

The Automation of Science

**Ross D. King,
University of Manchester & AIST, ross.king@manchester.ac.uk**



Scientific Discovery

Technology Drivers

- n Improved computer hardware:
 - faster processors, more processors, GPUs, ...
- n Improved data availability:
 - computers recording almost everything, deep data, ...
- n Improved computer software:
 - new machine learning methods, deep mining, ...

Artificial Intelligence (AI)

- n There have been multiple AI hype cycles, but this time it seems different.
- n AI, especially machine learning, is now the hottest technology on the planet. Speed of Advance has surprised me.
- n Machine Learning is the core technology of Google, Facebook, Amazon, ...Tencent, Alibaba, Baidu, ...

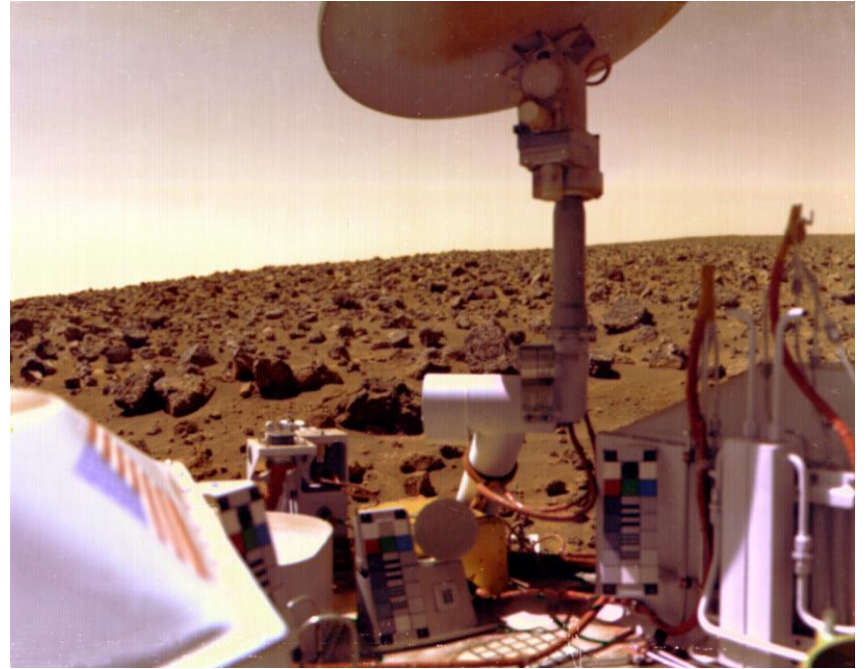
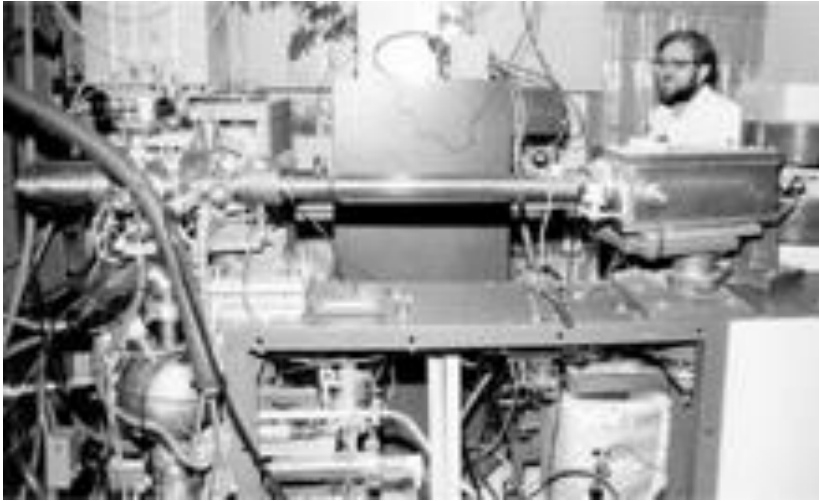
AI Systems have Superhuman Scientific Reasoning Powers

- n Flawlessly remember vast numbers of facts
- n Execute flawless logical reasoning
- n Execute optimal probabilistic reasoning,
- n Learn more rationally than humans
- n Learn from vast amounts of data
- n Extract information from millions of scientific papers.
- n Etc.

Scientific Discovery

- n Scientific problems are abstract, but involve the real-world.
- n Scientific problems are restricted in scope – no need to know about “Cabbages and Kings”.
- n Nature is honest – no malicious agents.
- n Nature is a worthy object of our study.
- n The generation of scientific knowledge is a public good.

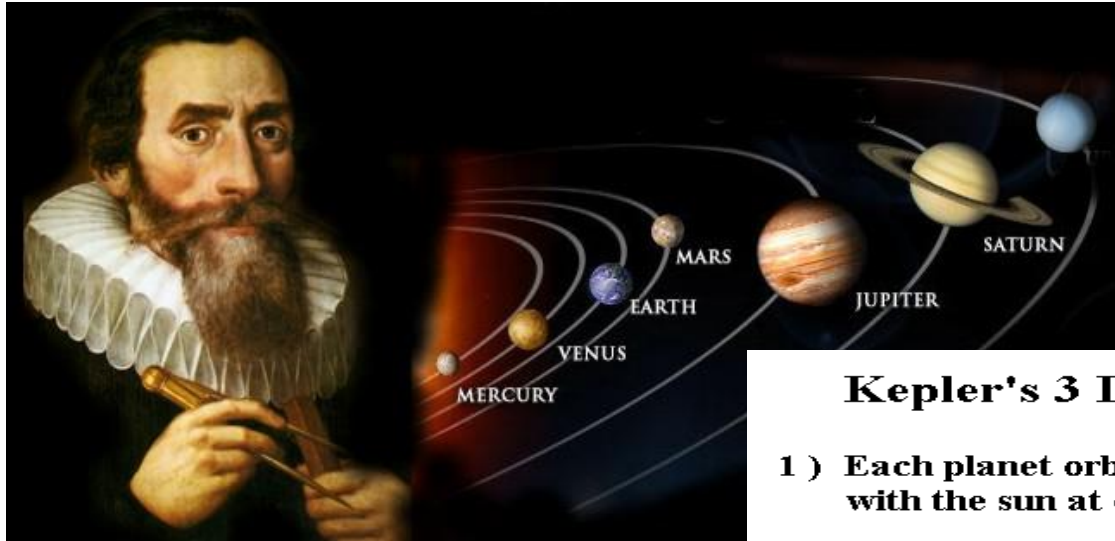
Meta-Dendral



Analysis of mass-spectrometry data.

Joshua Lederburg, Ed. Feigenbaum, Bruce Buchanan,
Karl Djerassi, *et al.* 1960-70s.

