

Automating Science using Robot Scientists

Ross D. King,
University of Manchester, ATI, ross.king@manchester.ac.uk



The State-of-the-Art - AI

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, ...

Computer Chess - 1997



Abstract world: 64 squares,
32 pieces.

Computers play chess better
than the best humans, and
computers can now make
strikingly beautiful moves.

Deep Blue v Kasparov

Famous game 2 – Machine had no fear

Jeopardy - 2011



- n Watson is a computer system capable of answering questions posed in natural language.
- n Watson easily beat Brad Rutter and Ken Jennings the best human players of Jeopardy.

Google Driverless Car - 2016



- n The system combines maps with sensors to automate driving.
- n Has driven long distances with no human intervention.
- n Laws changed in a number of states, countries.

Computer Go - 2016



Combination of brute-force search and machine learning (deep learning + reinforcement learning) to evaluate board positions.

AlphaGo v Lee Sedol

March 2016, won 4 : 1.

Scientific Discovery

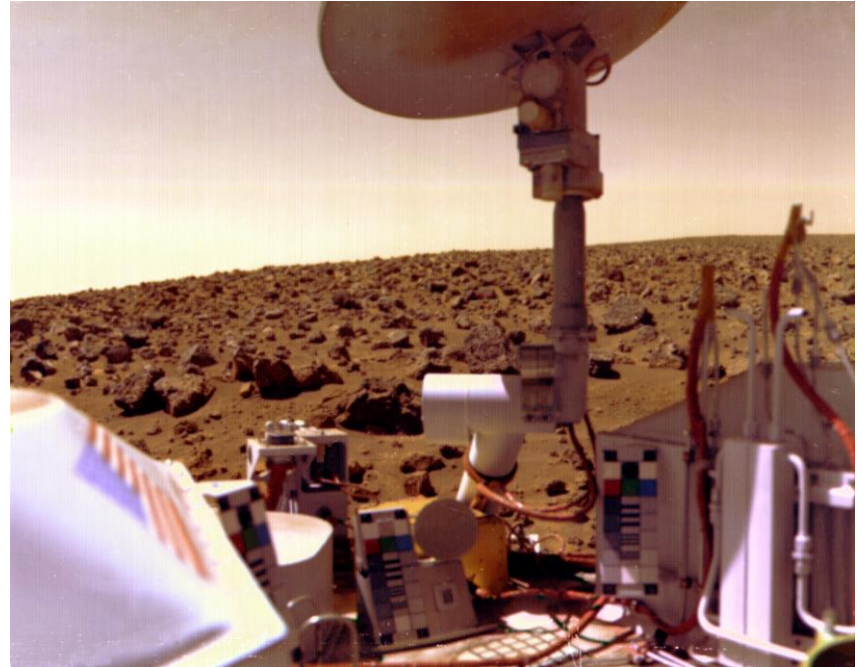
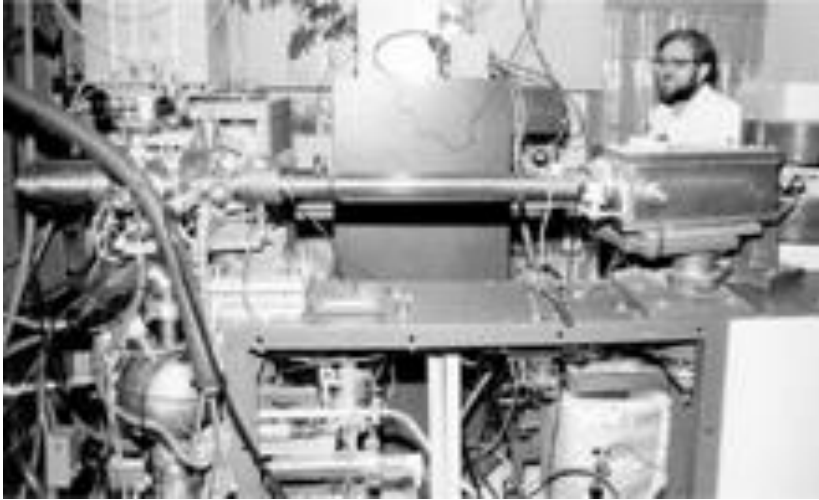
AI Systems have Superhuman Scientific Reasoning Powers

- n Flawlessly remember vast numbers of facts
- n Execute flawless logical reasoning
- n Execute near 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).

Into the Lab



n Automated discovery
in a chemistry
laboratory.

Zytkow, et al. (1990)

Jan Zytkow (1944-2001)

The State-of-the-Art
Laboratory Automation

Lab Automation

- n Laboratory Robotics technology is steadily advancing.
- n Almost everything that a human can do in the lab can now be automated.
- n Robots can work longer, faster, more accurately than humans.
- n Problems: high capital-costs, high running-costs, poor integration, poor software.

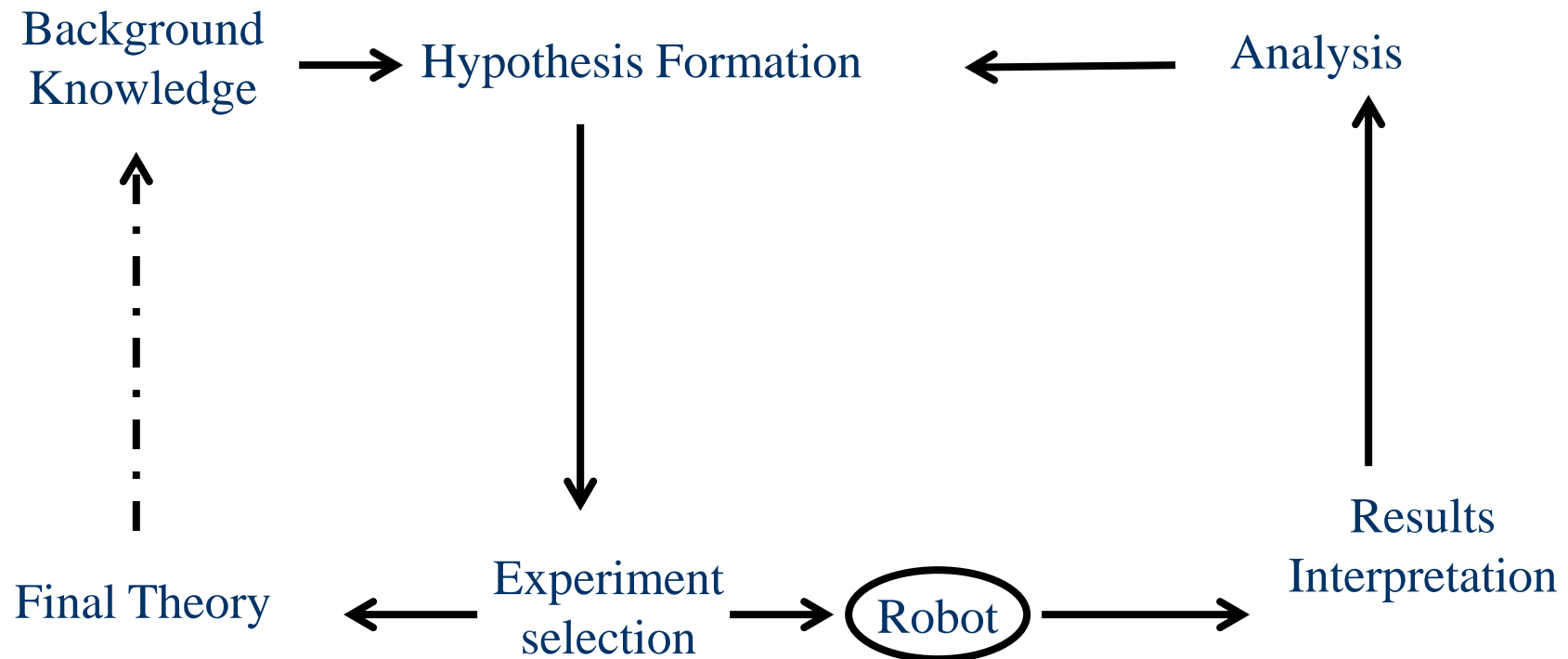
Cloud Laboratory Automation

- n Idea:
 - Laboratory automation is expensive and difficult to use.
 - Therefore put automation in the cloud.
 - Customers use an API to design the experiments.
 - Benefits of scale.
- n Two main companies: Emerald Cloud Services, Transcriptics.
- n I worry about control of the API and/or Balkanisation of APIs

Robot Scientists

The Concept of a Robot Scientist

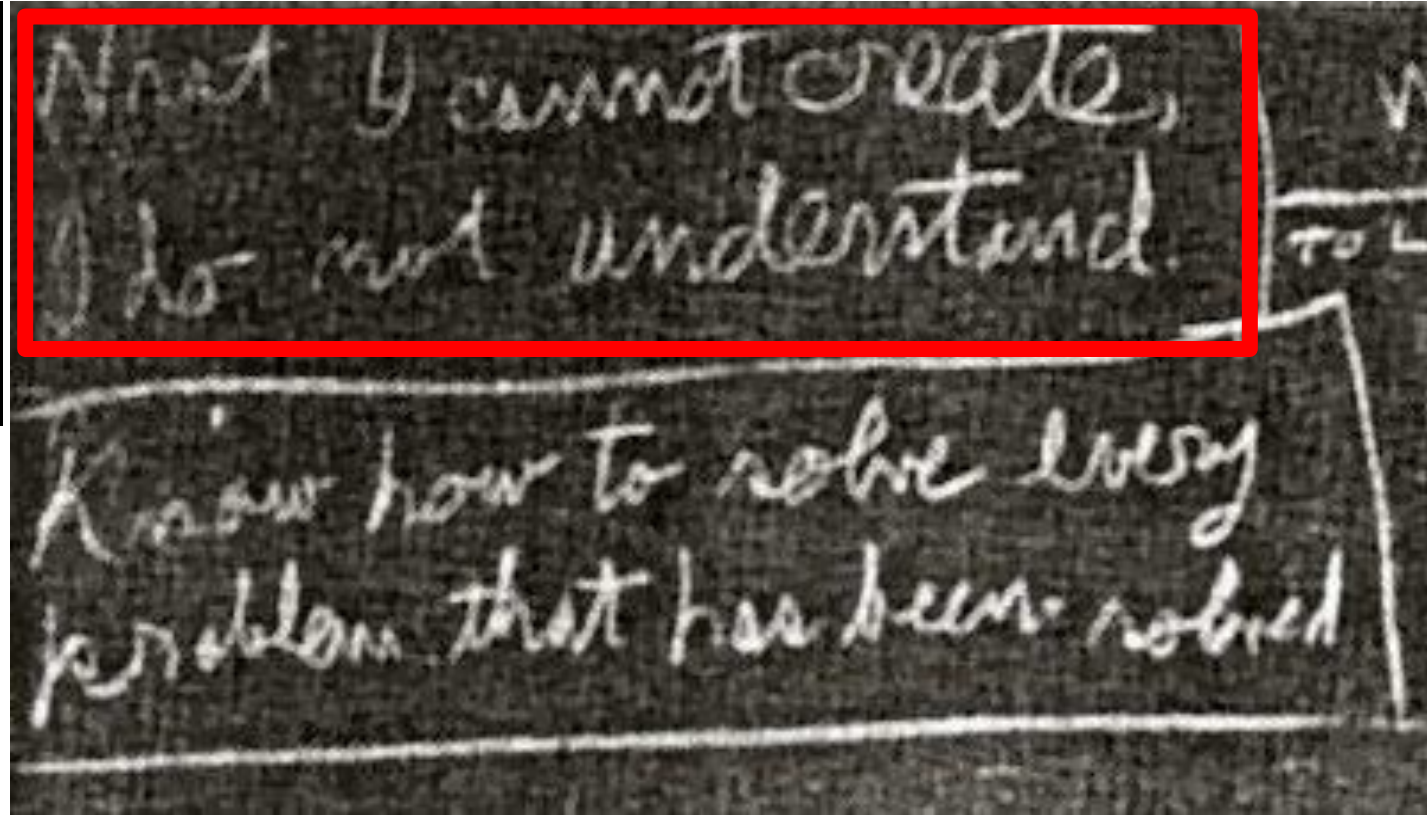
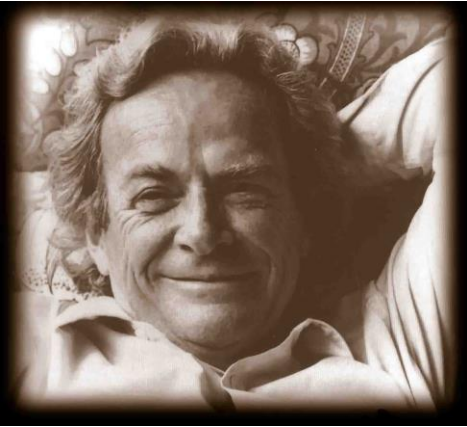
Computer systems capable of originating their own experiments, physically executing them, interpreting the results, and then repeating the cycle.



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.*

The Complexity of Biological Systems

- n Even simple “model” biological systems like that of *E. coli* and yeast are incredibly complicated.
- n Thousands of genes, proteins, small-molecules, interacting together in complicated spatial temporal ways.
- n Ockham's razor doesn't work - system evolved.

- n Not enough PhDs in the world to disentangle these systems.
- n Need help - Robot Scientists.

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

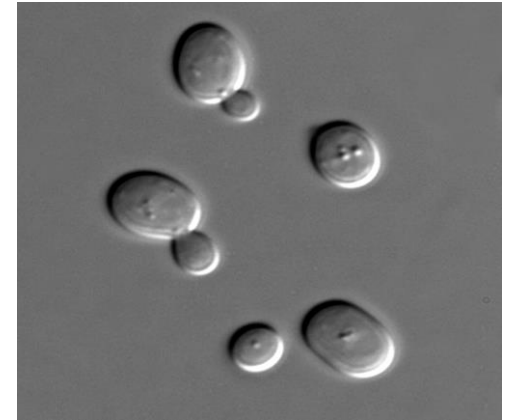
- n 2008-2015 Eve – Drug Design for Tropical Diseases
 - Sophisticated Laboratory Automation: Collaboration with Steve Oliver (Cambridge)

Williams et al. (2015) *Royal Society Interface*, DOI 10.1098/rsif.2014.1289

- n 2015-2019 Eve – Human cells - Cancer, Yeast - Aging
 - DARPA, CHIST-ERA, EPSRC.

Adam

The Application Domain

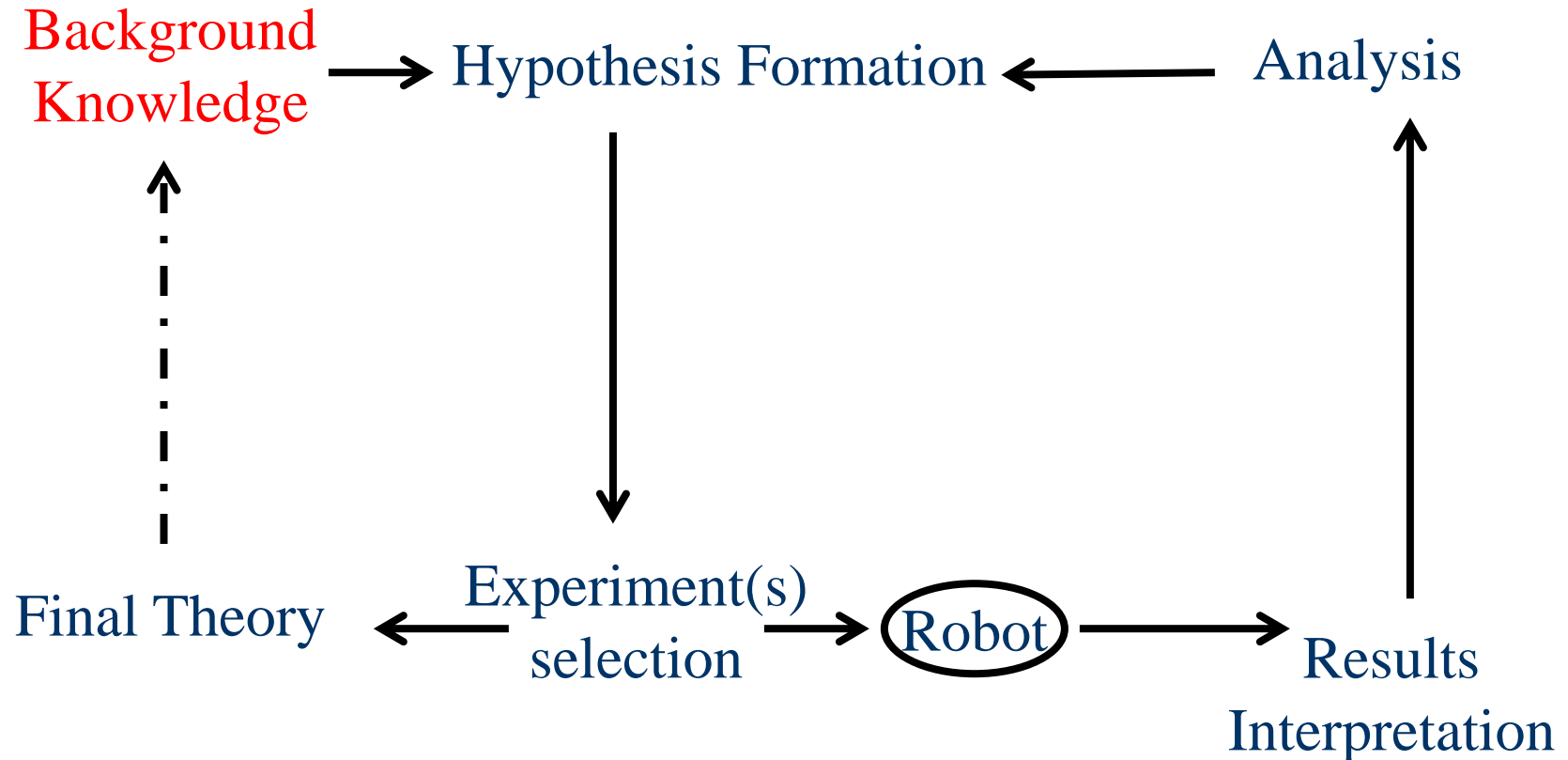


- n Functional genomics
- n In yeast (*S. cerevisiae*) ~15% of the 6,000 genes still have no known function.
- n EUROFAN 2 made all viable single deletant strains.
- n Task to determine the “function” of a gene by growth experiments.

