

# Becoming Dataware

Richard Mortier  
Horizon Digital Economy Research  
University of Nottingham

With thanks to James Goulding, Derek McAuley,  
Anil Madhavapeddy (CUCL)



# RCUK Digital Economy

- Digital Economy is defined by the Research Councils as:  
*The novel design or use of information and communication technologies to help transform the lives of individuals, society or business.*
- The Digital Economy programme is:
  - Cross-Research Council (EPSRC, ESRC, AHRC)
  - Funded 2008-2011 for:
    - £80m research – including 3 x £12m hubs
    - £36m training – 8 x DTCs
  - Aimed at realising the transformational impact of ICT for all aspects of Business, Society and Government



# What is Horizon?

- A Digital Economy Research Centre at the University of Nottingham comprising:
  - A Digital Economy Hub
    - £20m from RCUK and university
    - Spokes at Cambridge, Reading, Exeter, Brunel
  - A Doctoral Training Centre
    - £15m from RCUK and university
    - 20 PhD students per year for 5 years
  - Now 120+ partner companies, from 40 in initial bid
  - 3 TEDDI projects
  - ... + future Digital Economy projects

# Our Digital Footprints

“Every time we register for a new web service, or upload our photos and videos, we are enlarging our own digital footprints”

- Whether “informed” or not
  - Facebook, Google
- Growing digital footprint poses major societal challenges
- ...but also forms a key basis for the digital economy’s growth





