# The Transformative Potential of Al for Science: Will Al Write the Scientific Papers of the Future?

Yolanda Gil
Information Sciences Institute
and Department of Computer Science
University of Southern California

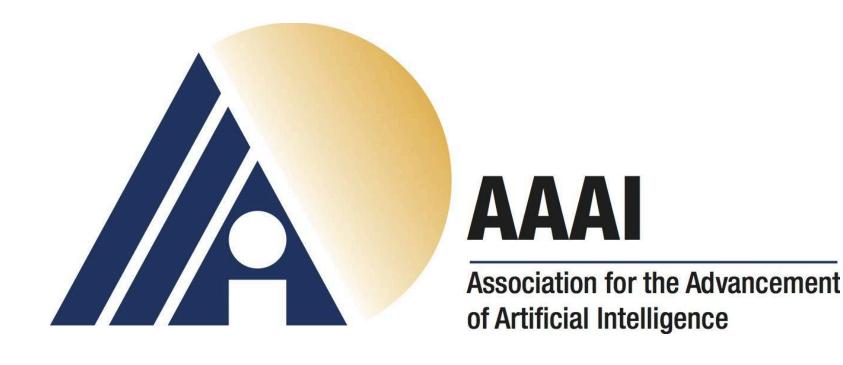
https://tinyurl.com/yolandagil-aaai2020



## A Personal Perspective on Al

## A Personal Perspective on Al

- The AI community has always been
  - Visionary
  - Broad
  - Inclusive
  - Interdisciplinary
  - Determined
- And dare I say
  - Successful



# A Personal View of Watershed Moments in AI: (I) 1980s

#### SOAR: AN ARCHITECTURE FOR GENERAL INTELLIGENCE

John E. Laird, Allen Newell and Paul S. Rosenbloom

University of Michigan Carnegie-Mellon University Stanford University



#### Integrating Planning and Learning: The PRODIGY Architecture

Manuela Veloso

Jaime Carbonell

Alicia Pérez

Daniel Borrajo

Eugene Fink

Jim Blythe

(veloso@cs.cmu.edu)

(carbonell@cs.cmu.edu)

(dborrajo@fi.upm.es)

(eugene@cs.cmu.edu)

(jblythe@cs.cmu.edu)

#### Theo: A Framework for Self-Improving Systems

Tom M. Mitchell, John Allen, Prasad Chalasani, John Cheng, Oren Etzioni, Marc Ringuette, Jeffrey C. Schlimmer

#### Learning representations by backpropagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

## Scripts Plans Goals and Understanding

Roger C. Schank and Robert P. Abelson

#### Elephants Don't Play Chess

Rodney A. Brooks

MIT Artificial Intelligence Laboratory, Cambridge, MA 02139, USA

Robotics and Autonomous Systems 6 (1990) 3-15

# A Personal View of Watershed Moments in Al: 1990s

1990: Sphynx shows speaker independent large vocabulary continuous speech recognition – used it to write my PhD thesis

1992: Soar flies helicopter teams in simulations and is indistinguishable to commanders from human-controlled

aircraft

1992: TD-Gammon autonomously learns to play backgammon at human player levels

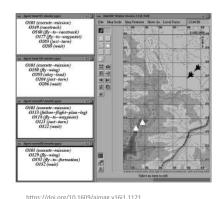
1995: SKICAT identifies five new quasars in the Second Palomar Sky Survey

1995: Navlab is the first trained car to drive autonomously on highways to cross the United

States

1997: Deep Blue defeats human chess world champion and gets grandmaster-level rating

1999: RAX flew a spacecraft autonomously, demonstrating planning, monitoring, and fault repair







https://flic.kr/p/9Dpww

# A Personal View of Watershed Moments in Al: 2000s

2000: The Gene Ontology is shown to describe over 15,000 gene products for drosophila, mouse, and yeast

**2000: Kismet demonstrates and recognizes emotions** 

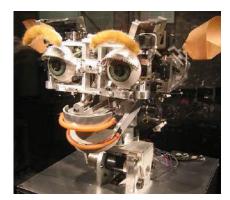
2003: USC ISI's statistical machine translation prototype beats hand-crafted commercial systems

2004: RDF semantic specifications become W3C recommendations for the GGG (Giant Global Graph)

2007: Stanley wins \$1M for first autonomous high-speed off-road driving

2007: First robot soccer team against human players in exhibition game at RoboCup

2009: Pragmatic Chaos ensemble learning wins \$1M competition to predict user film ratings



R D F



https://commons.wikimedia.org/w/index.php?curid=374949

https://commons.wikimedia.org/w/index.php?curid=2079835

# A Personal View of Watershed Moments in Al: 2010s

2010: Siri voice-activated personal assistant is released as a smart phone app

2011: CMDragons team of 10 soccer robots coordinate plays for routine passing, interception, and goal

scoring

2011: Watson takes first place in Jeopardy Q/A game defeating two human champions

2012: AlexNet scores improved 10 percentage points in the ImageNet visual recognition challenge

2013: Cognitive Tutor shows 8 percentile points average improvement in algebra in 25,000 students study

2016: Knowledge Graph used as semantic backbone in one third of 100B monthly searches

2019: Wikidata records 8 billion triples and over 880 billion edits, surpassing Wikipedia as most edited Wikimedia

site

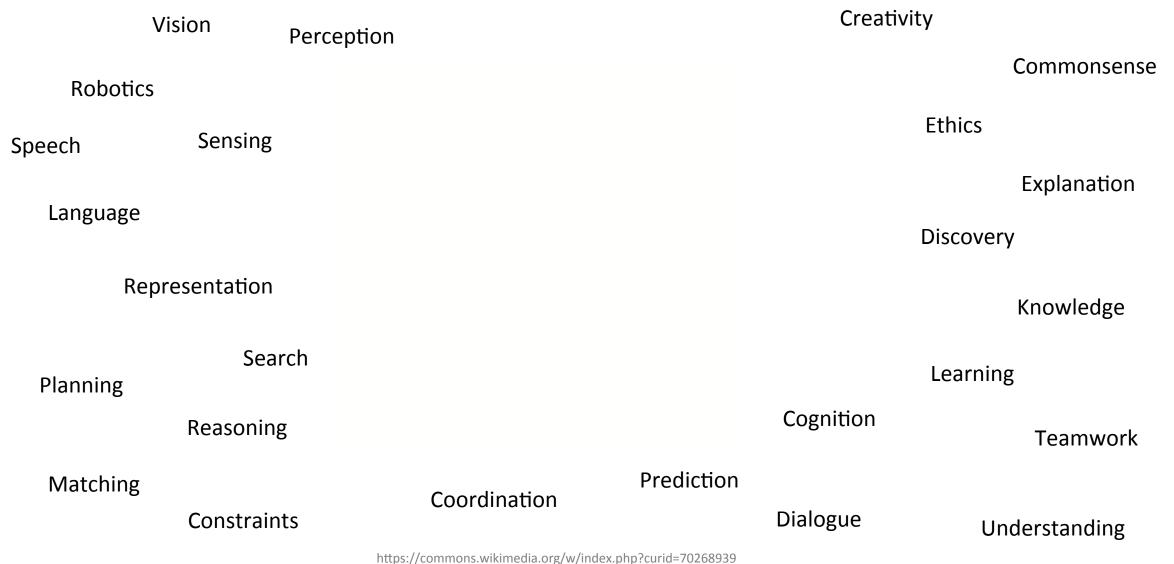






https://www.youtube.com/watch?v=4QtBSDSC2pk

## Diversity and Breadth of Advances in Al

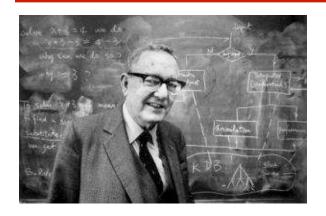


## Al to Address Major Future Challenges



# The Imperative for AI in Science

## Al and Scientific Discovery

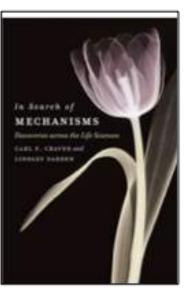


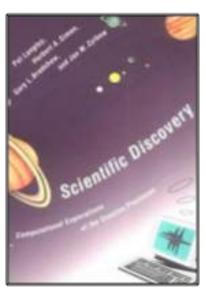


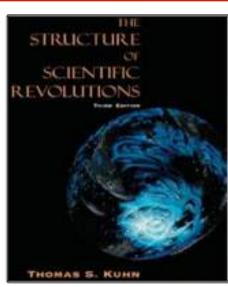


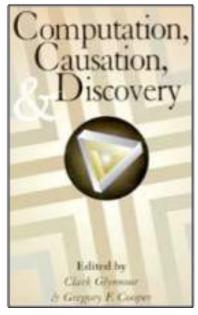


- [Lenat 1976]
- [Lindsay et al 1980]
- [Langley 1981]
- [Falkenhainer 1985]
- [Kulkarni and Simon 1988]
- [Cheeseman et al 1989]
- [Zytkow et al 1990]
- [Simon 1996]
- [Valdes-Perez 1997]
- [Todorovski et al 2000]
- [Schmidt and Lipson 2009]









## Human Limitations Curb Scientific Progress [Gil DSJ'17]

- Not systematic
  - e.g., [Peters et al PLOS 2014]

- Errors
  - e.g., [Herndon et al CJE 2013]

- Biases
  - e.g., [Anderson et al ACS 2014]

- Poor reporting
  - e.g., [Garijo et al PLOS 2013]

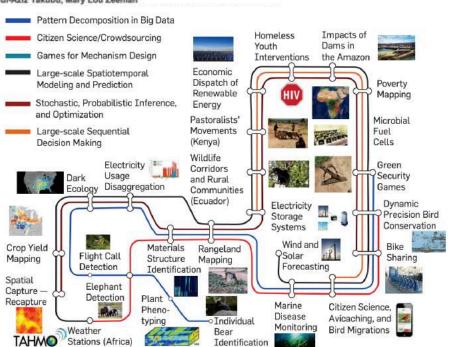


#### COMMUNICATIONS

https://doi.org/10.1145/3339399

#### Computational Sustainability: Computing for a Better World and a Sustainable Future

By Carla Gomes, Thomas Dietterich, Christopher Barrett, Jon Conrad, Bistra Dilkina, Stefano Ermon, Fei Fang, Andrew Farnsworth, Alan Fern, Xiaoli Fern, Daniel Fink, Douglas Fisher, Alexander Fiecker, Daniel Freund, Angela Fuller, John Gregoire, John Hopcroft, Steve Kelling, Zico Kolter, Warren Powell, Nicole Sintov, John Selker, Bart Selman, Daniel Sheldon, David Shmoys, Millind Tambe, Weng-Keen Wong, Christopher Wood, Xiaojian Wu, Yexiang Xue, Amulya Yadav, Abdul-Aziz Yakubu, Mary Lou Zeeman

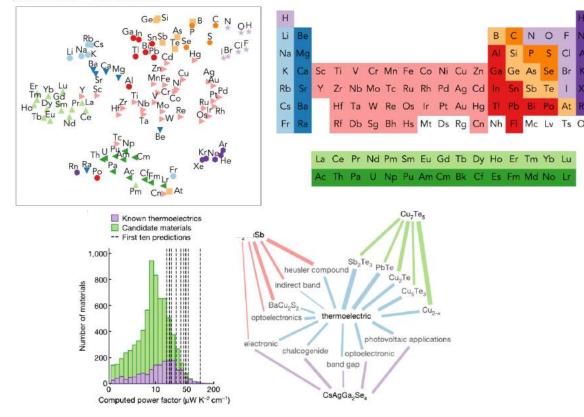


#### nature

https://doi.org/10.1038/s41586-019-1335-8

## Unsupervised word embeddings capture latent knowledge from materials science literature

 $\label{eq:Vahe Tshitoyan$^{1,3*}$, John Dagdelen$^{1,2}$, Leigh Weston$^{1}$, Alexander Dunn$^{1,2}$, Ziqin Rong$^{1}$, Olga Kononova$^{2}$, Kristin A. Persson$^{1,2}$, Gerbrand Ceder$^{1,2*}$, Anubhav Jain$^{1*}$.}$ 