Bacon



Kepler's 3 Laws of Planetary Motion

- 1) Each planet orbits the sun in an elliptical path with the sun at one focus**
- 2) The radius vector (from sun to planet) sweeps out equal areas in equal time intervals**
- 3) The square of the period is proportional to the cube of the semi-major axis of the orbit**

$$\text{i.e. } T^2 = k a^3 \quad \text{for some constant } k$$

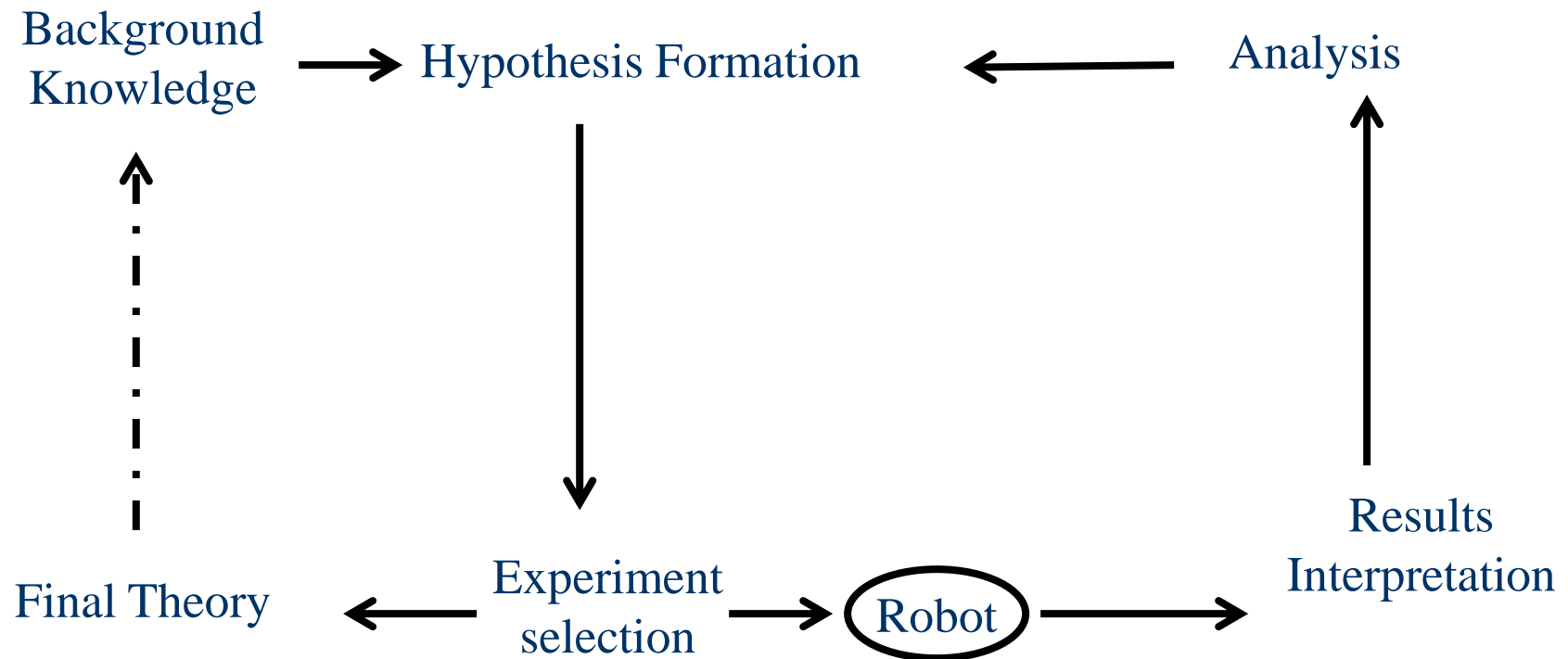
Figure 11.1

Rediscovering physics and chemistry: Langley, Bradshaw, Simon (1979).

Robot Scientists

The Concept of a Robot Scientist

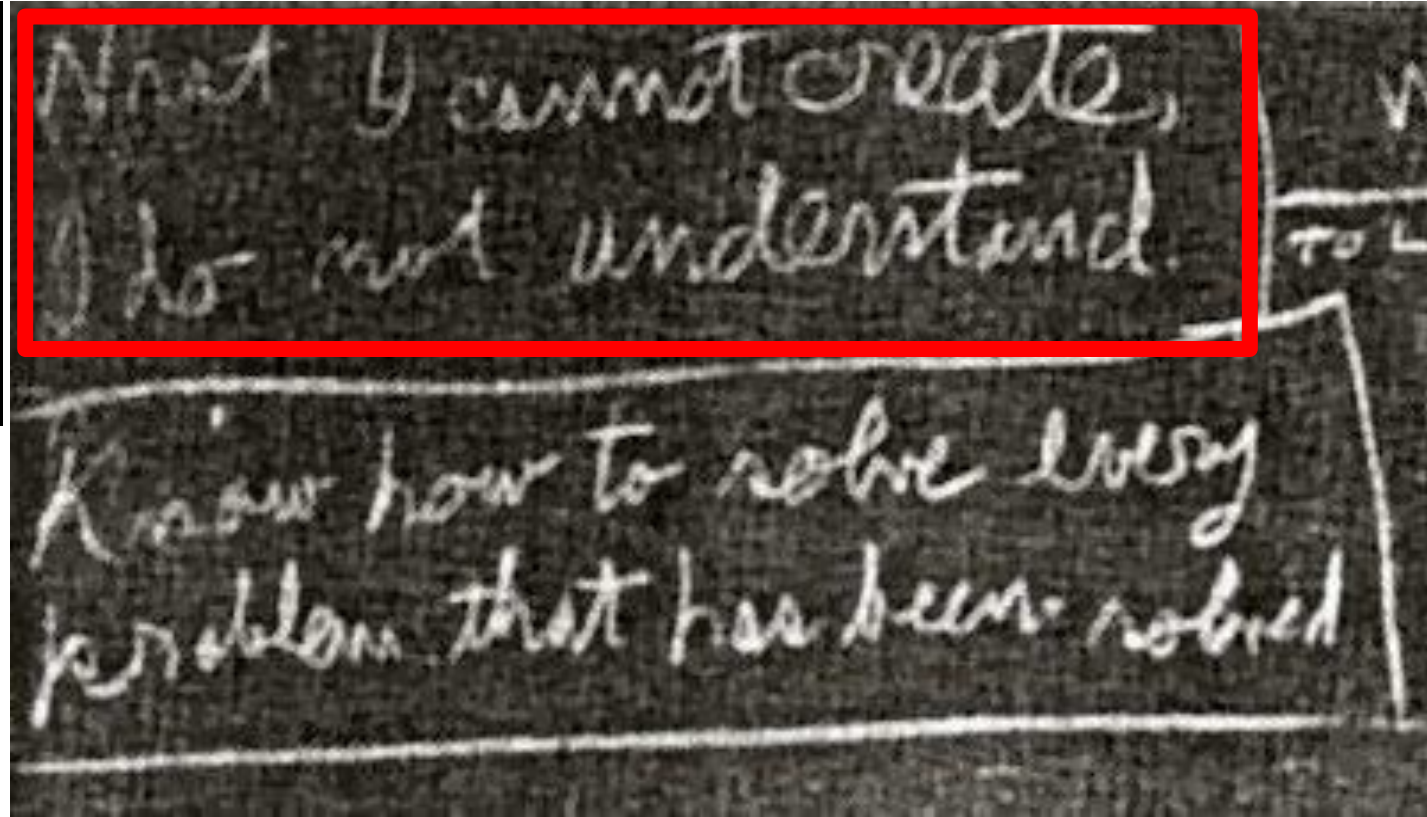
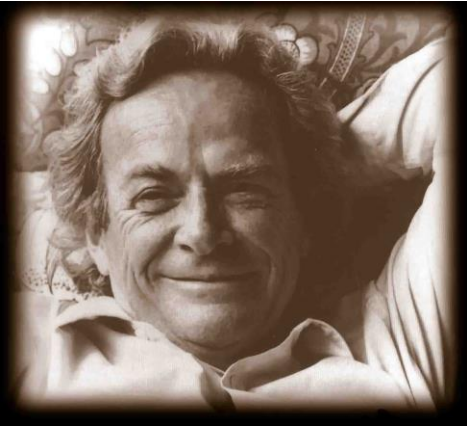
Computer systems capable of originating their own experiments, physically executing them, interpreting the results, and then repeating the cycle.