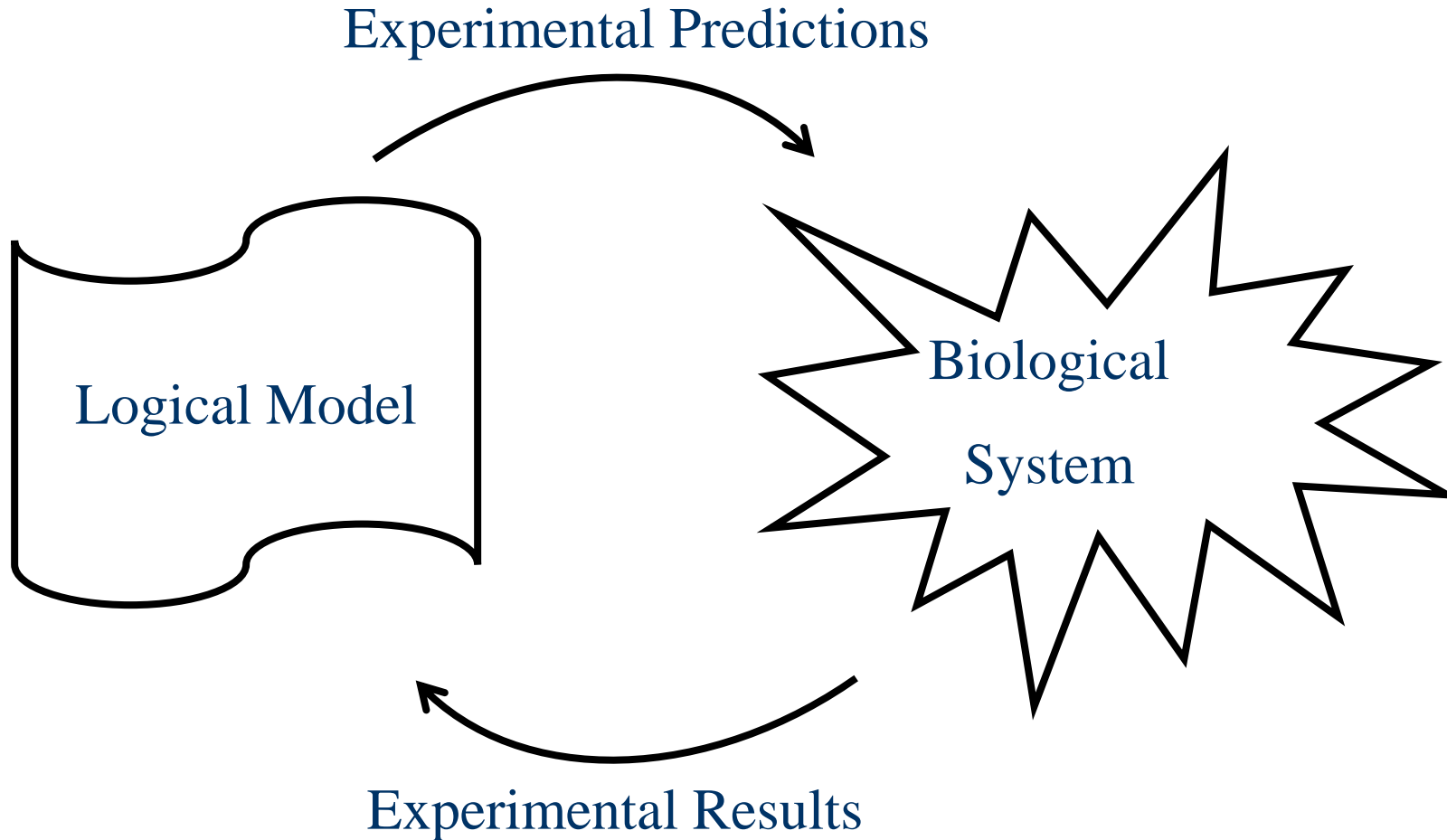
Formalising the Problem

- 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.

The Experimental Cycle



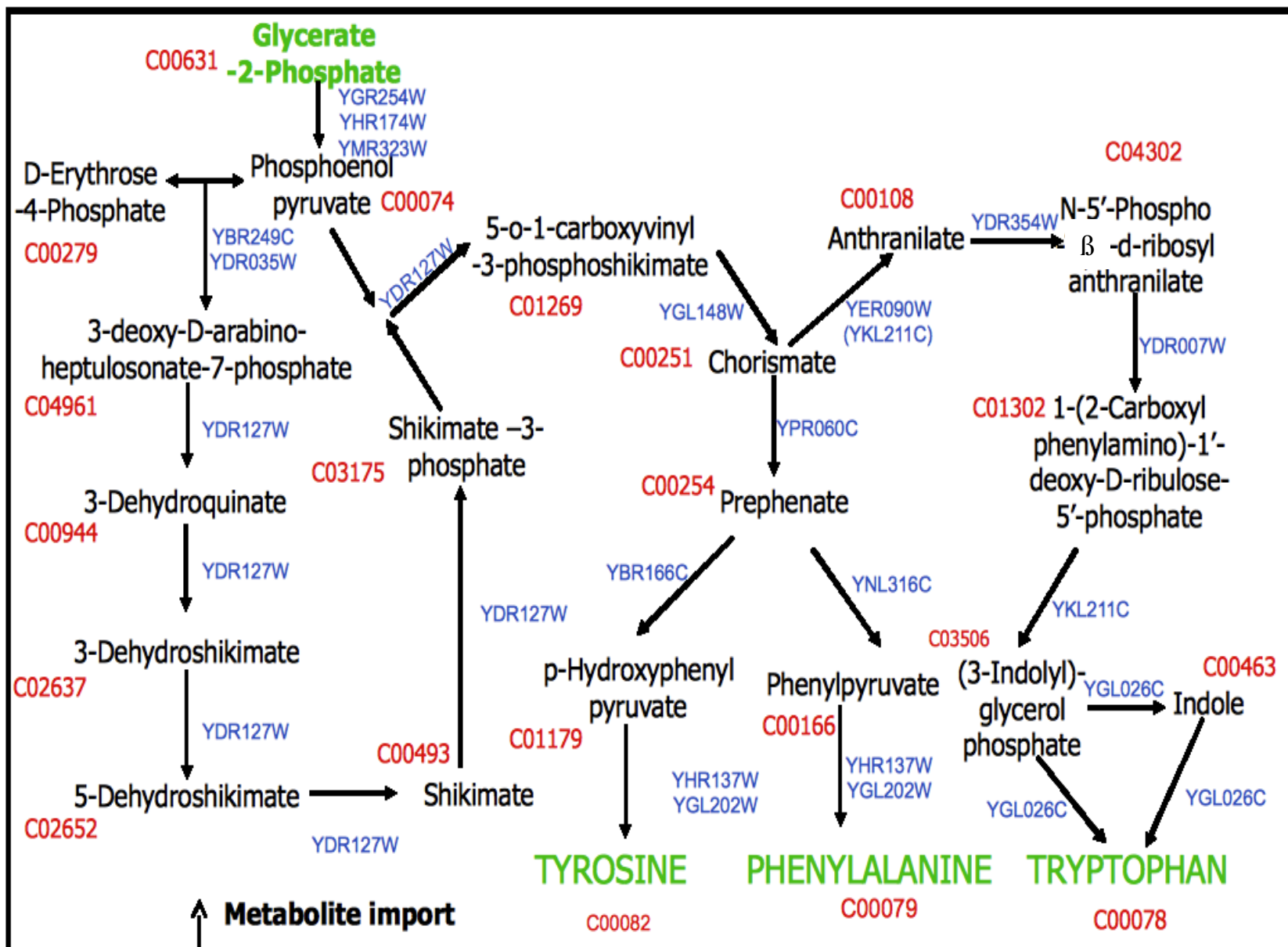
Model v Real-World



Logical Cell Model

- n We have developed a logical formalism for modelling metabolic pathways (encoded in Prolog). This is essentially a directed labeled hyper-graph: with metabolites as nodes and enzymes as arcs.
- n If a path can be found from cell inputs (metabolites in the growth medium) to all the cell outputs (essential compounds) then the cell can grow.

Phenylalanine, Tyrosine, and Tryptophan Pathways for *S. cerevisiae*



Noise

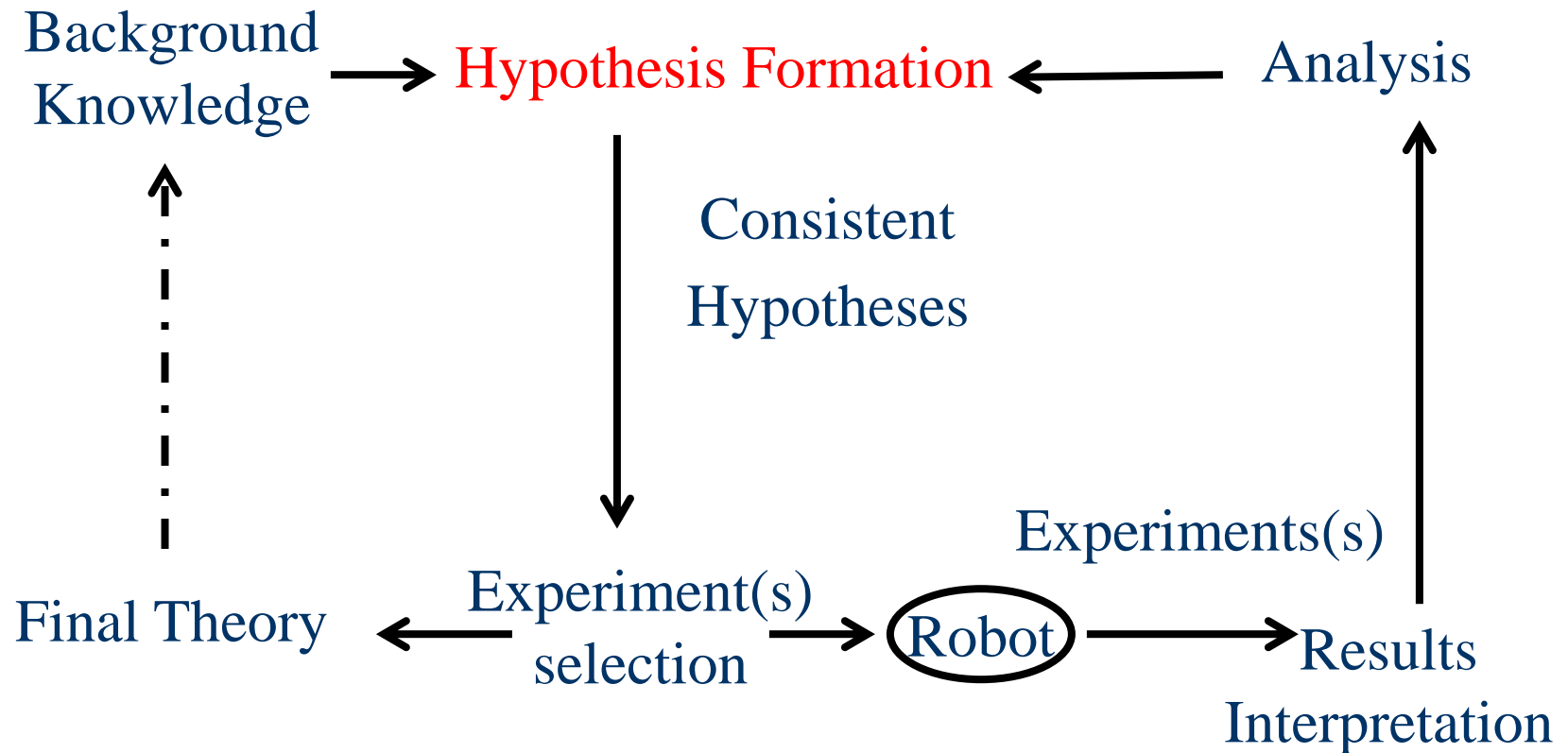
There were two types of noise in the physical experiments:

- n Experiment and measurement noise
(we estimate that ~25% of the experiments are noisy).
- n Noise due to errors in the background knowledge
(the model does not agree with ~2% of experimental results).

Genome Scale Model of Yeast Metabolism

- n We have extended our model of aromatic amino acid metabolism to cover most of what is known about yeast metabolism.
- n Includes 1,166 ORFs (940 known, 226 inferred)
- n Growth if path from growth medium to defined end-points.
- n State-of-the-art accuracy in predicting cell viability

The Experimental Cycle



Inferring Hypotheses

- n Science is based on the hypothetico-deductive method.
- n In the philosophy of science. It has often been argued that only humans can make the “leaps of imagination” necessary to form hypotheses.
- n In biology most hypothesis generation is abductive, not inductive.
- n Adam used abductive inference to infer missing arcs/labels in its metabolic graph - hypotheses. With these missing nodes Adam could then deductively infer (explain) the observed 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

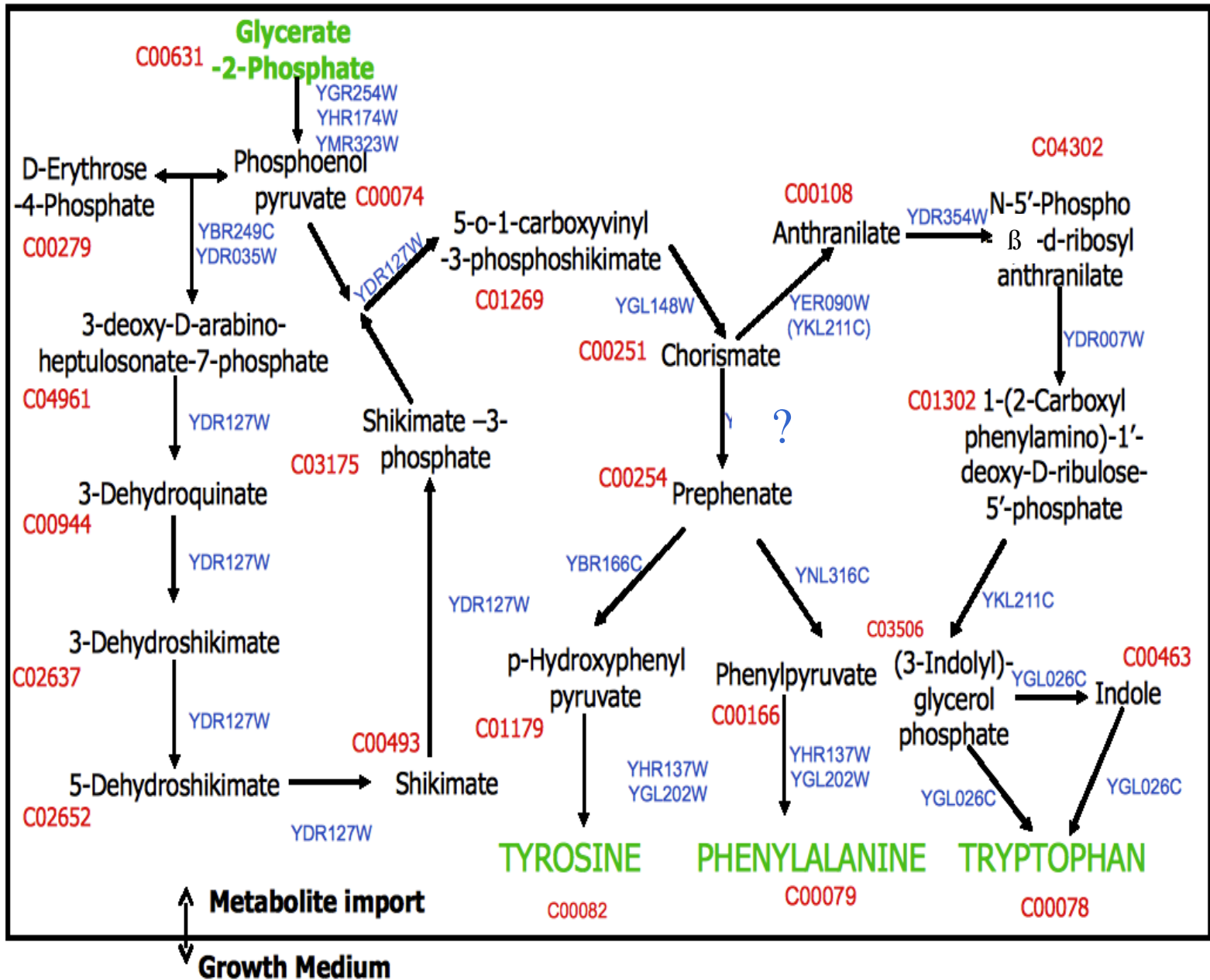
∴ All swans are white.

Bruce is a black swan.

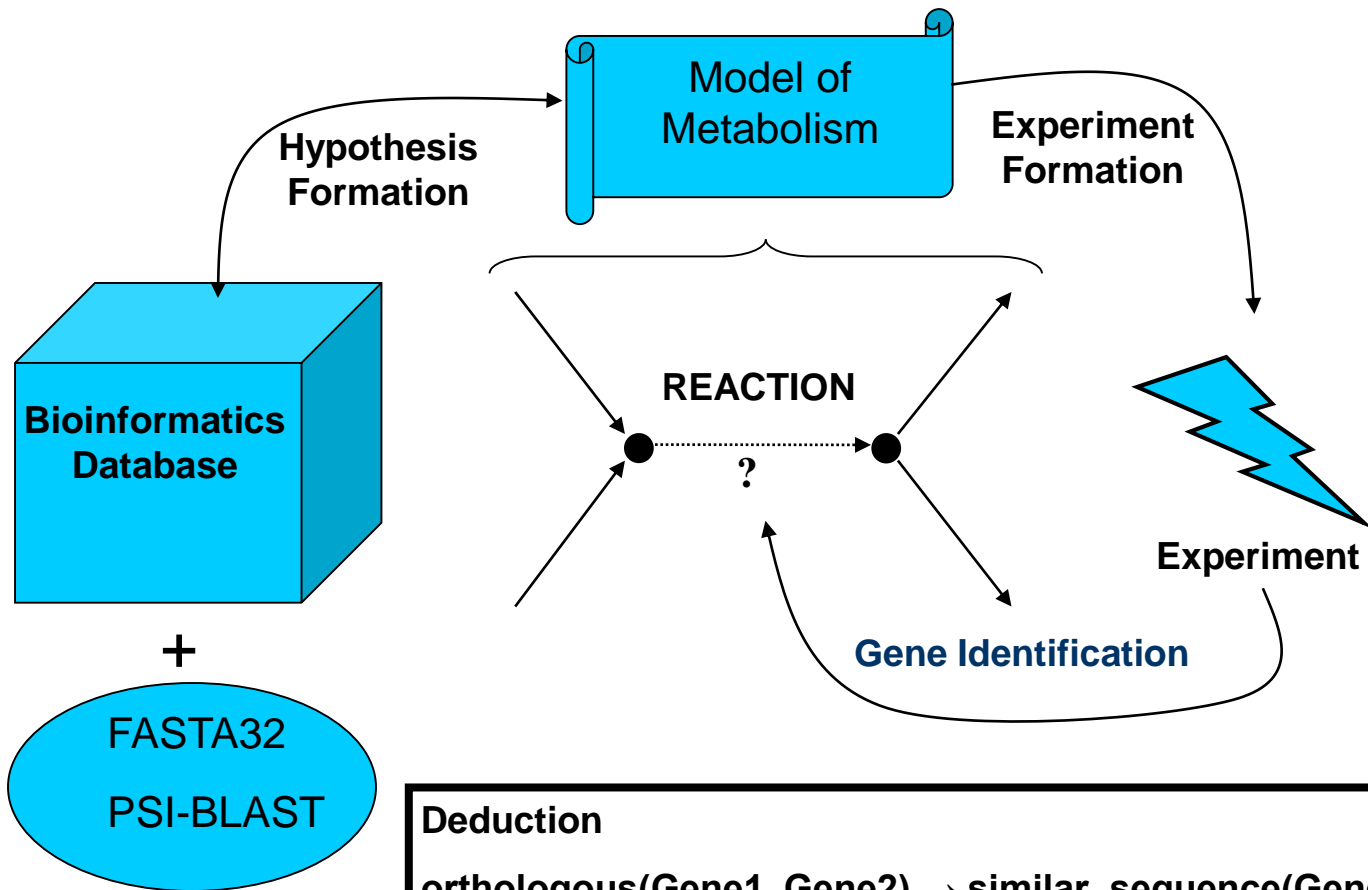
Orphan Enzymes

- n Our model of yeast metabolism has “orphan enzymes” enzymes which catalyse biochemical reactions known to be in yeast, but which do not have identified parent genes
- n We use bioinformatics to abduce genes which encode for these orphan enzymes.