#### Knowledge-based Biomedical Data Science 2019

Tiffany J. Callahan<sub>1</sub>, Harrison Pielke-Lombardo<sub>1</sub>, Ignacio J. Tripodi<sub>1,2</sub>, Lawrence E. Hunter<sub>1</sub>\*

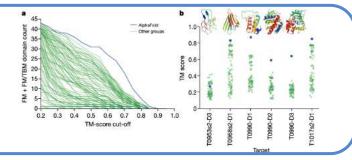
https://arxiv.org/pdf/1910.06710

https://doi.org/10.1038/s41586-019-1923-7

Improved protein structure prediction using

#### Improved protein structure prediction using potentials from deep learning

https://doi.org/10.1038/s41586-019-1923-7 Received: 2 April 2019 Andrew W. Senior<sup>14\*</sup>, Richard Evans<sup>14</sup>, John Jumper<sup>14</sup>, James Kirkpatrick<sup>14</sup>, Laurent Sifre<sup>14</sup>, Tim Green<sup>1</sup>, Chongli Giri, Augustin Zidek<sup>1</sup>, Alexander W. R. Nelson<sup>1</sup>, Alex Bridgland<sup>1</sup>, Hugo Penedones<sup>2</sup>, Stig Petersen<sup>1</sup>, Karen Simonyan<sup>1</sup>, Steve Crossan<sup>1</sup>, Pushmeet Kohli<sup>1</sup>, David T. Jones<sup>2,1</sup>, David Silver<sup>1</sup>, Koray Kavukcuoglu <sup>1</sup>8 Demis Hassabis<sup>1</sup>

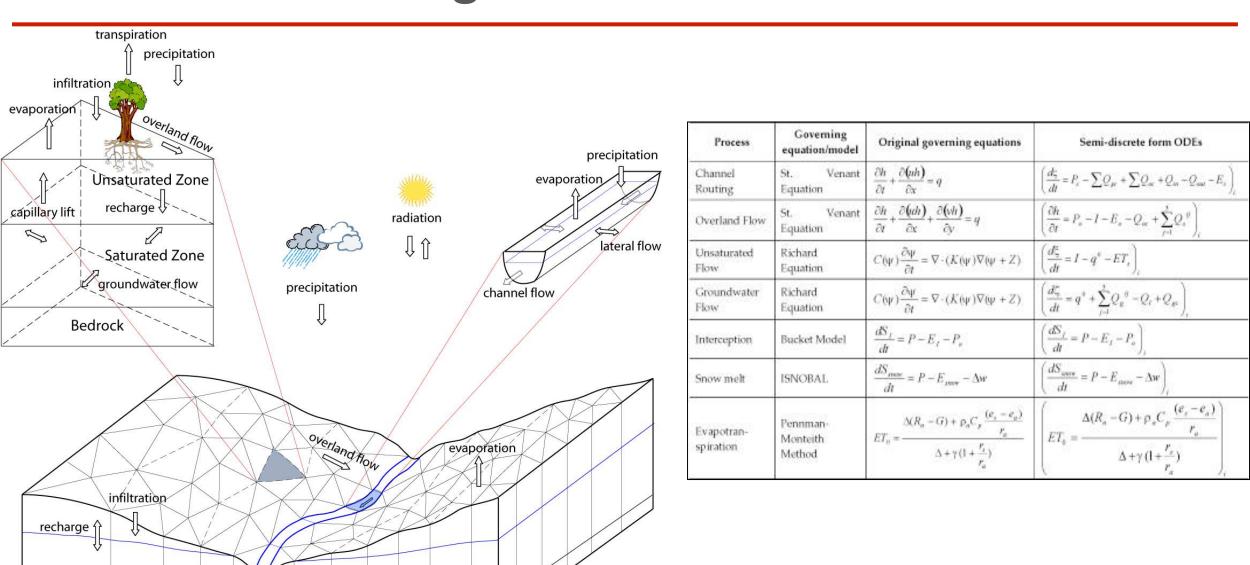




## Capturing Scientific Knowledge

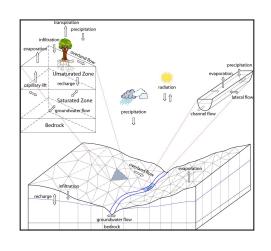
## Scientific Knowledge

groundwater flow bedrock



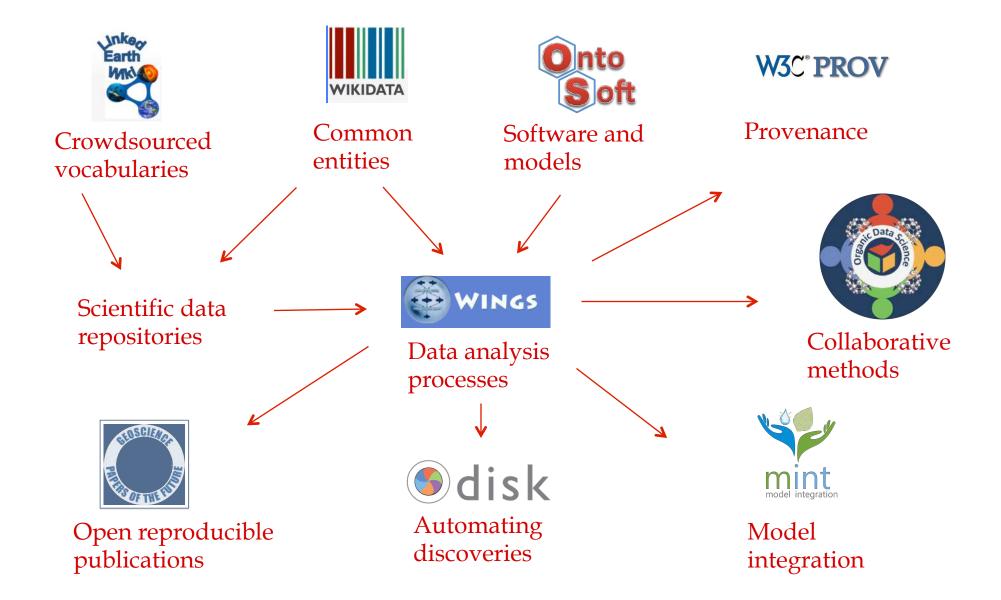
http://www.pihm.psu.edu/pihm home.html

## Supporting Compositionality of Scientific Knowledge

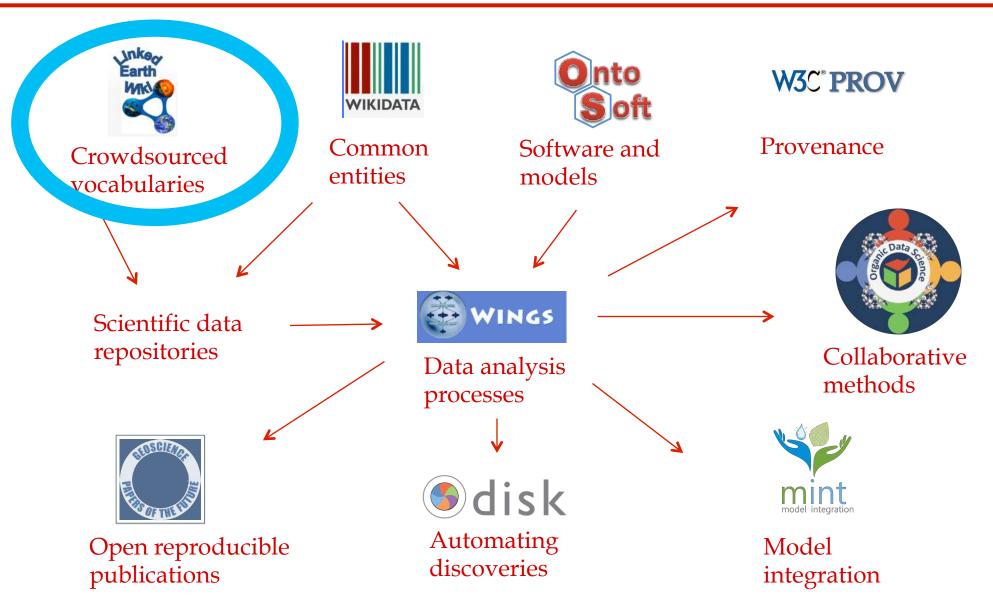


- Data formats
- Physical variables
- Constraints for use
- Adjustable parameters
- > Interventions

## Capturing Scientific Knowledge

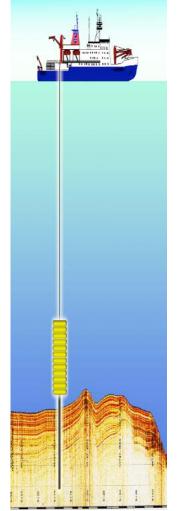


## Capturing Scientific Knowledge



## Low-Cost Creation of Scientific Vocabulary Standards

[Gil et al ISWC 2017; Khider et al PP 2019; Emile-Geay et al PAGES 2018]





https://commons.wikimedia.org/wiki/File:An\_ice\_core\_segment.jpg



https://commons.wikimedia.org/wiki/File:Gravity-corer\_hg.png

Work with Deborah Khider, Julien Emile-Geay, Daniel Garijo (USC); Nick McKay (NAU)

**Problem:** Diversity of requirements for metadata

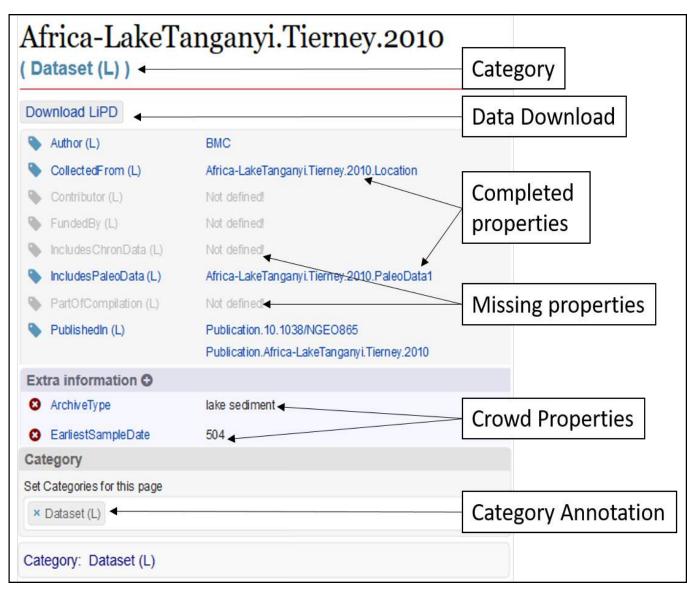
**Approach**: Semantic technologies used for controlled crowdsourcing facilitate creation of community standards to describe highly heterogeneous scientific data

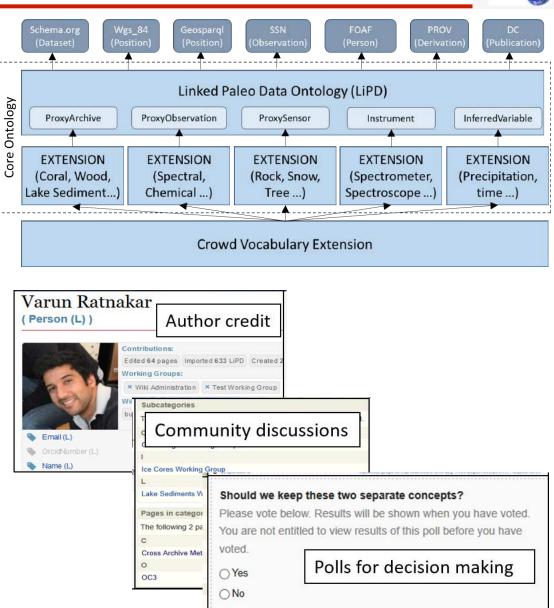
- Organic growth: As scientists annotate their datasets, they propose new metadata properties
- <u>Crowdsourcing</u>: Scientists proposed properties for reuse, vote on priorities
- Editorial oversight: Editors decide what properties will be in future versions

