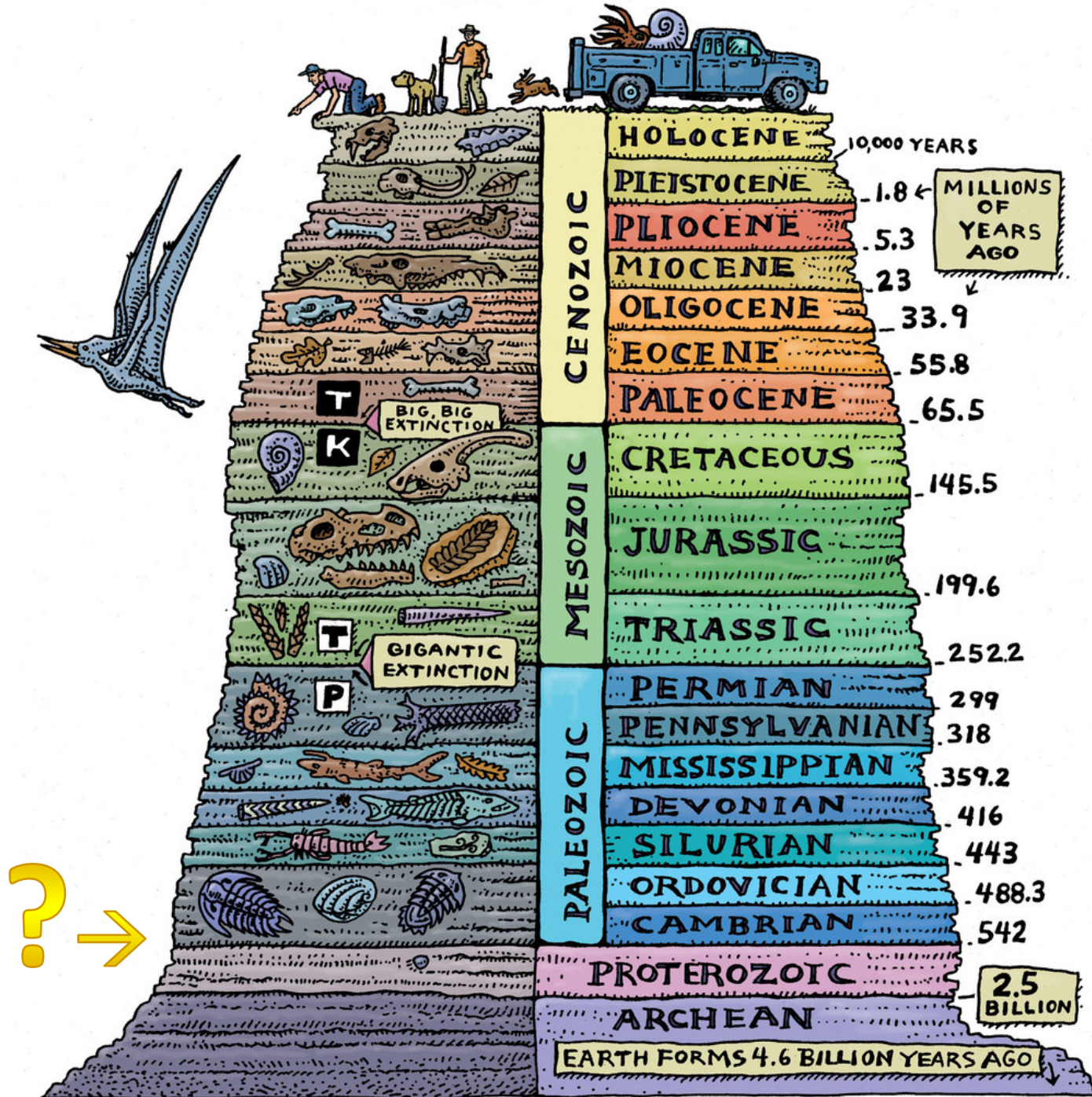


Flipping the Light Switch

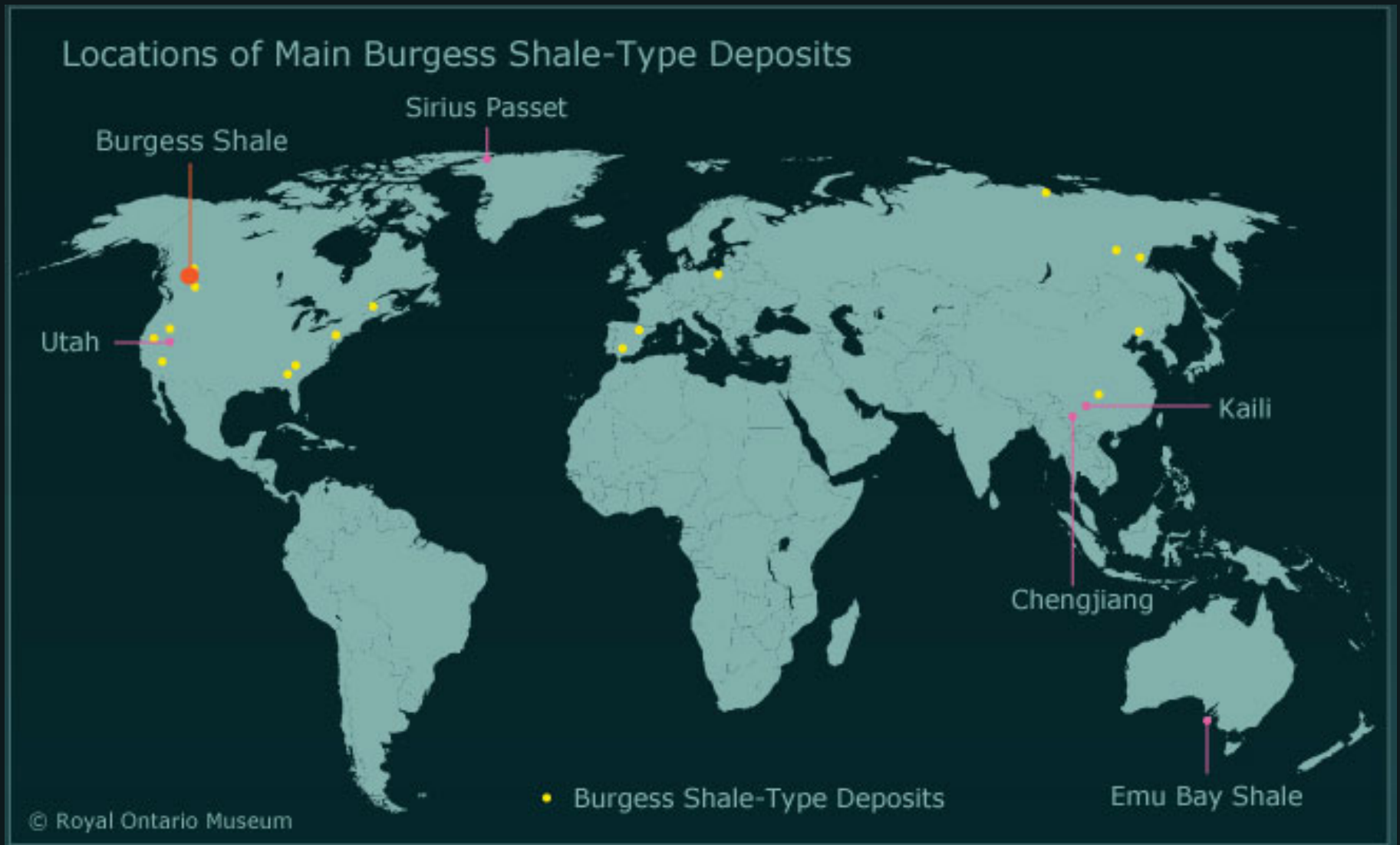
Using the Cambrian Explosion
as a Metaphor for Accelerated
Scientific Discovery

