# Grid computing and computational statistics
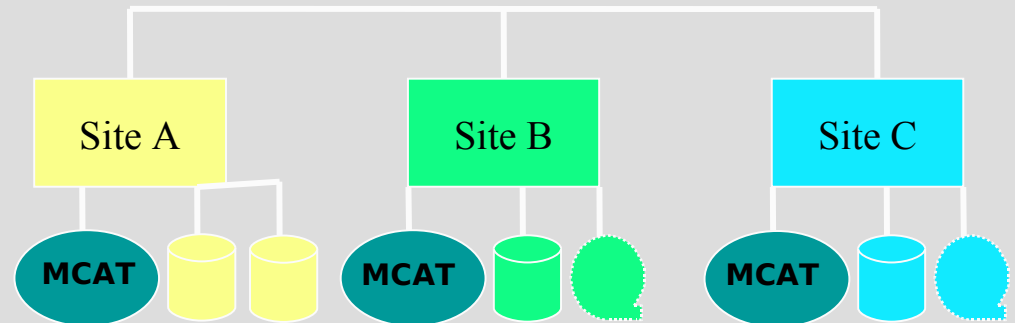
Nigel Sim
VeRG/SIDA

# Three focuses

- Data Grid library
  - PGL
- Computational infrastructure
  - Grid Agents
- Data mining:
  - Ensembles
  - Permutation testing

# The Data Grid

- Federation between sites

- Meta data directly attached to data

- Multiple data locations

- Virtualisation!

# Personal Grid Library

- Implementation of a web-based digital library
- Provides an interface to templated metadata manipulation
  - Define object classes (Images, reports, datasets…)
  - Define appropriate metadata tags
  - Define metadata display template
- Exists by placing a metadata template file in a directory – used by interface to render library

# Personal Grid Library cont...

# Grid Agents

- Cross platform data flow execution system
- Can utilise R, compiled libraries, legacy apps, etc
- Handles in-flight data conversion

# Grid Agents cont...

- Scalable work placement
  – Master-slave using workflows
- Utilise multiple CPUs on multiple machines



*R* monitoring interface

SOAP
RML

**MS** – Hybrid master-slave node
**M** – Master node
**S** – Slave node

Site B

# Data mining

- Mainly focusing on spectral data
  - NIR (Quality control)
  - Mass spec (blood proteins)

- Computationally intensive
  - Cross validation
  - Permutation test
  - Optimisation (EA)



```
Variable reduction
        ↓
   Build models
        ↓
  Validate models
        ↓
   Final models
```

# Permutation tests

- A non-parametric test
- Can be adapted to many situations
  - Model testing
  - Variable testing



Data → Build Model

*Randomise data labels*

Cross validate Model

NULL Distribution

Predictive performance

# Ensembles

- Linear combination of many models

$$F(x) = \sum_{i=1}^{N} w_i f_i(x)$$

$$f(x) = \text{Any model}$$

- Bootstrap using Lasso

- Optimise using evolutionary algorithms

•Stepwise regression
•Some coefficients may end up zero

|    | age    | sex     | bmi    | map    | tc      | ldl    | hdl     |
|----|--------|---------|--------|--------|---------|--------|---------|
| 0  | 0      | 0       | 0      | 0      | 0       | 0      | 0       |
| 1  | 0      | 0       | 60.12  | 0      | 0       | 0      | 0       |
| 2  | 0      | 0       | 361.89 | 0      | 0       | 0      | 0       |
| 3  | 0      | 0       | 434.76 | 79.24  | 0       | 0      | 0       |
| 4  | 0      | 0       | 505.66 | 191.27 | 0       | 0      | -114.1  |
| 5  | 0      | -74.92  | 511.35 | 234.15 | 0       | 0      | -169.71 |
|    |        | ...     | ...    | ...    | ...     | ...    |         |
| 11 | -7.01  | -237.1  | 521.08 | 321.55 | -580.44 | 313.86 | 0       |
| 12 | -10.01 | -239.82 | 519.84 | 324.39 | -792.18 | 476.75 | 101.04  |

# Combined use case

1) Ingest raw data files

2) Start a Grid Agents workflow

3) Results stored as metadata

4) Viewed using PGL