Motivation: Philosophical

- n What is Science?
- n The question whether it is possible to automate scientific discovery seems to me central to understanding science.
- n There is a strong philosophical position which holds that we do not fully understand a phenomenon unless we can make a machine which reproduces it.

Richard Feynman's Blackboard



“What I cannot create, I do not understand”

Motivation: Technological

- n Robot Scientists have the potential to increase the productivity of science. They can work cheaper, faster, more accurately, and longer than humans. They can also be easily multiplied.
 - *Enabling the high-throughput testing of hypotheses.*
- n Robot Scientists have the potential to improve the quality of science.
 - *by enabling the description of experiments in greater detail and semantic clarity.*

Robot Scientist Timeline

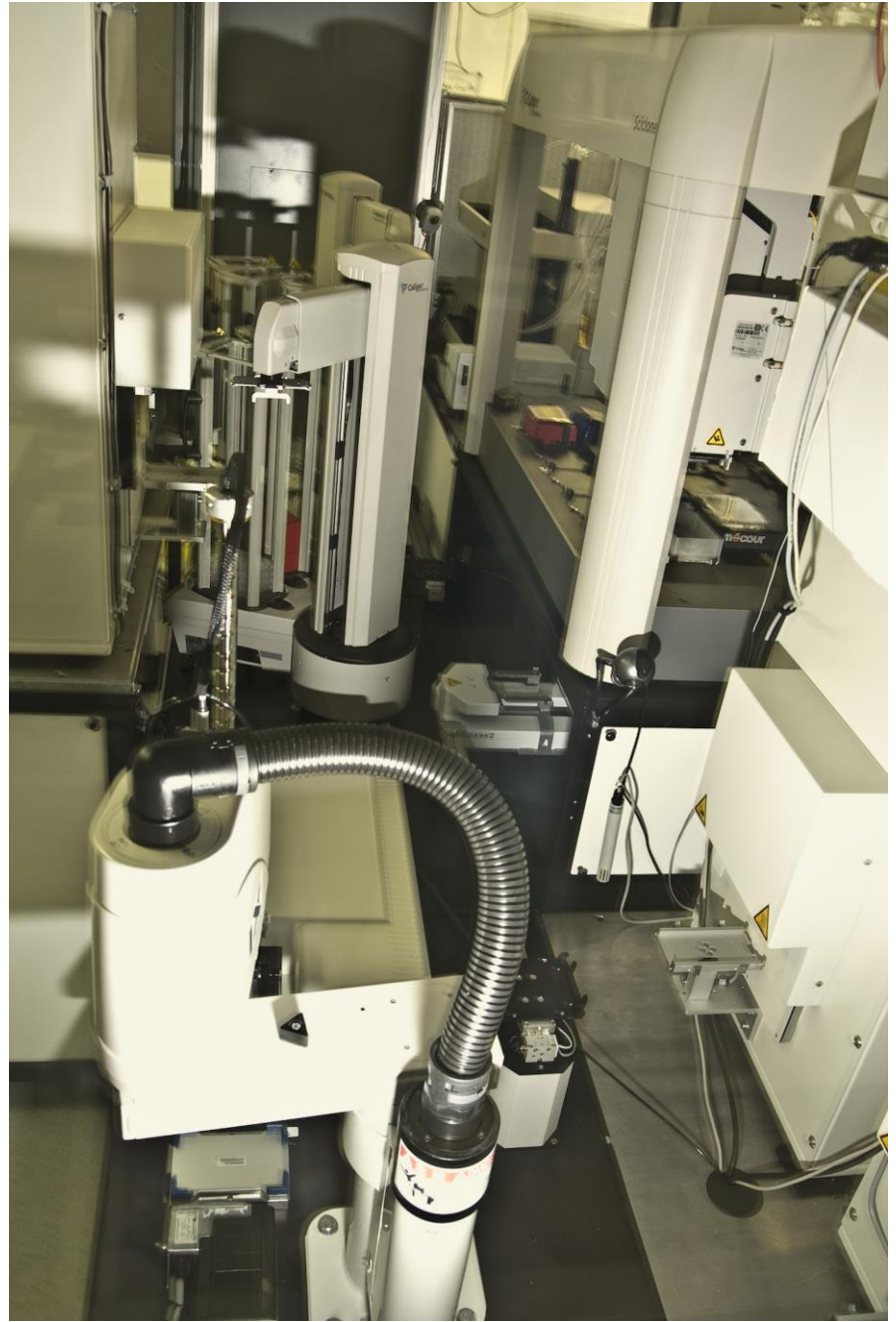
- n 1999-2004 Initial Robot Scientist Project
 - Limited Hardware: Collaboration with Douglas Kell (Aber Biology), Steve Oliver (Manchester), Stephen Muggleton (Imperial)

King et al. (2004) *Nature*, 427, 247-252
- n 2004-2011 Adam – Yeast Functional Genomics
 - Sophisticated Laboratory Automation: Collaboration with Steve Oliver (Cambridge).

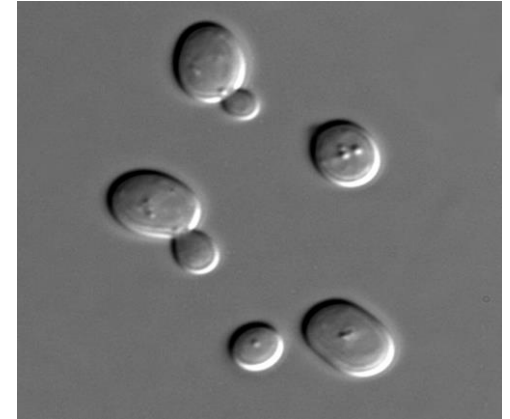
King et al. (2009) *Science*, 324, 85-89
- n 2008-2015 Eve – Drug Design for Tropical Diseases
 - Sophisticated Laboratory Automation: Collaboration with Steve Oliver (Cambridge)

Williams et al. (2015) *Royal Society Interface*, DOI 10.1098/rsif.2014.1289
- n 2015-2018 Eve – Human cells - Cancer, Yeast - Aging
 - DARPA, CHIST-ERA

Adam



The Application Domain



- n Functional genomics
- n In yeast (*S. cerevisiae*) ~15% of the 6,000 genes still have no known function.
- n EUROFAN 2 made all viable single deletant strains.
- n Task to determine the “function” of a gene by growth experiments.

Formalising the Problem

- n Use logic programming to represent background knowledge: metabolism modelled as a directed labeled hyper-graph.
- n Use abduction to infer new hypotheses:
 - Abductive logic programming.
 - Techniques from Bioinformatics.
- n Use active learning to decide efficient experiments: cost of compounds and time.
- n Use machine learning to decide meaning of experimental results.

Types of Logical Inference

Deduction

Rule: All swans are white.

Fact: Daffy is a swan.

∴ Daffy is white.

Abduction

Rule: All swans are white.