# In More Detail

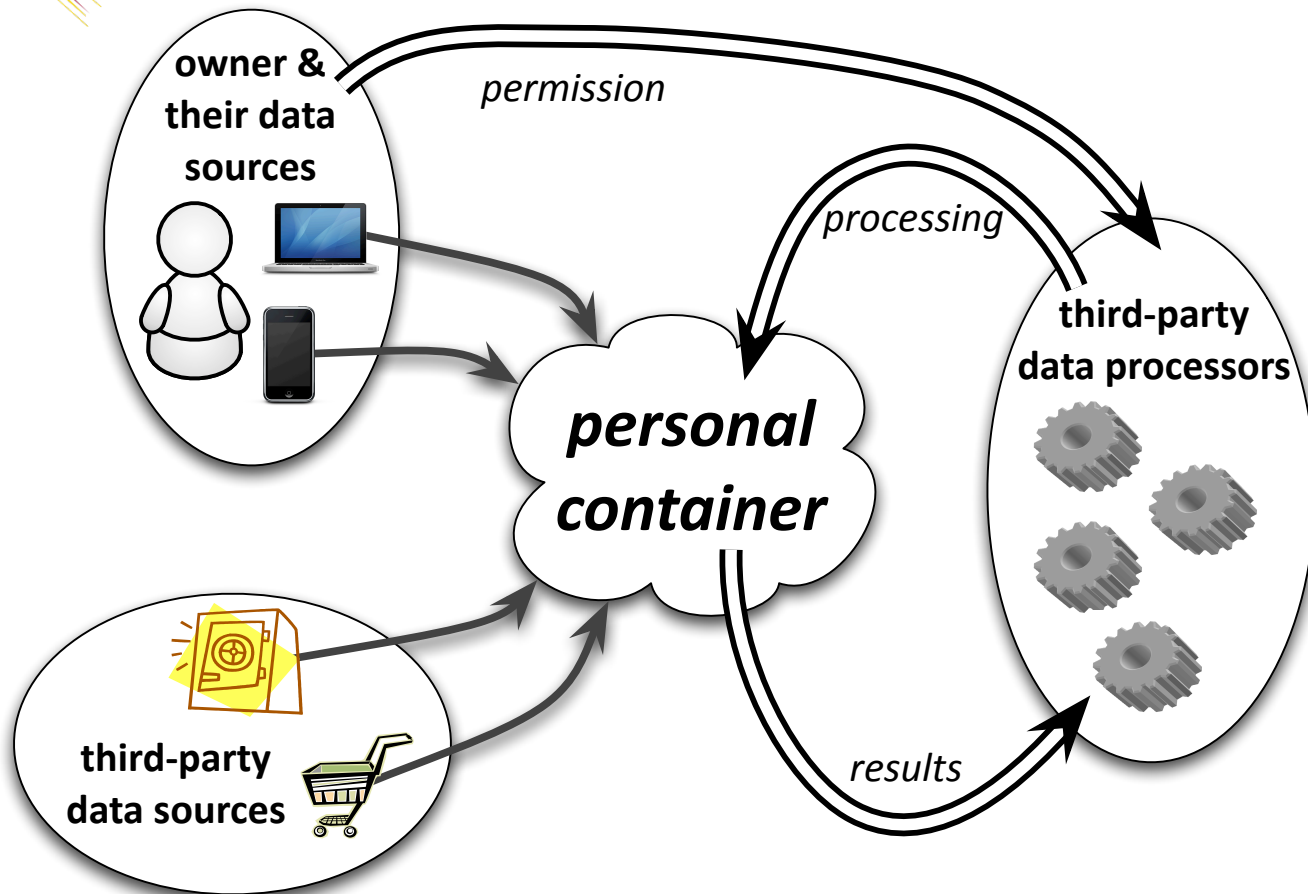
---

- We are making more and more information about our lives available, digitally
  - Whether or not we realise it
  - Often in very large, very rich data silos (e.g., Facebook)
  - “Contextual footprint”
