# Newton's ideas and methods are preserved forever: how about yours?

*Marco Roos*, Kristina Hettne, Jun Zhao, Mark Thompson

Cloud and Workflows for Reproducible Bioinformatics

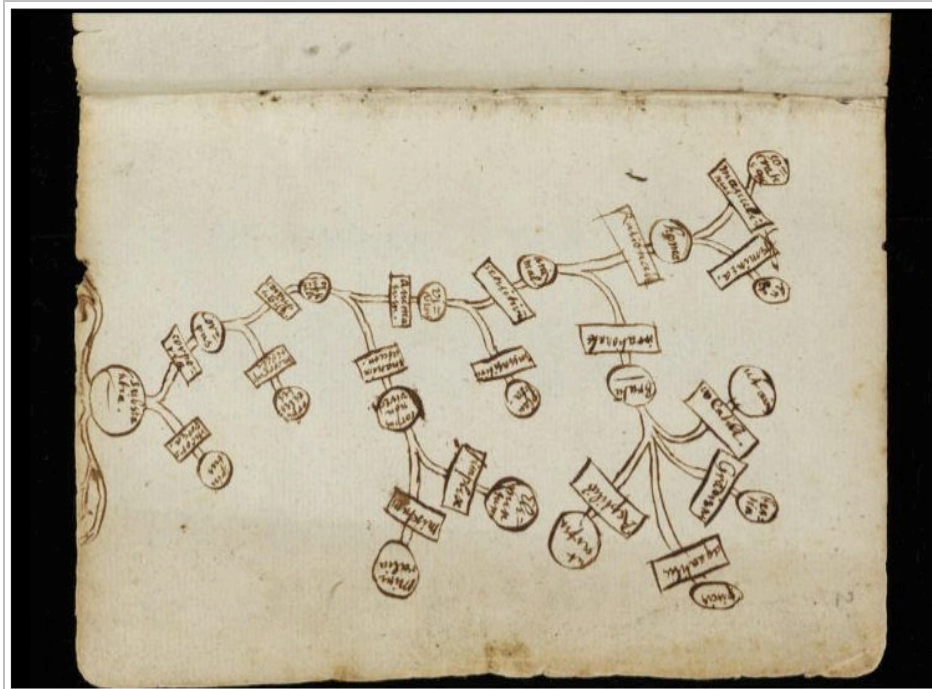Shenzhen, December 19, 2012

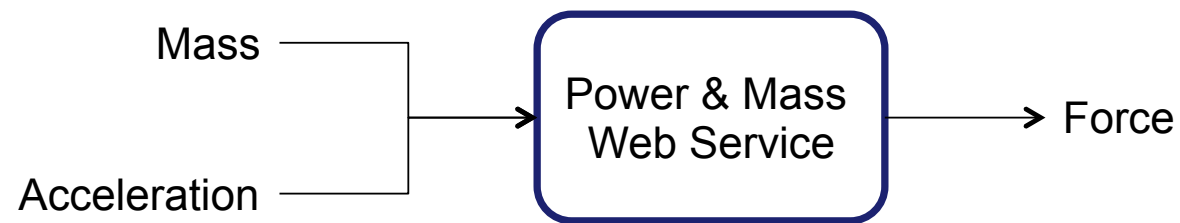**Newton's Books Scanned, Posted Online**

Officials at the Cambridge University, in the United Kingdom, announce that they recently made most of the books written by Sir Isaac Newton available online. These works, some of the most important scientific documents ever, are now available to the general public here. Newton was the Lucasian Ch... [ read more >> ]