**Results**: A new standard for paleoclimate (PaCTS 1.0) with one (!!) single initial face-to-face meeting

## Controlled Crowdsourcing: Continuous Ontology Growth





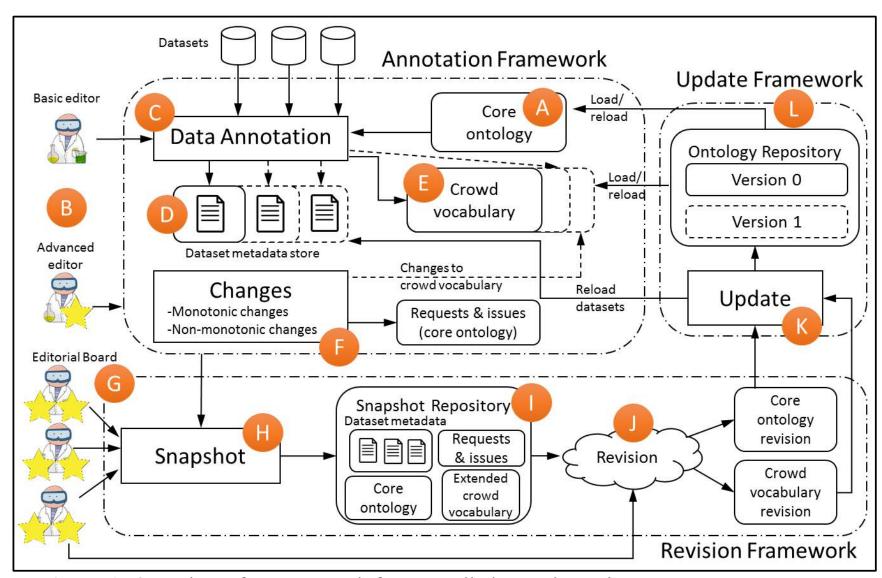


## Living with a Live Ontology



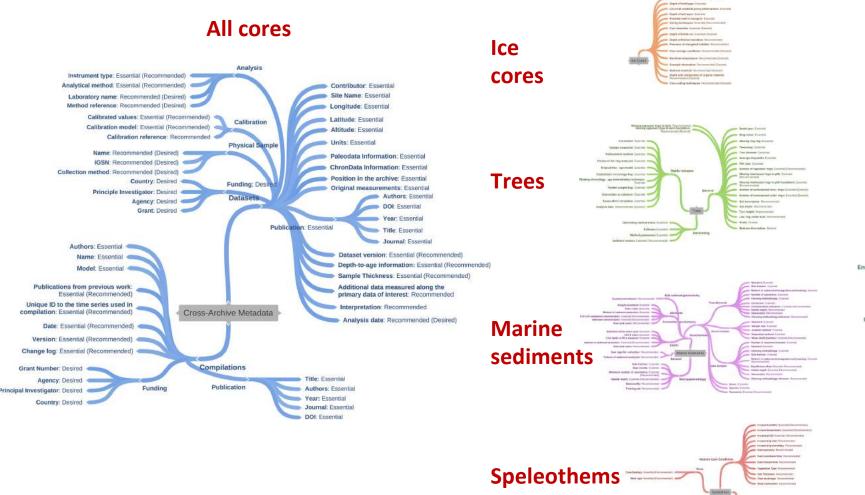
<u>Challenge</u>: Continuous revisions of ontologies + annotated data

- Maintain core + crowd ontology
- Editors decide when to update
  - Automated upgrades when monotonic changes
  - Semi-automated upgrades when non-monotonic changes

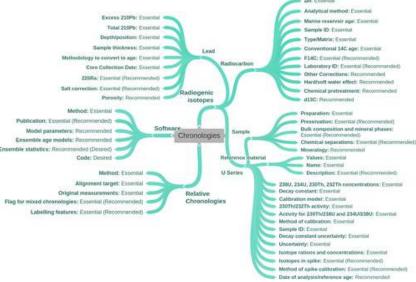


## The Paleoclimate Community reporTing Standard (PaCTS)

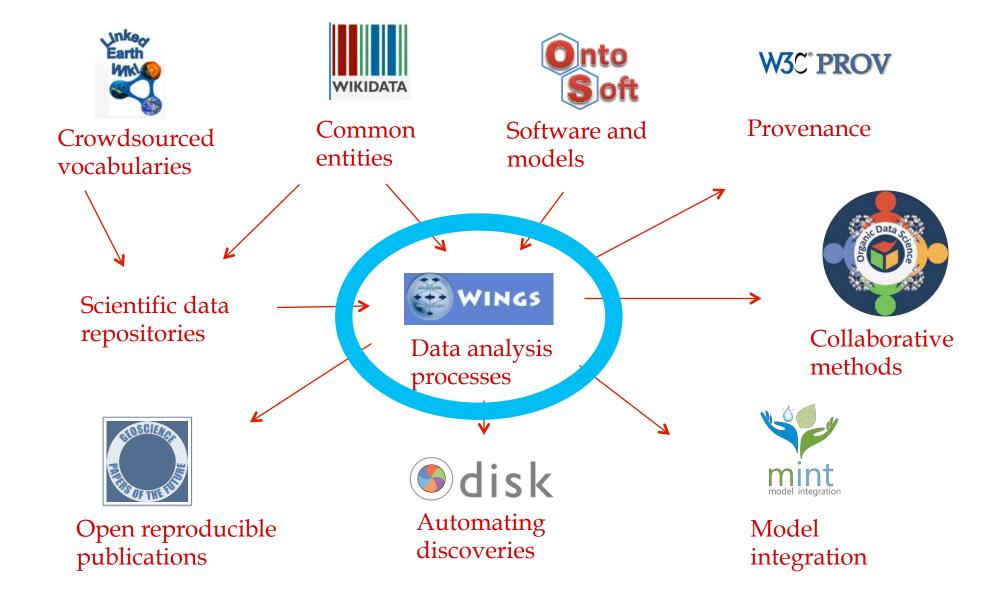




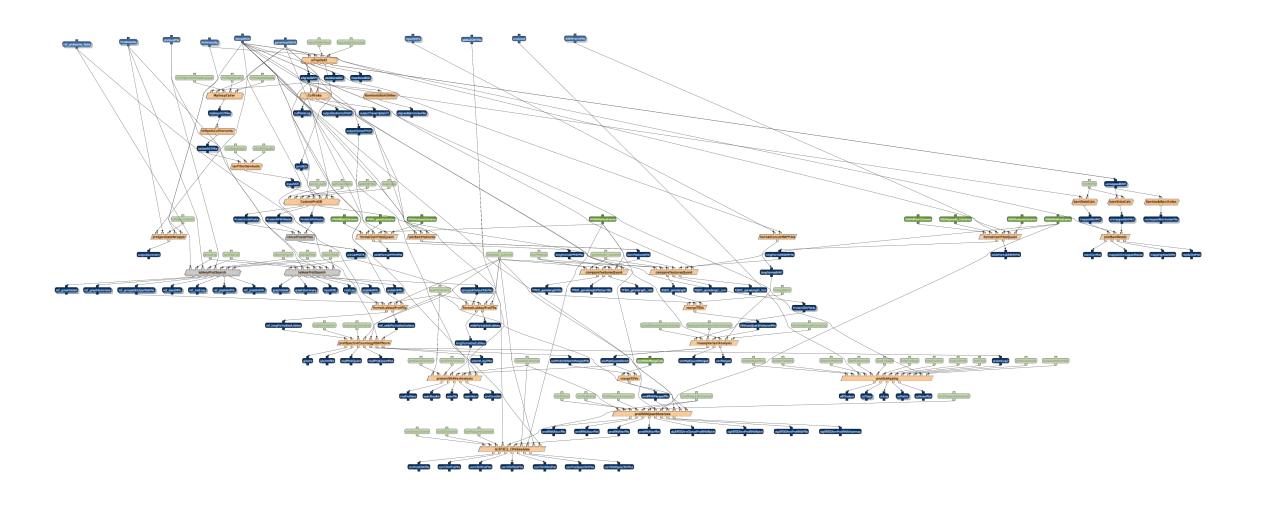
#### **Chronologies**



## Capturing Scientific Knowledge



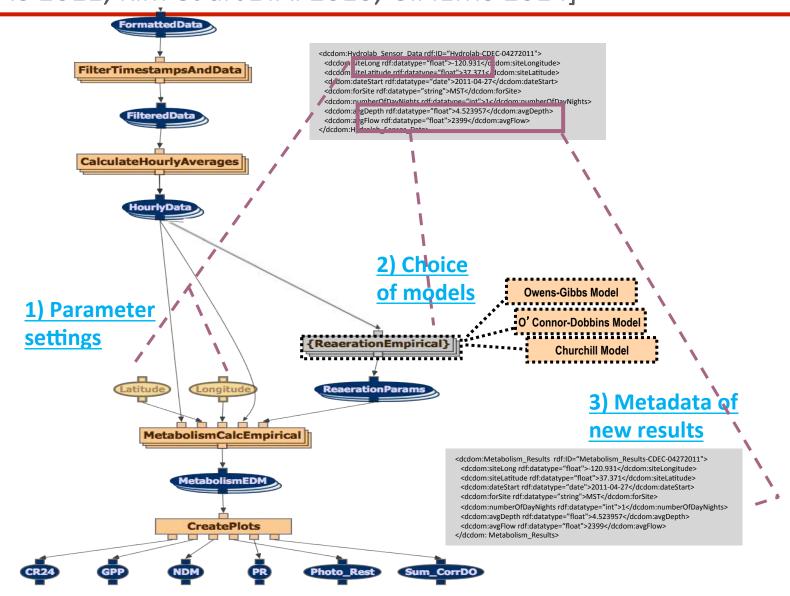
## Workflows



#### Semantic Workflows

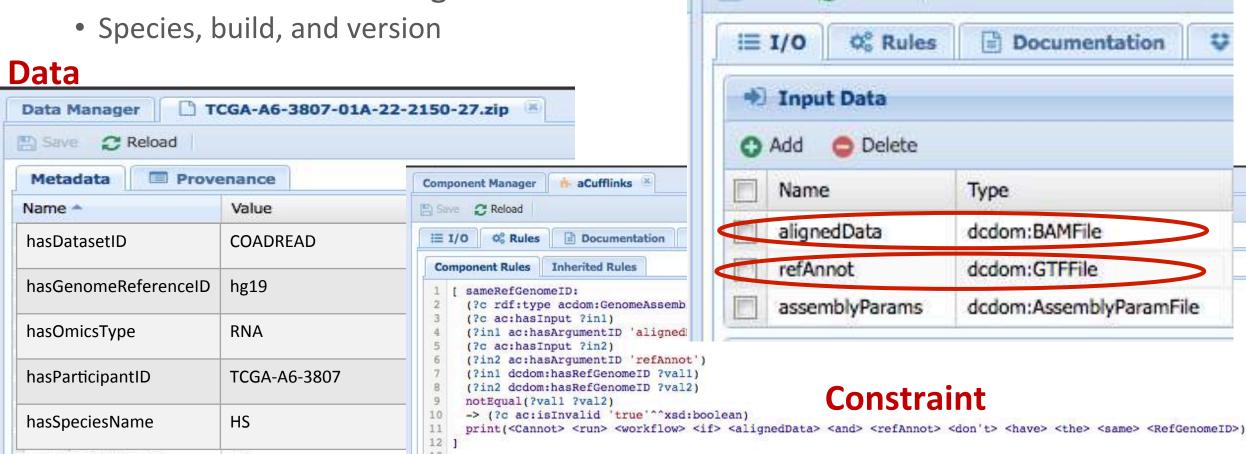
[Gil et al JETAI 2011; Gil et al IEEE IS 2011; Kim et al JETAI 2010; Gil IEMS 2014]