- Simultaneously, there are more and more opportunities for mutually beneficial exploitation of digital personal data
  - E.g., Shopping basket optimizer;  
Boots prescription conflict detector

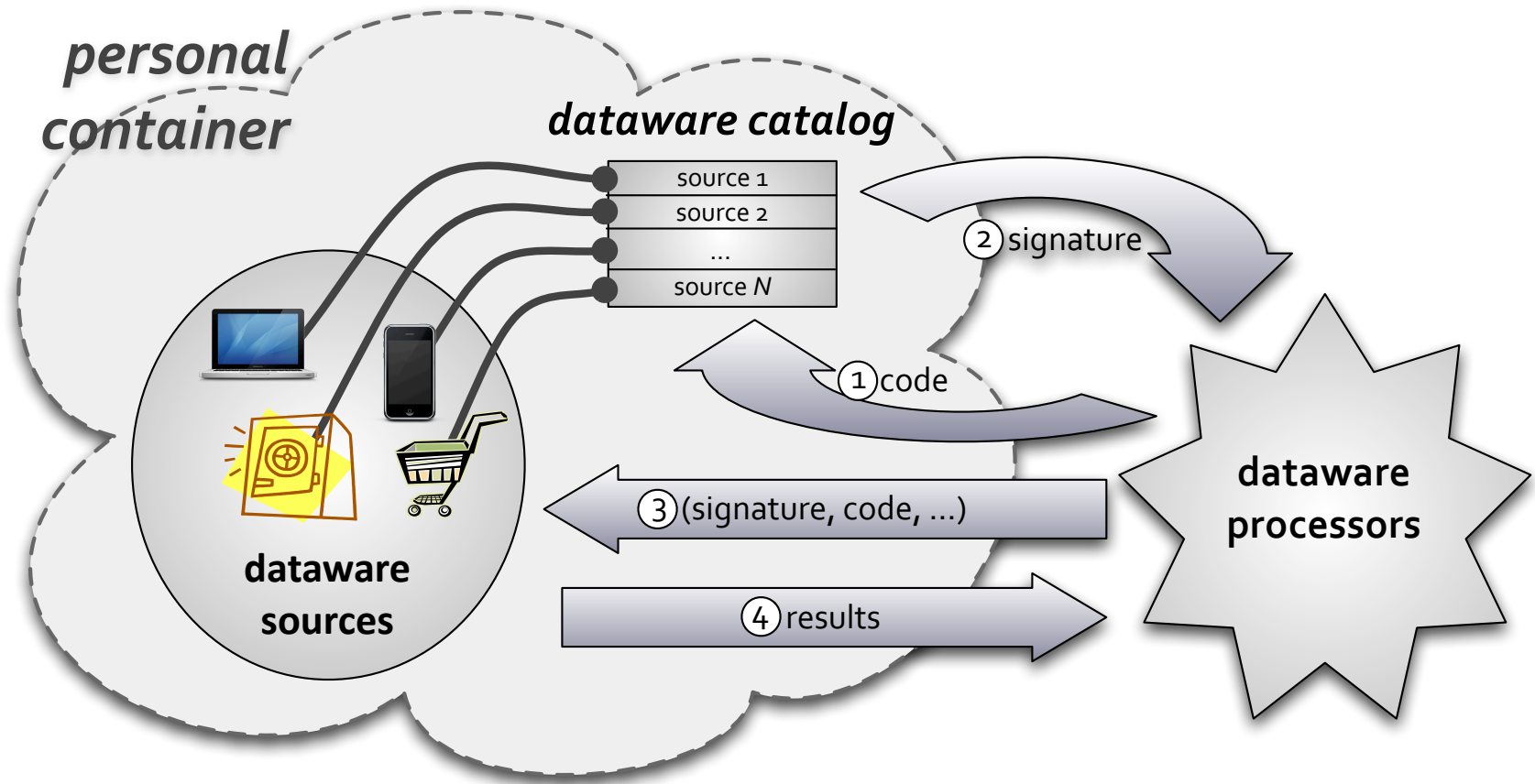
## Key Challenge:

How do we enable individuals to control collection and exploitation of both *their data* and *data about them*?

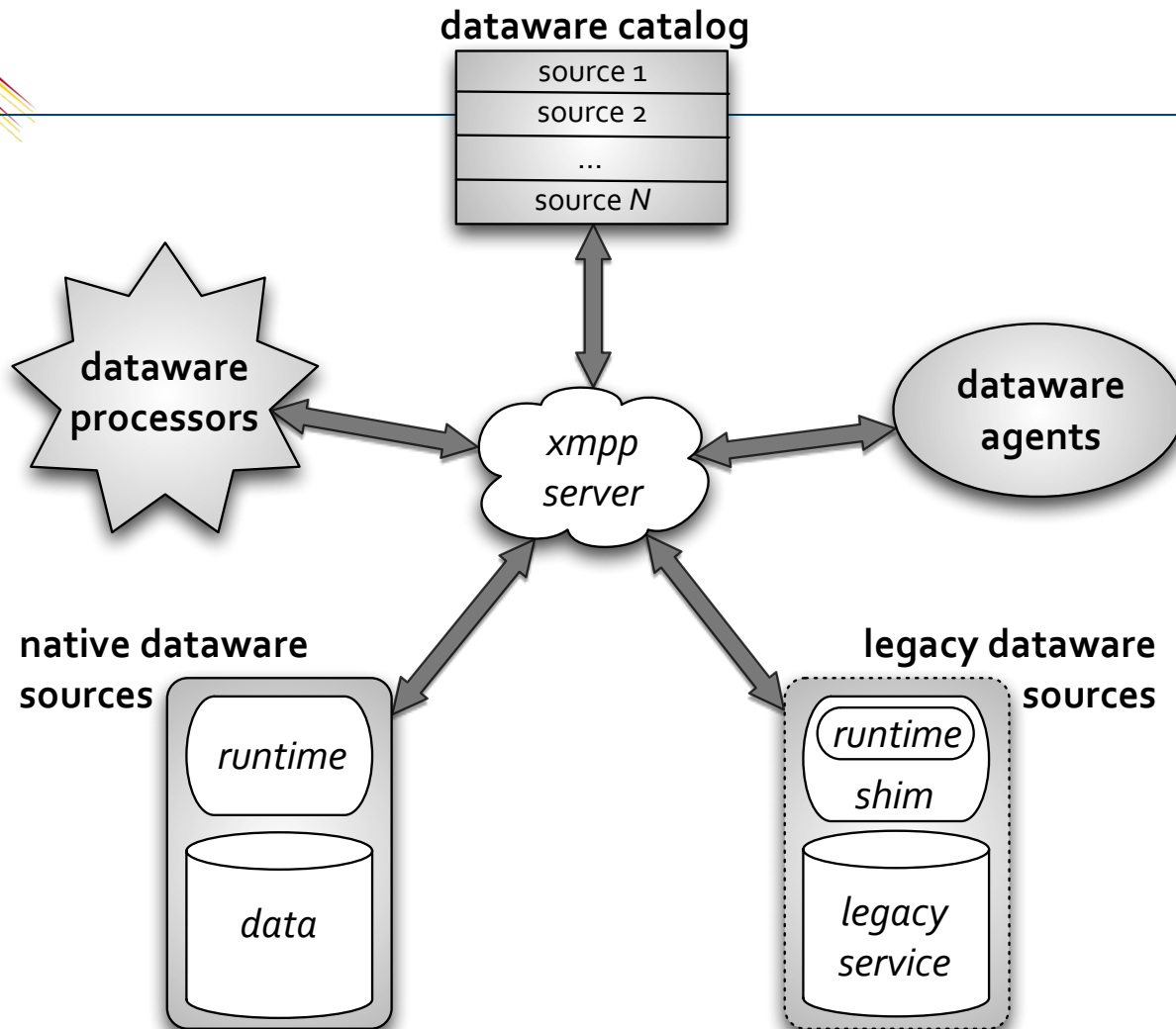
# Dataware



# Components



# Implementation





# Shim Login



my data sphere

Log in | Facebook - Mozilla Firefox

Facebook, Inc. (US) [https://ssl.facebook.com/login.php?api\\_key=1089099891703838&skip\\_ap](https://ssl.facebook.com/login.php?api_key=1089099891703838&skip_ap)

Facebook login

Log in to use your Facebook account with data-chant.

Email address:

Password:

☒ Keep me logged in

[Forgotten your password?](#)

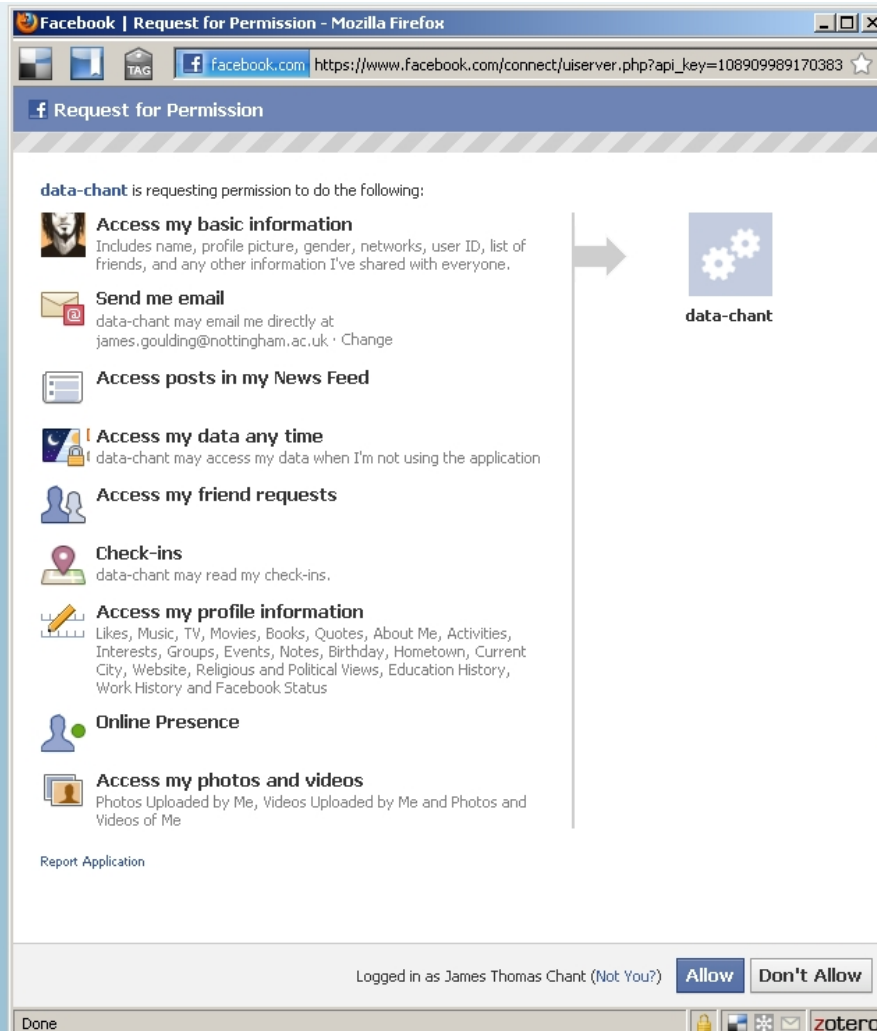
Sign up for Facebook

Log in Cancel

Done

login via  
facebook

# Shim Permission



data sphere

login via  
facebook

# Web Interface

[My Data](#)[Dataspheres](#)[History](#)[Statistics](#)[Datablog](#)[Login](#)[Register](#)