**Image comment:** This book includes many notes from Isaac Newton's studies and, increasingly, his own explorations into mathematics, physics and metaphysics
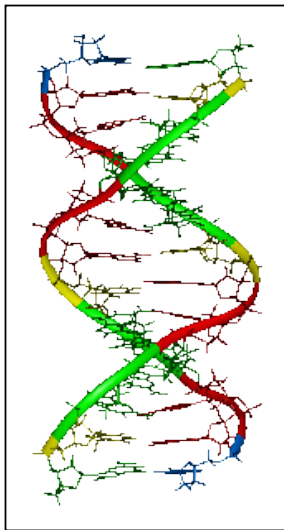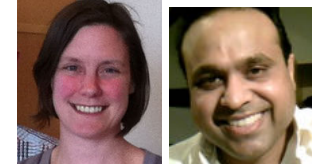**Image credits:** © Cambridge University Library

# Reproduced workflows

Mass ─────────────┐
                  │
                  ├──→ ┌─────────────────┐
                  │    │  Power & Mass   │ ──→ Force
Acceleration ─────┘    │  Web Service    │
                       └─────────────────┘

# Case study

Bioinformatics analysis of Metabolic Syndrome
Kristina Hettne, Harish Dharuri

**Genome Wide Association Studies**

*What is the genetic basis for the diseases associated with Metabolic Syndrome?*

# Reproducible Science

**Research article**

**Open Access**

## Mutant huntingtin activates Nrf2-responsive genes and impairs dopamine synthesis in a PC12 model of Huntington's disease

### Abstract

**Background:** Huntington's disease is a progressive autosomal dominant neurodegenerative disorder that is caused by a CAG repeat expansion in the HD or Huntington's disease gene. Although micro array studies on patient and animal tissue provide valuable information, the primary effect of mutant huntingtin will inevitably be masked by secondary processes in advanced stages of the disease. Thus, cell models are instrumental to study early, direct effects of mutant huntingtin. mRNA changes were studied in an inducible PC12 model of Huntington's disease, before and after aggregates became visible, to identify groups of genes that could play a role in the early pathology of Huntington's disease.

**Results:** Before aggregation, up-regulation of gene expression predominated, while after aggregates became visible, down-regulation and up-regulation occurred to the same extent. After aggregates became visible there was a down-regulation of dopamine biosynthesis genes accompanied by down-regulation of dopamine levels in culture, indicating the utility of this model to identify functionally relevant pathways. Furthermore, genes of the anti-oxidant Nrf2-ARE pathway were up-regulated, possibly as a protective mechanism. In parallel, we discovered alterations in genes which may result in increased oxidative stress and damage.

**Conclusion:** Up-regulation of gene expression may be more important in HD pathology than previously appreciated. In addition, given the pathogenic impact of oxidative stress and neuroinflammation, the Nrf2-ARE signaling pathway constitutes a new attractive therapeutic target for HD.

## Methods

### Cell culture

Inducible rat PC12 cell lines expressing an exon 1 fragment of huntingtin with 23 (Q23) or 74 (Q74) glutamine repeats fused to the Green Fluorescent Protein (GFP), [11,12] were cultured in standard high glucose Dulbecco's modified Eagle's medium (DMEM, Invitrogen Life Technologies, Carlsbad, USA) supplemented with 100 U/ml penicillin/streptomycin (Invitrogen Life Technologies), 2 mM L-glutamine (Invitrogen Life Technologies), 10% heat-inactivated horse serum (Invitrogen Life Technologies), 5% Tet-approved heat inactivated fetal bovine serum (Clontech, Palo Alto, USA), 100 µg/ml G418 (Invitrogen Life Technologies) and 75 µg/ml hygromycin (Invitrogen Life Technologies) at 37°C and 10% $CO_2$. Cells were induced with 1 µg/ml doxycycline (dox, Clontech) and harvested on day 0 (uninduced cells), 1 day (when only a few cells expressing mutant huntingtin contain aggregates) and 5 days (when nearly all cells expressing mutant huntingtin contain aggregates) [12]. The same culture conditions were used for PC12 cells without a construct, to eliminate the effect of doxycycline treatment on gene expression.