Fact: Daffy is white.

∴ Daffy is a swan.

Induction

Fact: Daffy is a swan and white.

Fact: Tweety is a swan and white

∴ All swans are white.

Types of Logical Inference

Deduction

Rule: All swans are white.

Fact: Daffy is a swan.

∴ Daffy is white.

Abduction

Rule: All swans are white.

Fact: Daffy is white.

∴ Daffy is a swan. Daffy is a duck.

Induction

Fact: Daffy is a swan and white.

Fact: Tweety is a swan and white

∴ All swans are white.

Types of Logical Inference

Deduction

Rule: All swans are white.

Fact: Daffy is a swan.

∴ Daffy is white.

Abduction

Rule: All swans are white.

Fact: Daffy is white.

∴ Daffy is a swan.

Induction

Fact: Daffy is a swan and white.

Fact: Tweety is a swan and white

∴ All swans are white.

Types of Logical Inference

Deduction

Rule: All swans are white.

Fact: Daffy is a swan.

∴ Daffy is white.

Abduction

Rule: All swans are white.

Fact: Daffy is white.

∴ Daffy is a swan.

Induction

Fact: Daffy is a swan and white.

Fact: Tweety is a swan and white

∴ All swans are white.

Bruce is a black swan.

Novel Science

- n Adam generated and confirmed novel functional-genomics hypotheses concerning the identify of genes encoding enzymes catalysing orphan reactions in the metabolic network of the yeast *S. cerevisiae*.
- n Adam's conclusions have been manually verified using bioinformatic and biochemical evidence.
- n Adam was the first machine to autonomously discover novel scientific knowledge: hypothesise, and experimentally confirm.

Novel Scientific Knowledge

Orphan Enzyme		Hypothesised Gene	Prob.	Acc.	No.	Existing Annotation	Dry	Wet
1	glucosamine-6-phosphate deaminase (3.5.99.6)	YHR163W (SOL3)	$<10^{-4}$	97	8	'6-phosphogluconolactonase' ida	-	-
2	glutaminase (3.5.1.2)	YIL033C (BCY1)	$<10^{-4}$	92	11	'cAMP-dependent protein kinase inhibitor' ida	x ?	-
3	L-threonine 3-dehydrogenase (1.1.1.103)	YDL168W (SFA1)	$<10^{-4}$	83	6	'alcohol dehydrogenase' ida	-	-
4	purine-nucleoside phosphorylase (2.4.2.1)	YLR209C (PNP1)	$<10^{-4}$	82	11	'purine-nucleoside phosphorylase' ida	✓	-
5	2-aminoadipate transaminase (2.6.1.39)	YGL202W (ARO8)	$<10^{-4}$	80	3	'aromatic-amino-acid transaminase' ida	✓	✓
6	5,10-methenyltetrahydrofolate synthetase (6.3.3.2)	YER183C (FAU1)	$<10^{-4}$	80	4	'5,10 formyltetrahydrofolate cyclo-ligase' ida	✓	-
7	glucosamine-6-phosphate deaminase (3.5.99.6)	YNR034W (SOL1)	$<10^{-4}$	79	2	'possible role in tRNA export'	-	-
8	pyridoxal kinase (2.7.1.35)	YPR121W (THI22)	$<10^{-4}$	78	1	'phosphomethylpyrimidine kinase' iss	-	-
9	mannitol-1-phosphate 5-dehydrogenase (1.1.1.17)	YNR073C	$<10^{-4}$	78	6	'putative mannitol dehydrogenase' iss	-	-
10	1-acylglycerol-3-phosphate O-acyltransferase (2.3.1.51)	YDL052C (SLC1)	0.0001	80	6	'1-acylglycerol-3-phosphate O-acyltransferase' ida	✓	-
11	glucosamine-6-phosphate deaminase (3.5.99.6)	YGR248W (SOL4)	0.0002	78	2	'6-phosphogluconolactonase' ida	-	-
12	maleylacetoacetate isomerase (5.2.1.2)	YLL060C (GTT2)	0.0003	76	3	'glutathione S-transferase' ida	-	-
13	serine O-acetyltransferase (2.3.1.30)	YJL218W	0.0005	78	2	'unknown function'	-	-
14	L-threonine 3-dehydrogenase (1.1.1.103)	YLR070C (XYL2)	0.0052	75	6	'xylitol dehydrogenase' ida	-	-
15	2-aminoadipate transaminase (2.6.1.39)	YJL060W (BNA3)	0.0084	73	3	'kynurenine aminotransferase' ida	-	✓
16	pyridoxal kinase (2.7.1.35)	YNR027W	0.0259	76	2	'involved in bud-site selection' iss	-	-
17	polyamine oxidase (1.5.3.11)	YMR020W (FMS1)	0.0289	78	4	'polyamine oxidase' ida	✓	-
18	2-aminoadipate transaminase (2.6.1.39)	YER152C	0.0332	74	3	'uncharacterized'	-	✓
19	L-aspartate oxidase (1.4.3.16)	YJL045W	0.1300	72	1	'succinate dehydrogenase isozyme' iss	-	-
20	purine-nucleoside phosphorylase (2.4.2.1)	YLR017W (MEU1)	0.1421	72	6	'methylthioadenosine phosphorylase' ida	✓	-

Formalising Science

Formalisation of Science

- n The goal of science is to increase our knowledge of the natural world through the performance of experiments.
- n This knowledge should be expressed in formal logical languages.
- n Formal languages promote semantic clarity, which in turn supports the free exchange of scientific knowledge and simplifies scientific reasoning.

Robot Scientist & Formalisation

- n Robot Scientists provide excellent test-beds for the development of methodologies for formalising science.
- n Using them it is possible to completely capture and digitally curate all aspects of the scientific process.
- n The ontology LABORS is designed to enable the open access of the Robot Scientist experimental data and metadata to the scientific community.