## my datasphere

### Your data

- » Live Feed
- » Full History
- » Statistics
- » Your Details

### Datasphere

- » Your Datasphere
- » Activity Summary
- » Recent Activity
- » Available

### Sharing

- » Accounts Summary
- » Your Requests
- » Sharing Stats

### Datasphere Source summary:



common name:  
**Facebook Datasphere Shim**

namespace:  
**ds:facebook**

subscription status:  
**ACTIVE**

[unsubscribe](#)

1 2 next »

### Fri 19 November



New Facebook "status" post: "I think we are alive alive o" 7:38pm

### Thu 18 November



In your Facebook bio, your "quotes" has changed to "Quotation changed again so there! And again!" 9:47am



In your Facebook bio, your "birthday" has changed to "03/10/1995" 9:47am

### Wed 17 November



New Facebook "status" post: "mighty booooooshka" 5:05pm

### Tue 16 November



New Facebook "status" post: "test" 2:10pm



In your Facebook bio, your "bio" has been removed. 2:02pm



In your Facebook bio, your "hometown" has changed to Olton, Solihull, United Kingdom (107013716004816) 2:01pm

### Your datasphere details:

name:

**James Goulding**

datasphere id:

**james.goulding@gmail.com**

### Datasphere you are running:



### Datasphere requests:

### Available datasphere:



[see all...](#)

```

ec2-user@ip-10-48-138-30:~$ java -jar datasphere.jar -t 8080 -p critic41
2010-11-25 10:57:40 [CONFIG] --- DSCatalog: loading configuration file... [SUCCESS]
2010-11-25 10:57:40 [FINE] --- DSDataManager: admin password has been specified via the command line... [SUCCESS]
2010-11-25 10:57:40 [INFO] --- DSCatalog: Establishing database connection and consistency...
2010-11-25 10:57:40 [INFO] --- DSDatabaseManager: Connecting to database for persistence... [SUCCESS]
2010-11-25 10:57:40 [INFO] --- DSDataManager: Checking System Table integrity... [SUCCESS]
2010-11-25 10:57:40 [INFO] --- DSCatalog: Attempting to start server components...
2010-11-25 10:57:41 [INFO] --- DSCatalog: HTTP server setup...[SUCCESS]
2010-11-25 10:57:41 [INFO] --- DSChatServer: Starting the internal XMPP server on port 5222... [SUCCESS]
2010-11-25 10:57:41 [INFO] --- DSChatServer: Creating XMPP bots for 2 clients...
2010-11-25 10:57:43 [FINE] --- DSClientBot: connecting bot for client [testreceiver@jabber.org]... [SUCCESS]
2010-11-25 10:57:43 [FINER] --- DSPresenceListener: [testreceiver@jabber.org] Dataware available for (testreceiver@jabber.org/datasphere)
2010-11-25 10:57:44 [FINE] --- DSClientBot: connecting bot for client [james.goulding@gmail.com]... [SUCCESS]
2010-11-25 10:57:44 [FINE] --- DSChatServer: 2 XMPP bots have been connected... [SUCCESS]
2010-11-25 10:57:44 [INFO] --- DSCatalog: XMPP server setup...[SUCCESS]
2010-11-25 10:57:44 [INFO] Datasphere setup and ready for FULL service...
2010-11-25 10:57:44 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware available for (j-eliza@appspot.com/bot)
2010-11-25 10:57:44 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware available for (data-chant@appspot.com/bot)
2010-11-25 10:57:44 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware unavailable for (testsender@jabber.org)
2010-11-25 10:59:25 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware available for (james.goulding@gmail.com/Talk.v10498203F13)
2010-11-25 10:59:25 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware available for (james.goulding@gmail.com/Talk.v10498203F13)
2010-11-25 10:59:25 [FINER] --- DSPresenceListener: [james.goulding@gmail.com] Dataware available for (james.goulding@gmail.com/Talk.v10498203F13)
2010-11-25 11:01:53 [FINE] --- DSUpdateListener: [james.goulding@gmail.com] <DSUpdate><namespace>ds:facebook</namespace><primaryTag>ds:fb:post</primaryTag>
<description>New Facebook "status" post: "what the world needs now..."</description><crud>create</crud><ctime>1290682887000</ctime><ftime>1290682913014</ftime>
<rtime>null</rtime><total>0</total><tag>ds:communication</tag><meta><message>what the world needs now..."</message><category>status</category></meta></DSUpdate>
e>