Gully A. Burns

Intelligent Systems Division, Information Sciences
Institute



Paleontological Treasure Troves

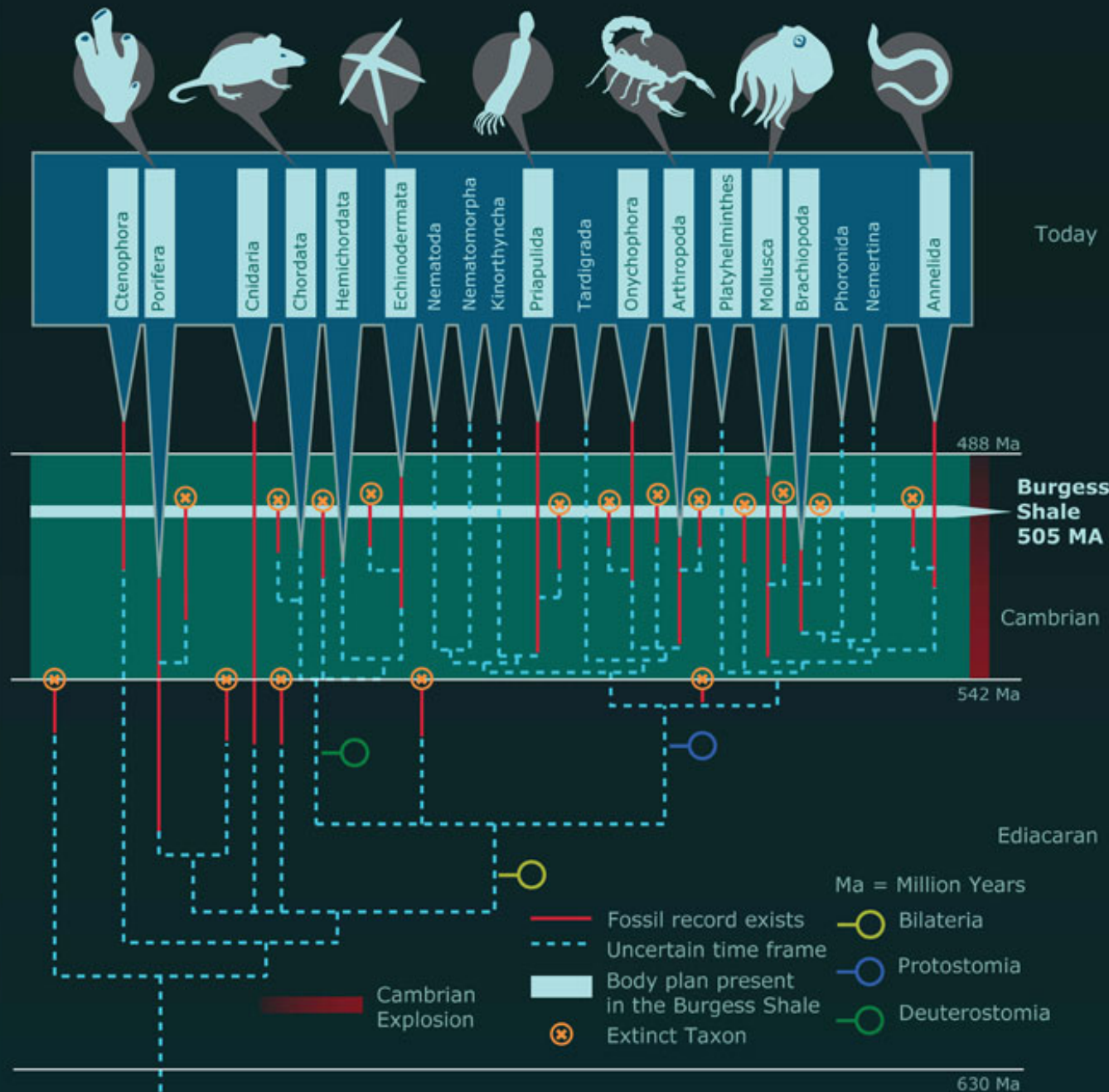


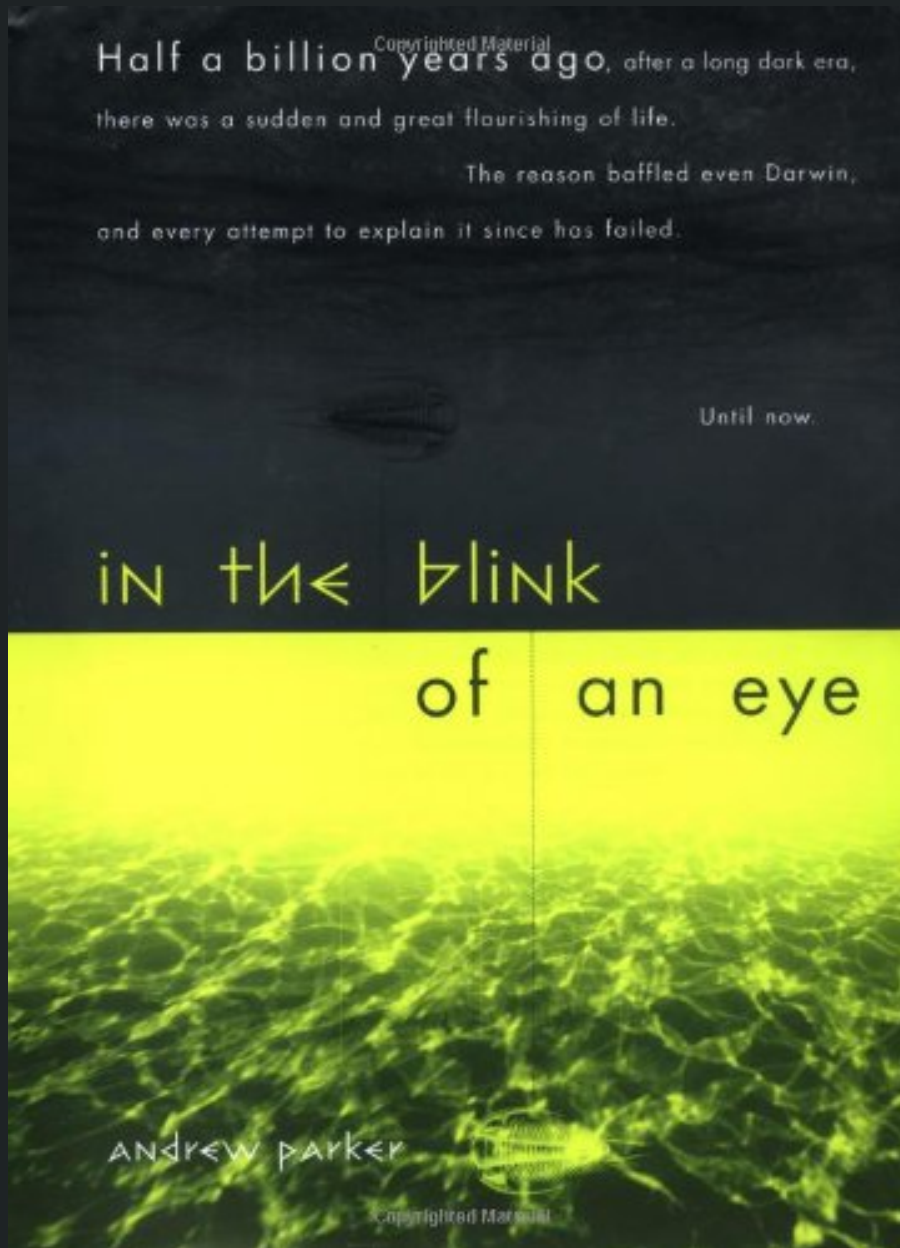
(from <http://burgess-shale.rom.on.ca/>)

The Cambrian Explosion

Every metazoan phylum came into existence during this period

Evolutionary Tree (from <http://burgess-shale.rom.on.ca/>)





The 'Light Switch' Hypothesis

Andrew Parker hypothesizes that the cause of the dramatic evolutionary radiation was the evolution of the early eye and visual system

This made predatory behavior much more successful, leading to defensive specializations and transformation from precambrian forms (without skeletons) to armored and more highly evolved Cambrian animals (skeletons, spikes, teeth, swimming capabilities and armor).

(not everyone agrees: see Morris (2003) 'On the First Day, God Said . . .' for a scathing book review)

Using the 'Light Switch' Metaphor for Scientific Discovery

Consistent with this hypothesis, we might attempt to look for:

- **'Light'** – a ubiquitous data signal with rich substructure
- **'Evolution'**– Accelerated scientific discovery
- **'Eyes'** – Experimental methods + data gathering
- **'Visual System'** - Data analysis and statistics
- **'Intelligence'** – Scientific Theory

‘Cancer in the Age of Algorithms’ – a success story

Dr. Shirley Pepke had late-stage ovarian cancer and used advanced unsupervised learning to help her fight her cancer (with a colleague at ISI: Dr. Greg ver Steeg)

“It was suddenly like looking at the dictionary of tumor biology. It was suddenly pulling out all this information other algorithms couldn’t. It looked really beautiful and informative.”

