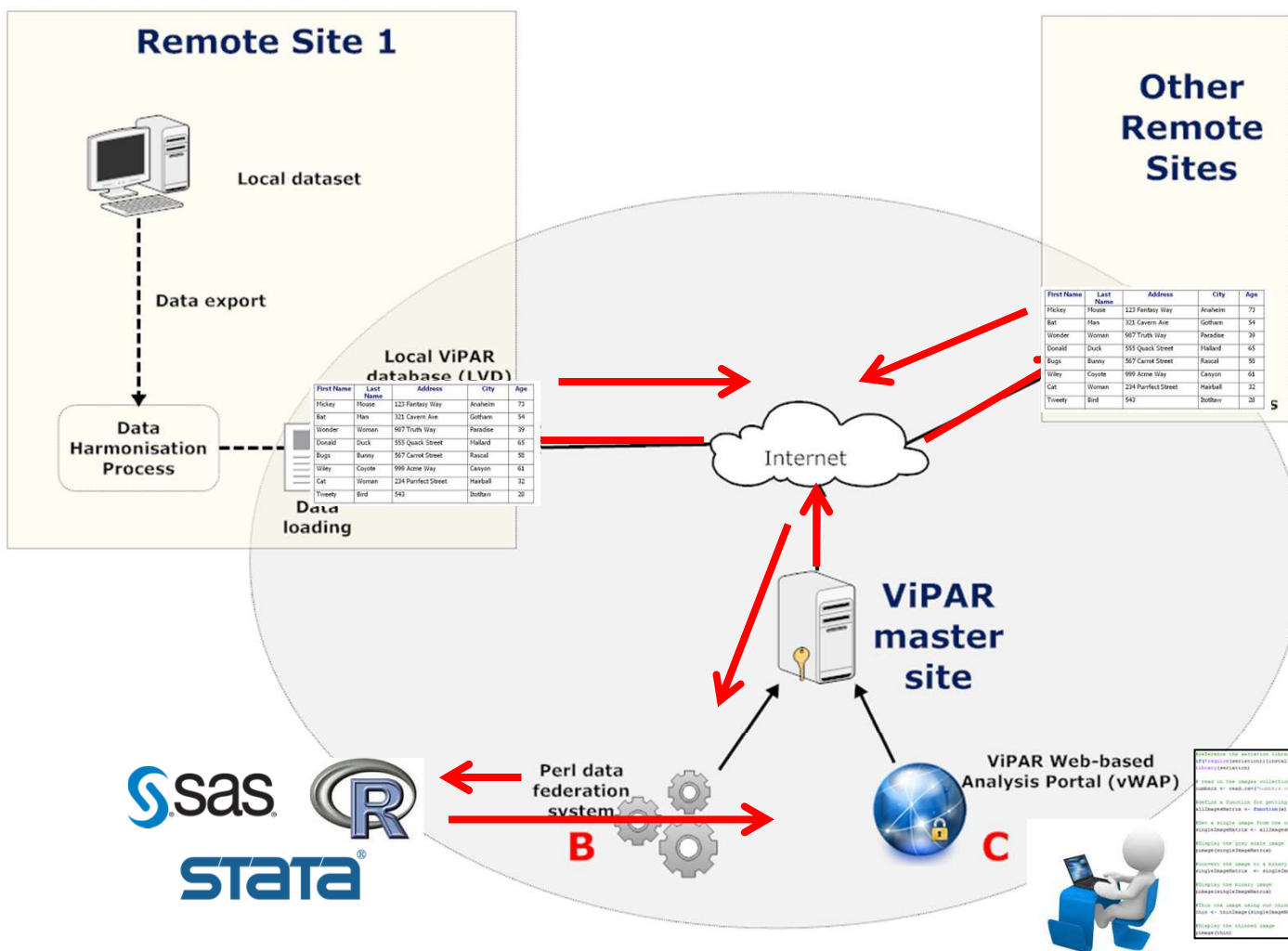TELETHON KIDS INSTITUTE
Discover. Prevent. Cure.

# ViPAR: overview

# ViPAR – a custom data federation platform

**Open-source stack – Linux, Apache, Mysql, OpenSSL**
**Available as a pre-built VM images for easy access**

http://bioinformatics.childhealthresearch.org.au/software/vipar/
(pre-built VM images for easy use – Vbox and Vmware)
https://gitlab.com/kim.carter/ViPAR (source code)
https://groups.google.com/forum/#!forum/vipar (support forum)

**ViPAR Daemon – 1500 lines of Perl code**
-Handles multiple clients simultaneously, and can perform eg parallel data retrieval (all sites at the same time)
-Implements its own federation (and analysis) engine