```

ec2-46-51-150-243.eu-west-1.compute.amazonaws.com:8080/user\_history

Docs Webmail 31 Calendar Gmail Facebook Outlook Jambon SIGWEB Flickr

# my datasphere

**Your data**

- » Live Feed
- » Full History
- » Statistics
- » Your Details

**Dataware**

- » Your Dataware
- » Activity Summary
- » Recent Activity

**Thu 25 November**

- New Facebook "status" post: "what the world needs now..." 11:01am
- New Facebook "status" post: "Another glorious sunny morning..." 9:25am

**Fri 19 November**

- New Facebook "status" post: "I think we are alive alive o" 7:38pm

**Thu 18 November**

data-chant@appspot.com

Send voicemail Call Send Files Email

data-chant: {"namespace":"ds:facebook","primaryTag":"ds:fb:post","tags":["ds:communication"],"ctime":"1290682887000","ftime":"1290682913014","crud":"create","description":"New Facebook \"status\" post: \"what the world needs now...\", \"meta\": {\"message\":\"what the world needs now...\", \"category\":\"status\"}}

data-chant is online.  
Sent at 11:01 AM on Thursday

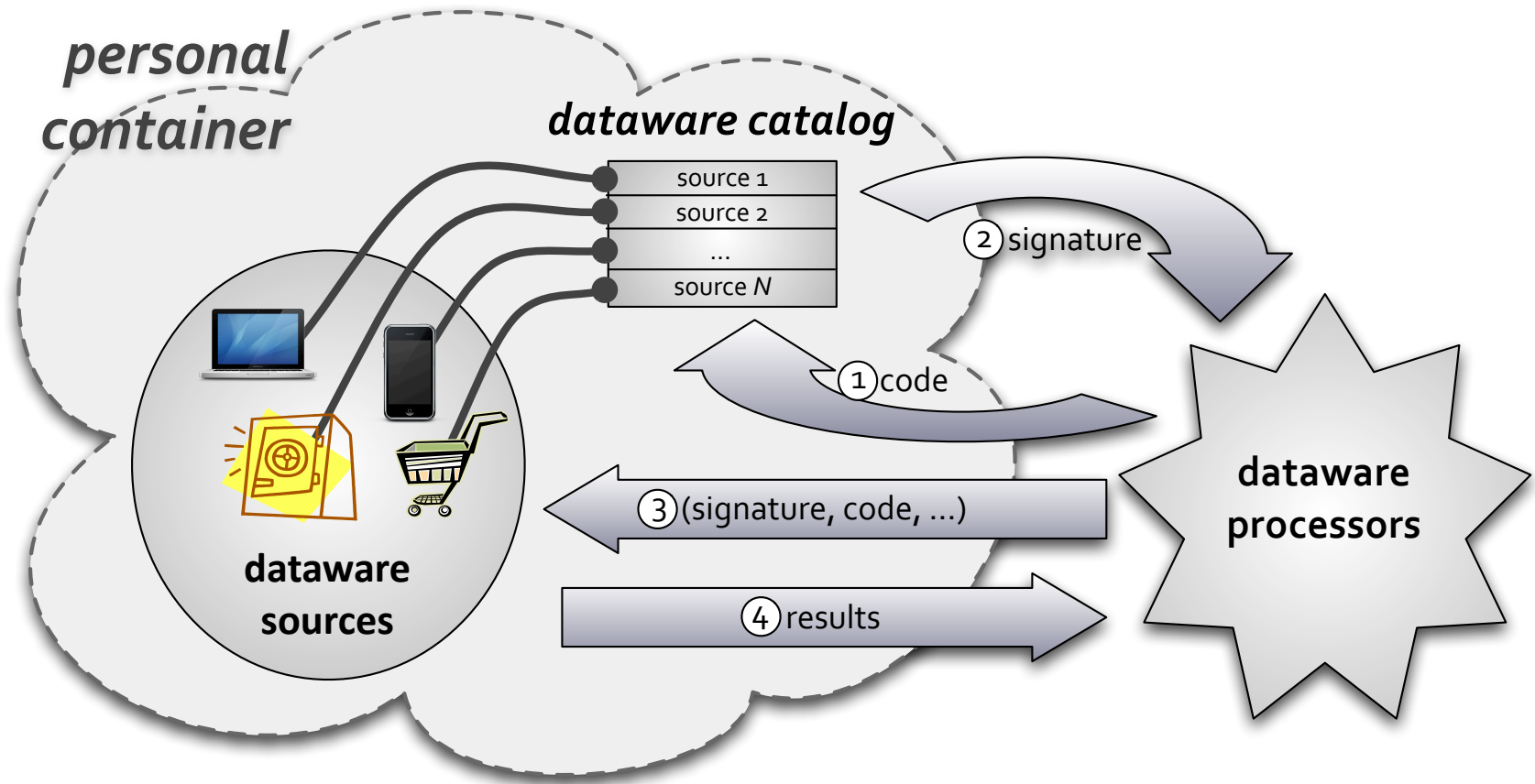
Your datasphere details:

name:  
**James Goulding**  
datasphere id:  
**james.goulding@gmail.com**

Dataware you are running:



# Processing





# Applications

---

- An application is a list of source schemas and signed components
  - Each element defines processing occurring on some data
- Each component starts out as
  - *Code*, describing the computation
  - *Schedule*, against which the code is executed
  - *Result schema*, describing the results
  - *Result destination*, JID to which results are sent



# Application Signing

---

- Processor presents application to catalog
- Catalog applies user policy to determine whether to sign
  - Static policy concerning data to be processed, source, &c
  - Or interacting directly with user via web, chat, &c
- To sign, replace each component with one or more:
  - *Result schema, schedule, code*, as before
  - *Source JID*, specifying data source
  - *Destination JID*, to which results are sent
  - *Signature*, an HMAC for the component
- ...And then return to the submitting processor



# Application Execution

---

- Processor distributes signed components to given sources
- Sources then:
  - Install components, notifying catalog
  - Execute components in runtime, as per schedule
  - Send results to specified destination
- During execution, the runtime monitors code
  - Enables periodic reporting back to user
  - Exceptions raised if code operates outside boundaries





# JOINing Data

---

- A clear issue is how to deal with data JOINS
  - Recall: we're avoiding giving access to raw data
  - We must also avoid giving access to intermediate results
- Setup *trusted third-party sources*
  - Have no data of their own
  - Act as destination for other code
  - May host subsequent application components
- Makes sense that the catalog serves as at least one such



# Summary

---

- So conceptually, a *dataware application* is:
  - A network of running components,
  - Each processing your personal data,
  - In a manner acceptable to you
- Your *dataware* implements your *personal container*:
  - Defining APIs to your data,
  - Enabling third-parties to process your data,
  - While ensuring *you* retain control,
  - And they *don't* get copies of your data
    - (unless you want them to!)



# Evolution or Revolution?

- One way to put this in context is via Van Jacobson's *content centric networking*:
  - Telephones care about building paths (not calls)
  - Internet cared about connections (not data)
  - Content centric networking cares about data (not results)
- ...so, finally,
- *Dataware*, cares about results
  - Computation is as mobile as data, if not more so
  - (Anil) cf. datacentre computing, map-reduce, &c



# Status

---

- First version of catalog built
- First legacy source built
  - Facebook, using Graph API
  - Twitter, Google Gdata on their way
- Now prototyping
  - Dataware Processor and first applications, and
  - Native source
- Challenge: how to build protocols for trusted aggregation
  - Would like to guarantee *at-least-N* contributors