- Workflow constituents (step, input data, results, parameters) are assigned an identifier that can be referenced in constraints
- Input datasets have metadata that can be referenced in the constraints
- Constraints are used to customize a workflow to a given dataset:
  - Set parameters
  - Generate new metadata
  - Choose steps
  - Validation



## Example: Semantic Constraints on Workflows

 The alignment step (TopHat) and the assembly step (Cufflinks) must be done with the same reference genome



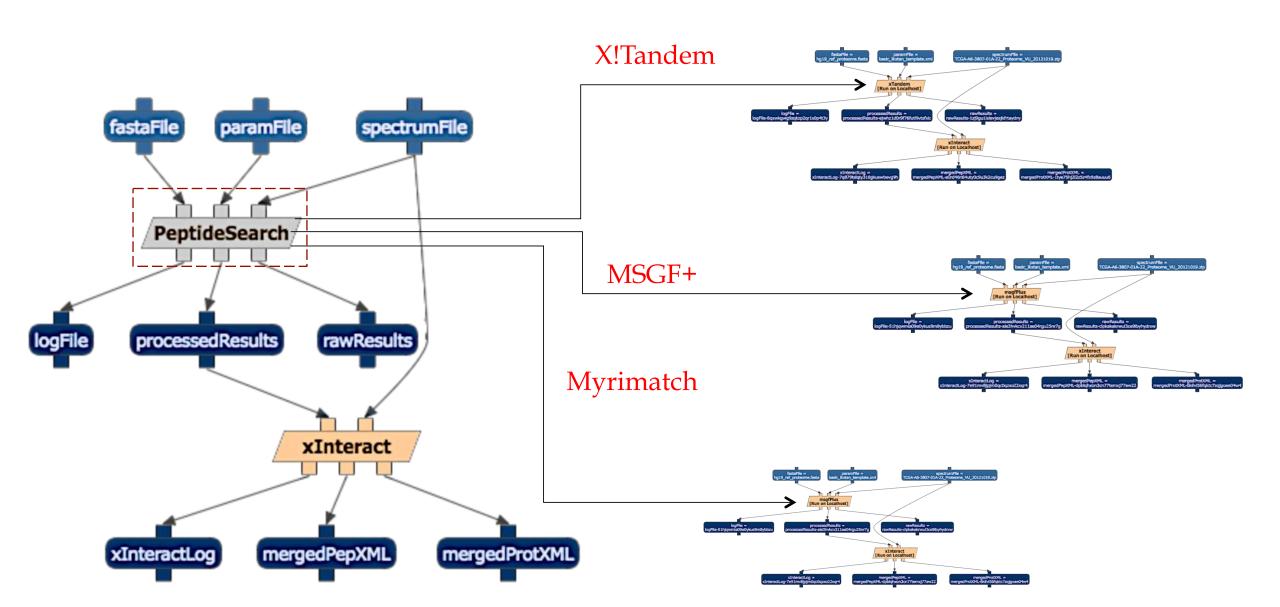
Component Manager

Reload

Step

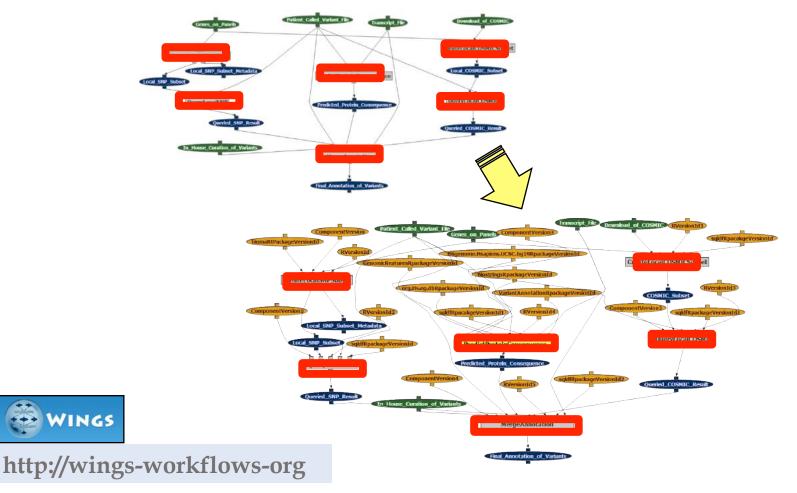
aCufflinks |

#### Abstractions in Semantic Workflows



## WINGS: Workflow Composition

Workflow reasoning algorithms use constraint-based planning to generate executable workflows from high-level workflows

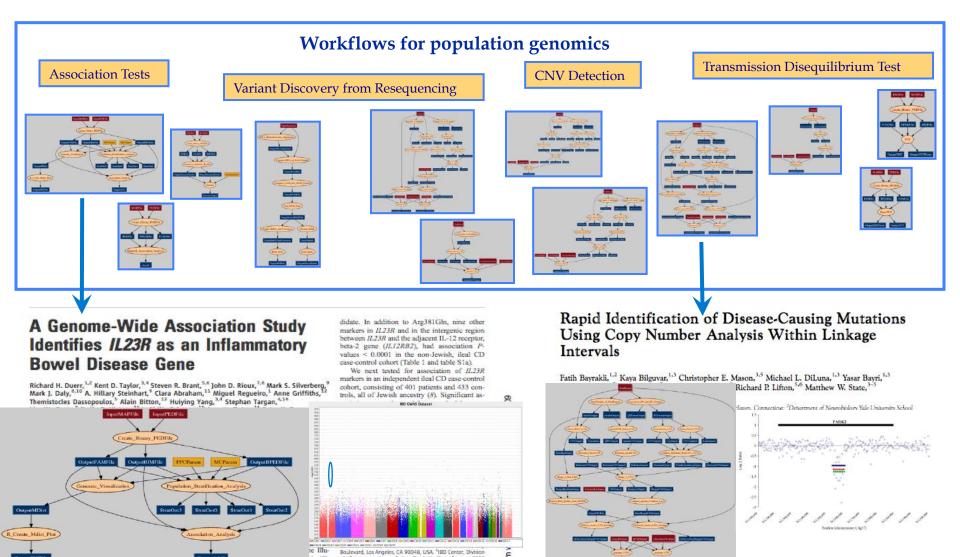


unified well-formed reg. Seed workflow from request seeded workflows Find input data requirements binding-ready workflows Data source selection bound workflows Parameter selection configured workflows Workflow instantiation workflow instances Workflow grounding ground workflows Workflow ranking top-k workflows Workflow mapping executable workflows

## Reproducing Work in Population Genomics

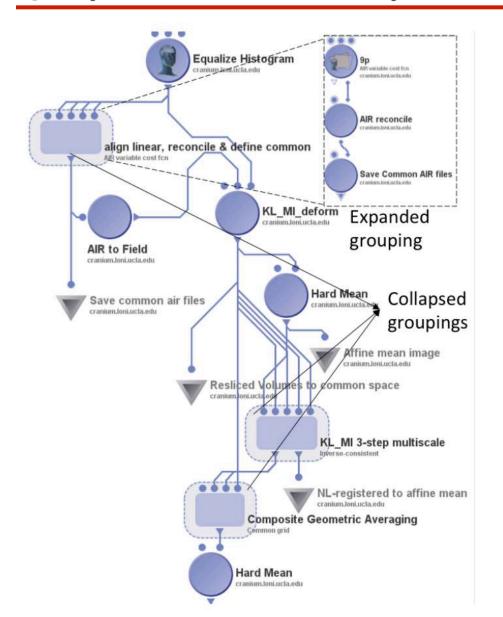
[Gil et al 2012]

#### Work with Christopher Mason (Cornell University)

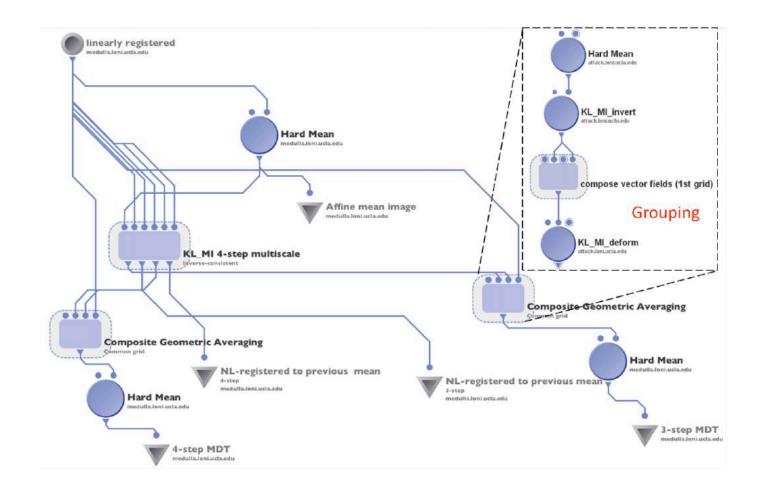


## Learning Reusable Workflow Fragments

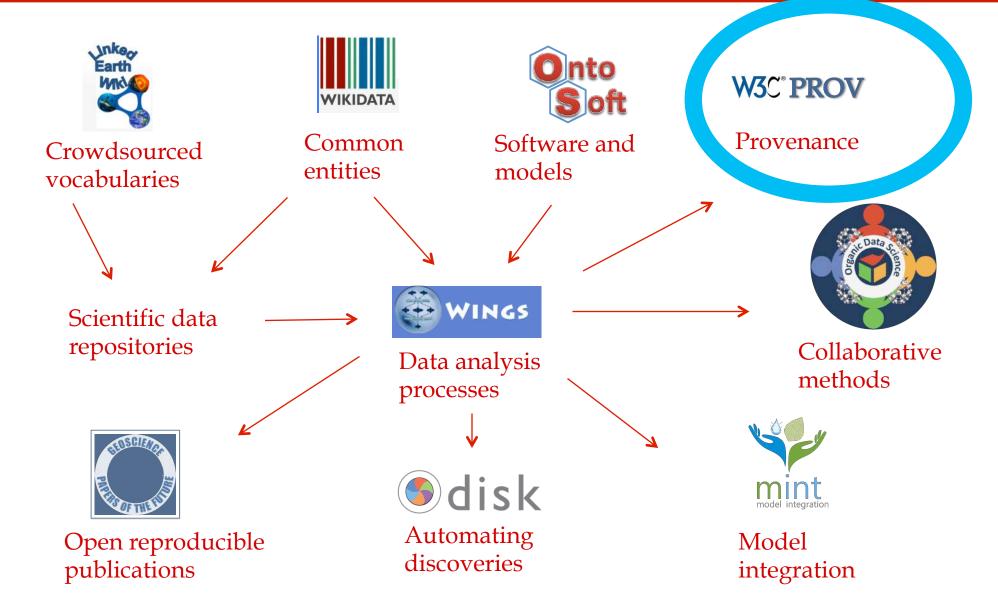
[Garijo et al FCGS 2017; Garijo et al eScience 2014]



Work with Daniel Garijo (USC) and Oscar Corcho (UPM)