**VWAP – (ViPAR Web-based Analysis Portal)**
-8000 lines of Perl code
-1200 lines of Javascript

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Behind the scenes

# Behind the scenes

**MySQL**
Open-source database management system, with powerful features including database federation
Widely used, especially for large databases

**SSH**
Powerful, free solution to securely exchange data between 2 networked devices over an insecure channel (such as the Internet)
Enterprise reliability and security
"tunnels" allow us to transfer of programs services between remote computers dynamically

**SSL**
Client-side and Server-side SSL certificates used (along with user accounts and roles) to provide layered, strong security model

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Data model

**Study**
Overarching conglomeration of analysis projects, data resources, data dictionary and variables and users

**Data dictionary**
A version controlled set of Variables and Missing data fields

**Variable**
A data field of a given type (eg continuous, categorical, date) with a missing data type

**Sites, Servers and Resources**
A site is a nominal geo-location housing a ViPAR LVD server accessible on a specified SSH port, with one or more (certified) harmonised dataset resources on the server

# Data model (cont)

**Analysis "projects"**

Defined research question(s) (eg agreed by AOC)
- Specified list of Resources to be included
- Specified list of Variables available for analysis
- Specified list of Users with access @ analyst or guest level (read-only)
- Results/outputs can be tagged as "sharing" accessible to enable non-project users within ViPAR to see specific results

**Users**
- Can be designated as "Admin" (see all projects and get extra admin menu options"
- Can be member of multiple analysis projects
- Unique user/pass, along with client SSL certificate (recommended)

# Behind the scenes: data model

# How does ViPAR protect my data?

**Network level**
- Private, encrypted network connecting each data site (LVD) to the master site (VMS). Protected by hardware & software firewalls
- Only the VMS can initiate connections to the LVDs
- Sites maintain their own data (in the LVD) – no data is saved at the VMS

**User-level**
- Unique user pass & cert combo (with client SSL certificates) – you can't get to the site login page without
- Only users who run code get to "pool" the data (within a defined project)
- Implicit "trust" to give max flexibility in analysis -> that users won't try to access individual data

**Portal-level**
- Everything is logged, all run code, all interactions including "deletions".
- No direct access to the data

**Data-level design**
- Design your dictionary to minimise potentially identifying fields

# ViPAR Portal (for iCARE)

# Running an analysis

# Accessing outputs

# Management

# Has VιPAR been successful?

**Currently**:

- 7 remote sites (Nor,Fin,Den,Swe,Isr,US,WA), containing approx 9million records across 60 harmonised fields
- Passed multiple ethics approvals (at each site)
- Data retrieval with ViPAR: approx 2minutes in serial (one site at a time); approx 30-40 seconds in parallel

**The International Collaboration for Autism Registry Epidemiology (iCARE): multinational registry-based investigations of autism risk factors and t**

Schendel DE[1], Bresnahan M, Carter K
Parner ET, Reichenberg A, Sandin S, ...

+ Author information

**Autism risk associated with parental age and with increasing difference in age between the parents.**

Sandin S[1], Schendel D[2], Magnusson P[3], Hultman C[3], Surén P[4], Susser E[5], Grønborg T[6], Gissler M[7], Gunnes N[4], Gross R[8], Henning M[9], Bresnahan M[5], Sourander A[10], Hornig M[11], Carter K[12], Francis R[12], Parner E[6], Leonard H[12], Rosanoff M[13], Stoltenberg C[14], Reichenberg A[15].

+ Author information

Open/close author information list

## Abstract

The International Collaboration for A
Norway, Sweden, USA) to promote
iCARE devised solutions to challeng
integrating existing national or state
Analyses are performed using datal
unprecedented resource in autism r
course of autism.

## Abstract

Advancing paternal and maternal age have both been associated with risk for autism spectrum disorders (ASD). However, the shape of the association remains unclear, and results on the joint associations is lacking. This study tests if advancing paternal and maternal ages are independently associated with ASD risk and estimates the functional form of the associations. In a population-based cohort study from five countries (Denmark, Israel, Norway, Sweden and Western Australia) comprising 5 766 794 children born 1985-2004 and followed up to the end of 2004-2009, the relative risk (RR) of ASD was estimated by using logistic regression and splines. Our analyses included 30 902 cases of ASD. Advancing paternal and maternal age were each associated with increased RR of ASD after adjusting for confounding and the other parent's age (mothers 40-49 years vs 20-29 years, RR=1.15 (95% confidence interval (CI): 1.06-1.24), P-value<0.001; fathers≥50 years vs 20-29 years, RR=1.66 (95% CI: 1.49-1.85), P-value<0.001). Younger maternal age was also associated with increased risk for ASD (mothers <20 years vs 20-29 years, RR=1.18 (95% CI: 1.08-1.29), P-value<0.001). There was a joint effect of maternal and paternal age with increasing risk of ASD for couples with increasing differences in parental ages. We did not find any support for a modifying effect by the sex of the offspring. In conclusion, as shown in multiple geographic regions, increases in ASD was not only limited to advancing paternal or maternal age alone but also to differences parental age including younger or older similarly aged parents as well as disparately aged parents.Molecular Psychiatry advance online publication, 9 June 2015; doi:10.1038/mp.2015.70.