Ontologies

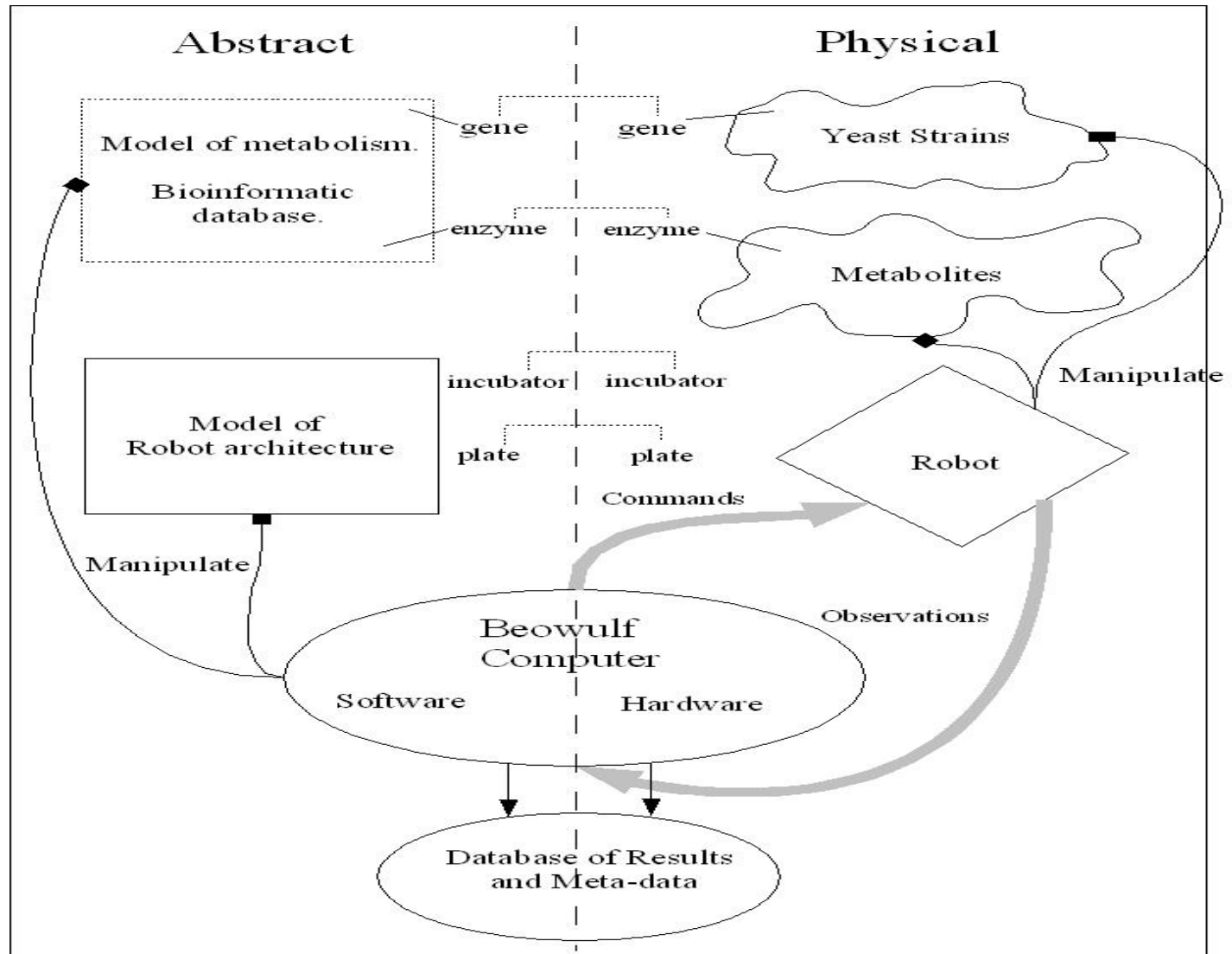
An ontology is “a concise and unambiguous description of what principal entities are relevant to an application domain and the relationship between them”*.

*Schulze-Kremer, S., 2001, Computer and Information Sci. 6(21)

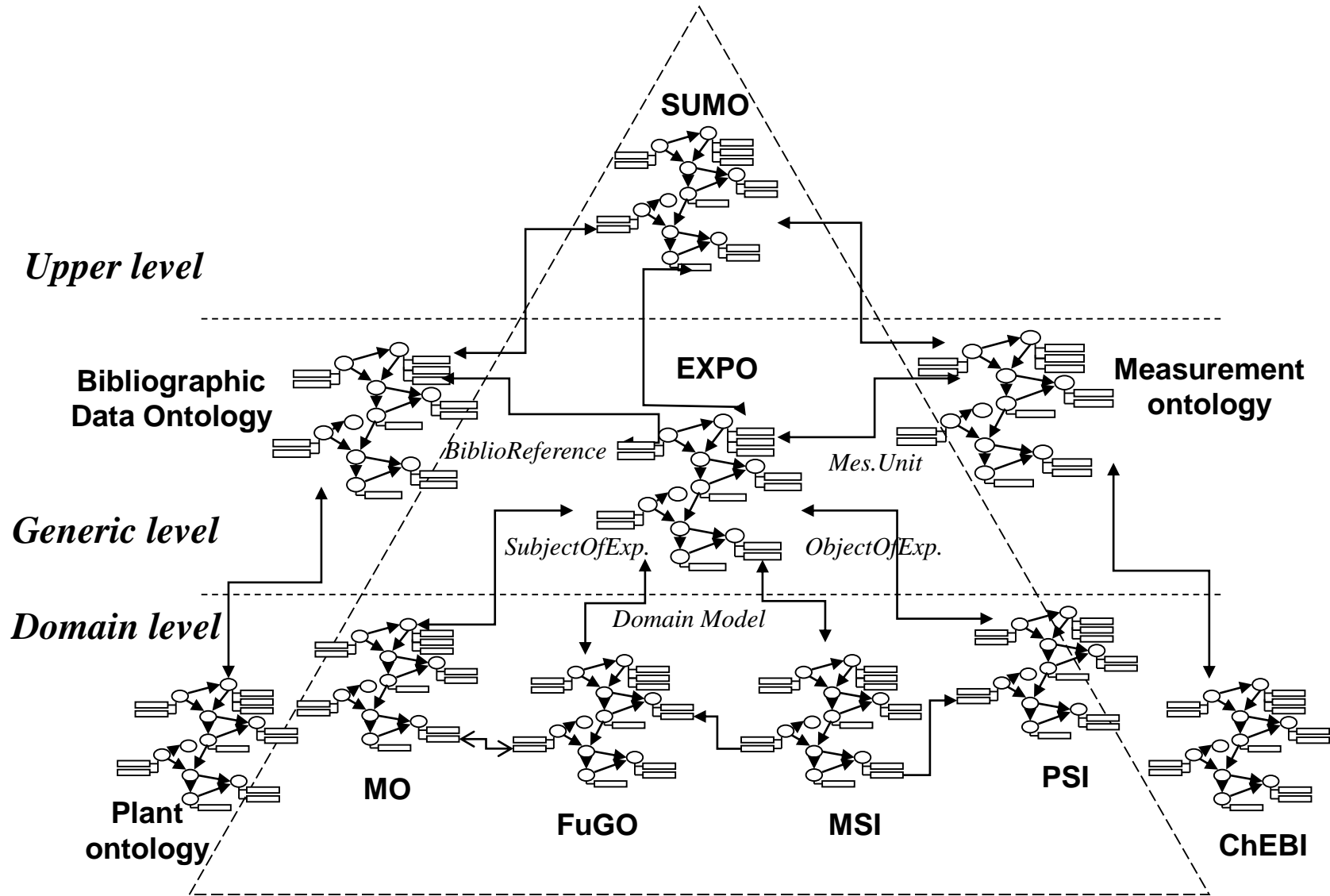
Dualism

- n The most fundamental ontological division in our design of Adam is between <abstract> and <physical> objects
- n We argue for this ontological division because it makes explicit the separation between models and reality.
- n All the objects which Adam deals with computationally are <abstract>, and all the objects it deals with physically are <physical>.