Dr. Pepke switched treatments and her cancer is in remission

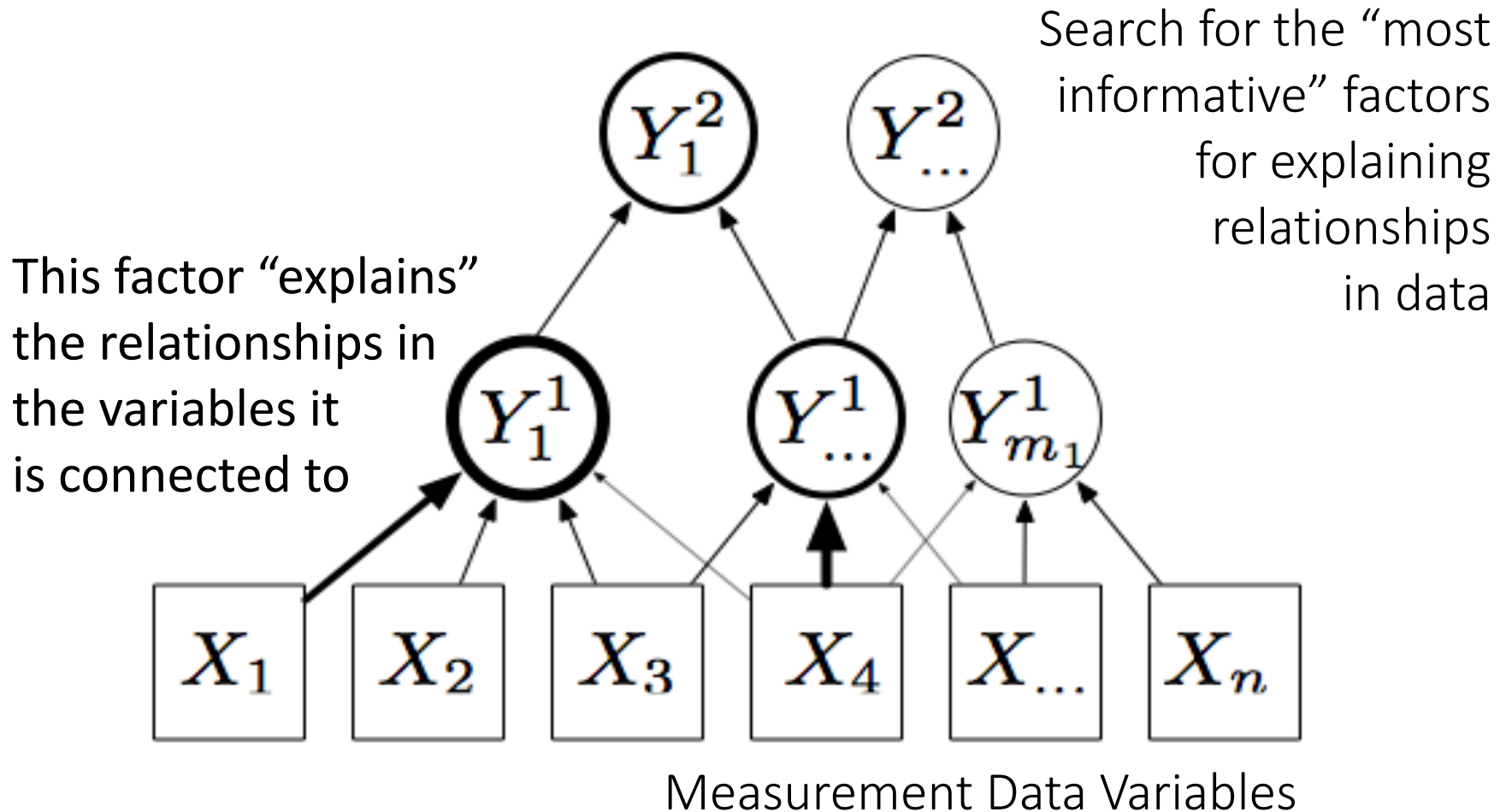


Dr. Shirley Pepke and Dr. Greg Ver Steeg

Washington Post: <http://gul.ly/9u24h>

Soundcloud Link: <http://gul.ly/9u23o>

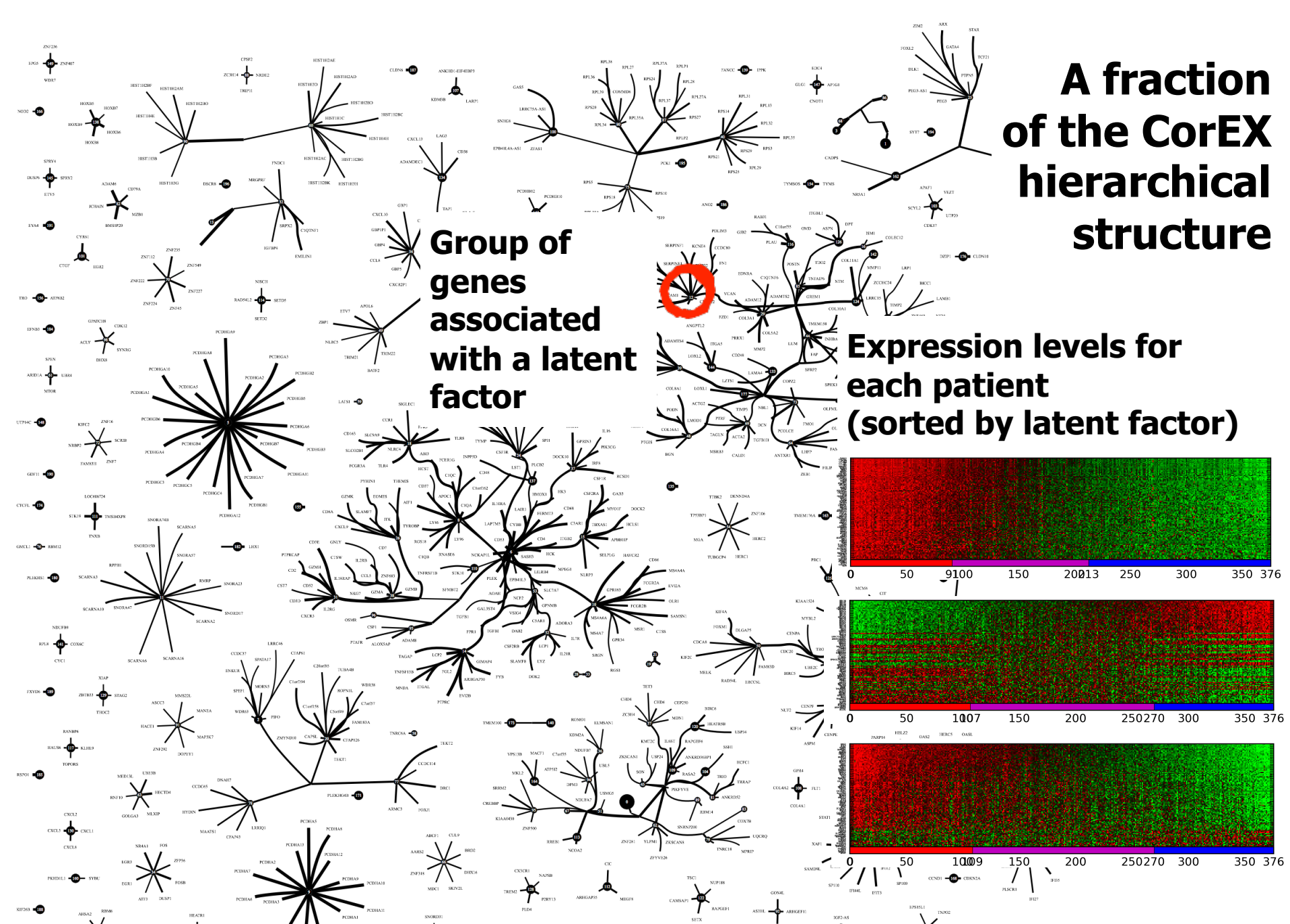
Visual System: Greg ver Steeg's 'Correlation Explanation' (CorEx)