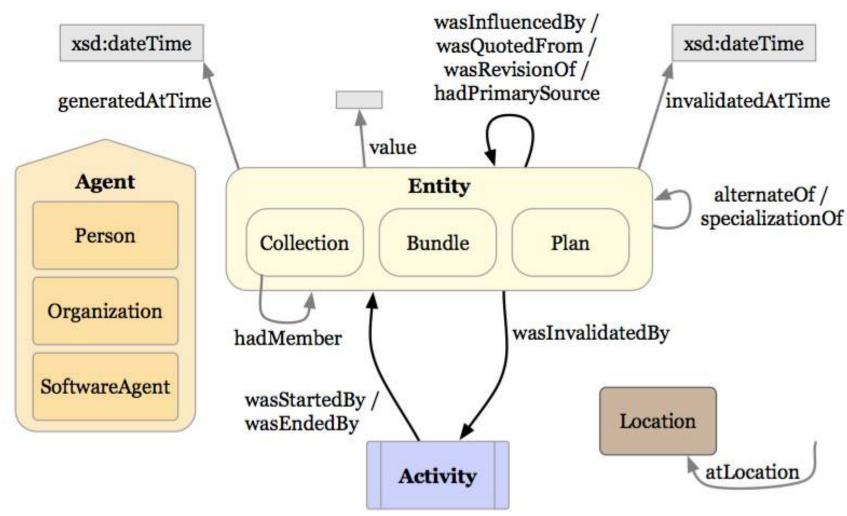


## Capturing Scientific Knowledge

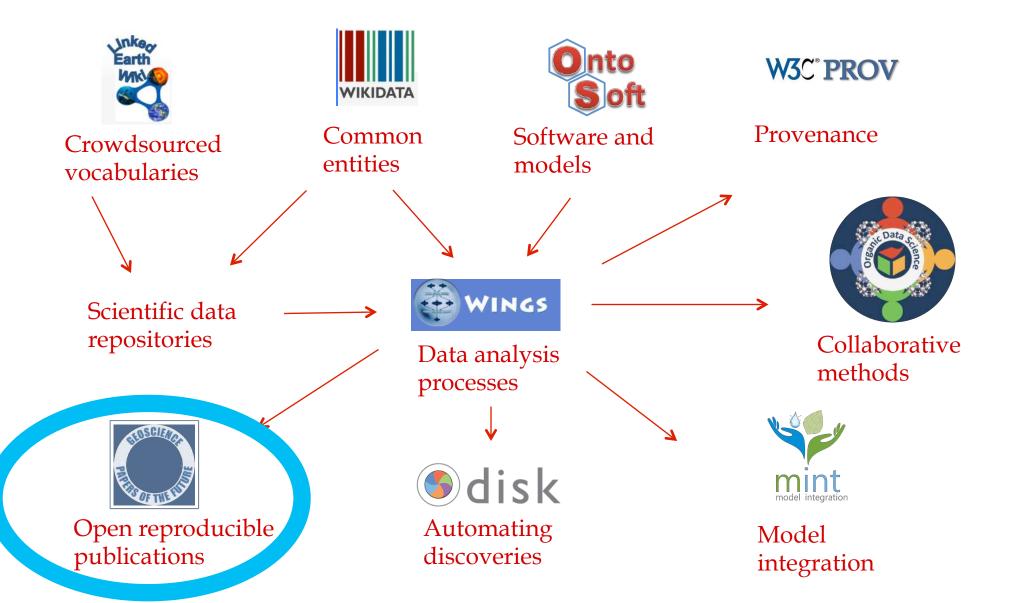


#### The W3C PROV Provenance Standard

[Gil and Miles 2013; Groth and Moreau 2013; Moreau et al 2014]



## Capturing Scientific Knowledge



## www.scientificpaperofthefuture.org

[Gil et al ESS 2016; Essawy et al EMS 2017; Goodman et al PLOS CB 2014]

#### Scientific Paper of the Future

#### **Modern Paper**

#### Text:

Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

#### Data:

Include data as supplementary materials and pointers to data repositories

#### Reproducible Publication

#### **Software:**

For data preparation, data analysis, and visualization

#### **Provenance and methods:**

Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

#### **Open Science**

#### **Sharing:**

Deposit data and software (and provenance/workflow) in publicly shared repositories

#### **Open licenses:**

Open source licenses for data and software (and provenance/workflow)

#### Metadata:

Structured descriptions of the characteristics of data and software (and provenance/workflow)

#### **Digital Scholarship**

#### **Persistent identifiers:**

For data, software, and authors (and provenance/workflow)

#### **Citations:**

Citations for data and software (and provenance/workflow)





**Special Section: Geoscience Papers of the Future** 

#### NATURE REVIEWS | NEUROSCIENCE

Scanning the horizon: towards transparent and reproducible neuroimaging research

Russell A. Poldrack<sup>1</sup>, Chris I. Baker<sup>2</sup>, Joke Durnez<sup>1,3</sup>, Krzysztof J. Gorgolewski<sup>1</sup>, Paul M. Matthews<sup>4</sup>, Marcus R. Munafò<sup>5,6</sup>, Thomas E. Nichols<sup>7</sup>, Jean-Baptiste Poline<sup>8</sup>, Edward Vul<sup>9</sup> and Tal Yarkoni<sup>10</sup>

#### Towards the neuroimaging paper of the future

In this Analysis article, we have outlined a number of problems with current practice and made suggestions for improvements. Here, we outline what we would like to see in the neuroimaging paper of the future, inspired by related work in the geosciences<sup>71</sup>.



GEOPHYSICS Call for Papers
Reproducible Research:
Geophysics Papers of the Future



On Reproducible AI:
Towards Reproducible
Research, Open Science, and
Digital Scholarship in AI
Publications

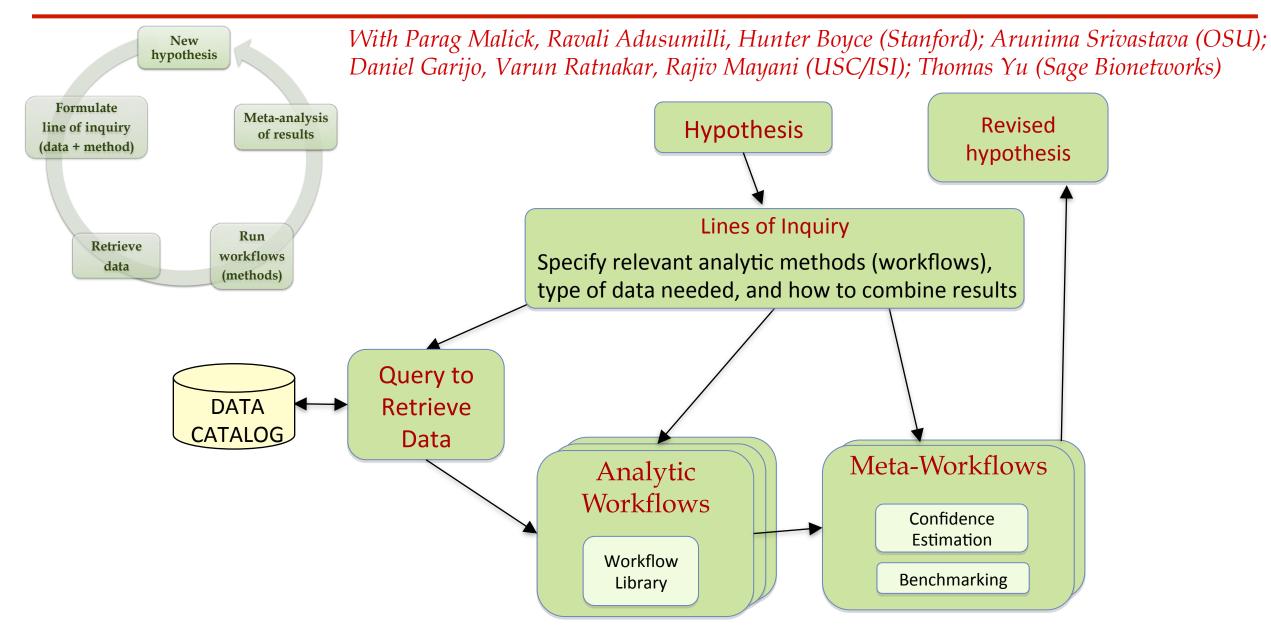
Odd Erik Gundersen, Yolanda Gil, David W. Aha

# Al for Systematic Scientific Data Analysis

### Automated Discovery: Hypotheses and Lines of Inquiry



[Gil et al ACS 2016; Gil et al AAAI 2017; Garijo et al 2017; Srivastava et al PSB 2019; Mallick et al 2020]



### Automated Discovery in DISK

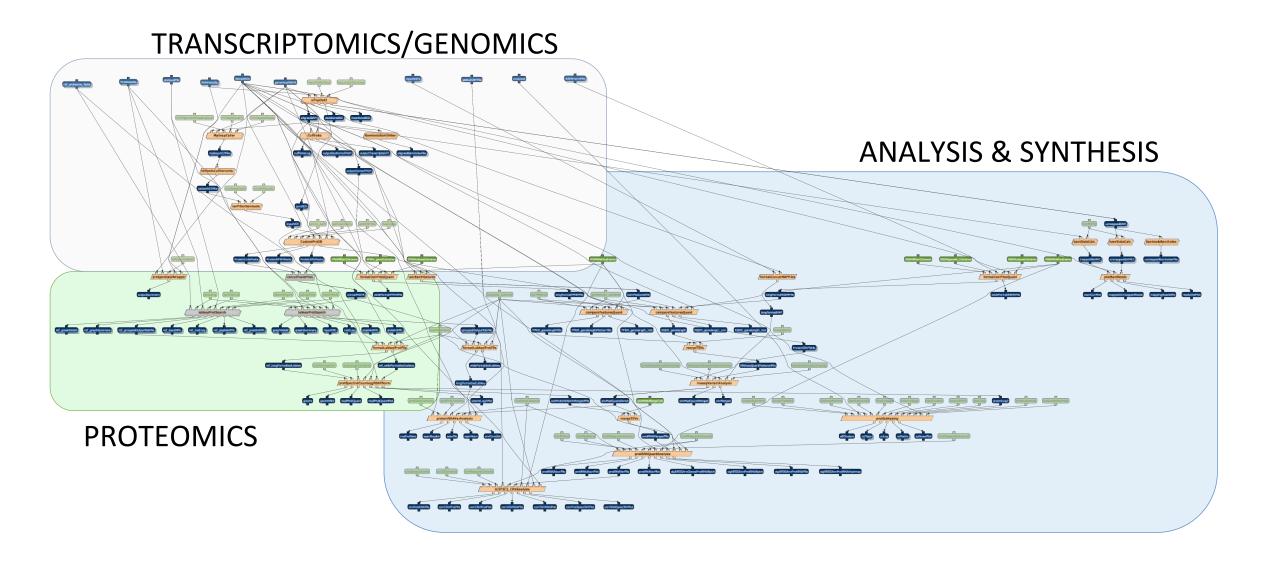


hypothesis Protein PRKCDBP is expressed in samples of patient P36

in samples of patient P36 lines of inquiry revision Lines of Inquiry data PRKCDBP mutation Protein association with patient is expressed in P36 This line of inquiry is used for protein->patient association Query (Ctrl-Space for suggestions) 1 ?x :expressedIn ?sample 3 ?sample :collectedFrom ?p 4 ?p a :Patient 5 ?sample a :Sample 6 ?e1 :experimentedOn ?sample 7 ?e2 :experimentedOn ?sample workflows 8 FILTER(?e1 != ?e2) meta-9 ?e1 :produceData ?d1 10 ?e2 :producedData ?d2 workflows 11 ?d1 a :MassSpecData 12 ?d2 a :RNASeq Workflows to Run + proteogenomic\_analysisBasic Variable Bindings: InputFASTQ = ?d2, input-file = ?d1 proteomics\_analysis Variable Bindings: input-file = ?d1 W3C° PROV Meta-Workflows to Run + Protein\_Diff\_WF Workflow repository