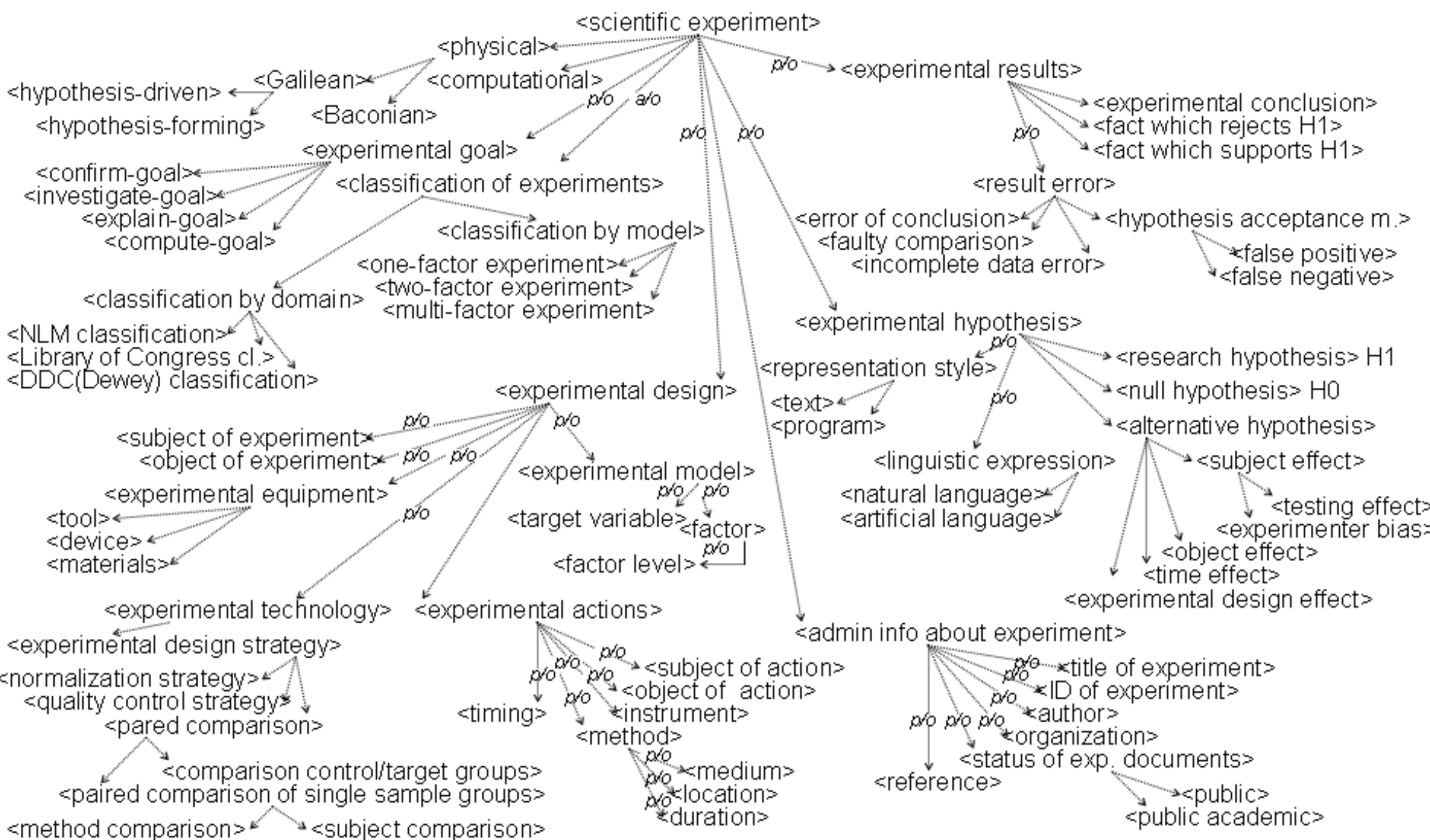
Overall View of the Universe



The Position of EXPO

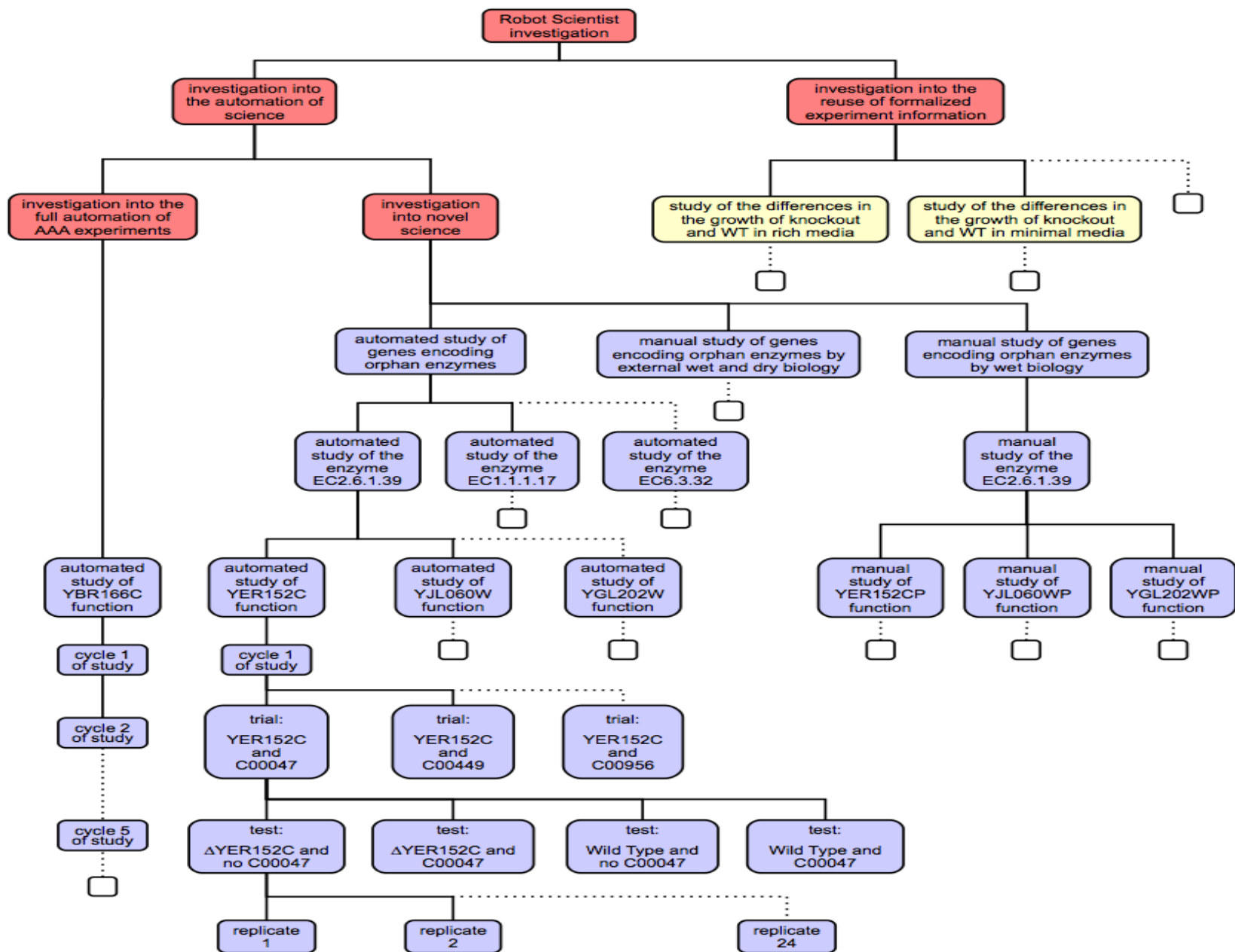


Small Section of EXPO



Adam's Investigations

- n This formalisation involves >10,000 different research units in a nested tree-like structure 11 levels deep.
- n It logically connects >6.6 million OD600_{nm} measurements to hypotheses, experimental goals, results, etc.
- n No previous large-scale experimental work has been so comprehensively described and recorded.