### Hybridization design

For each construct, we performed duplicate experiments with 2 independent cell lines for each construct (biological replicates). Furthermore, from each cell line, two separate RNA isolations were performed (technical
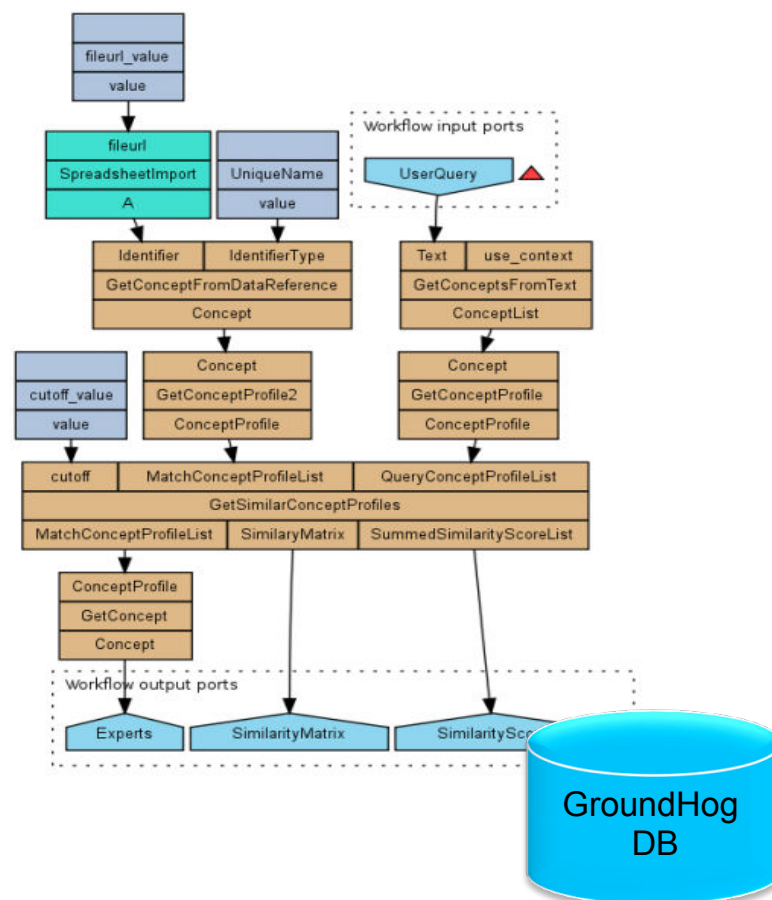
## Preservation for the wet laboratory scientist

# Reproducible Science?

What is the digital equivalent?

Is it equally good?

Can we do better?
- *or worse?*



GroundHog DB

# What is our incentive?

## Nobility

Good Reproducible Science



## Greater Good

Serve the public

# What is our incentive?



**Fame and Glory**

Getting on with it…

I'll be the first in Nature

## CHALLENGE

## Stimulate preservation and reproducibility while speeding up the research process

**Workflow 4Ever**

Enhance the research cycle
**What slows us down?**

SW

**Research Question**

**Find** Methods and Data, + their Owners

**Get** Methods and Data

**Understand** Methods and Data

**Format** (Align) Data

**Design** the Analysis

**Compute**

**Interpret** Results

**Publish**

nbic    Open PHACTS    Concept Web Alliance    LUMC

# Bottlenecks

- Loosing track of what you did

- Messy storage

- Preparing material for a publication

- Understanding the computational procedure

- Communication with (non-technical) colleagues

- Keeping tools working

- Getting credit for digital results outside of traditional publications

# Getting on with workflows

**Workflow 4Ever**

~~Monolithic Tool →~~

Web Services → Workflows → (Web) Tool

Example: Anni 2.0 → Anni workflows

http://workflow.biosemantics.org/t2web/workflow/2725

![Workflow 4Ever logo]

# Digital Repository
## myExperiment.org

# The recipes store

- Find workflows
- Share workflows & files
- Find people
- Build communities
- Publish packages
- Tag workflows
- Score, rate, comment

# Instructions for workflow authors
## 10 Best Practices for creating workflows

1. Make a sketch workflow
2. Use modules
3. Think about the output
4. Provide example inputs and outputs
5. Annotate
6. Test execution from outside local environment
7. Choose services carefully
8. Reuse existing workflows
9. Advertise
10. Maintain

**10/10**

Reproducible Science
**Is a workflow sufficient?**

**Useful Preservation
=
Understandable Objects**

Reproduce, Reuse, Repurpose, Repair, ...