Centrepiece of a successful NIH grant bid (2013-2017) $5.5M

**TELETHON KIDS INSTITUTE**
Discover. Prevent. Cure.

# For more information

Original article

## ViPAR: a software platform for the Virtual Pooling and Analysis of Research Data

Kim W. Carter,*[†] Richard W. Francis[†] and the International Collaboration for Autism Registry Epidemiology

Telethon Kids Institute, Centre for Child Health Research, University of Western Australia, Perth, WA, Australia

*Corresponding author. Telethon Kids Institute, The University of Western Australia, 100 Roberts Road, Subiaco, Perth, Western Australia, 6008. E-mail: Kim.Carter@telethonkids.org.au
[†]These authors contributed equally.
A full list of authors and affiliations appears at the end of the paper.

## Abstract

**Background:** Research studies exploring the determinants of disease require sufficient statistical power to detect meaningful effects. Sample size is often increased through centralized pooling of disparately located datasets, though ethical, privacy and data ownership issues can often hamper this process. Methods that facilitate the sharing of research data that are sympathetic with these issues and which allow flexible and detailed statistical analyses are therefore in critical need. We have created a software platform for the Virtual Pooling and Analysis of Research data (ViPAR), which employs free and open

*Int J Epi* 2015; doi: 10.1093/ije/dyv193

# VıPAR vs DaтaSHIELD

# ViPAR vs DataSHIELD

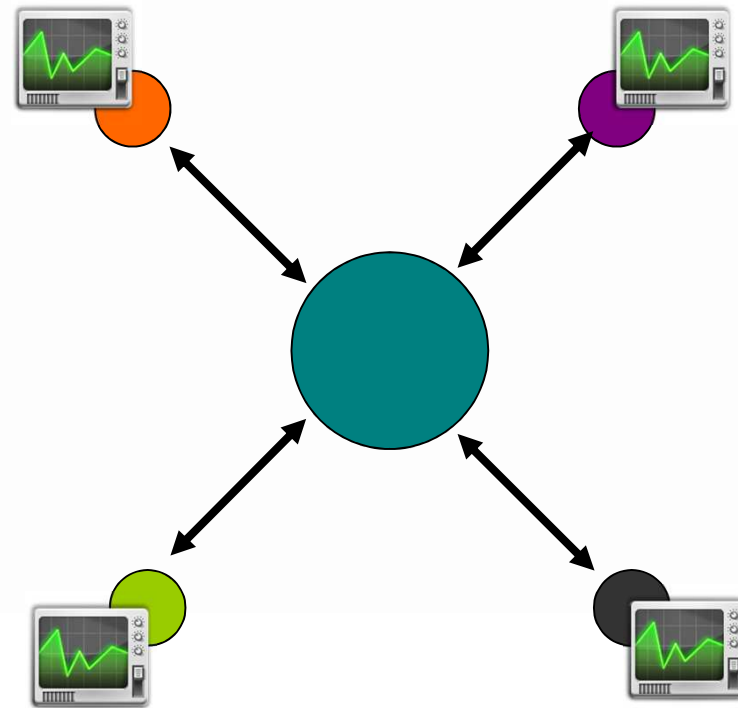- Data remains at study site
- Harmonisation is required before copies of data are sent 'to the cloud'
- Generation of pooled dataset occurs 'in the cloud'
- Results generated from a combined virtual dataset
- Pooled virtual dataset not saved/stored – only exists while analysis is running
- Straightforward analysis syntax in multiple languages

- Data remains at the study site
- Harmonisation is required before sending
- No data is sent anywhere
- No pooled dataset is generated
- Statistical methods sent to each study site, where methods are run locally
- Results are pooled back at analysis centre, to produce final result
- Requires specific analysis syntax and software (R) to make this happen

TELETHON
KIDS
INSTITUTE
Discover. Prevent. Cure.