Levels in the Formalisation

Investigation into the automation of Science

Investigation into the automation of novel science

Investigation into the automated discovery of genes encoding orphan enzymes

Automated study of E.C.2.6.1.39 encoding

Cycle 1 of automated study of YER152C function

YER152C and Lysine automated trial

Experiment 1 (wild-type no metabolite)

Replicate 1 (well)

Observation 1

automated study of yer152c function

has text representation:

automated study: automated study of yer152c_function

has domain of study: functional genomics

has investigator

has goal: 'To test the hypothesis that the

with enzyme class

has organism class

has ncbi taxonomy ID

has hypothesis

has research hypothesis

has negative hypothesis

has cycle 1 of study

has study result

encodes(yer152c)

highest

proportion

has study conclusion

has datalog representation:

```
a:automated_study(X) :- a:automated_study(X),
a:hypotheses-set(X) :- a:research_hypothesis(X),
a:cycle_of_study(X) :- a:cycle_1_of_study(X),
a:hypotheses-set(X) :- a:negative_hypothesis(X),
a:domain_of_study(Y) :- a:automated_study(X),
a:investigator(Y) :- a:automated_study(X), a:investigator(Y),
a:goal(Y) :- a:automated_study(X), a:has_goal(Y),
a:organism_of_study(Y) :- a:automated_study(X),
a:hypotheses-set(Y) :- a:automated_study(X),
a:cycle_of_study(Y) :- a:automated_study(X),
a:study_result(Y) :- a:automated_study(X),
a:study_conclusion(Y) :- a:automated_study(X),
a:domain_of_study(X) :- a:functional_genomics(X),
a:investigator(X) :- a:adam(X),
a:goal(X) :- a:to_test_the_hypothesis_that_g(X),
a:organism_of_study(X) :- a:saccharomyces(X),
a:study_result(X) :- a:the_strength_of_evidence(X),
a:study_conclusion(X) :- a:hypothesis_1_conclusion(X)
```

has OWL representation:

```
<?xml version="1.0"?>
<rdf:RDF
  xmlns="http://www.owl-ontologies.com/Ontology1204198571.owl#"
  >
  <owl:Class rdf:ID="goal"/>
  <owl:Class rdf:ID="study_result"/>
  <owl:Class rdf:ID="ncbi_taxonomy_ID"/>
  <owl:Class rdf:ID="cycle_of_study"/>
  <owl:Class rdf:ID="negative_hypothesis">
    <rdfs:subClassOf>
      <owl:Class rdf:ID="hypotheses-set"/>
    </rdfs:subClassOf>
  </owl:Class>
  <owl:Class rdf:ID="domain_of_study"/>
  <owl:Class rdf:ID="organism_of_study"/>
  <owl:Class rdf:ID="cycle_1_of_study">
    <rdfs:subClassOf rdf:resource="#cycle_of_study"/>
  </owl:Class>
  <owl:Class rdf:ID="automated_study">
    <rdfs:subClassOf>
      <owl:Restriction>
        <owl:someValuesFrom rdf:resource="#goal"/>
        <owl:onProperty>
          <owl:ObjectProperty rdf:ID="has_goal"/>
        </owl:onProperty>
      </owl:Restriction>
    </rdfs:subClassOf>
    <rdfs:subClassOf>
      <owl:Restriction>
        <owl:someValuesFrom rdf:resource="#organism_of_study"/>
        <owl:onProperty>
          <owl:ObjectProperty rdf:ID="has_organism_of_study"/>
        </owl:onProperty>
      </owl:Restriction>
    </rdfs:subClassOf>
  </owl:Class>
</rdf:RDF>
```

Eve



Drug Design

The Application Domain



Malaria



Shistosomiasis



Leishmania

Chagas



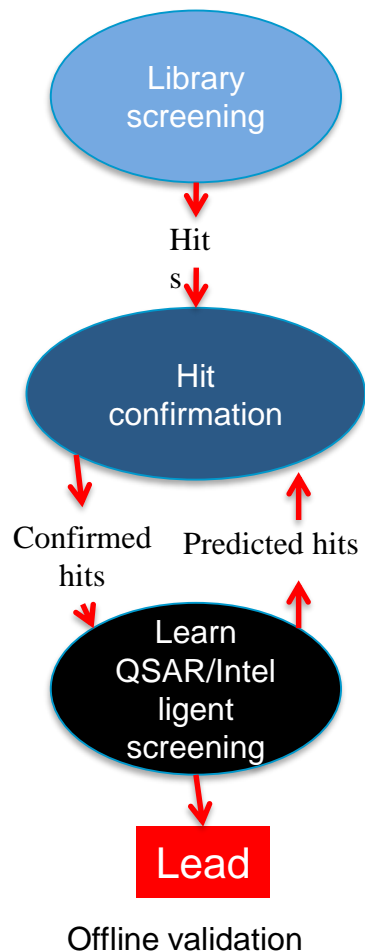
Why Tropical Diseases?

- n Millions of people die of these diseases, and hundreds of millions of people suffer infection.
- n It is clear how to cure these diseases – kill the parasites.
- n They are “neglected”, so avoid competition from the Pharmaceutical industry.

Formalising the Problem

- n Use graphs and standard chemoinformatic methods to represent background knowledge - the use of relations is planned.
- n Uses induction (quantitative structure activity relationship – QSAR learning) to infer new hypotheses.
- n Use active learning to decide efficient experiments, and econometric model to decide what compounds to test.

Eve's Automation of Pipeline

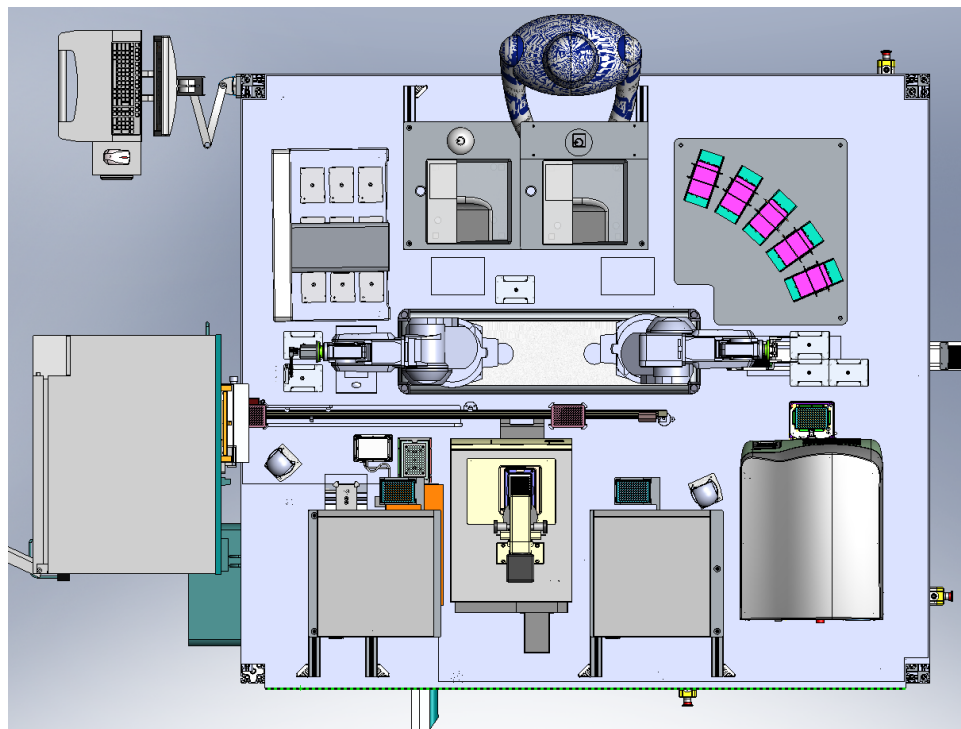


- Standard library screening is brute force:
- Eve uses intelligent screening
- In the standard “pipeline” the 3 processes are not integrated.
- In Eve automated and integrated.