*What is this doing?*

# myExperiment Packs



"Experiment sketch"

"Workflow to display supporting documents"

MedLine abstracts until 2009

"SNPs from GWAS"

Hypothesis document

my experiment

**Research Object Model**
Aggregation and Annotation Model for Digital Methods

http://wf4ever.github.com/ro/

# Research Object (RO) Model

RO = ORE + AO + vocabularies

Object Re-use and Exchange (OAI-ORE)

    Describes aggregations of resources:

    data, metadata, papers, *etc.*

Annotation Ontology (AO)

    Associates RDF metadata descriptions with resources

Generic and domain-specific vocabularies

    Used in annotation bodies to provide information about resources (types, dependencies, descriptions, *etc.*)

Builds on RDF, leading to RDF as a natural implementation choice

Model specification: http://wf4ever.github.com/ro/

# Research Object Model

# Research Object: "Hello World"



https://github.com/wf4ever/ro-catalogue/tree/master/v0.1/HelloWorld

# Help organize the materials and methods of computational analysis
## Research Object Portal

*Materials & Methods of Metabolic Syndrome Analysis*
Kristina Hettne
Harish Dharuri

http://sandbox.wf4ever-project.org/portal

# Expected on myExperiment

**Research Objects inside!**

- Packs more prominent

- Start a pack when you upload a workflow

- Upload wizards, pack management, export

- Checklists, automated star ratings

- Add workflow runs and example data

- Sticky annotations



RO-enabled myExperiment mockup

# Nanopub.org

# Examples

# Examples in RDF format

# Validator

# Example: LOVD

# Nanopublications of Genetic Variations visualized on the genome

Zuotian Tatum, Jesse van Dam

**Other Sources**

**Nanopublication Store**

**Other Tools**

31

# Summary (1/2)

- Preservation under the hood of digital research tools

- Research Object Model: annotated aggregates

- Nanopublication: fine-grained digital credit
  **Check Nanopub.org to stay updated**

# Summary (2/2)

- Semantic Web for exchange and interoperability

- In progress: RO-enabling myExperiment
  **Watch myExperiment.org in 2013!**

- Plans to RO-enable Taverna, Galaxy, GenomeSpace

# Acknowledgements

35

# Thank you for your attention

http://biosemantics.org

# Reproducible Science

Research article

**Open Access**

**Mutant huntingtin activates Nrf2-responsive genes and impairs dopamine synthesis in a PC12 model of Huntington's disease**

Willeke MC van Roon-Mom*[1], Barry A Pepers[1,2], Peter AC 't Hoen[1], Carola ACM Verwijmeren[1], Johan T den Dunnen[1,3], Josephine C Dorsman and GertJan B van Ommen[1]
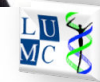
Preserved materials and methods for the 'wet laboratory' scientist