Phenylalanine, Tyrosine, and Tryptophan Pathways for *S. cerevisiae*



Automated Model Completion



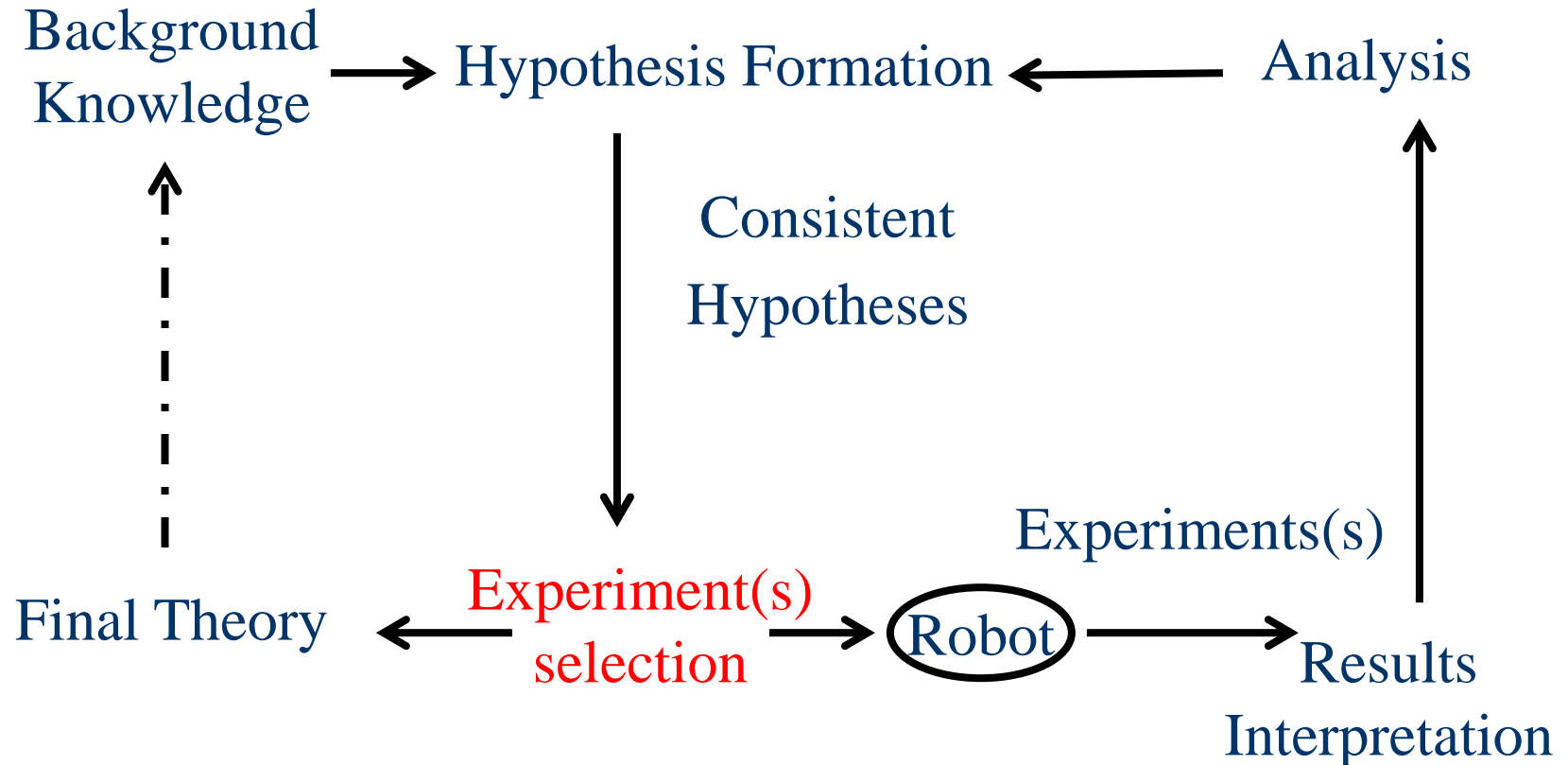
Deduction

$\text{orthologous}(\text{Gene1}, \text{Gene2}) \rightarrow \text{similar_sequence}(\text{Gene1}, \text{Gene2}).$

Abduction

$\text{similar_sequence}(\text{Gene1}, \text{Gene2}) \rightarrow \text{orthologous}(\text{Gene1}, \text{Gene2}).$

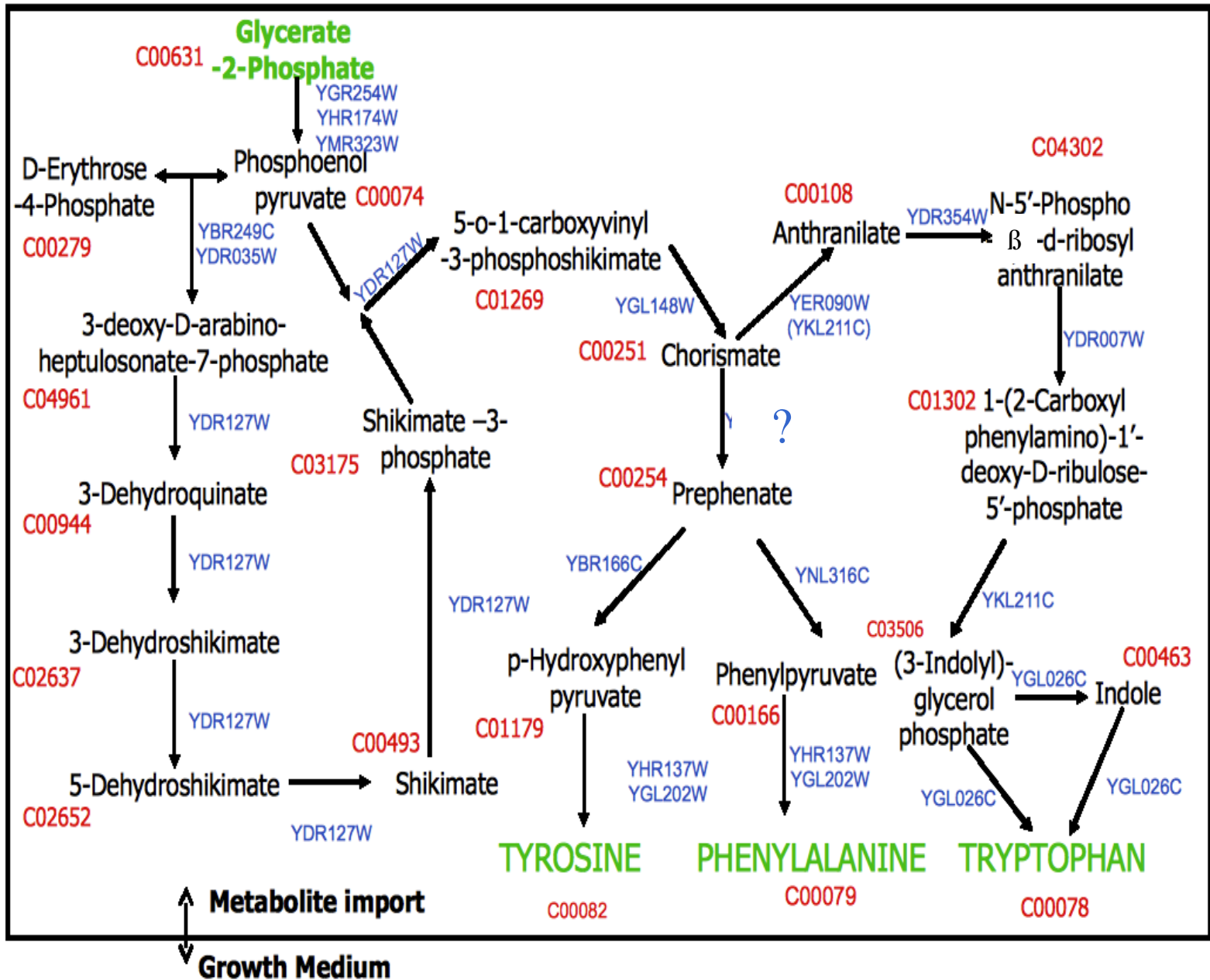
The Experimental Cycle



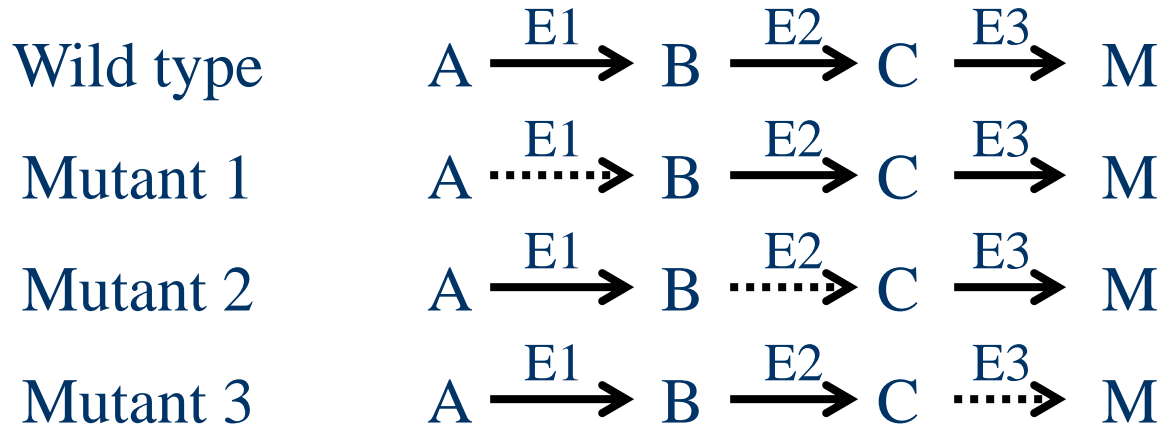
A Discriminating Experiment

- n Hypothesis 1: YDR060C codes for the enzyme the reaction: chorismate → prephenate.
- n Hypothesis 2: YDR060C codes for the enzyme the reaction: chorismate → anthranilate.
- n These can be distinguished by growing the knockout YDR060C on prephenate or anthranilate.

Phenylalanine, Tyrosine, and Tryptophan Pathways for *S. cerevisiae*



Auxotrophic Experiments



| | A | B | C | M |
|-----------|---|---|---|---|
| Wild type | + | + | + | + |
| Mutant 1 | - | + | + | + |
| Mutant 2 | - | - | + | + |
| Mutant 3 | - | - | - | + |

Experimental Methodology

- n Experiments consist of making particular growth media and testing if the mutants can grow (add metabolites to a basic defined medium).
- n A mutant is auxotrophic if cannot grow on a defined medium that the wild type can grow on.
- n By observing the pattern of chemicals that recover growth the function of the knocked out mutant can be inferred.

Inferring Experiments

Given a set of hypotheses we wish to infer an experiment that will efficiently discriminate between them

Assume:

- n Every experiment has an associated cost.
- n Each hypothesis has a probability of being correct.

The task:

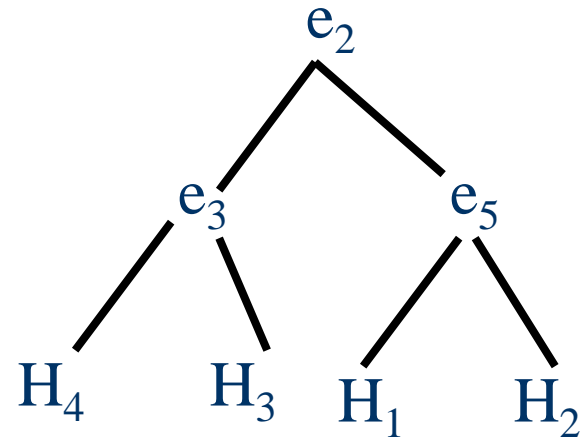
- n To choose a series of experiments which minimise the expected cost of eliminating all but one hypothesis.

Active Learning

- n In the 1972 Fedorov (Theory of optimal experiments) showed that this problem is in general intractable (NP complete).
- n However, it can be shown that the problem is the same as finding an optimal decision tree; and it is known that this problem can be solved “nearly” optimally in polynomial time.

How to choose the best experiment

| | e_1 | e_2 | ... | e_m |
|-------|-------|-------|-----|-------|
| H_1 | T | F | ... | T |
| ... | F | T | ... | F |
| H_n | F | T | ... | T |



Choosing the best experiment is equivalent to choosing the best node in a decision tree.

Recurrence Formula

$EC(H, T)$ denote the minimum expected cost of experimentation given the set of candidate hypotheses H and the set of candidate trials T :

$$EC(\emptyset, T) = 0$$

$$EC(\{h\}, T) = 0$$

$$EC(H, T) \gg \min_{t \in T} [C_t + p(t)(\text{mean}_{t \in (T-t)} C_t) J_{H[t]} + (1 - p(t)) \text{mean}_{t \in (T-t)} C_t J_{H[\bar{t}]}]$$

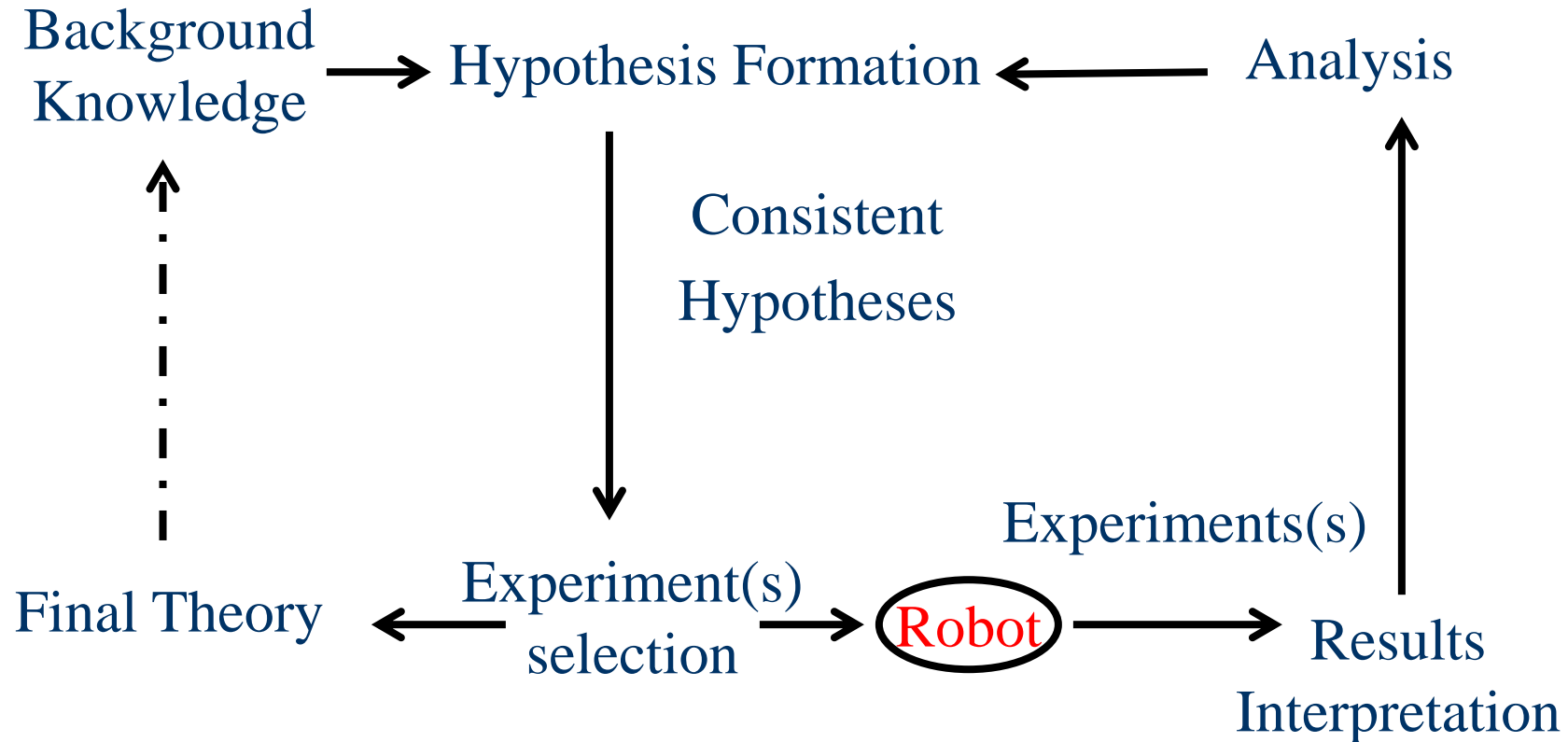
$$J_H = -\sum_{h \in H} p(h) \log_2(p(h))$$

C_t is the monetary price of the trial t

$p(t)$ is the probability that the outcome of the trial t is positive

$p(t)$ can be computed as the sum of the probabilities of the hypotheses (h) which are consistent with a positive outcome of t

The Experimental Cycle



Adam

- n Designed to fully automate yeast growth experiments.
- n Had a -20C freezer, 3 incubators, 2 readers, 3 liquid handlers, 3 robotic arms, 2 robot tracks, a centrifuge, a washer, an environmental control system, etc.
- n Was capable of initiating ~1,000 new experiments and >200,000 observations per day in a continuous cycle.