### 





### Representing Hypotheses

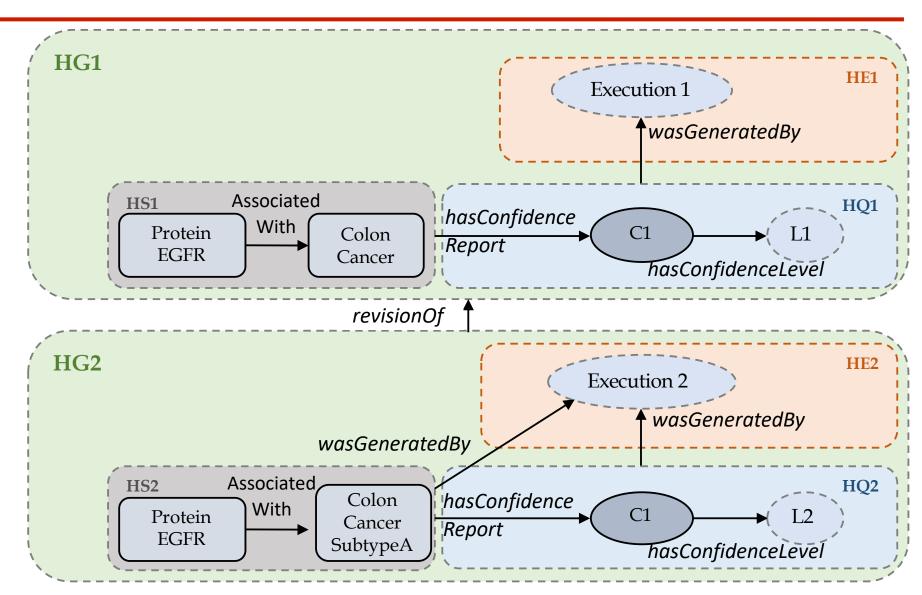


Statement

Qualifier

**Evidence** 

**Evolution** 



DISK Hypothesis Ontology:

http://disk-project.org/ontology/disk

### Reproducing a Seminal Cancer Omics Paper

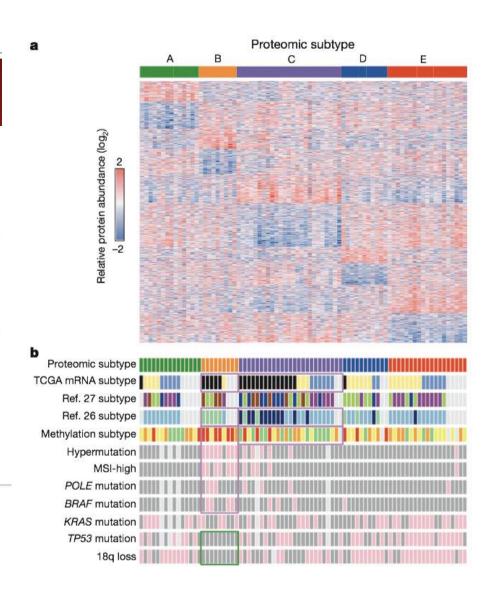


### Proteogenomic characterization of human colon and rectal cancer

Bing Zhang, Jing Wang, Xiaojing Wang, Jing Zhu, Qi Liu, Zhiao Shi, Matthew C. Chambers, Lisa J. Zimmerman, Kent F. Shaddox, Sangtae Kim, Sherri R. Davies, Sean Wang, Pei Wang, Christopher R. Kinsinger, Robert C. Rivers, Henry Rodriguez, R. Reid Townsend, Matthew J. C. Ellis, Steven A. Carr, David L. Tabb, Robert J. Coffey, Robbert J. C. Slebos, Daniel C. Liebler & the NCI CPTAC

#### Affiliations | Contributions | Corresponding author

Nature **513**, 382–387 (18 September 2014) | doi:10.1038/nature13438 Received 19 September 2013 | Accepted 02 May 2014 | Published online 20 July 2014

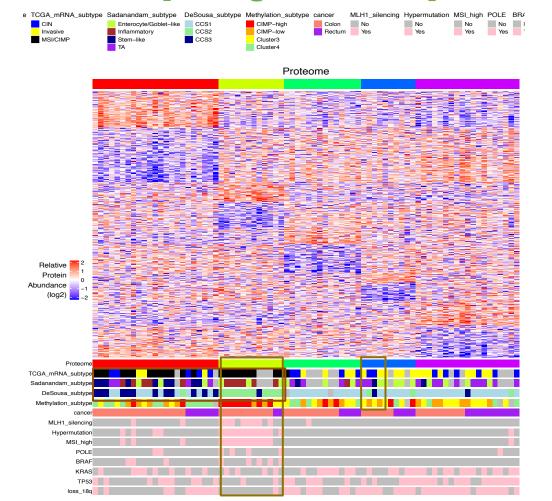


http://www.nature.com/nature/journal/v513/n7518/full/nature13438.html

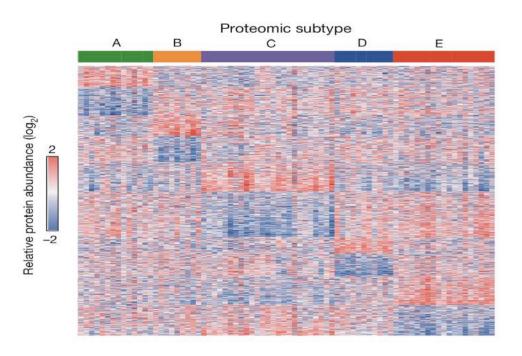
### Patient Proteomic Subtyping



### Original [Zhang et al 2014]



### Reanalysis [Mallick et al 2020]



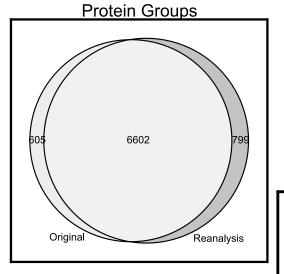
### Large Differences in Protein Identification

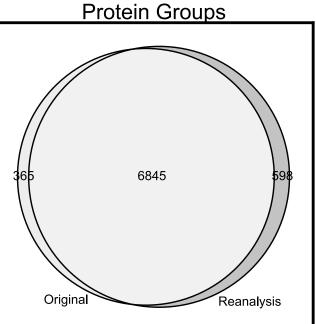


#### 10% Difference

### **Peptides** 1b867 12976 113035 Original Reanalysis

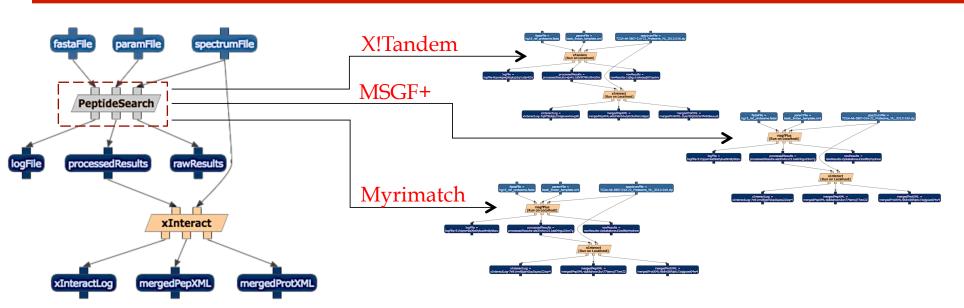
#### 10% to 5% Difference





### Multi-Omics: Sensitivity to Methods/Software

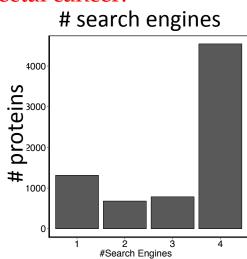


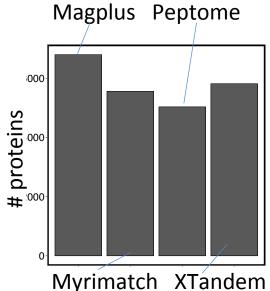


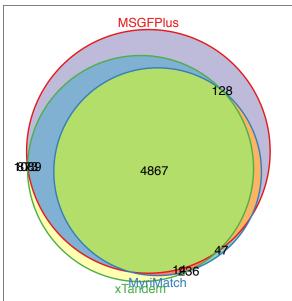
35% of protein identifications are not robust to changing just one analysis step

#### Detection of important proteins in colorectal cancer:

Protein	MSGF+	Myrimatch	X!Tandem
APC	Found	Not found	Found
PIK3CA	Found	Not found	Not found
VTI1A	Found	Not found	Found
PMS1	Found	Not found	Not found
NTHL1	Found	Found	Not found