## Methods

### Cell culture

Inducible rat PC12 cell lines expressing an exon 1 fragment of huntingtin with 23 (Q23) or 74 (Q74) glutamine repeats fused to the Green Fluorescent Protein (GFP), [11,12] were cultured in standard high glucose Dulbecco's modified Eagle's medium (DMEM, Invitrogen Life Technologies, Carlsbad, USA) supplemented with 100 U/ml penicillin/streptomycin (Invitrogen Life Technologies), 2 mM L-glutamine (Invitrogen Life Technologies), 10% heat-inactivated horse serum (Invitrogen Life Technologies), 5% Tet-approved heat inactivated fetal bovine serum (Clontech, Palo Alto, USA), 100 µg/ml G418 (Invitrogen Life Technologies) and 75 µg/ml hygromycin (Invitrogen Life Technologies) at 37°C and 10% $CO_2$. Cells were induced with 1 µg/ml doxycycline (dox, Clontech) and harvested on day 0 (uninduced cells), 1 day (when only a few cells expressing mutant huntingtin contain aggregates) and 5 days (when nearly all cells expressing mutant huntingtin contain aggregates) [12]. The same culture conditions were used for PC12 cells without a construct, to eliminate the effect of doxycycline treatment on gene expression.
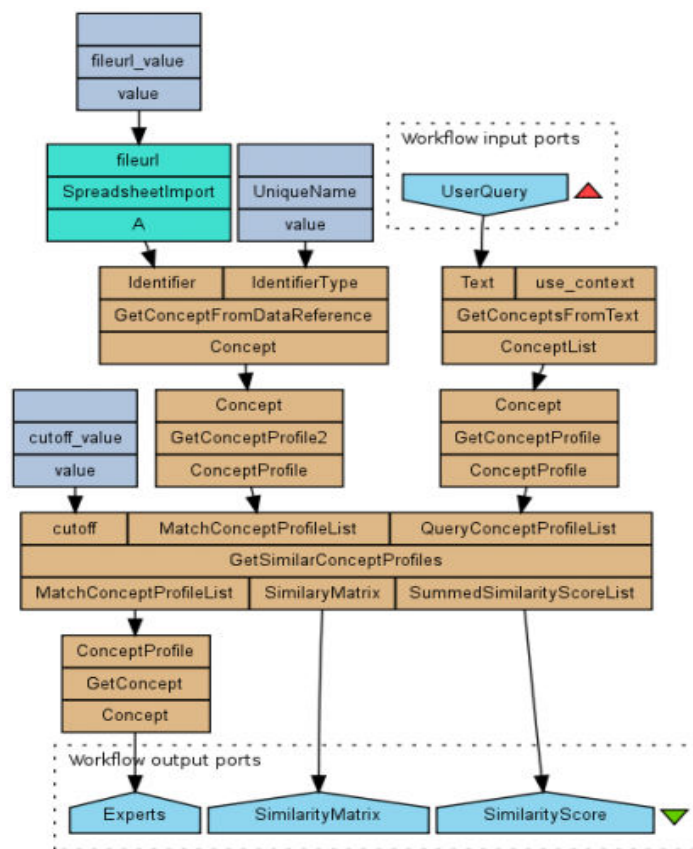
### Hybridization design

For each construct, we performed duplicate experiments with 2 independent cell lines for each construct (biological replicates). Furthermore, from each cell line, two separate RNA isolations were performed (technical

# Reproducible Science?

What is the digital
  equivalent?

Is it equally good?

Can we do better?
  *- or worse?*

# Reproducible Science

**What is the digital equivalent?**

Is it equally good?

Can we do better?
*– or worse?*

*Can you tell what this is doing?*

# What is our incentive?

## Nobility

Good Reproducible Science



## Greater Good

Serve the public

# What is our incentive?

**Fame and Glory**

Getting on with it…
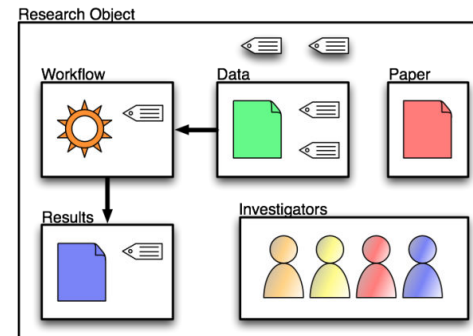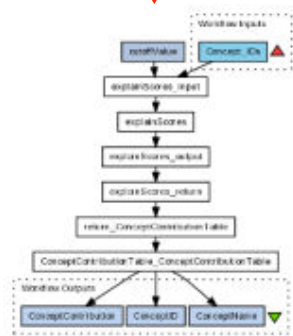
# Our aim