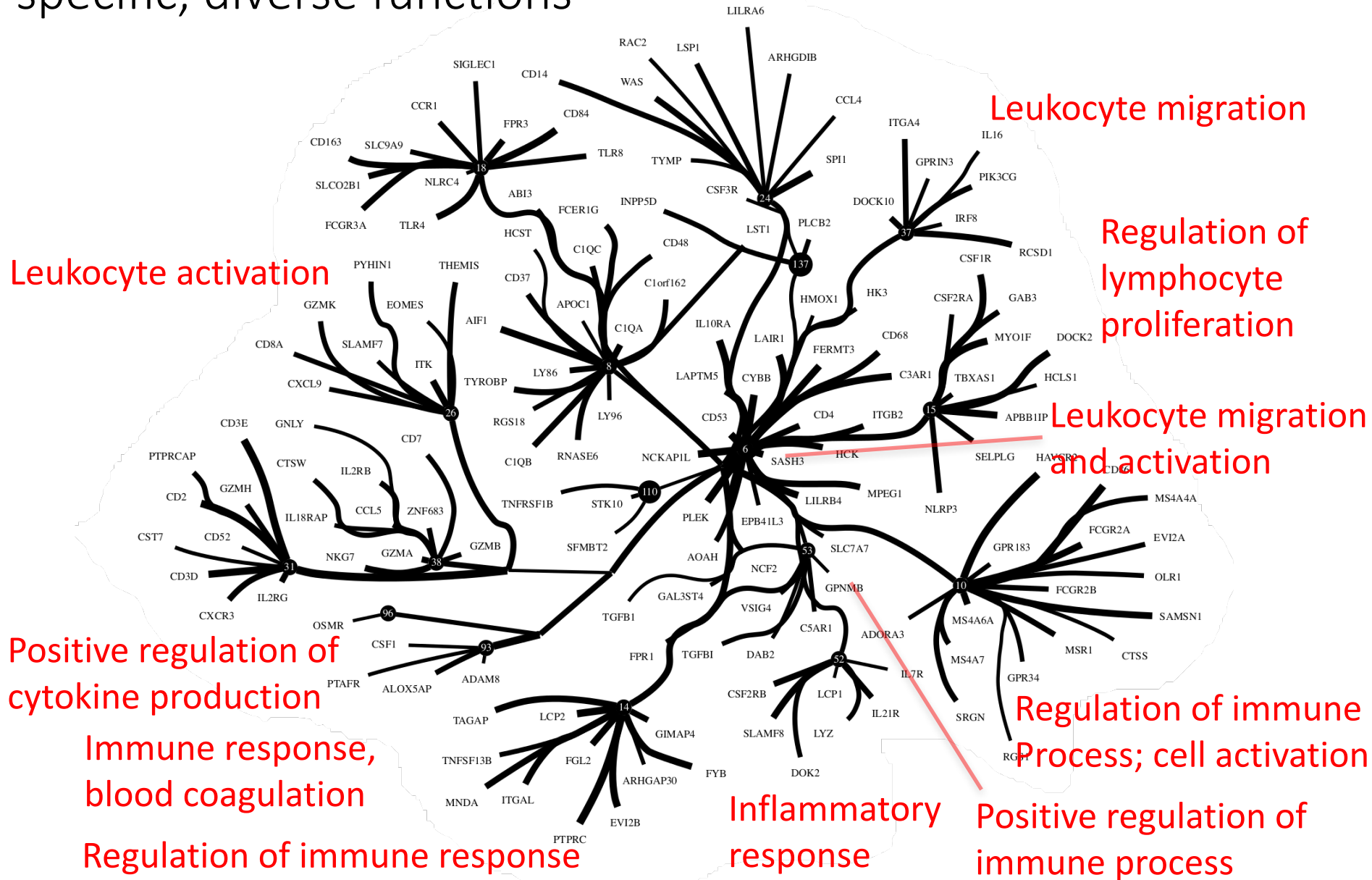
A fraction of the CorEX hierarchical structure

Group of genes associated with a latent factor

Expression levels for each patient (sorted by latent factor)



Gene Ontology Annotations show that groups reflect strong, specific, diverse functions



Where's the Light Switch?

- 'Light' => Gene Expression
- 'Evolution' => Precision Medicine
- 'Eyes' => Gene Chip Cancer Assays
- 'Visual System' => CorEx (Unsupervised Learning)
- 'Intelligence' => Linkage to Gene Ontology (requires human curated models + interpretation)

How does this idea generalize in the broader context of accelerating science?

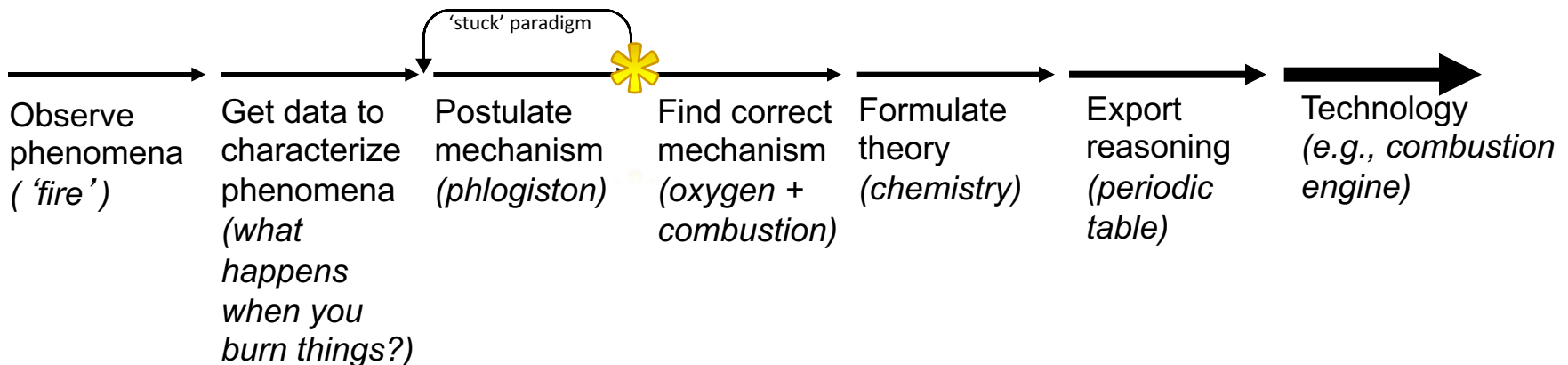
Scientific Discovery: From 'Elements' to the Elements