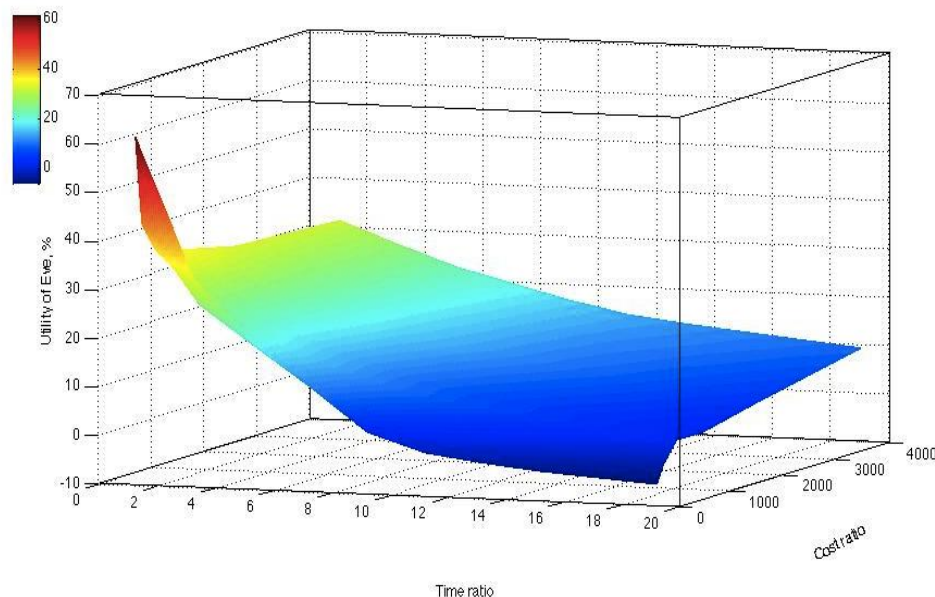
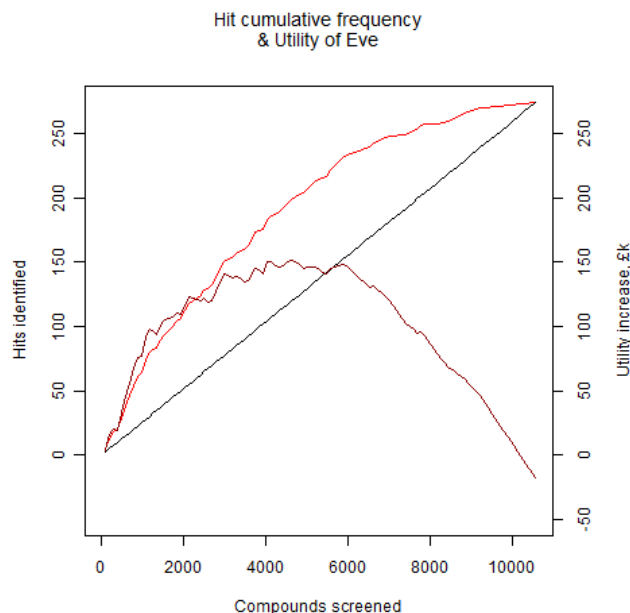
Eve's Hardware

Highlights of Eve's hardware:

- Acoustic liquid handling
- High throughput 384 well plates
- Two industrial robot arms
- Automated 60x microscope
- Liquid handlers, fluorescence readers, barcode scanners, dry store, incubator, tube decapper ...



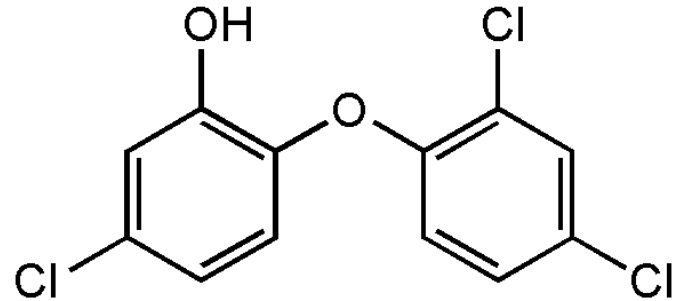
The Economics of Intelligent Screening



$$\Delta \text{Utility of Eve} = \sum_{1}^{Nm} (Tm + Cm) + \sum_{1}^{Nx} (Tc + Cc - Uh) + \sum_{1}^{Ne} (Tm - Tc + Cm - Cc)$$

Nm	-	Number of compounds not assayed by Eve
Tm	-	Cost of the time to screen a compound using the mass screening assay
Cm	-	Cost of the loss of a compound in the mass screening assay
Nx	-	Number of hits missed by Eve
Tc	-	Cost of the time to screen a compound using a cherry-picking (confirmation or intelligent) assay
Cc	-	Cost of the loss of a compound in a cherry-picking assay
Uh	-	Utility of a hit
Ne	-	Number of compounds assayed by Eve

Triclosan Repositioned for Malaria



- n Simple compound
- n Known to be safe – used in toothpaste.
- n Targets both DHFR and FAS-II – well established targets.
- n Demonstrated activity using multiple wet experimental techniques.
- n Works against wild-type and drug-resistant *Plasmodium falciparum*, and *Plasmodium vivax*.

Future Prospects

The Future?

- n In Chess/Go there is a continuum of ability from novices up to Grandmasters.
- n I argue that this is also true in science, from the simple research of Eve, through what most human scientists can achieve, up to the ability of a Newton or Einstein.
- n If you accept this, then just as in Chess/Go, it is likely that advances in computer hardware and software will drive the development of ever smarter Robot Scientists.
- n In favour of this argument are the ongoing development of AI and laboratory robotics.

Vision

- n The collaboration between Human and Robot Scientists will produce better science than either can alone – human/computer teams still play better chess than either alone.
- n Scientific knowledge will be primarily expressed in logic with associated probabilities and published using the Semantic Web.
- n The improved productivity of science leads to societal benefits: better food security, better medicines, etc.
- n The Physics Nobel Frank Wilczek is on record as saying that in 100 years' time the best physicist will be a machine

Conclusions

- n Science is a wonderful application area for AI.
- n Automation is becoming increasingly important in scientific research e.g. DNA sequencing, drug design.
- n The Robot Scientist concept is the logical next step in scientific automation.
- n The Robot Scientist Adam was the first machine to have discovered novel scientific knowledge.
- n The Robot Scientist Eve has found new lead compounds for neglected tropical diseases.
- n The Robot Scientist Eve can accelerate systems biology modelling.

Acknowledgments

- n The Robot Scientist team: Manchester, Aberystwyth, Cambridge, Brunel, Leuven, Thailand. (BBSRC)
- n Chicago Big Mechanism consortium. (DARPA)
- n AdaLab consortium. (CHIST-ERA)
- n AIST (Tokyo)