# Future Directions

- Transparent integration of ViPAR with other complementary techniques eg DataSHIELD
- Web 2.0 interface
- Building capabilities to handle complex and large data types eg imaging data, whole genome seq
  - NoSQL, Object and other DBs
- Further security enhancements
  - on the fly privacy checks
  - using containers (Docker) for running analyses for greater separation

# Acknowledgements

- Richard Francis, TKI



- Becca Wilson, Univ of Bristol
- Paul Burton, Univ. of Bristol

# A QUICK PLUG

## COINSTAC: A privacy enabled model and prototype for leveraging and processing decentralized brain imaging data

Sergey M. Plis[1*], Anand Sarwate[2], Dylan Wood[1], Christopher Dieringer[1], Drew Landis[1], Cory Reed[1], Sandeep R. Panta[1], Jessica A. Turner[1, 3], Jody Shoemaker[1], Kim Carter[4], Paul Thompson[5], Kent Hutchinson[6] and Vince D. Calhoun[1, 7]

[1]The Mind Research Network, USA
[2]Dept. of Electrical and Computer Engineering, Rutgers, The State University of New Jersey, USA
[3]Dept. of Psychology and Neuroscience Institute, Georgia State University, USA
[4]Telethon Kids Institute, the University of Western Australia, Australia
[5]Imaging Genetics Center, ENIGMA Center for Worldwide Medicine, Imaging, and Genomics, Departments of Neurology, Psychiatry, Engineering, Radiology, and Pediatrics, University of Southern California, USA
[6]Department of Psychology and Neuroscience, University of Colorado Boulder, USA
[7]Dept. of Electrical and Computer Engineering, University of New Mexico, USA

The field of neuroimaging has embraced the need for sharing and collaboration. Data sharing mandates from public funding agencies and major journal publishers have spurred the development of data repositories and neuroinformatics consortia.
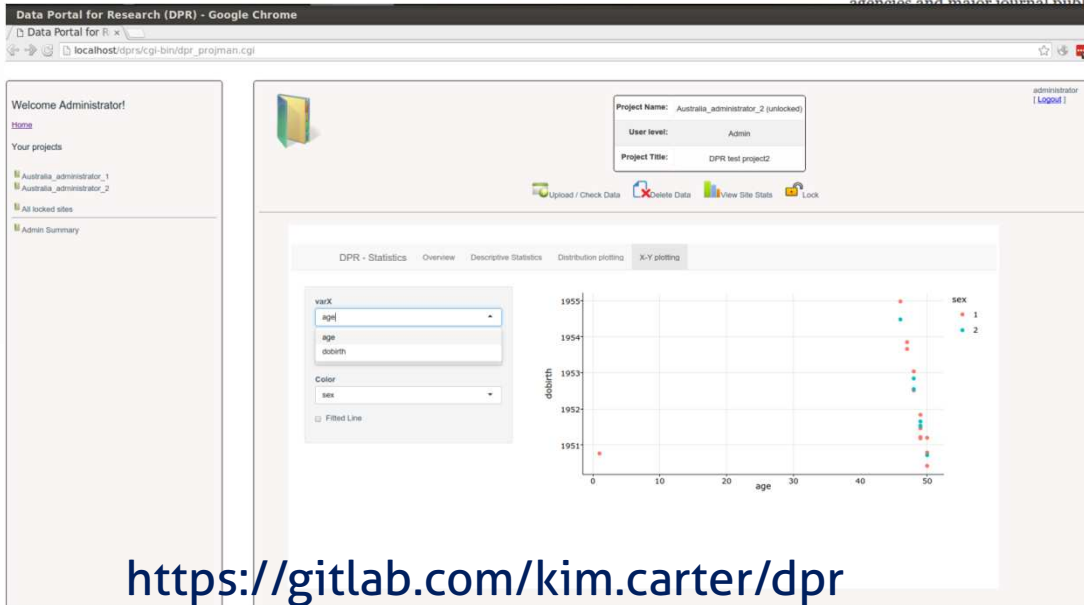...ta sharing still faces several hurdles. For example, open data sharing is on the rise but is not ...not easily shared, such as genetics. Current approaches can be cumbersome (such as negotiating ... There are also significant data transfer, organization and computational challenges. Centralized ... the issues. We propose a dynamic, decentralized platform for large scale analyses called the ...roimaging Suite Toolkit for Anonymous Computation (COINSTAC). The COINSTAC solution ...ral repositories, allows pooling of both open and ``closed" repositories by developing privacy-... algorithms, and incorporates the tools within an easy-to-use platform enabling distributed ... prototype system which we demonstrate on two multi-site data sets, without aggregating the ...ss sites, the COINSTAC model enables meta-analytic solutions to converge to ``pooled-data" ...vere in hand). More advanced approaches such as feature generation, matrix factorization ...incorporated into such a model. In sum, COINSTAC enables access to the many currently ...ly privacy enabled interface for decentralized analysis, and a powerful solution that ... solutions.

...rivacy, brain imaging, data sharing, Decentralized algorithms

...Dieringer C, Landis D, Reed C, Panta SR, Turner JA, Shoemaker J, Carter K, Thompson P, Hutchinson K and ...cy enabled model and prototype for leveraging and processing decentralized brain imaging data. *Front.* ...016.00365
...Jul 2016

https://gitlab.com/kim.carter/dpr

# Questions & Demo time?



http://xkcd.com/257/