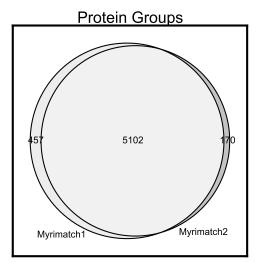


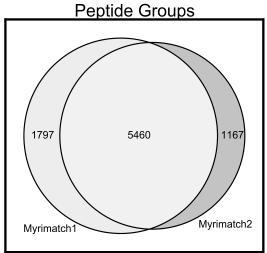


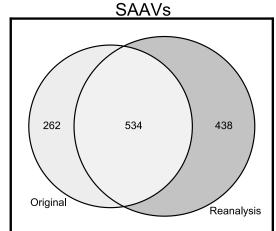
### Multi-Omics: Method Sensitivity

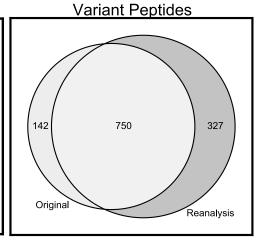


#### **Different Parameters**

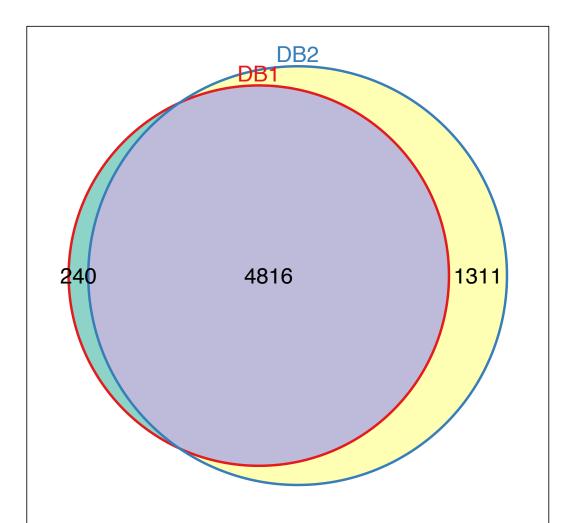






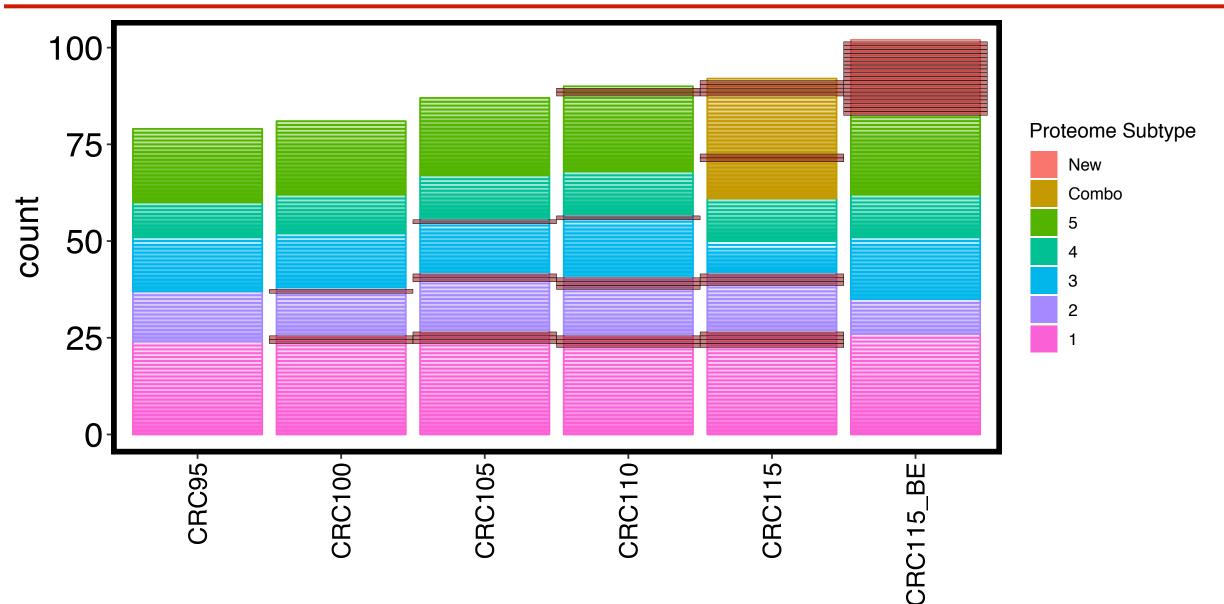


#### **NCBI vs UCSC Reference DB**



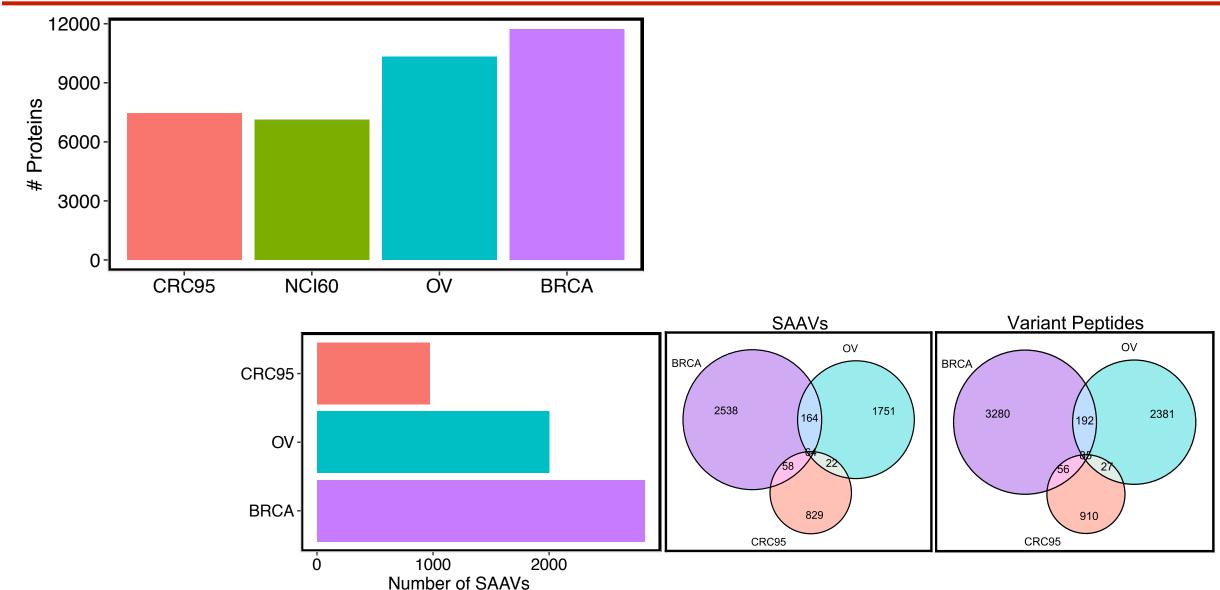
### Revising Hypotheses with Additional Data





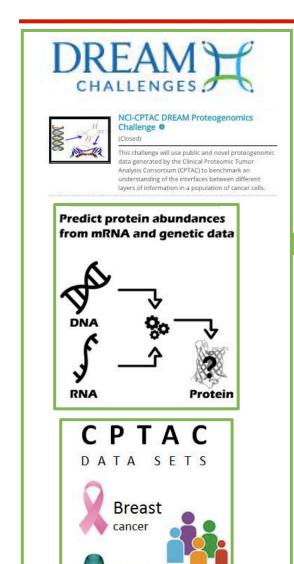
### Analyzing and Comparing Different Datasets disk



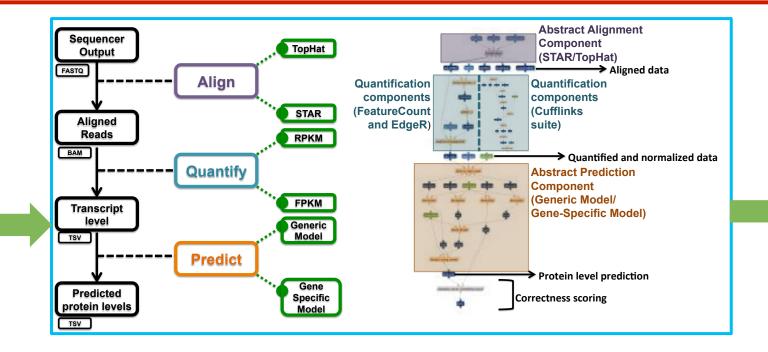


### Benchmarking for DREAM Challenges

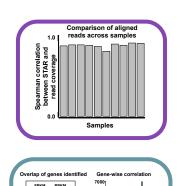


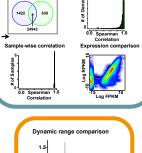


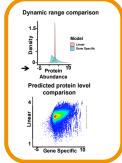
Ovarian



Alignment	Quantification	Predictive Model	Correctness Score	Time Taken
STAR	FPKM	Linear	0.2161	~29 hrs
STAR	RPKM	Linear	0.2155	~20 hrs
STAR	FPKM	Gene-Specific	0.9064	~29 hrs
STAR	RPKM	Gene-Specific	0.9124	~20 hrs
TopHat	RPKM	Linear	0.2053	~103 hrs
TopHat	RPKM	Gene-Specific	0.9080	~103 hrs







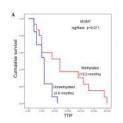
### Al for Automating Discovery

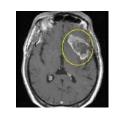
#### **Projection of Data Volumes in 2025**

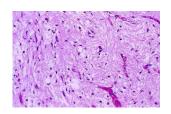
Data Phase	Astronomy	Twitter	YouTube	Genomics
Acquisition	25 zetta-bytes/year	0.5–15 billion tweets/year	500-900 million hours/year	1 zetta-bases/year
Storage	1 EB/year	1-17 PB/year	1–2 EB/year	2-40 EB/year
Analysis	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
Distribution	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

doi:10.1371/journal.pbio.1002195.t001

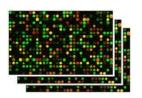
#### **Analytic Complexity for Each Type of Data**





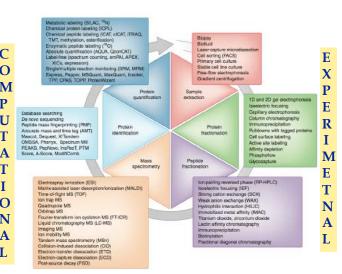








### Siloed **Expertise**



doi: 10.1038/nbt.1658

# Al for Interdisciplinary Science Frontiers

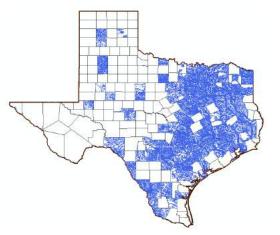
### MINT: Model INTegration

[Gil et al IEMS 2018; Garijo et al eScience 2019]



Collaboration with Daniel Garijo, Deborah Khider, Craig Knoblock, Ewa Deelman, Rafael Ferreira (USC/ISI), Vipin Kumar (UM), Scott Peckham (CU), Chris Duffy & Armen Kemanian (PSU), Kelly Cobourn (VT), Suzanne Pierce (UT)













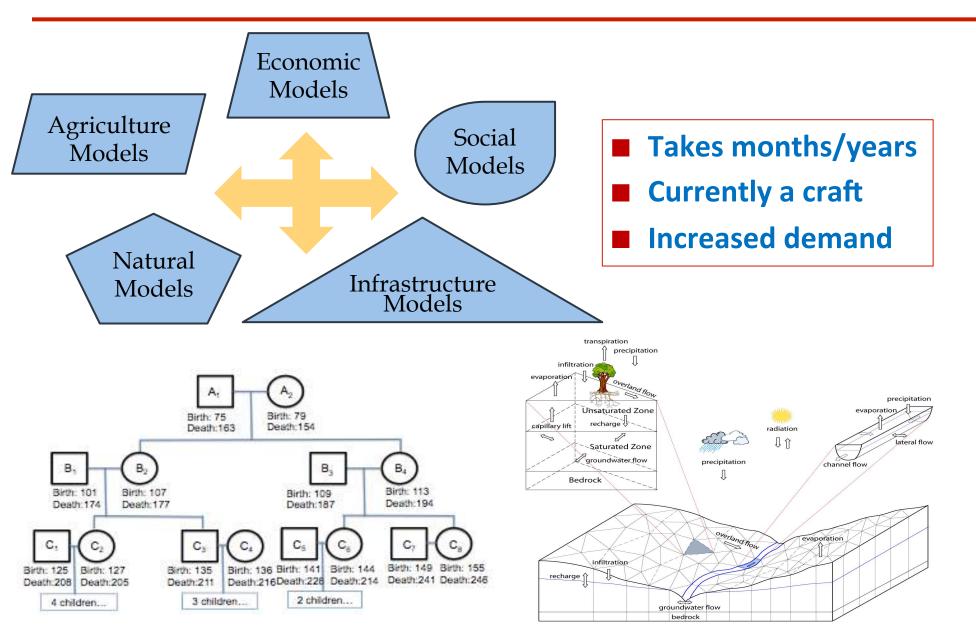




https://www.pexels.com/photo/bird-birds-blue-sky-drought-1178291/

### Integrated Modeling





### Mediation at many levels

Modeling scope

Model set up

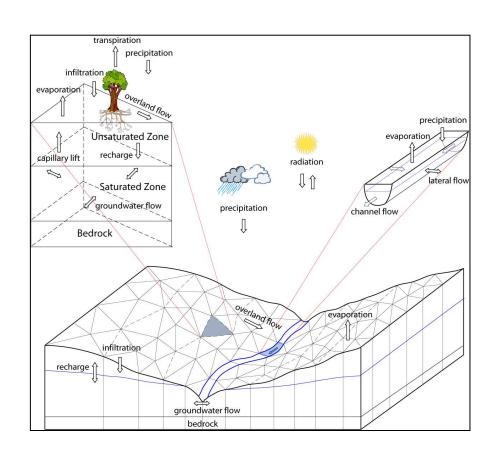
Variable mapping

Data ingestion

Spatial gridding

### Supporting Compositionality of Scientific Knowledge





- Data formats
- Physical variables
- Constraints for use
- Adjustable parameters
- Interventions

### Ontologies for Unambiguous Physical Variables



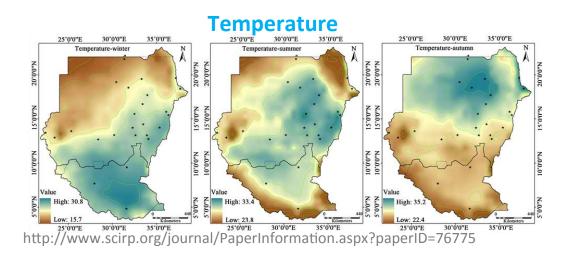
Work by Scott Peckham and Maria Stoica (CU)

- Ontology of standard scientific names
  - Eg SSN: watershed\_outlet\_water\_\_volume\_outflow\_rate is more precise than "streamflow" or "discharge"