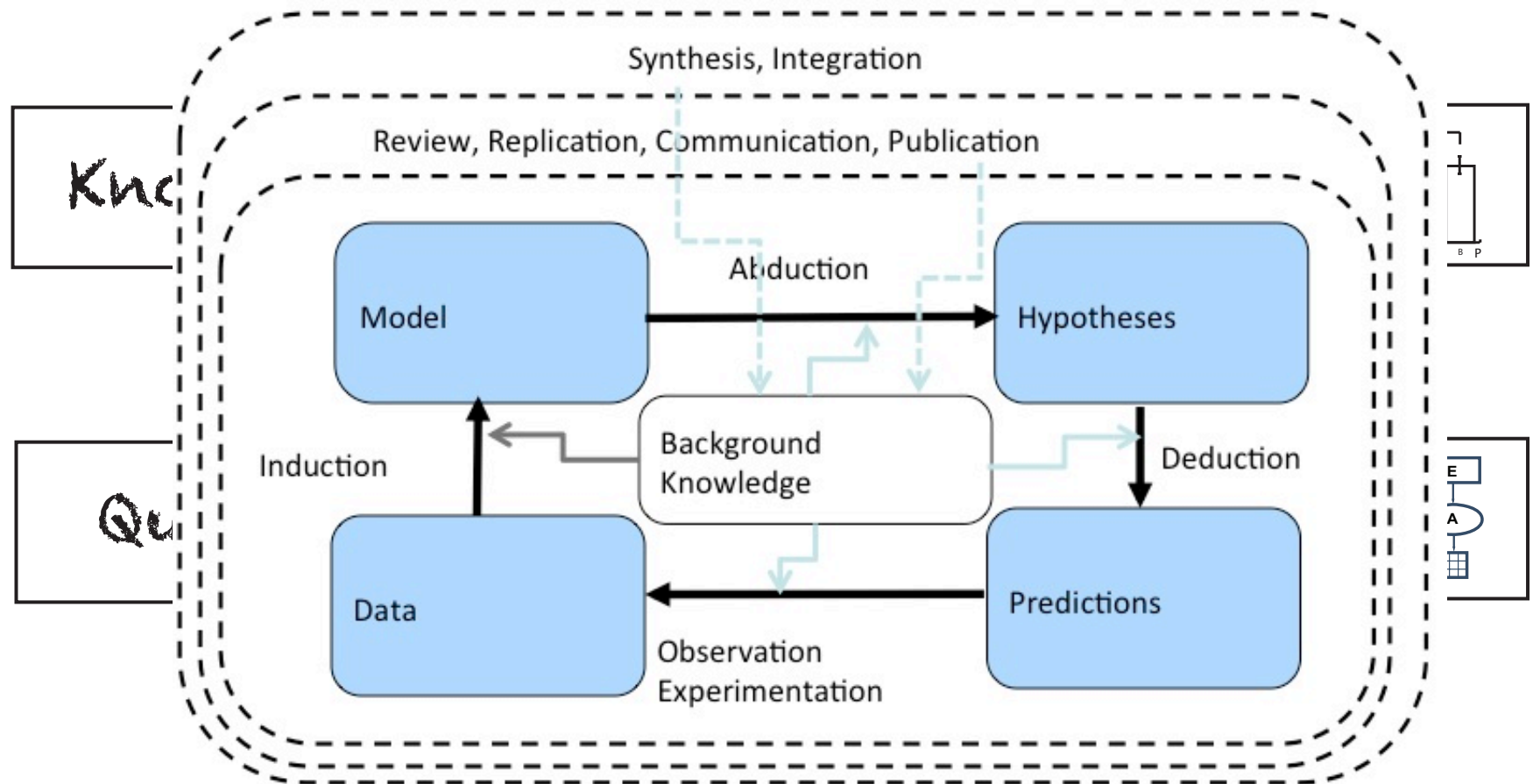
EUREKA!

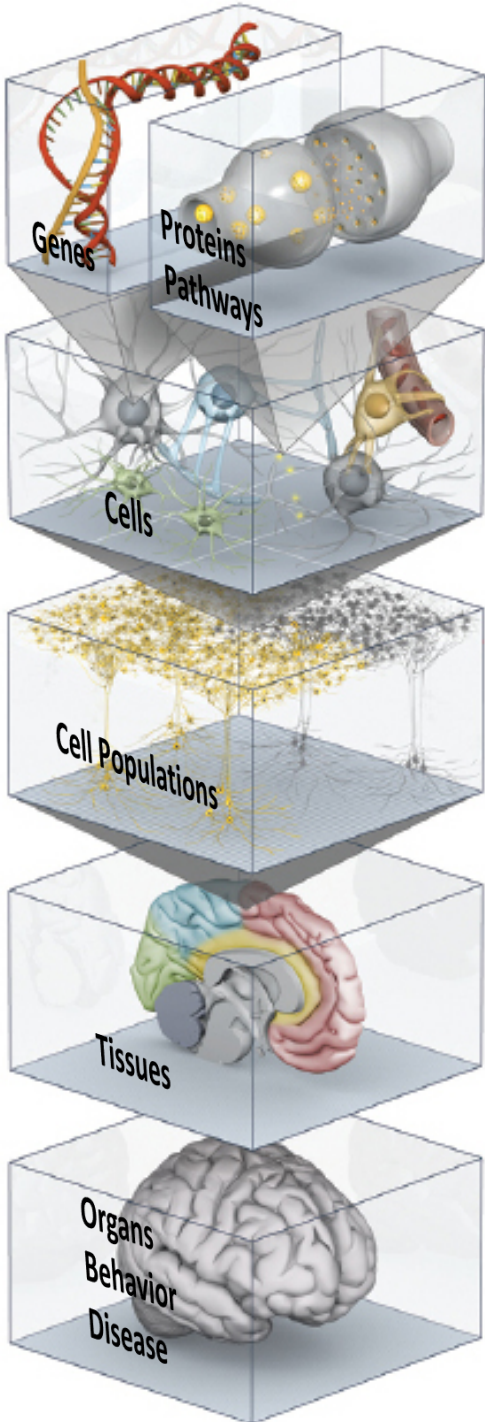


2 8 18 10 2	41 Nb Niobium 92.9063	42 Mo Molybdenum 95.96	43 Tc Technetium [98]	44 Ru Ruthenium 101.07	45 Rh Rhodium 102.9055	46 Pd Palladium 106.42	47 Ag Silver 107.8682	48 Cd Cadmium 112.411	49 In Indium 114.818	50 Sn Tin 118.710
2 8 18 32 10 2	73 Ta Tantalum 180.94788	74 W Tungsten 183.84	75 Re Rhenium 186.207	76 Os Osmium 190.23	77 Ir Iridium 192.217	78 Pt Platinum 195.084	79 Au Gold 196.966569	80 Hg Mercury 200.59	81 Tl Thallium 204.3833	82 Pb Lead 207.2
2 8 18 32 10 2	105	106	107 Bh Bohrium	108 Hs Hassium	109 Mt Meitnerium [276]	110 Ds Darmstadtium [281]	111 Rg Roentgenium [280]	112 Cn Copernicium [285]	113 Nh Nihonium [286]	114 Fl Flerovium [289]



Cycles of Investigation (‘KQED’ Model)





The Complexity of Data Spaces

Neuroscience is an example of a complex multi-level domain

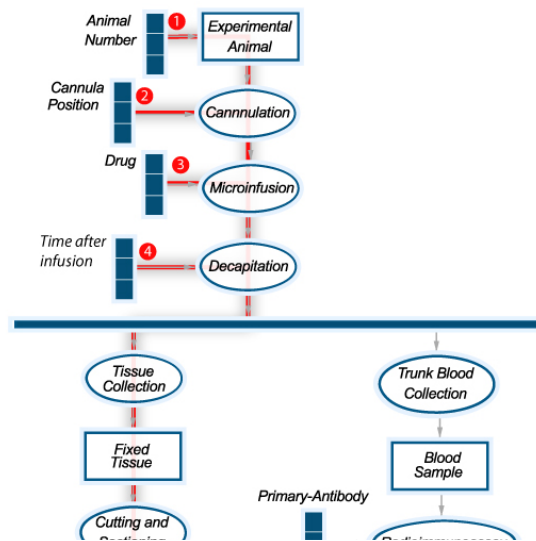
A variety of variables are needed to describe models across scales and systems.

Complex experiments have to be designed and executed by human experts to investigate and test these models.