Plan of Adam

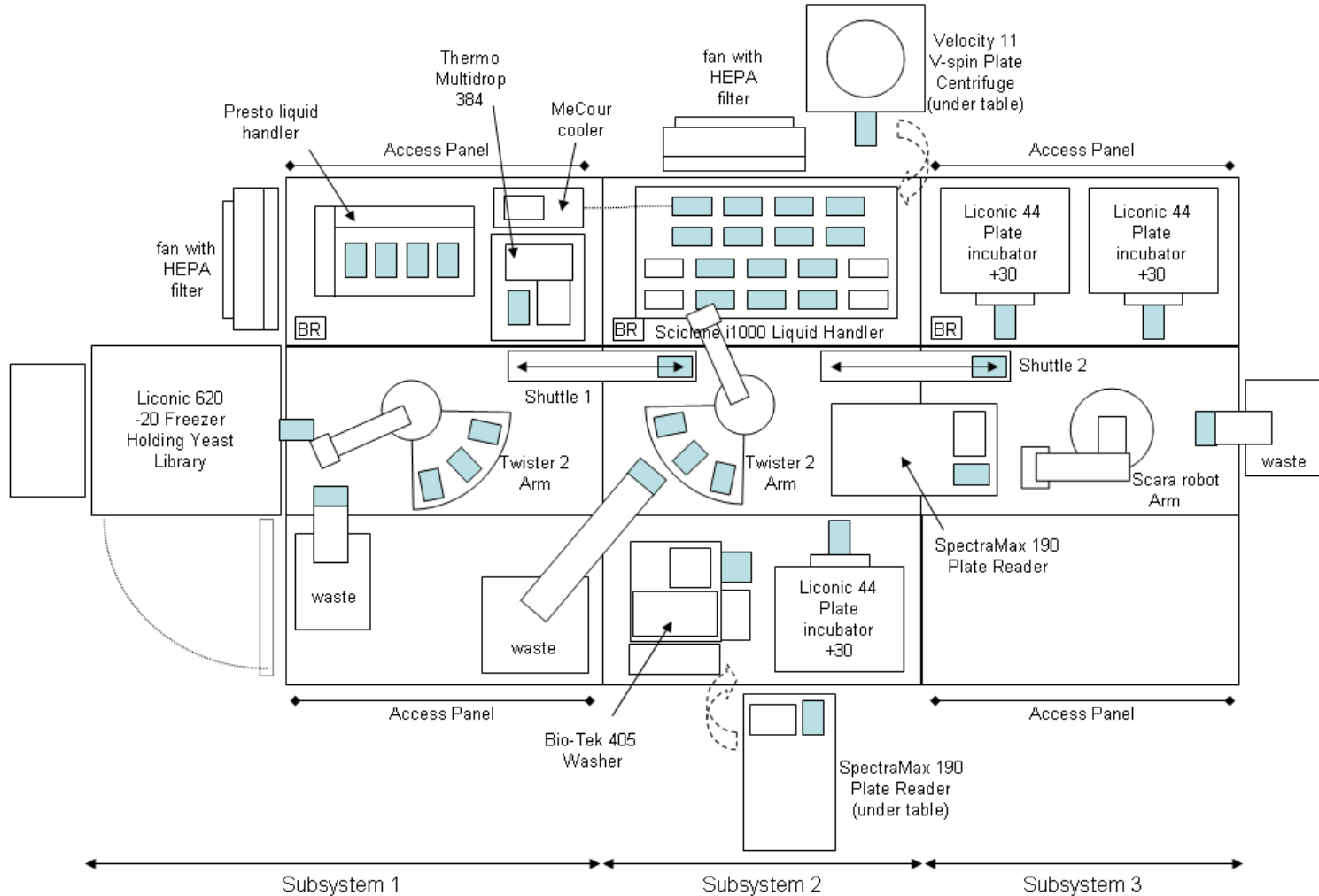
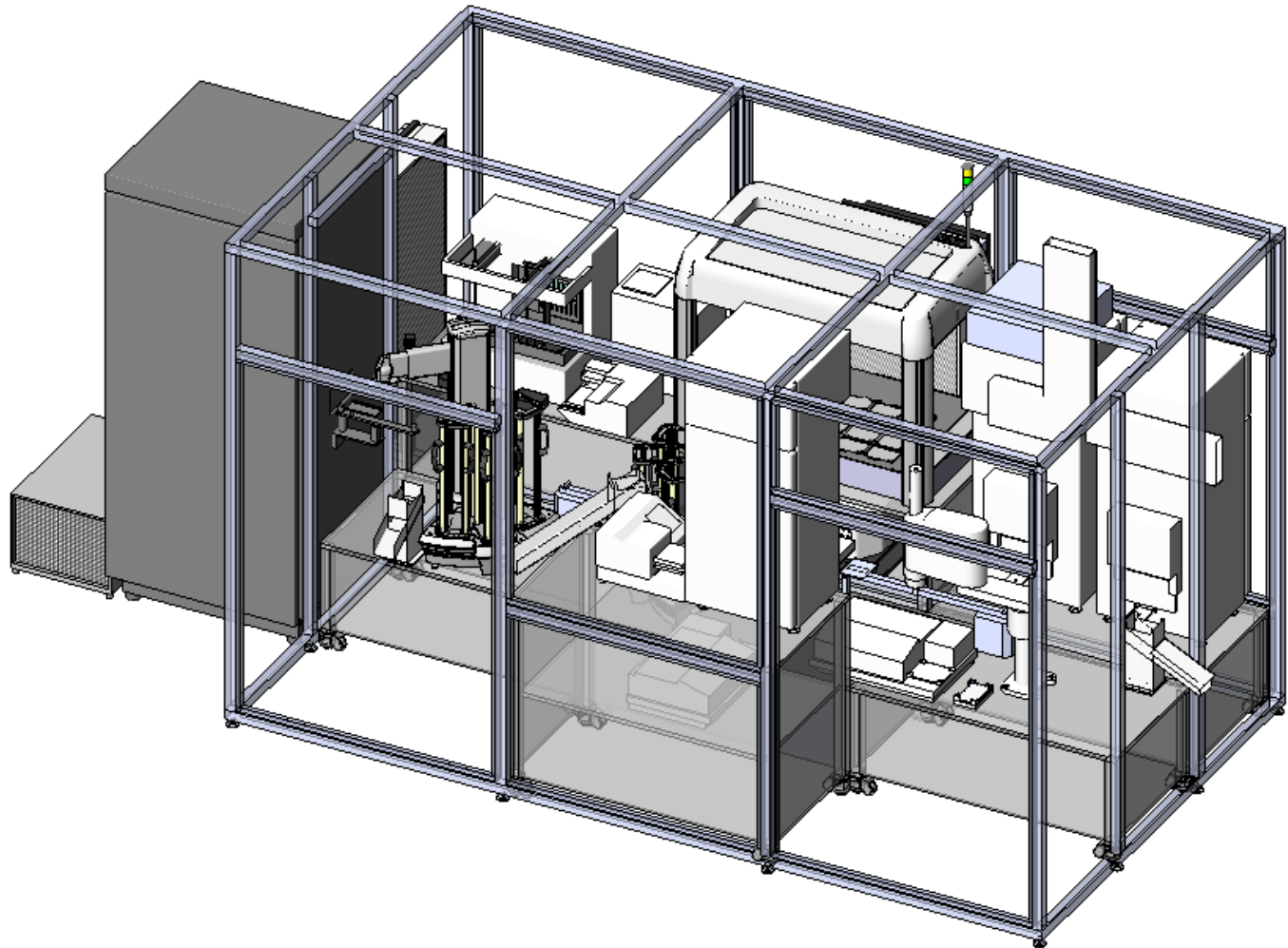
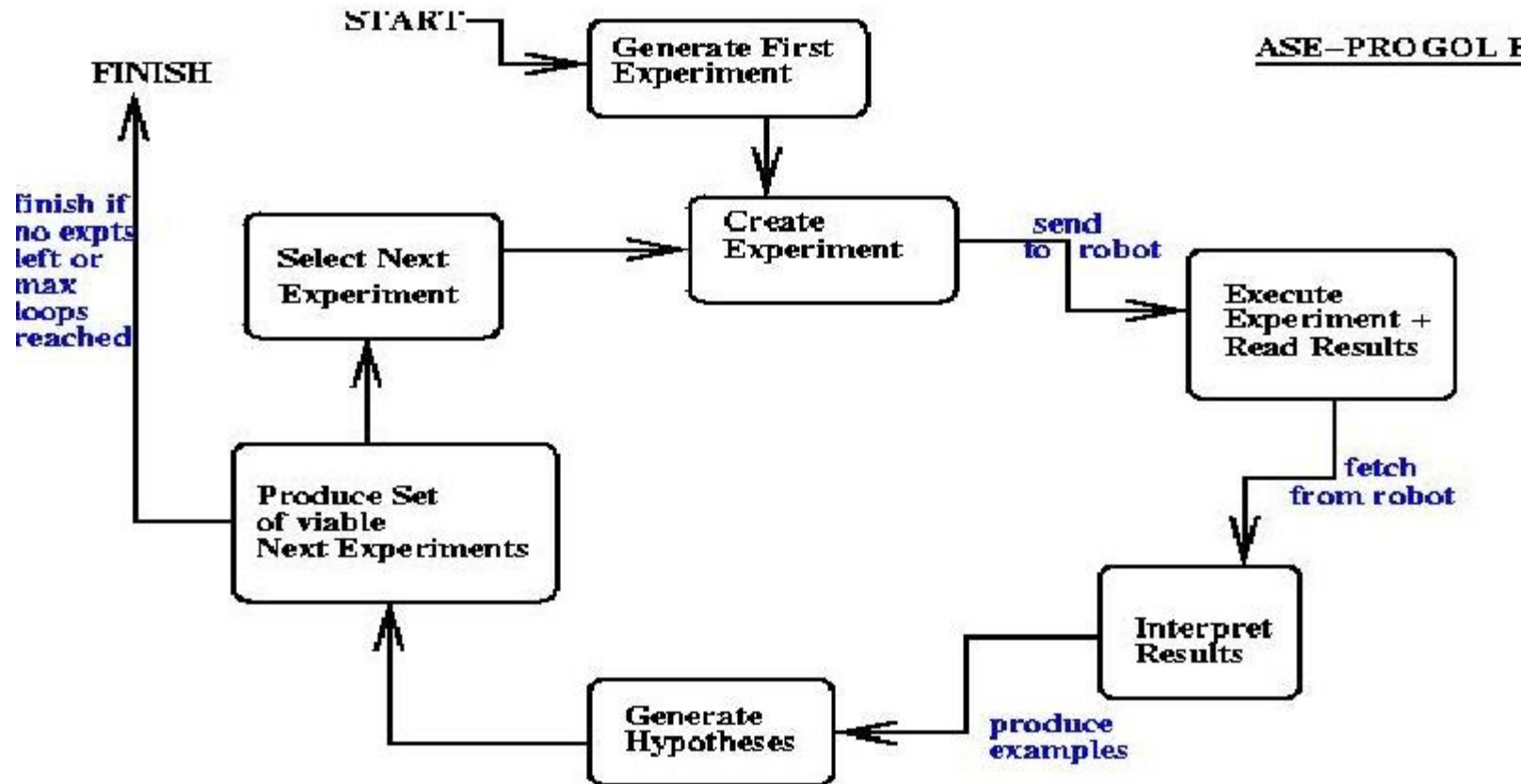