*'Useful' preservation*

**Support reproducibility
in tools and by guidelines that**
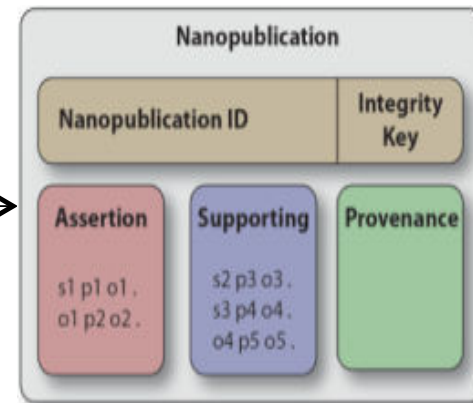
**speed up your research**
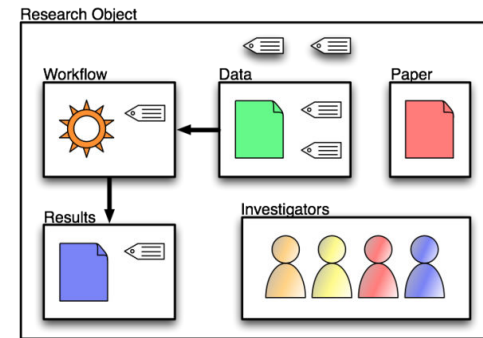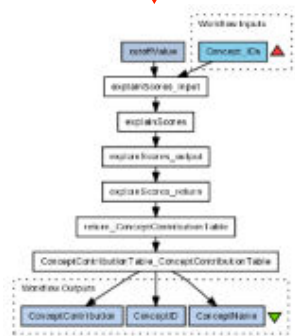
**get you acknowledgement**
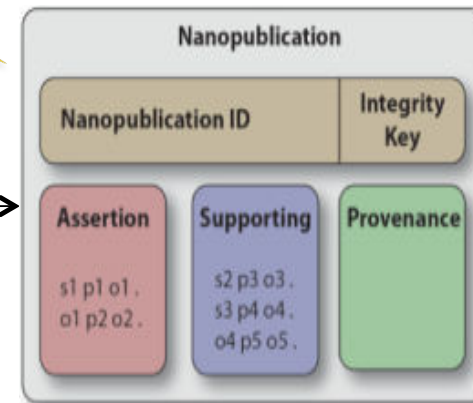
# Preservation

# Preservation

**Valuable for scientists**

**Digital Value**

**What?**

**How?**

**Research Results**

# Acknowledgements

**http://biosemantics.org/**

- Erik Schultes
- Andrew Gibson
- Reinout van Schouwen
- Kostas Karasavvas
- Kristina Hettne
- Harish Dharuri
- Eleni Mina
- Jesse van Dam
- Herman van Haagen
- Zuotian Tatum
- Johan den Dunnen
- Peter-Bram 't Hoen
- Barend Mons
- Gert-Jan van Ommen

- Paul Groth
- Frank van Harmelen

**Erasmus MC**
University Medical Center Rotterdam

- Erik van Mulligen
- Bharat Singh
- Jan Kors

- Christine Chichester
- Kees Burger - NBIC
- Spyros Kotoulas - VU
- Antonis Loizou - VU
- Valery Tkachenko - RSC
- Andra Waagmeester - Maastricht
- Sune Askjaer - Lundbeck
- Steve Pettifer - Manchester
- Lee Harland - Pfizer/CD
- Carina Haupt - Fraunhofer
- Colin Batchelor - RSC
- Miguel Vazquez - CNIO
- José María Fernández - CNIO
- Jahn Saito - Maastricht
- Andrew Gibson (Outside Expert) - Amsterdam
- Louis Wich - DTU