How to apply the Cambrian Metaphor here?

Knowledge Engineering from Experimental Design

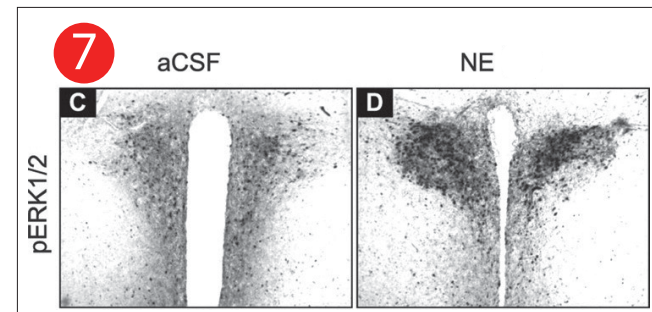
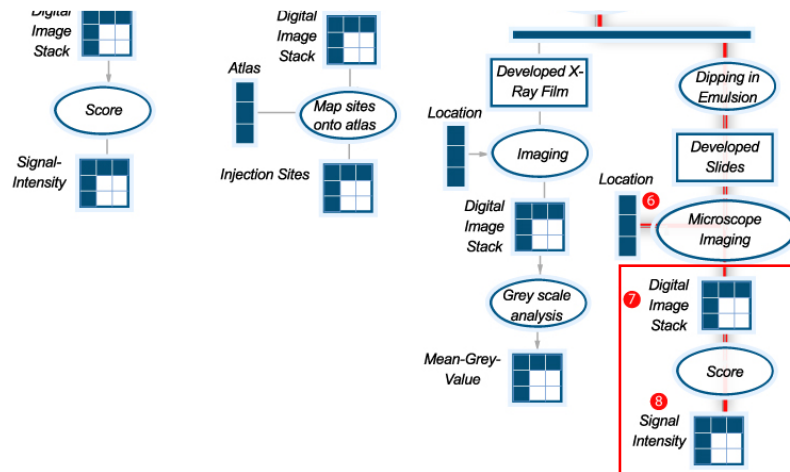
Gene Expression Lab Study
Khan et al. (2007)



Indexing Parameters
(‘Independent Variables’)

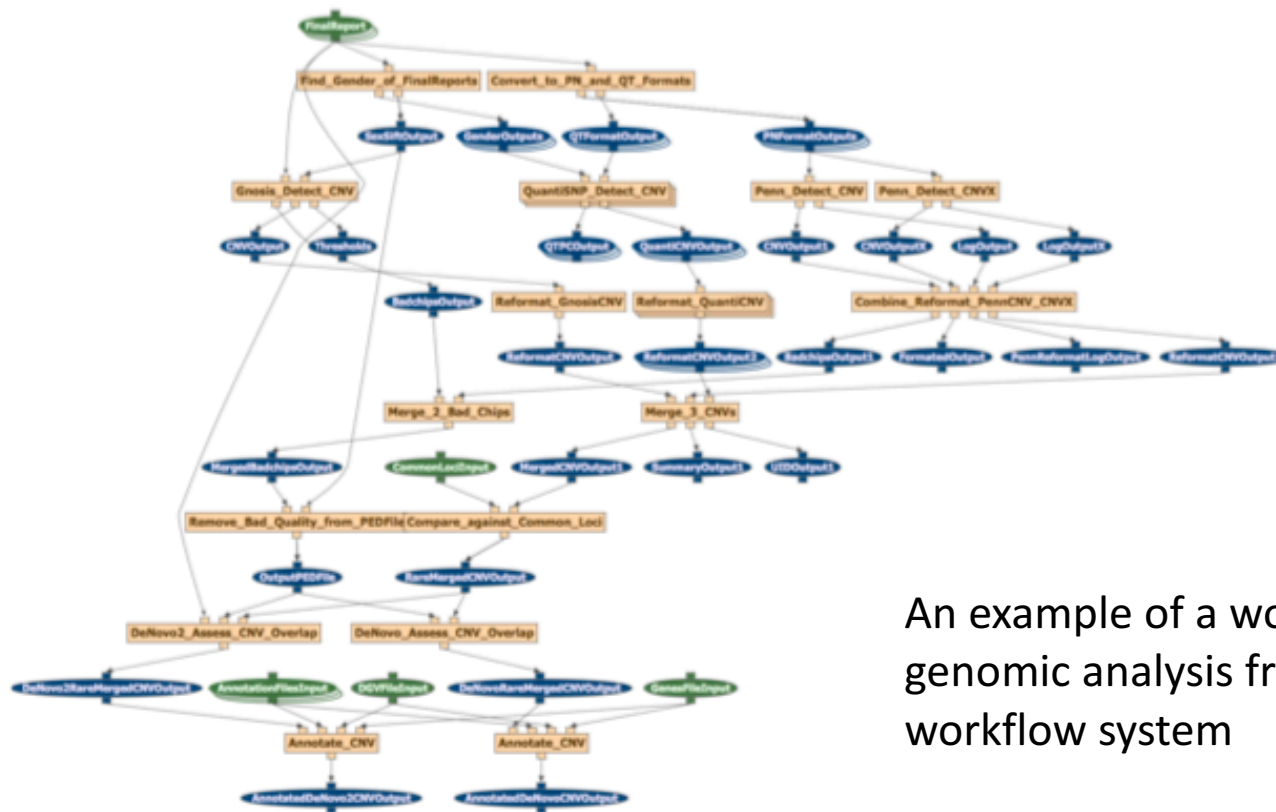
Measurements
(‘Dependent Variables’)

1	2	3	4	5	6	8
subject-id	Cannula Position	Drug	Time after infusion	probe	Location	Signal Intensity
RO2-164	PVH proximity	NE2	10 mins	pERK1/2	PVHmpd	high
RO2-175	PVH proximity	aCSF	10 mins	pERK1/2	PVHmpd	low