Diagram of Adam

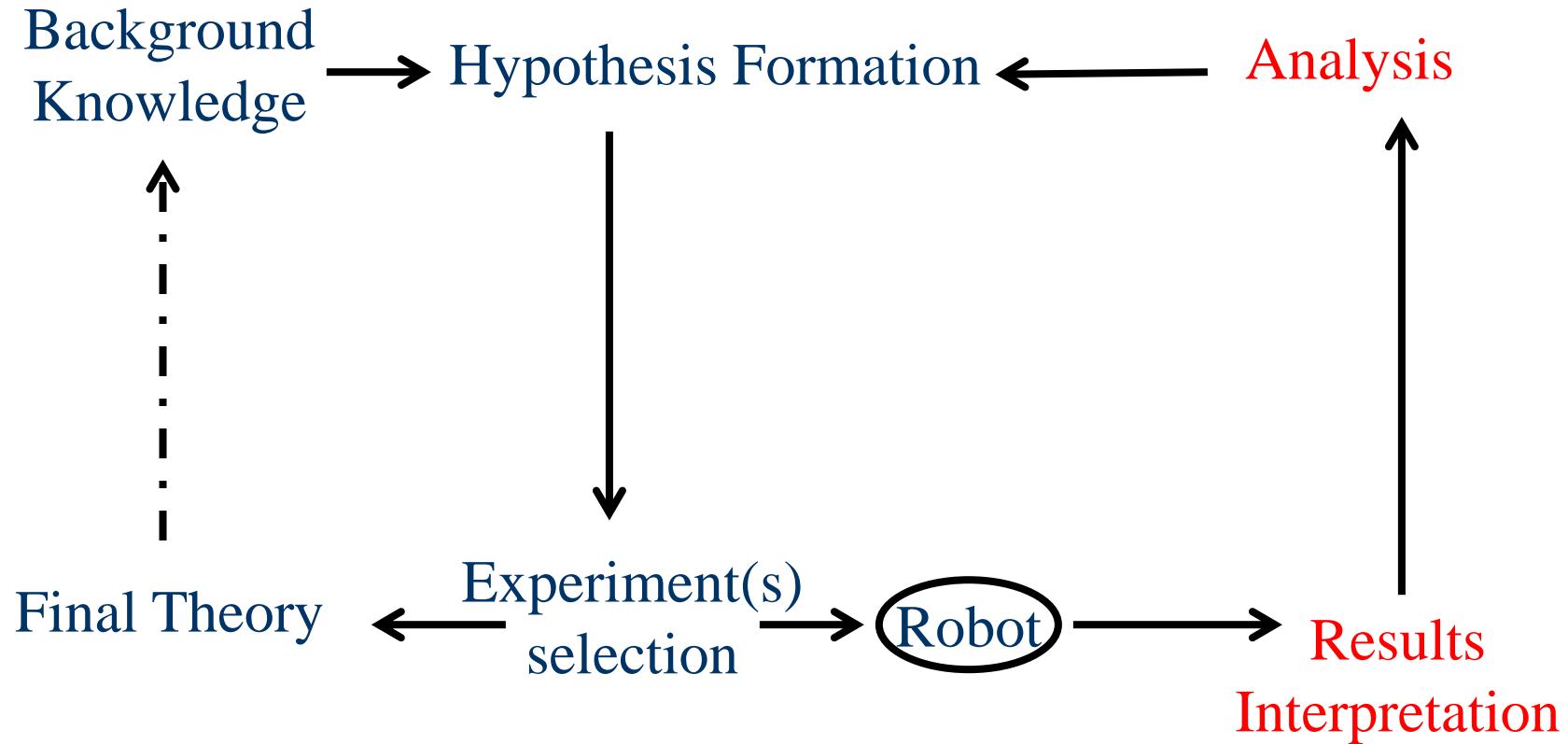


LIMS Setup

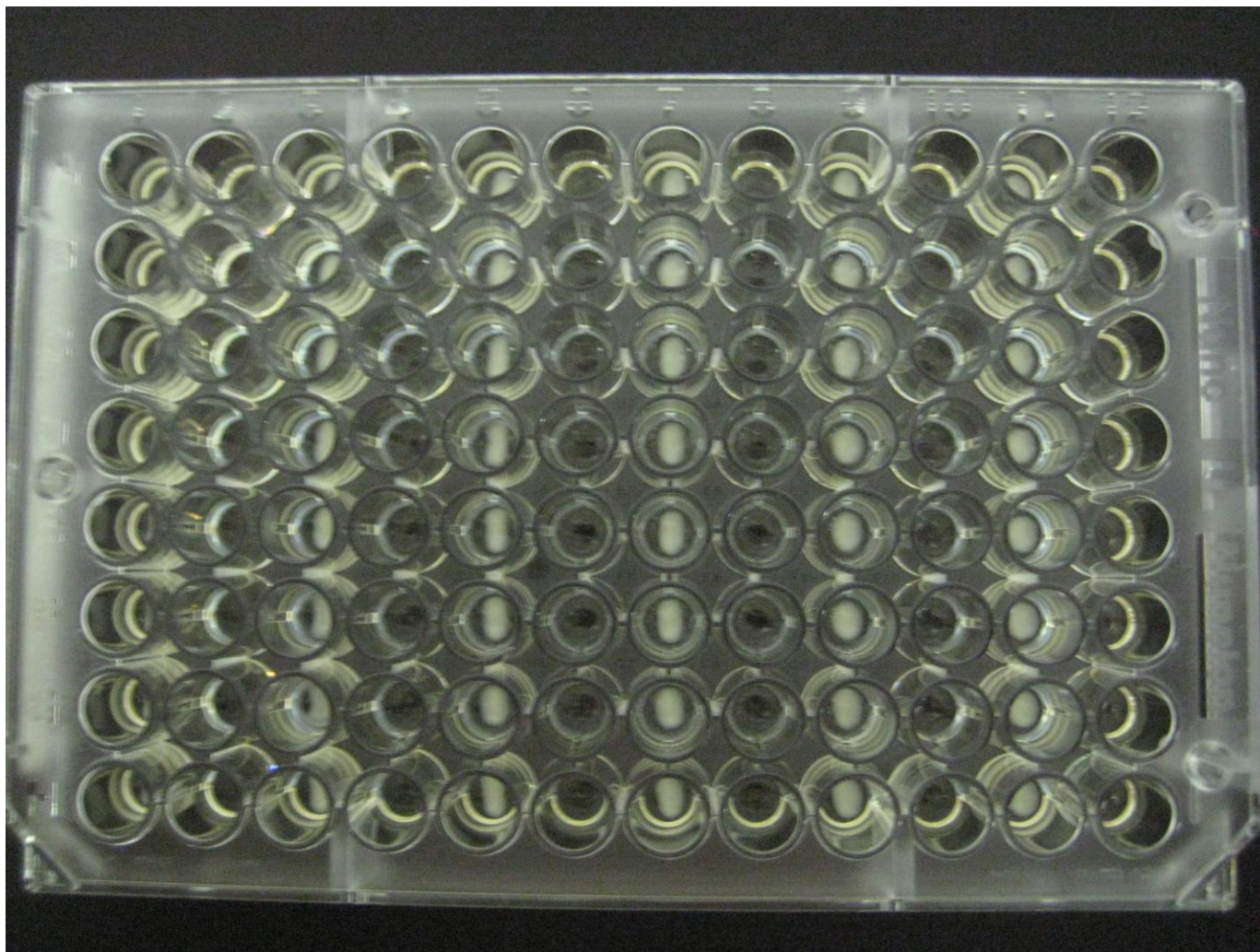


Adam in Action

The Experimental Cycle



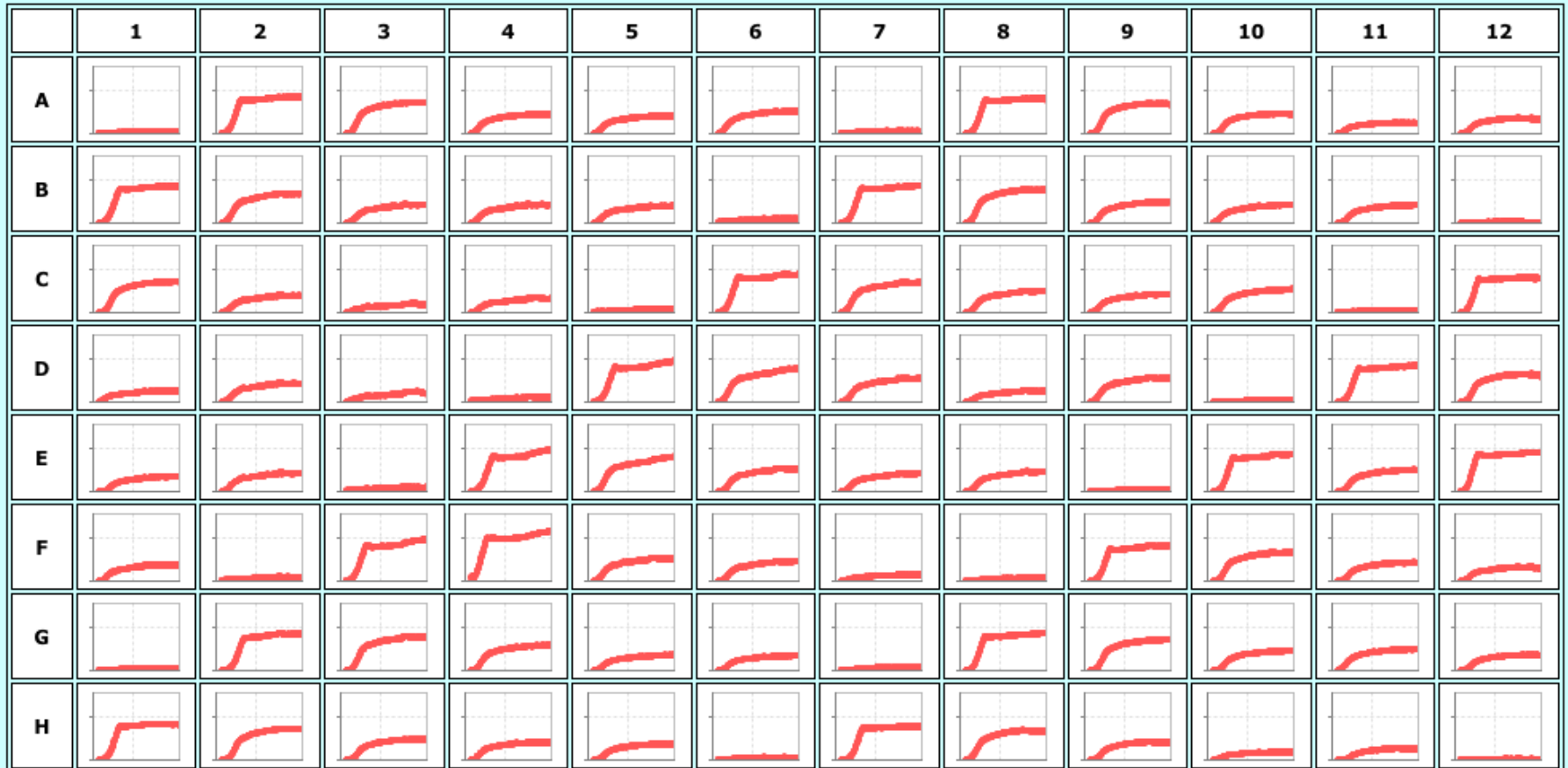
Growth plates



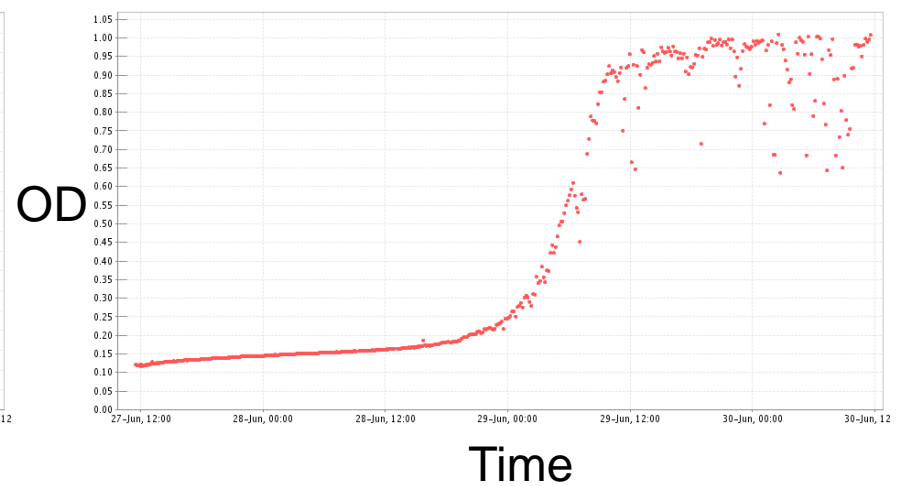
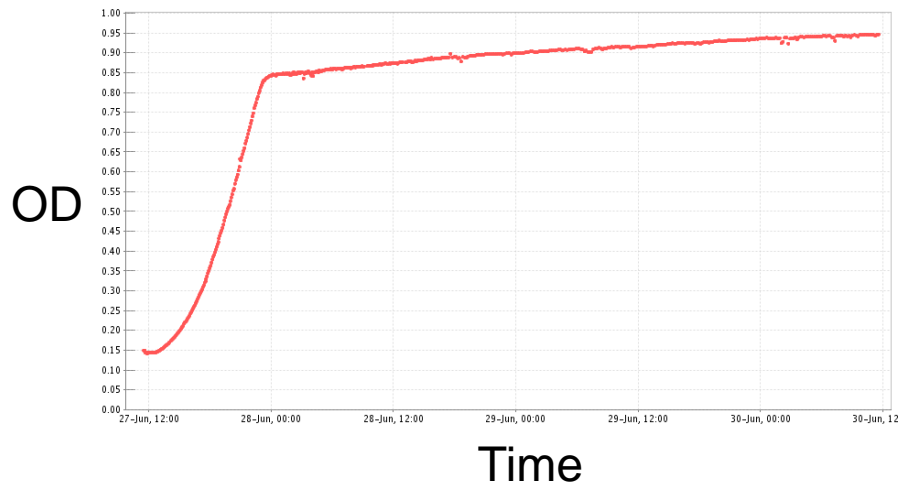
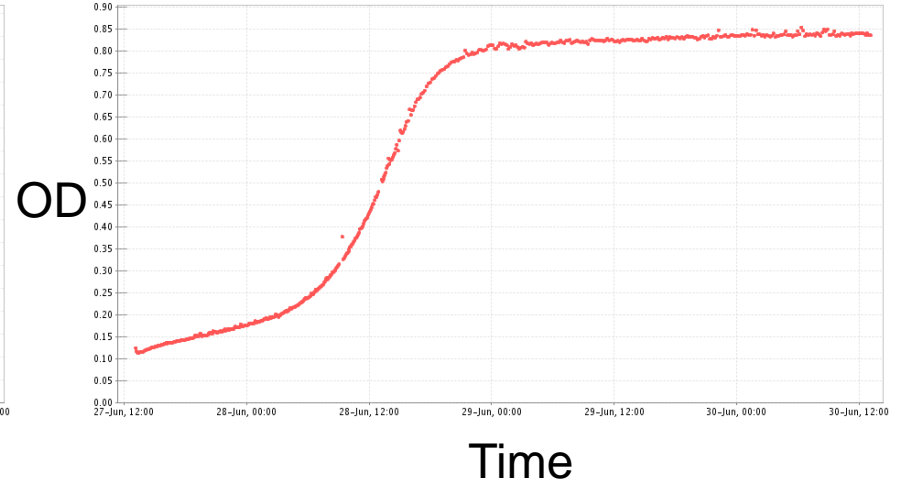
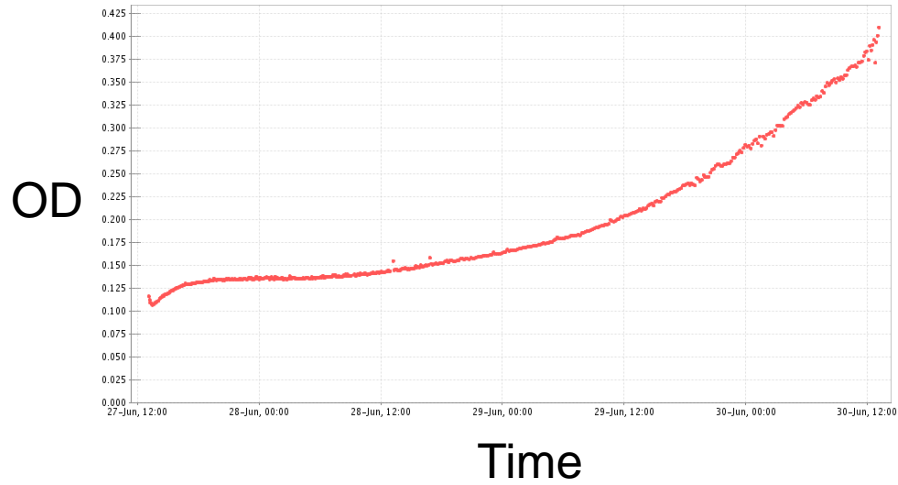
Example Growth Curves

Plate ID: jun09-titr7-expt001-plate001 Barcode 16248

Arginine as N source



Growth curves



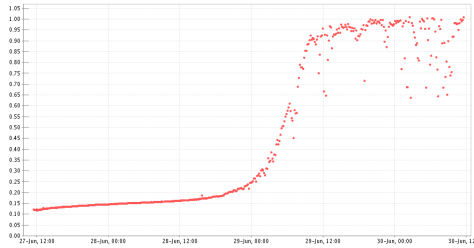
Qualitative to Quantitative

- n The functions of most genes in *S. cerevisiae* that when deleted result in auxotrophy (no growth) have already been discovered.
- n Most genes of unknown function only affect growth quantitatively.
- n They may have slower growth (bradytrophs), faster growth, higher/lower biomass yield, etc.

Test for Growth

- n We needed a reliable and automatic method to decide whether growth of a mutant strain had occurred or not.
- n Controls - positive (wild type) and negative (no yeast)
- n Trained a decision tree to fit predictions of model, attributes aggregates of measurements.
- n Simple tree with only a few leafs

Results



Spline based
curve fitting,
smoothing

Curve parameters

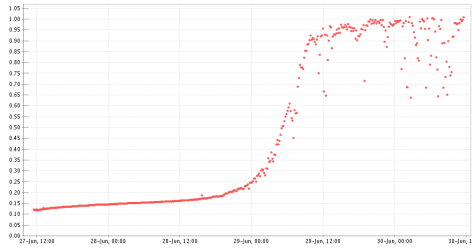
Lag time: 3.5

Maximum OD: 1.4

Growth rate: 0.44

Random forests,
Statistical tests

Results



Spline based
curve fitting,
smoothing

Curve parameters

Lag time: 3.5

Maximum OD: 1.4

Growth rate: 0.44

Two factor design:

Wildtype

Deletant

Wildtype

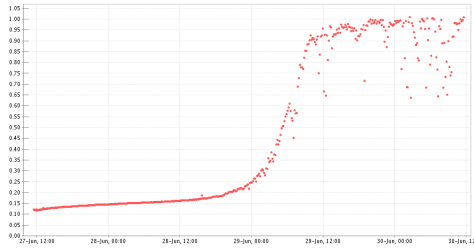
Deletant

with metabolite

with metabolite

Random forests,
Statistical tests

Results



Spline based
curve fitting,
smoothing

Curve parameters

Lag time: 3.5

Maximum OD: 1.4

Growth rate: 0.44

Two factor design:

Wildtype



Wildtype
with metabolite

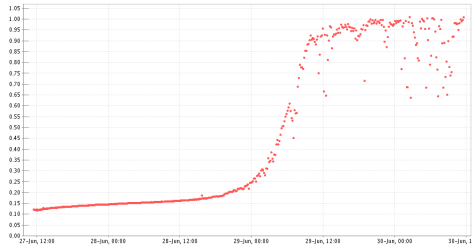
Deletant



Deletant
with metabolite

Random forests,
Statistical tests

Results



Spline based
curve fitting,
smoothing

Curve parameters

Lag time: 3.5

Maximum OD: 1.4

Growth rate: 0.44

Two factor design:

Wildtype

Deletant

Wildtype

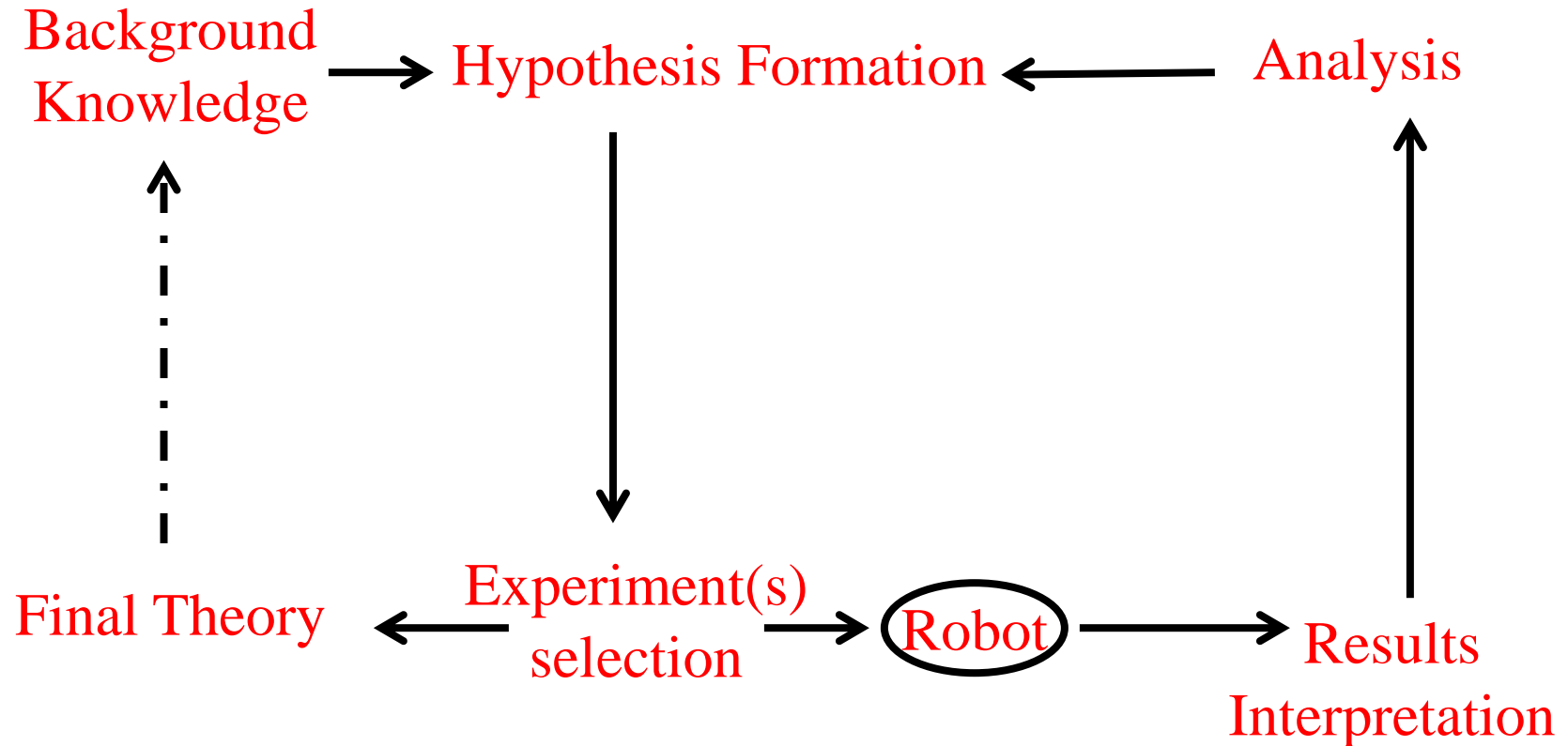
Deletant

with metabolite

with metabolite

Random forests,
Statistical tests

The Experimental Cycle



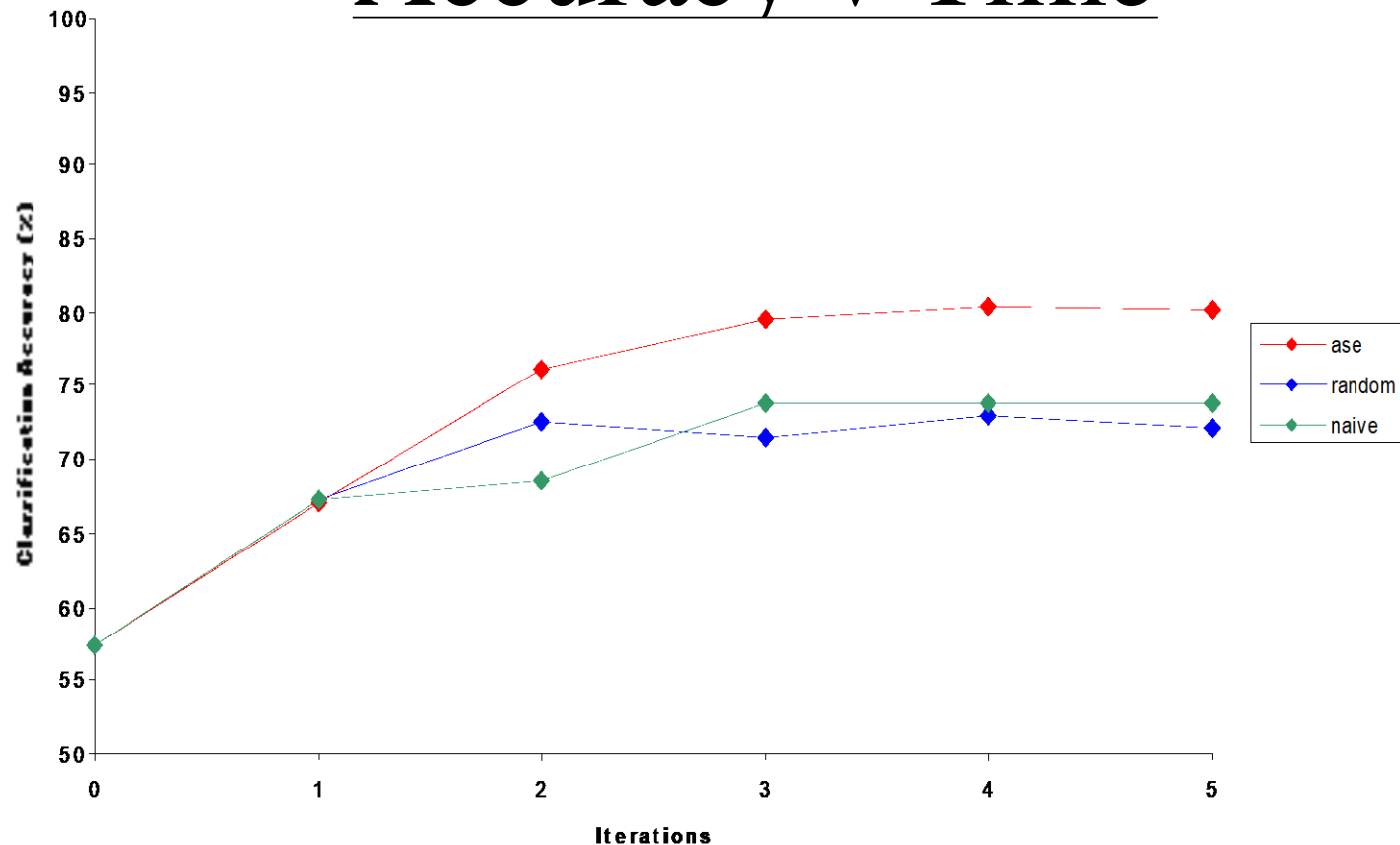
Closing the Loop

- n We physically implemented all aspects of Adam.
- n To the best of our knowledge Adam was the first AI system that can both explicitly form hypotheses and experiments, and physically do the experiments.

Time and Money

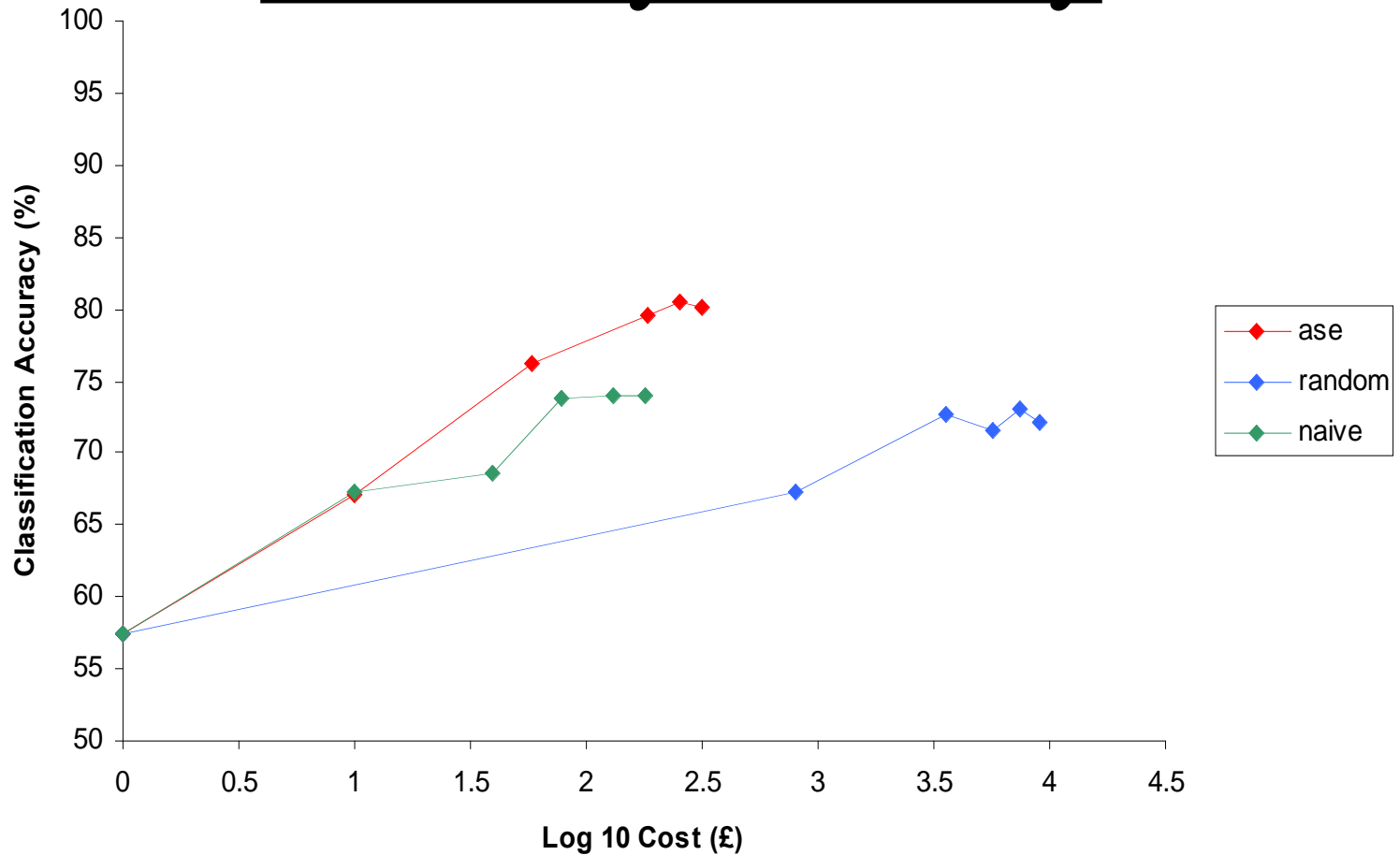
- n “Cost” is a positive function of time & money.
- n ASE dominates for both, therefore ASE dominates for any reasonable cost function.
- n For example: to achieve an accuracy of ~70%, ASE requires fewer trial iterations, and a hundredth of the price, of Random; and almost half the number of iterations, and a third of the price of Naïve.

Accuracy v Time



At the end of the 5th iteration: ASE 80.1%, Naïve 74.0%, Random 72.2%. ASE was significantly more accurate than either Naïve ($p < 0.05$) or Random ($p < 0.07$) using a paired t-test.

Accuracy v Money



Given a spend of $\leq \text{£}10^{2.26}$, ASE 79.5%, Naïve 73.9%, Random 57.4%. ASE was significantly more accurate than either Naïve ($p < 0.05$) or Random ($p < 0.001$).

Human Comparisons

- n We were interested to compare the performance of the Robot Scientist with that of humans.
- n We adopted the simulator to allow humans to choose and interpret the results of cycles of experimentation.
- n Compared nine graduate computer scientists and biologists.
- n No significant difference between the best humans and the Robot

Discovery of Novel Science

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.

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 | ✓ | - |

A 50 Year Old Puzzle

- n The enzyme 2-aminoadipate: 2-oxoglutarate aminotransferase was a locally orphan enzyme.
- n It is in the lysine biosynthesis pathway which has been studied for 50 years in fungi: target for antibiotics, and on path to penicillin.
- n Adam formed three hypotheses for the gene to encode this enzyme: YER152C, YJL060W, and YGL202W (in that order of probability).
- n Currently KEGG states that YGL202W is the gene.
- n Evidence from 1960's that at least 2 isoenzymes are involved.

Confirmed New Knowledge

- n Adam's differential growth experiments were consistent with all three genes encoding 2-oxoglutarate aminotransferase.
- n Manual experiments: purified protein + enzyme assays are consistent.
 - YGL202W literature confirmed.
 - YJL060W (was annotated as an arylformamidase, new (08) annotation kynurenine aminotransferase)
 - YER152C (currently not annotated)
- n YGL202W & YJL060W double knockout is lethal

Formalising Science

Formalisation of Science

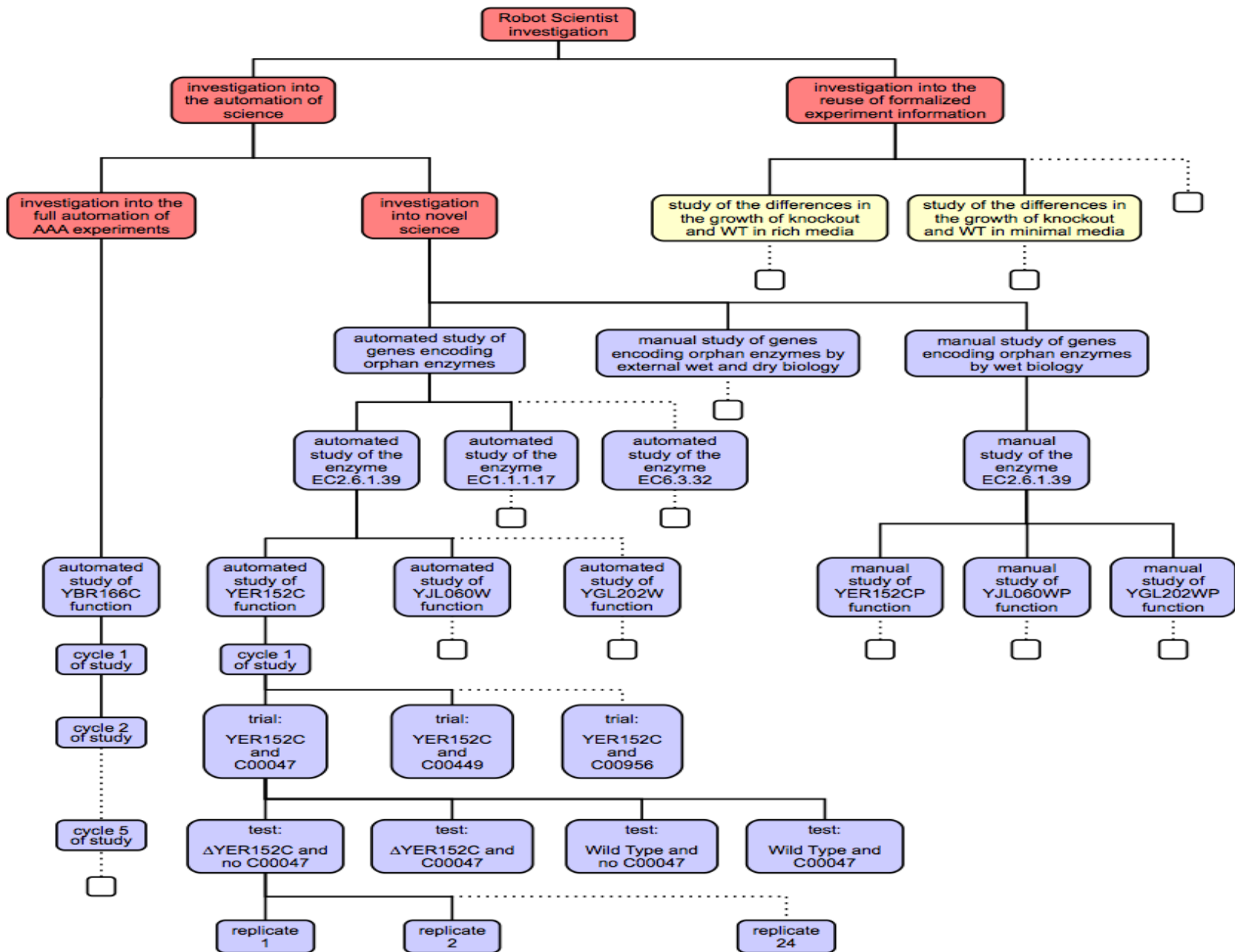
- n The goal of science is to increase our knowledge of the natural world through the performance of experiments.
- n This knowledge should be expressed in formal logical languages.
- n Formal languages promote semantic clarity, which in turn supports the free exchange of scientific knowledge and simplifies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

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

with enzyme class

has organism class

has ncbi taxonomy

has hypothesis

has research

has negative

has cycle 1 of

has study result

encodes(yer152c)