Russ et al (2011), BMC Bioinformatics

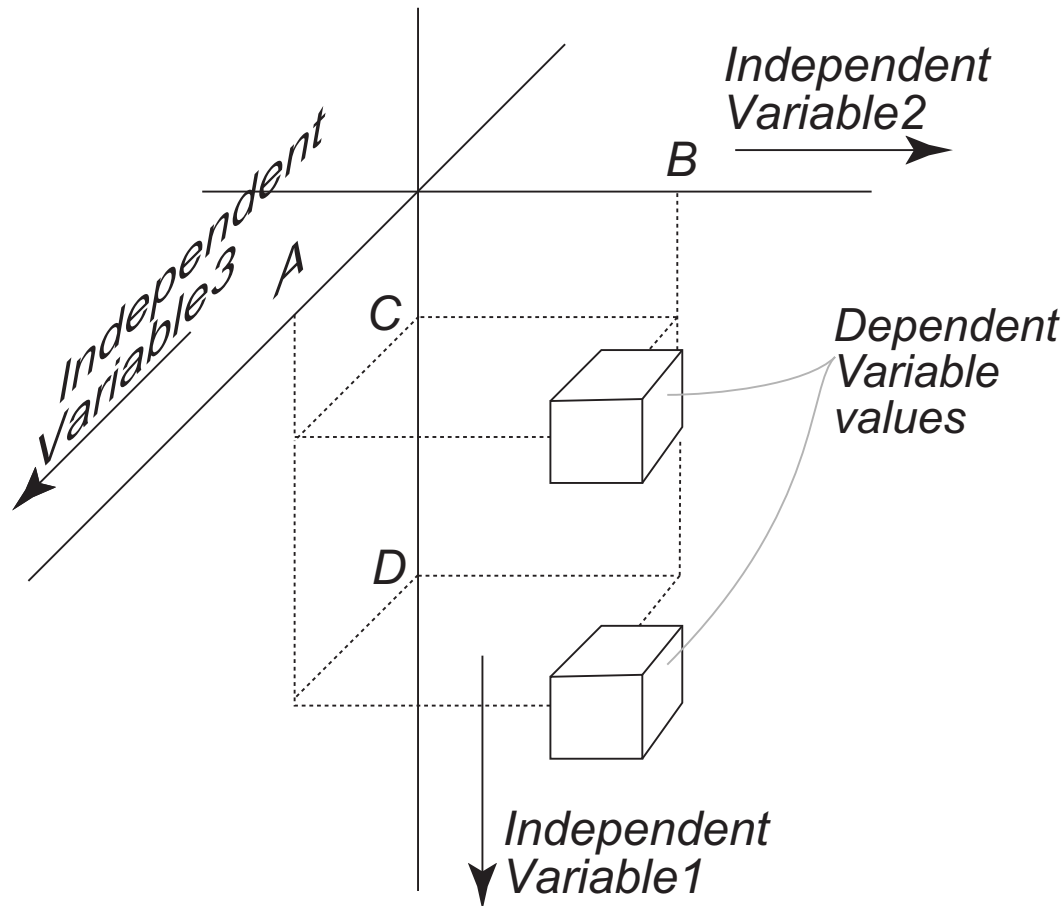
E-Science Workflows



An example of a workflow for genomic analysis from the WINGS workflow system

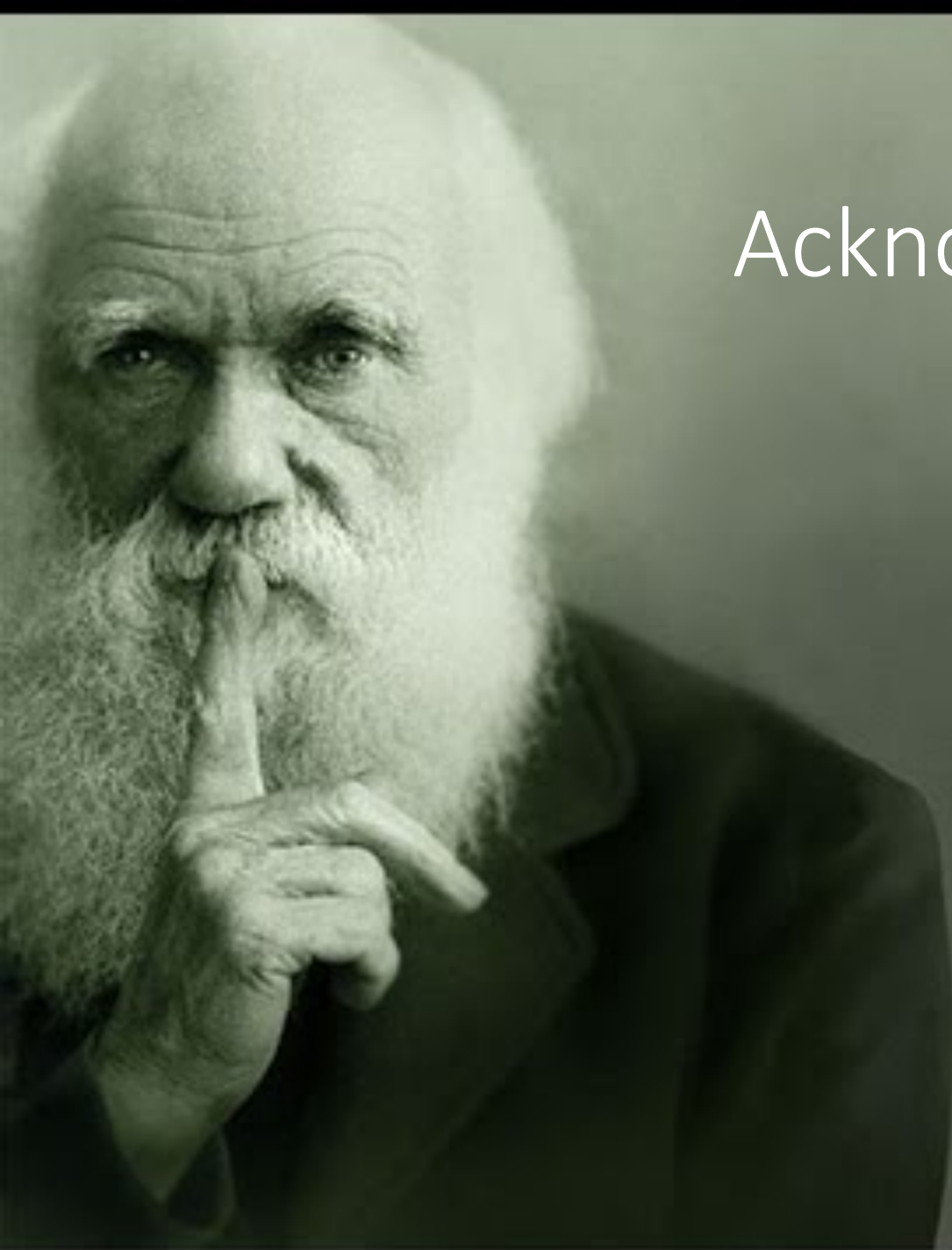
<http://www.wings-workflows.org/>

Measurements and Metadata Define the Data Space



Could we apply tools like CorEx to examine the structure of this space?

Would this be like suddenly flipping the switch and turning on the lights?



Acknowledgements

Greg ver Steeg
Oren Etzioni
Alan Watts
Prem Natarajan
Yolanda Gil
Ed Hovy

NIH, DARPA,
NSF, IARPA