highest

proportion

has study condition

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: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, a:goal(X) :- a:to_test_the_hypothesis_that_g(X), a:_encodes_an_enzyme_with_enzyme_class(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>

```

Eve



Drug Design

Parasitic Diseases targeted



Malaria



Shistosomiasis



Leishmania

Chagas



Why Tropical Diseases?

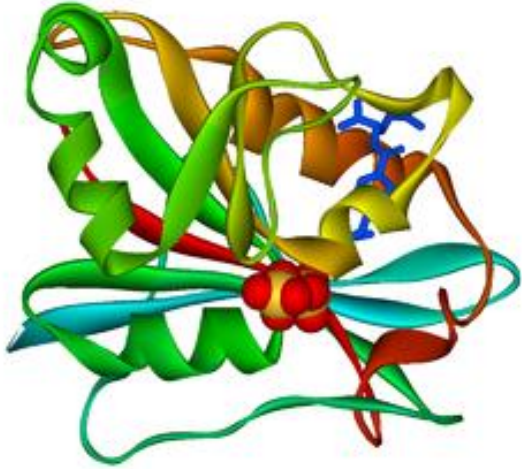
- n Millions of people die of these diseases, and hundreds of millions of people suffer infection.
- n It is clear how to cure these diseases – kill the parasites.
- n They are “neglected”, so avoid competition from the Pharmaceutical industry.

Parasites targeted

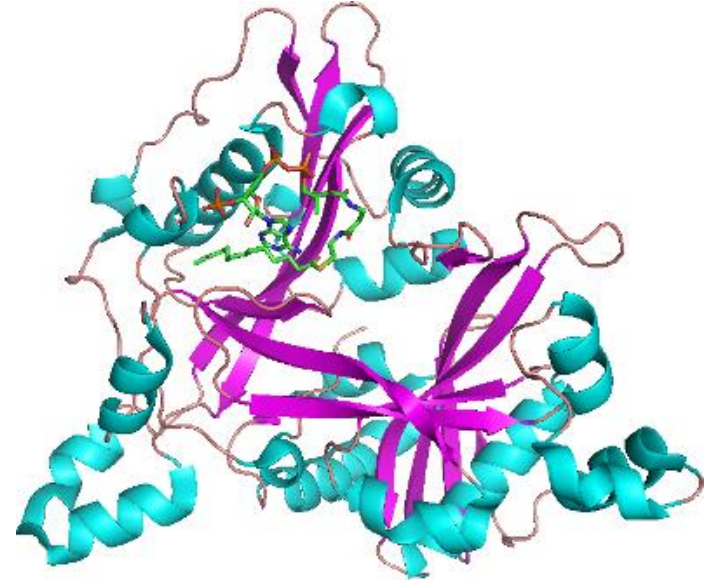
n Organism/disease:

- *Plasmodium falciparum* (malaria),
- *Plasmodium vivax* (malaria),
- *Trypanosoma brucei* (sleeping sickness),
- *Trypanosoma cruzi* (Chagas),
- *Leishmania major* (leishmania),
- *Schistosoma mansoni* (shistosomiasis),
- Staphylococcus aureus
- ...

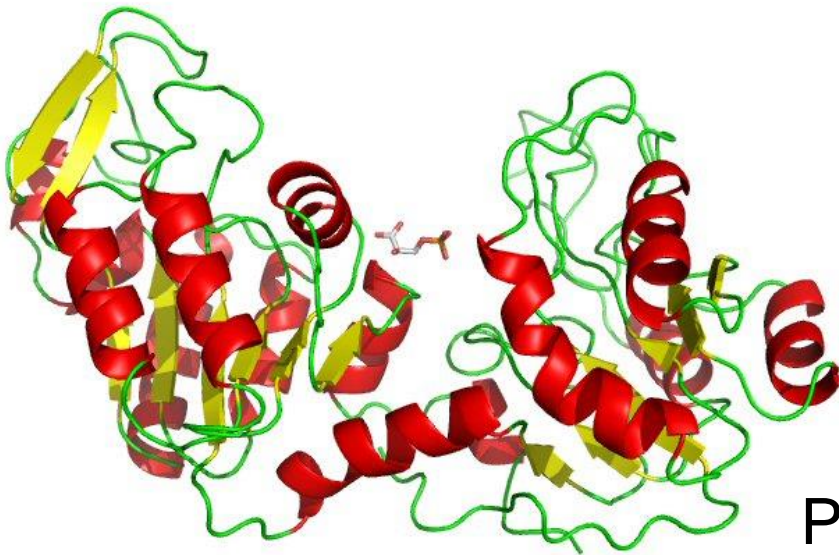
Enzymes Targeted



Dihydrofolate Reductase (DHFR)



N-myristoyl transferase

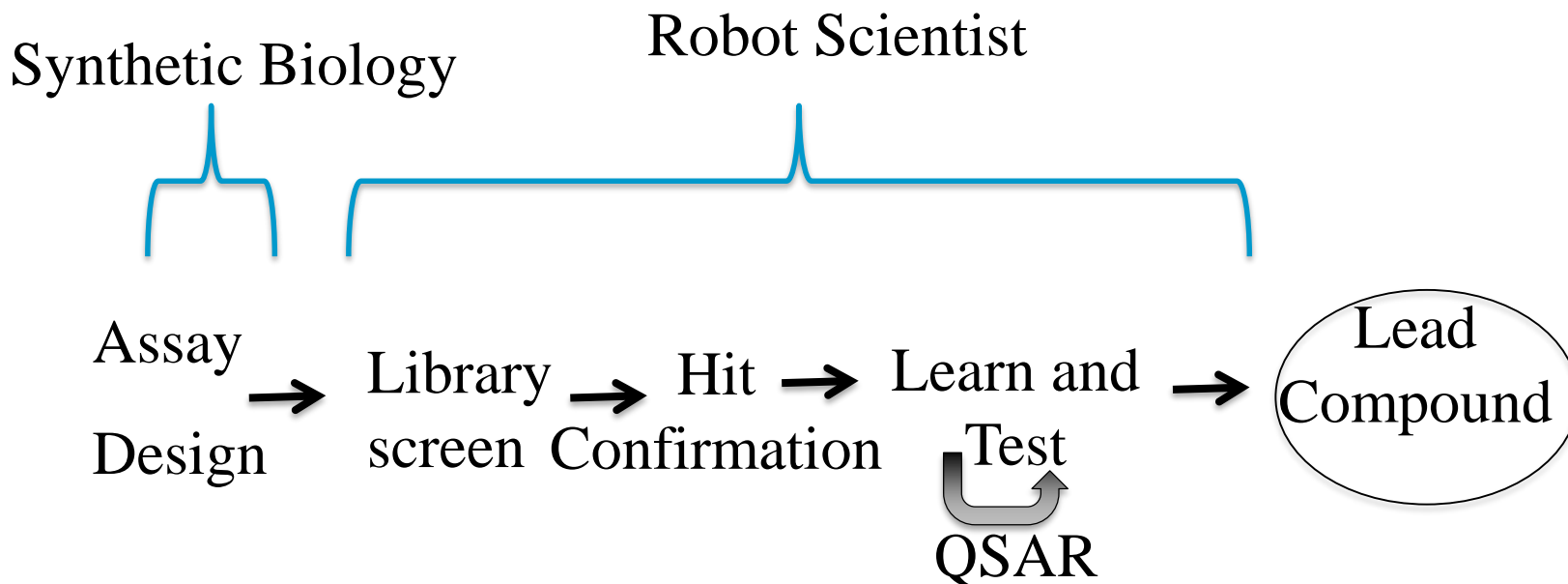


Phosphoglycerate kinase

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.

Automating Early Drug Development



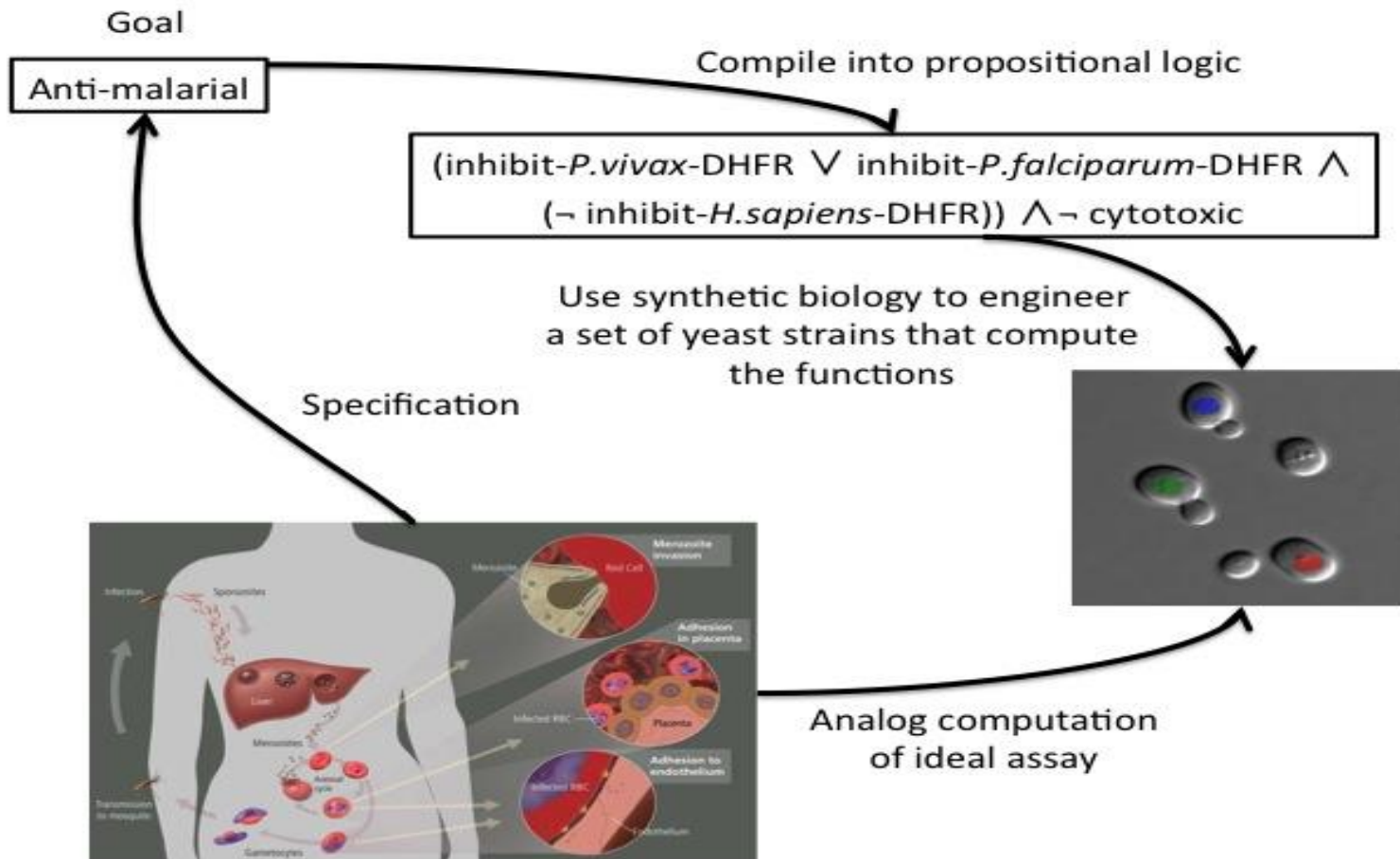
Design

- n Novel design features: Combines: screening, hit confirmation, and QSAR learning.
 - Uses synthetic biology based yeast assays.
 - During the standard screening process Eve is able to decide to switch to QSAR mode.
 - Use cycles of active learning to improve QSARs.

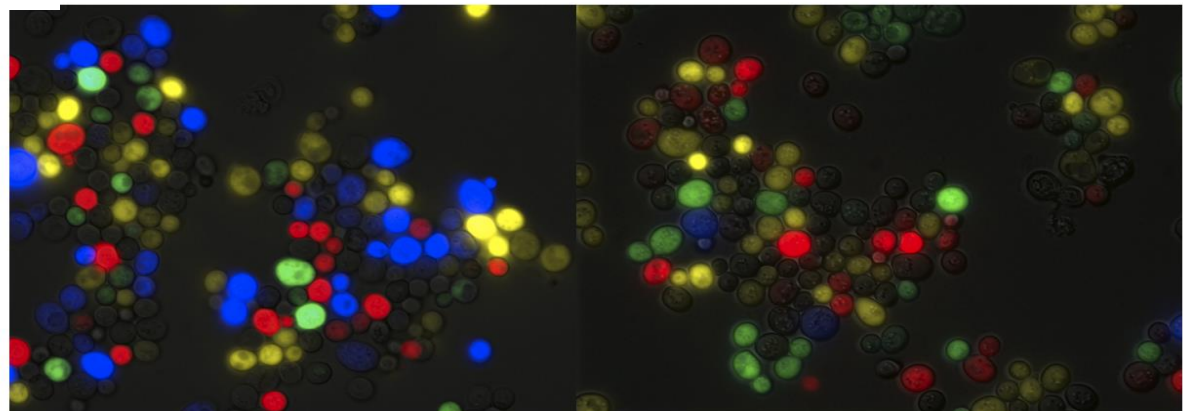
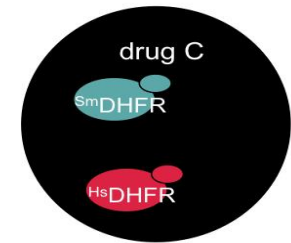
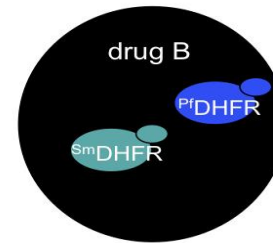
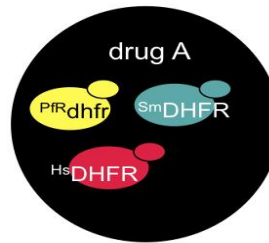
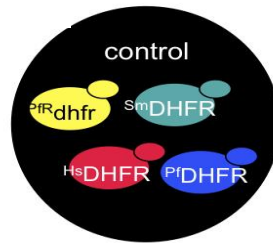
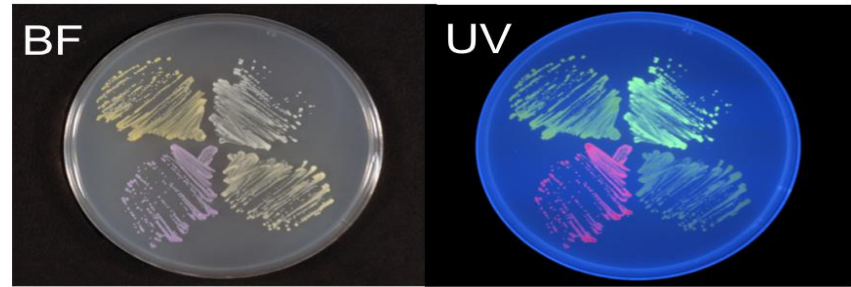
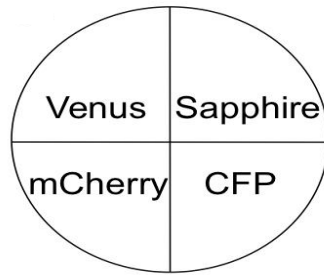
Synthetic Biology based Assays

- n Our idea is to engineer cells to be analog computers.
- n These computers will accurately estimate a biological function that corresponds to the set of desired assay properties.
- n The function estimated is the utility of a compound against a disease.
- n E.g. ((inhibit *P. vivax* DHFR) \wedge (\neg inhibit *H. sapiens* DHFR) \wedge (\neg cytotoxic)).

Synthetic Biology Workflow



Yeast Strains



No drug

Pyrimethamine

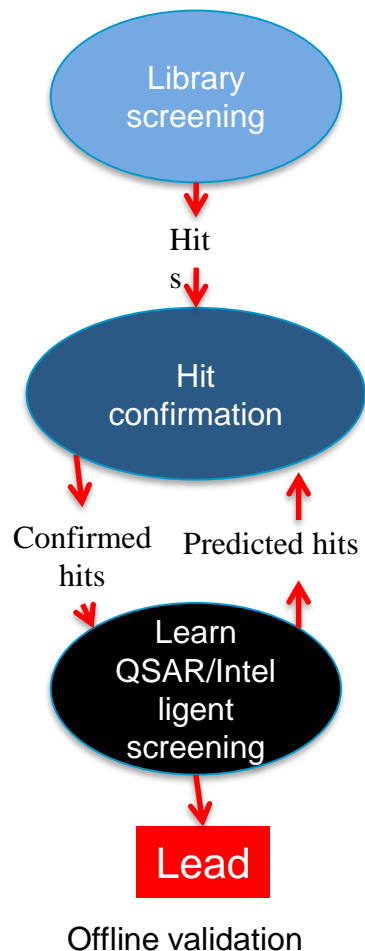
● y^{Hs}DHFR+yEp mCherry
● y^{Pf}DHFR+yEp Sapphire

● y^{PFR}dhfr+yEp Venus
● ySmDHFR+yEp CFP

Learning QSARs

- n Almost every form of statistical and machine learning method you can think of has been applied to QSAR learning.
- n Leading methods are logistic regression, support vector machines, random forests. ...
- n Eve currently uses Gaussian process models. Has the advantages of being generative and outputting probabilities – helps active learning.

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.

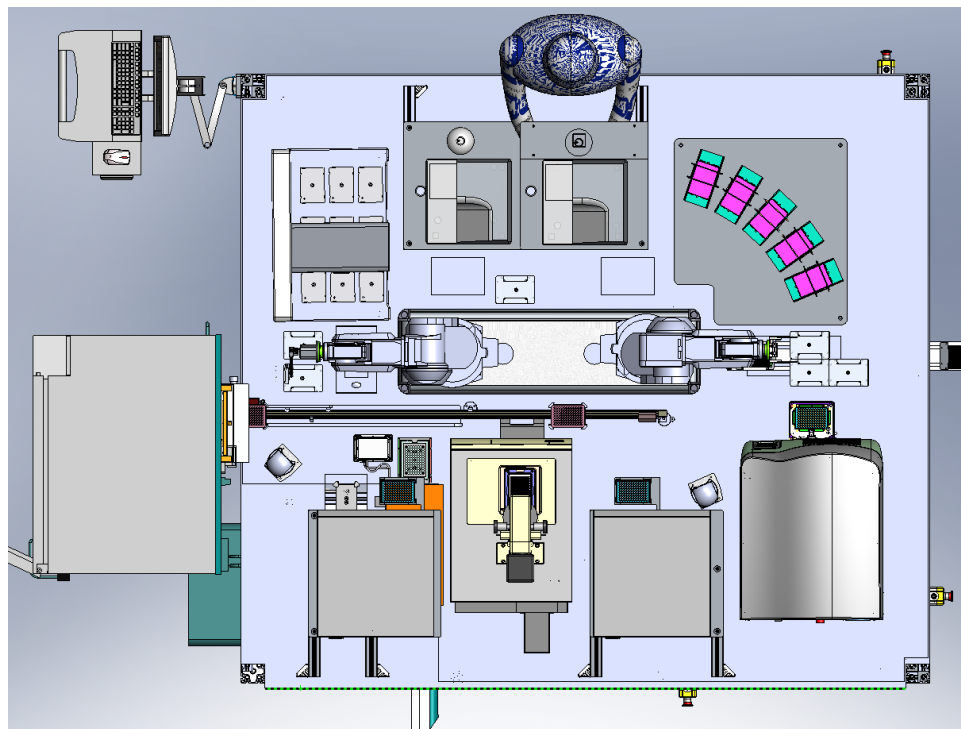
Intelligent v Brute-force Screening

- n We wished to compare our AI based screening against the standard brute-force approach: “begin at the beginning and go on till you come to the end: then stop” (Lewis Carroll).
- n While simple to automate standard screening is slow and wasteful of resources, since every compound in the library is tested. It is also unintelligent, as it makes no use of what is learnt during screening.
- n Use money to decide.

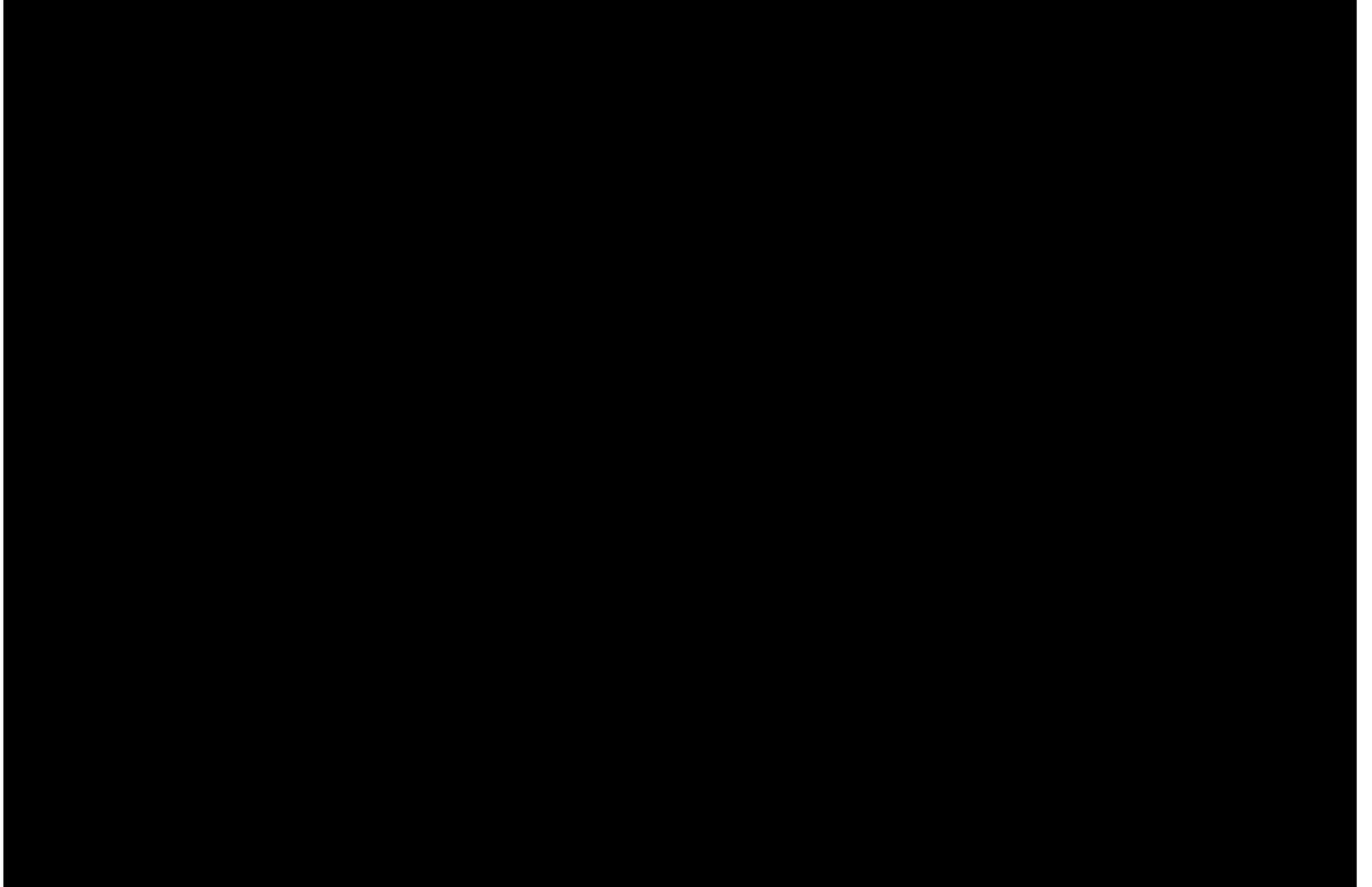
Eve's Hardware

Highlights of Eve's hardware:

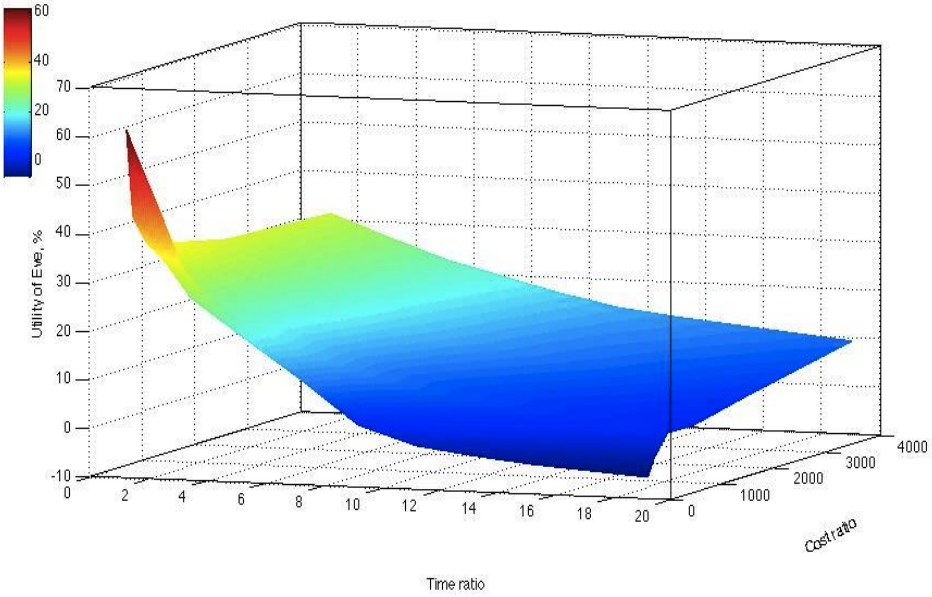
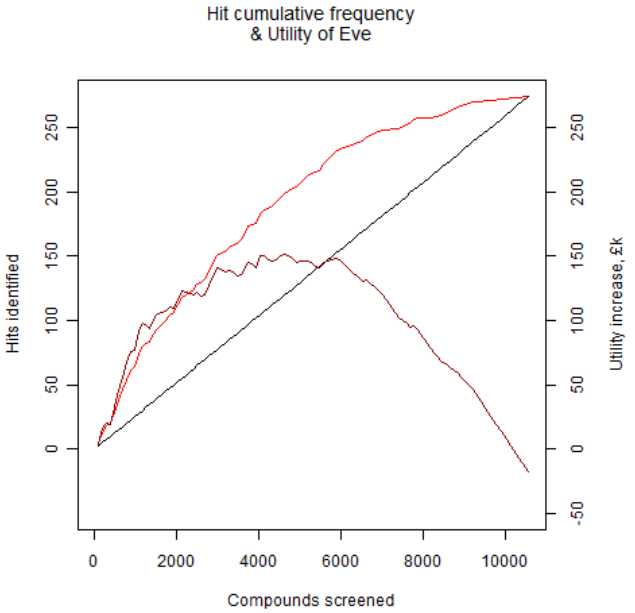
- Acoustic liquid handling
- High throughput 384 well plates
- Two industrial robot arms
- Automated 60x microscope
- Liquid handlers, fluorescence readers, barcode scanners, dry store, incubator, tube decapper ...



Eve



The Economics of Intelligent Screening

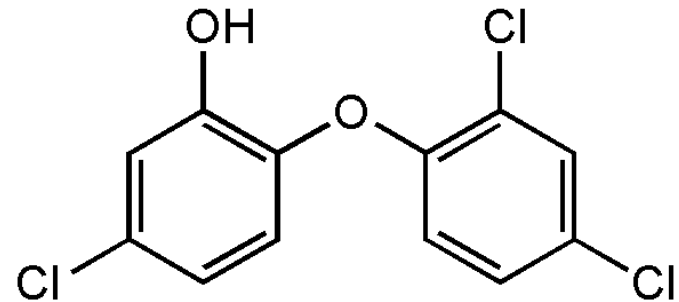


$$\Delta \text{Utility of Eve} = \sum_1^{Nm} (Tm + Cm) + \sum_1^{Nx} (Tc + Cc - Uh) + \sum_1^{Ne} (Tm - Tc + Cm - Cc)$$

- 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

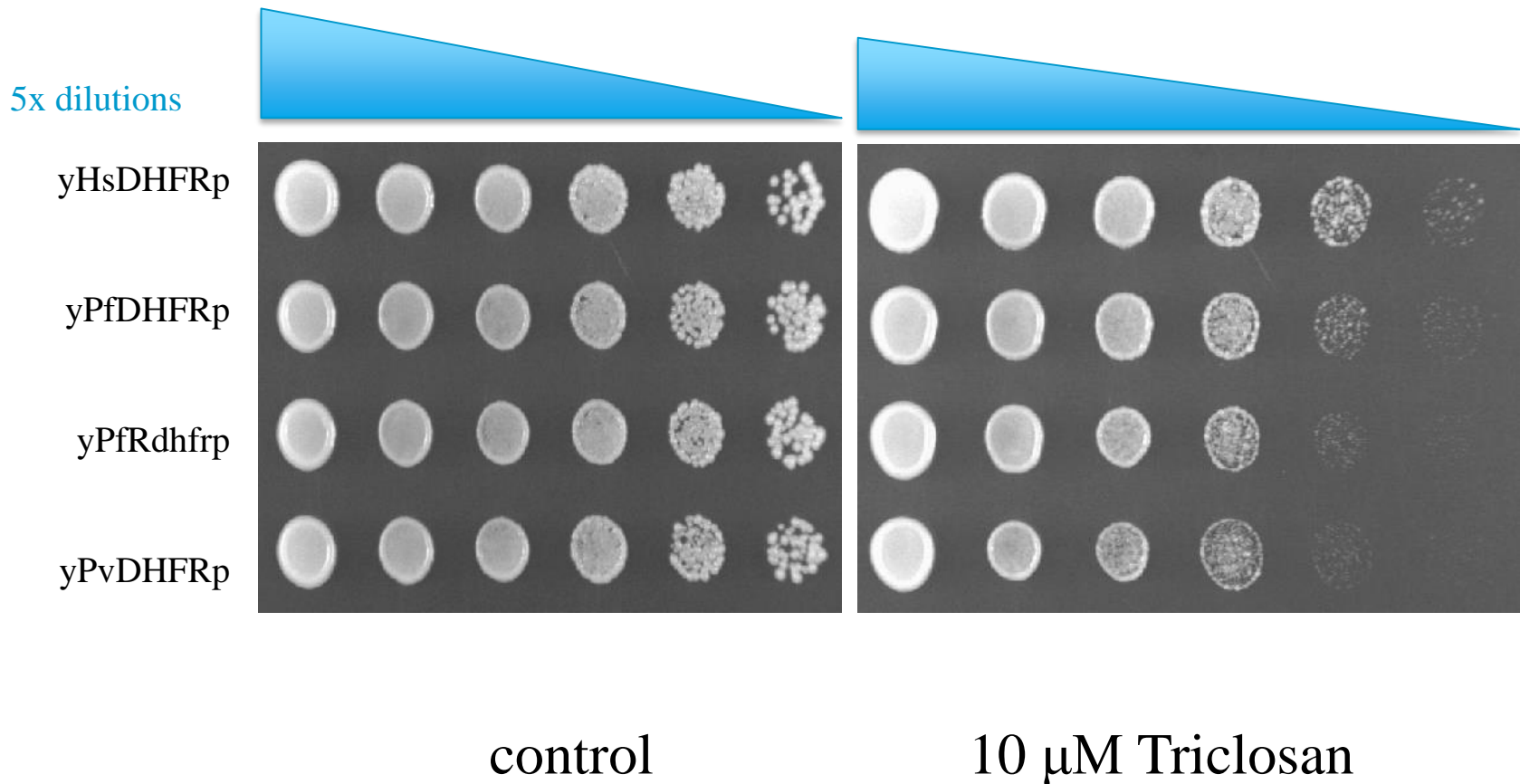
Triclosan

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*.

Triclosan specifically inhibits *Plasmodium* DHFRs



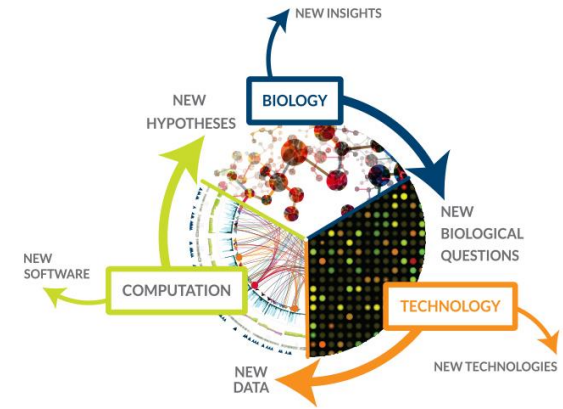
Eve



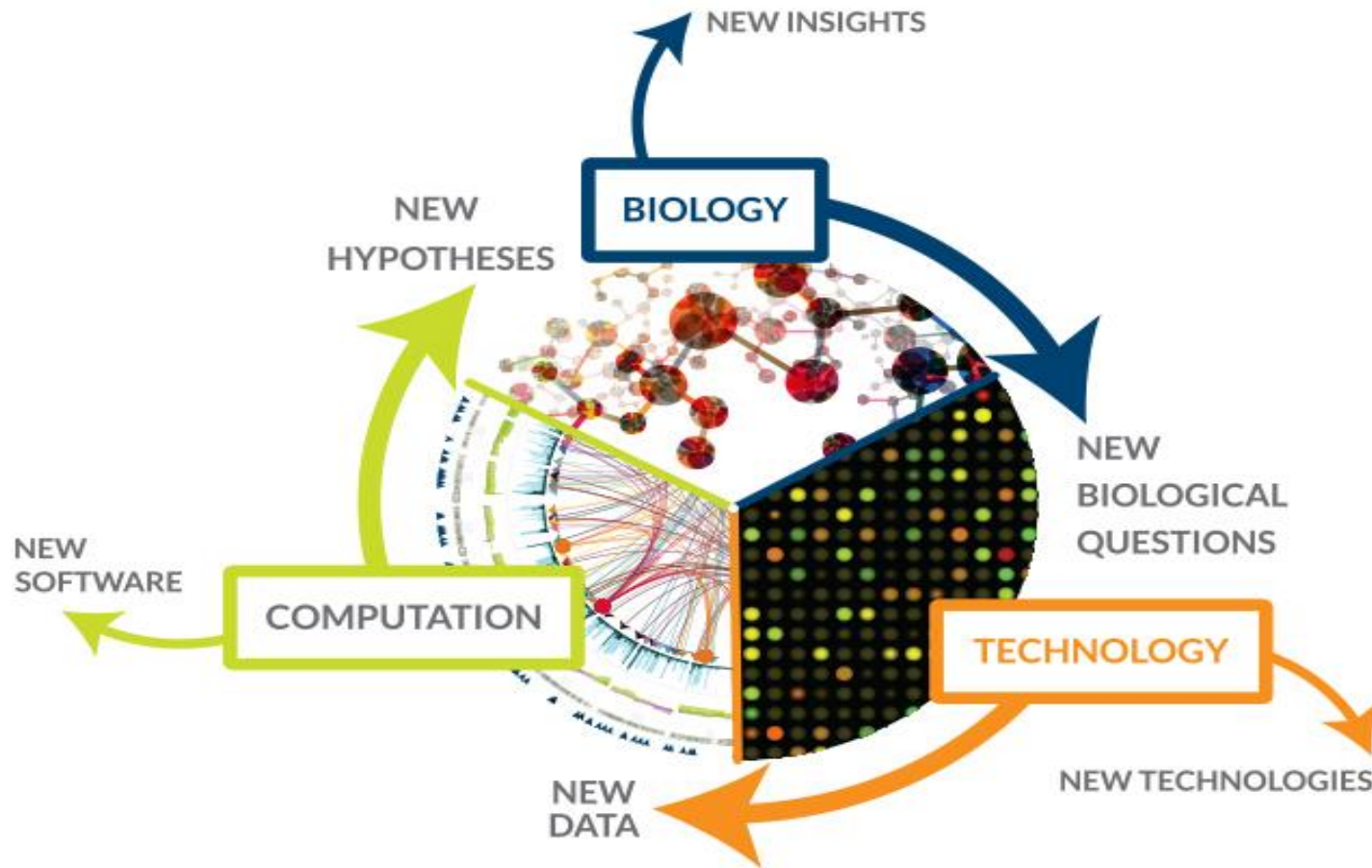
Systems Biology

The Application Domain

- n Systems Biology.
- n The modelling of Biological Systems.
- n Diauxic shift in yeast (*S. cerevisiae*).
- n Glucose → Ethanol → CO₂
- n A model of Biological Switching.
- n Understanding important for cancer – Warburg effect.
- n Understanding important for aging.



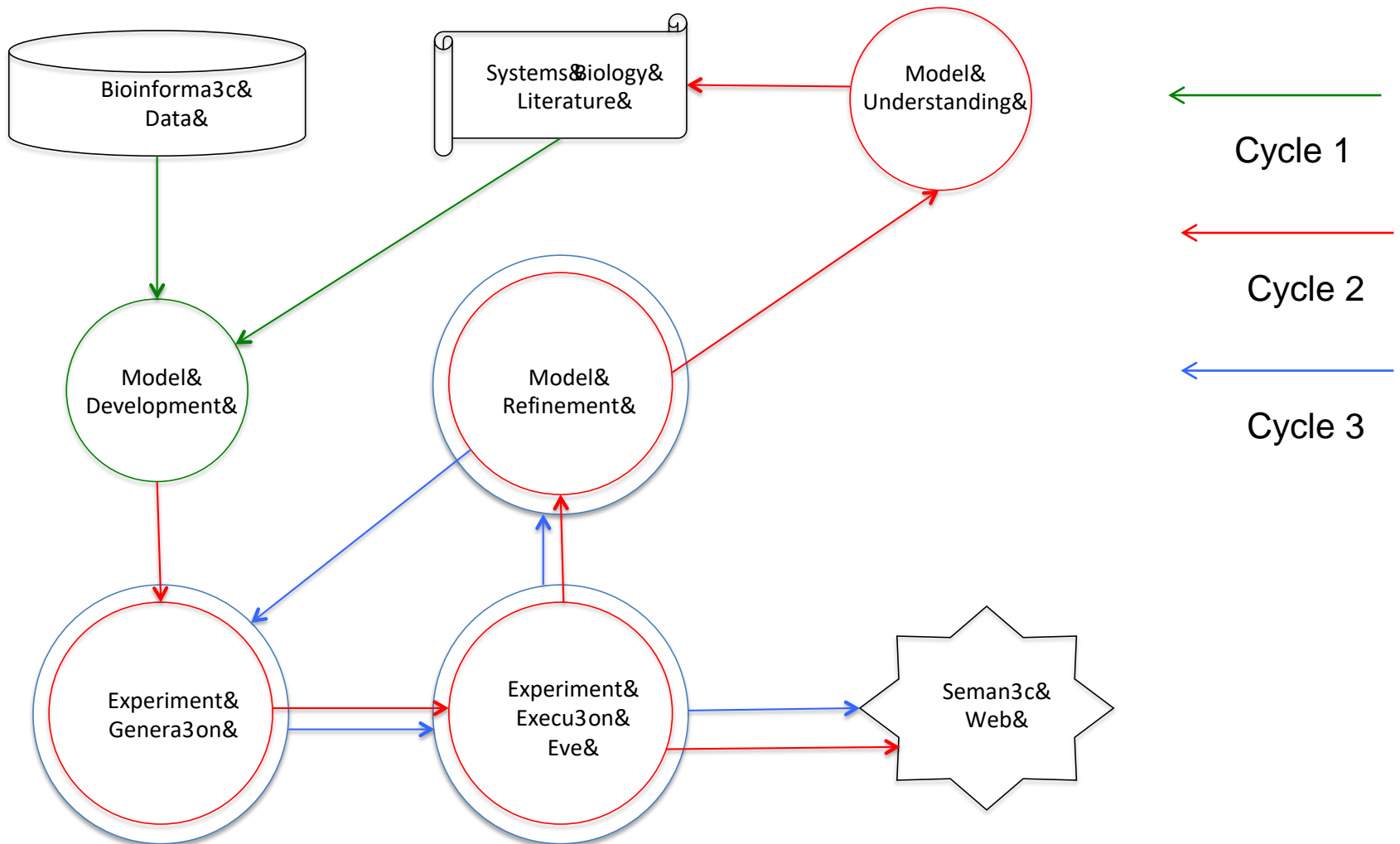
Systems Biology

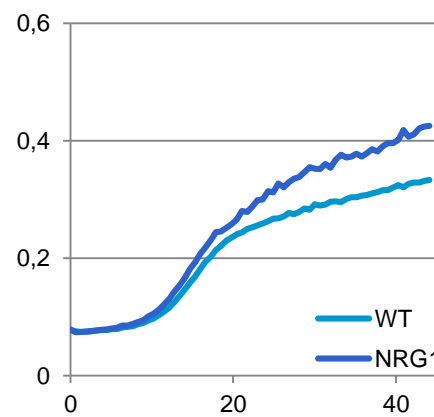
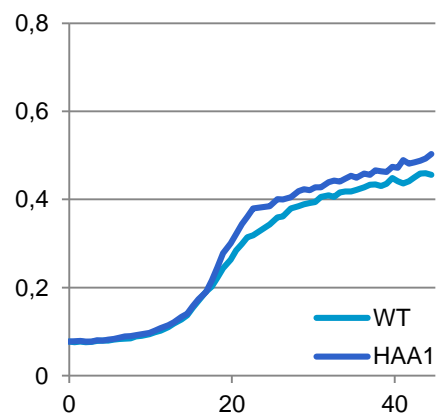
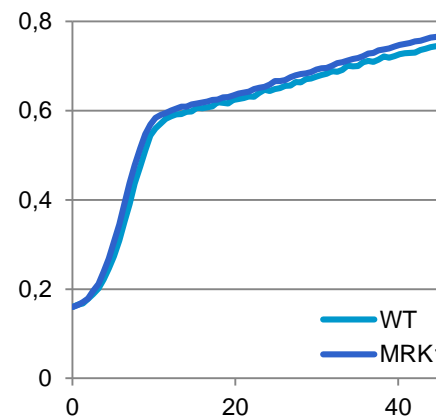
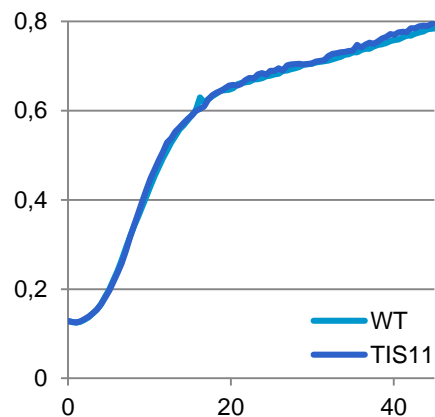
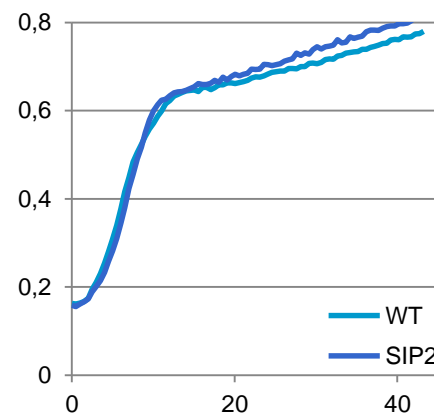
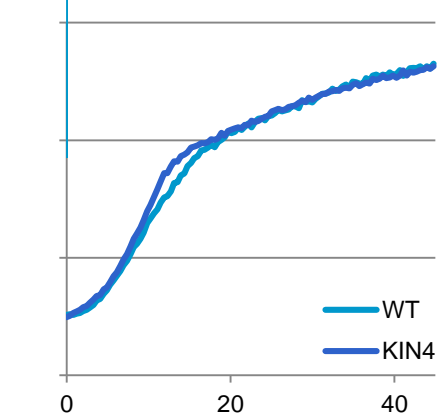
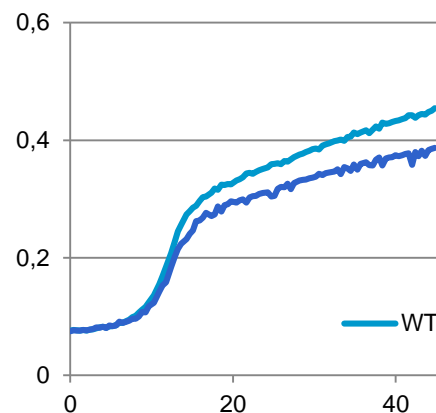
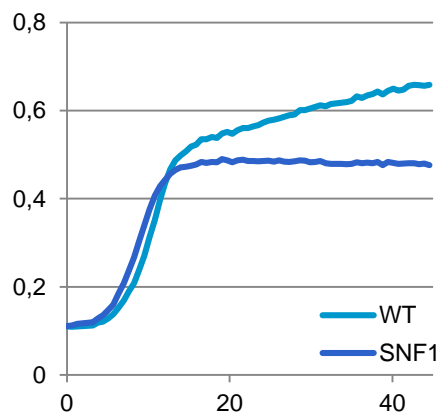
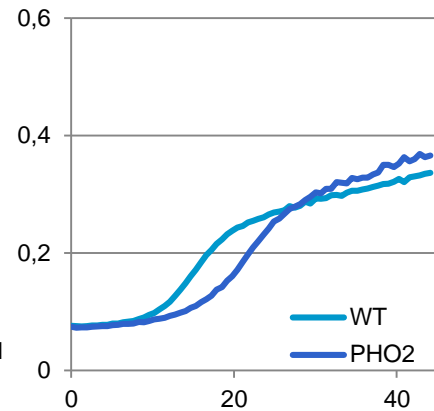
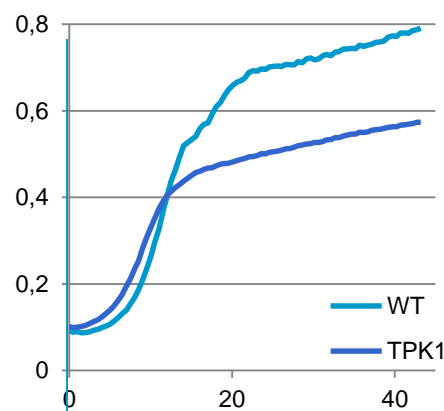
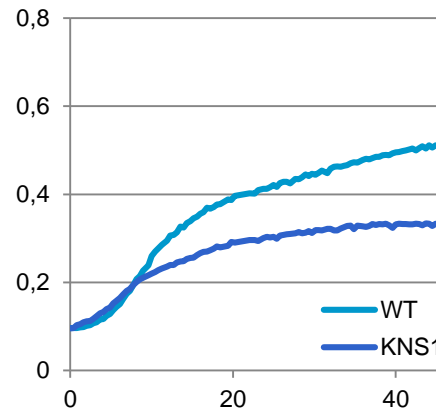
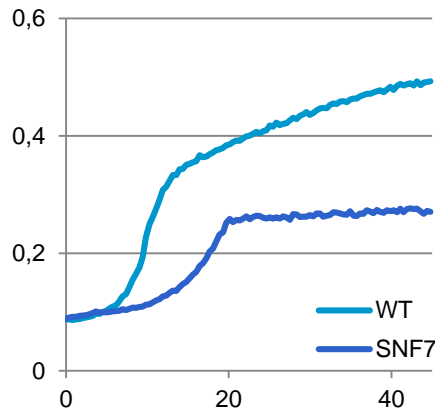


AdaLab Software

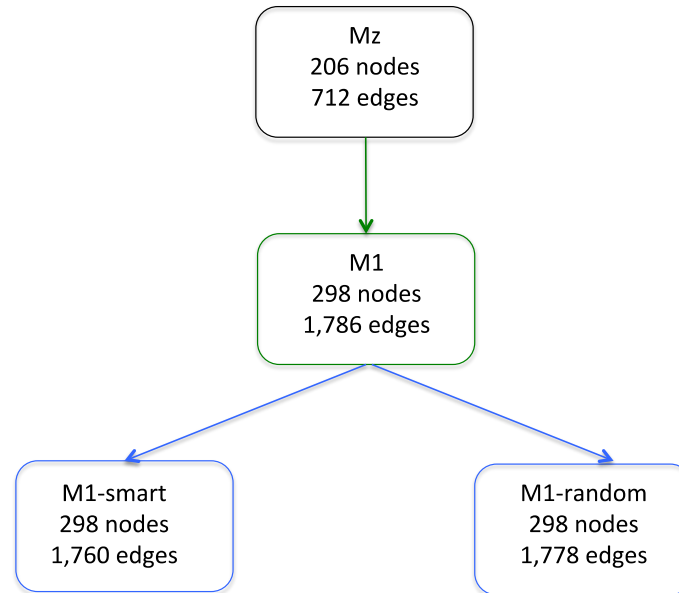
| Software Type | Name | Details |
|--|-------------|---|
| <i>AI Tools</i> | CoRegNet | Reconstruction and integrated analysis of co-regulatory networks. |
| | CoRegMine | Experiment design |
| | ELSA | Ensemble Learning of Spanning Arborescences. |
| | AdactiveFB | Active learning based on forward & backward simulation. |
| | MinerLC* | Graph mining tool. |
| | Adarev | Experiment design |
| | Adana | Data analysis |
| <i>Semantic Web Tools and Ontologies</i> | AdaLab-meta | An ontology for the description of metadata about datasets. |
| | AdaLab | A domain ontology to represent relevant system biology biological entities. |
| | UNO | An ontology of uncertainties. |
| | Eve-CV | Eve experiments control vocabulary |
| | RDF | RDF Knowledgebase |
| <i>Bioinformatic Resources</i> | Brauer | Microarray data of the yeast diauxic shift. |
| | Yeastract | A curated repository of regulatory associations between transcription factors and target genes in yeast. |
| | YeastKinome | A yeast kinase and phosphatase interactome resource. |
| <i>Systems Biology Models</i> | iMM904 | Model of yeast metabolism: a Flux Balance Analysis (FBA) model. |
| | Mz | Diauxic shift model derived from the literature. |
| | M1 | Diauxic shift model enhanced using bioinformatic data. |
| | M1-random | Diauxic shift model enhanced using bioinformatic data and high-throughput experiments. |
| | M1-smart | Diauxic shift model enhanced using bioinformatic data and hypothesis led experiments. |
| <i>Systems Biology Simulation</i> | DBN | Simulation of yeast cell signaling: a dynamic two-time slice Bayesian network, with linear Gaussian parameters. |
| | DFBA | Simulation of yeast metabolism: Dynamic Flux Balance Analysis. |
| <i>Statistics</i> | Yeast-stats | Yeast growth parameter estimation. |
| <i>Laboratory Robotics</i> | Overlord | Laboratory automation control. |

AdaLab Project: 3 Cycles





AdaLab Models



M1-smart > M1-random > M1 > Mz

| Cycle | Test Strains | Mz error | M1 error | M1-random error | M1-smart error | Wins | Ratio | Significance |
|-------|--------------|----------|------------------------|------------------------|------------------------|------|-------|-------------------------|
| 1 | 192 | 0.17 | 0.0033 | - | - | 192 | 98% | $< 2.2 \times 10^{-16}$ |
| 2 | 281 | - | 5.533×10^{-3} | 3.892×10^{-3} | - | 259 | 30% | $< 2.2 \times 10^{-16}$ |
| 2 | 281 | - | 5.533×10^{-3} | - | 2.158×10^{-3} | 269 | 74% | $< 2.2 \times 10^{-16}$ |
| 3 | 81 | - | - | 1.757×10^{-3} | 7.231×10^{-4} | 79 | 58.4 | $< 7.6 \times 10^{-9}$ |

Future Prospects

Next Generation Robot Scientist

- n Domain Systems Biology.
- n One of the most challenging tasks in modern science.
- n Ability to do ~10,000 individual hypothesis-led experiments in parallel.
- n AI able to automatically design, plan, execute experiments, and modify systems biology model.

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.