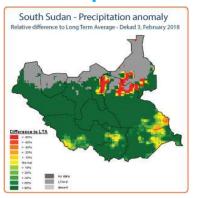
Label	Long Name	Description	Standard Name	Units
RV	rock volume	Rock volume expressed as volume over total volume	soil_rockvolume_fraction	m3 m-3
DZ	soil layer thickness	Soil layer thickness	soil_layerthickness	m
SLOPE	slope of the field	Average slope of field of interest	land_surfaceslope	m m-1
LAYER	soil layer number	Soil layer number from 1 to an integer that is user defined (or database defined)		
SOM	soil organic matter	Soil organic matter per unit of non-rock soil expressed as percentage over total mass	soil_matter~organicmass_fraction	kg kg-1 x- 100
SILT	silt percentage	Silt mass per unit of soil mass (no rocks) and expressed as a percentage	soil~no-rock_siltmass_fraction	%
BD	bulk density	Soil mass dry and wihtout rock divided by the sampled volume	soil~no-rock~drymass-per-volume_density	Mg m-3
CLAY	clay percentage	clay particle size fraction size fraction of each soil layer in %.	soil_clay_particlevolume_fraction	%

### Spatial Datasets of Varying Quality

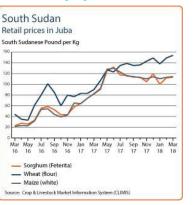




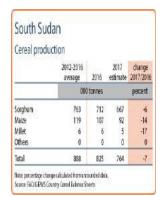
#### **Precipitation**



#### **Crop prices**

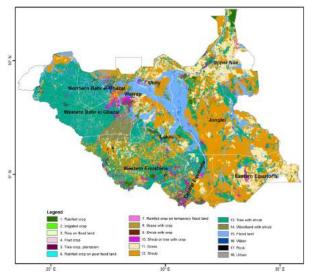


#### **Production**



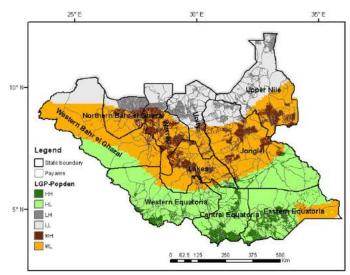
http://www.fao.org/giews/countrybrief/country.jsp?code=SSD

#### Land use



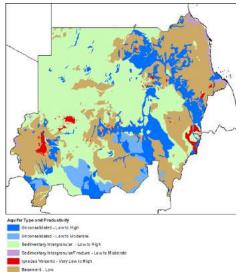
Diao, X. et al. DOI: 10.5772/47938

#### **Agricultural potential**



Diao et al DOI:10.5772/47938

#### **Aquifer productivity and recharge**



http://earthwise.bgs.ac.uk/index.php/Hydrogeology\_of\_Sudan

### **Automated Data Transformations**



Create a saturation file for each mesh cell

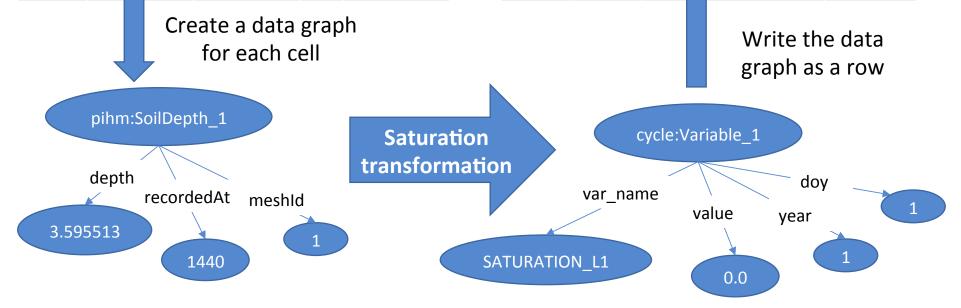
Work by Craig Knoblock, Yao-Yi Chiang Jay Pujara (USC)

Cycles.REINIT file

					CGCIIII
Time	Mesh	1	Mesh 2	•••	Mesh n
1440	3.595	513	6.534754		3.771523
2880	3.59	509	6.534728		3.771488
4320	3.59	505	6.534702	•••	3.771453
•••					

pg.gw file

ROT	DOY	VARIABLE	VALUE
1	1	INFILTRATION	0.0
1	1	SATURATION_L1	0.0
1	1	SAT TION_L2	0.012



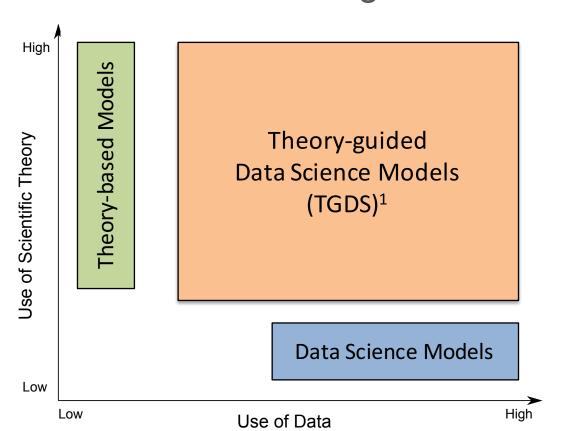
### Creating Virtual Gauges When No Data Available

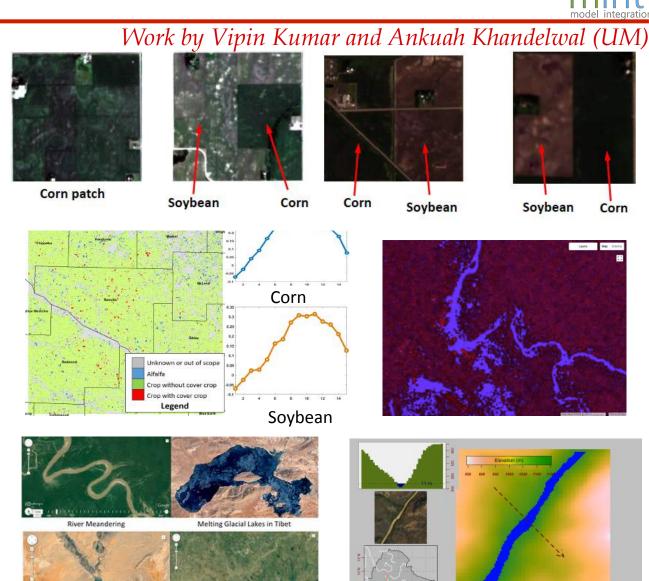


Corn

Soybean

Theory-guided data science incorporates biophysical laws into machine learning from remote sensing data





Shrinking Lake Mead

### Mapping Models to Interventions and Decisions



# Drivers and adjustable parameters

Weather

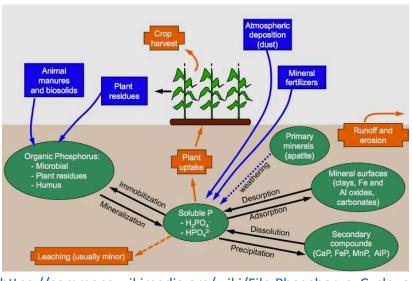
Crop types

Fertilizer

Planting dates

Weed factor

Land use



https://commons.wikimedia.org/wiki/File:Phosphorus Cycle copy.jpg

### Model parameters

Nitrogen stress

Soils

Solar radiation

Work by Armen Kemanian and Yuning Shi (PSU)

## Response variables

Crop yield

### Rich Representations of Models

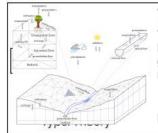




#### **Model Catalog**

Q Search models

Search on Full text



#### The Soil & Water Assessment Tool (SWAT)

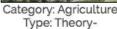
The Soil & Water Assessment Tool (SWAT) is a small watershed to river basin-scale model used to simulate the quality and quantity of surface and ground water and predict the environmental impact of land use; land management practices and climate change (https://swat.tamu.edu/2019)

Keywords: Soil, watershed, surface water, ground water, en... More details

2 versions

#### Cycles 🔼

Cycles simulates the productivity and the water-carbon and nitrogen balance of soil-crop systems subject to climate conditions and a large array of management constraints



Keywords: agriculture, cycles, crop growth, weather, soil, cr... More details



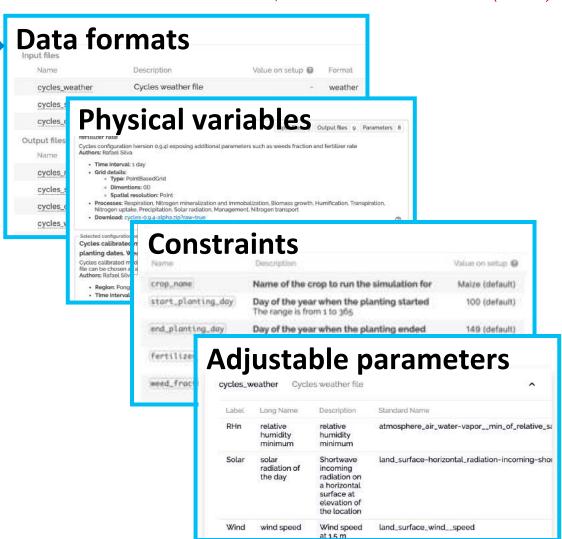
#### Economic aggregate crop supply response model (EACS)

The Aggregate crop supply response model (EACS) describes the aggregate crop supply response model for the country of South Sudan. This is a regional-scale aggregate model of agricultural supply for a specified set of crops (cassava; groundnuts; maize; sesame seed; and sorghum).

Category: Economy Type: Theory-Keywords: economy, land use, crop production, fertilizer c...

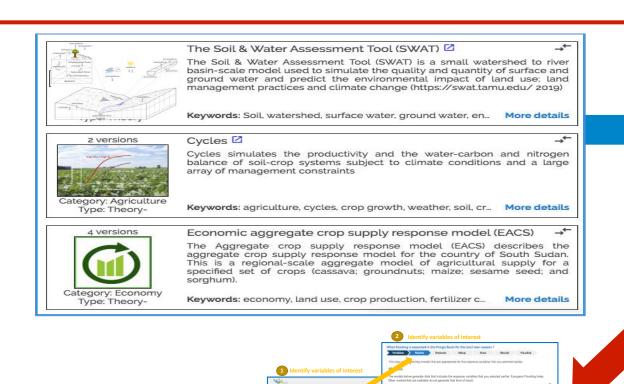
More details

Work with Daniel Garijo, Deborah Khider, Varun Ratnakar, Maximiliano Osorio (USC)

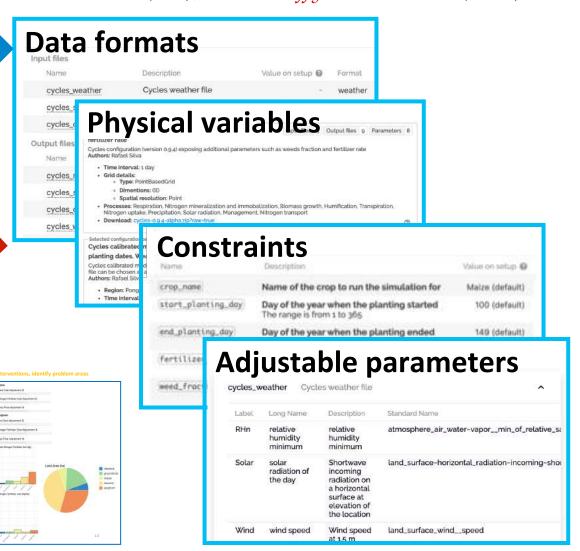


### MINT: From Models to Solutions





Work with Deborah Khider (USC); Suzanne Pierce and Lissa Pearson(UT); Chris Duffy and Lele Shu (PSU)



### Model Portability: Data Scarce Regions

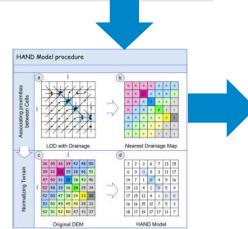


#### Height Above Natural Drainage Model (HAND)

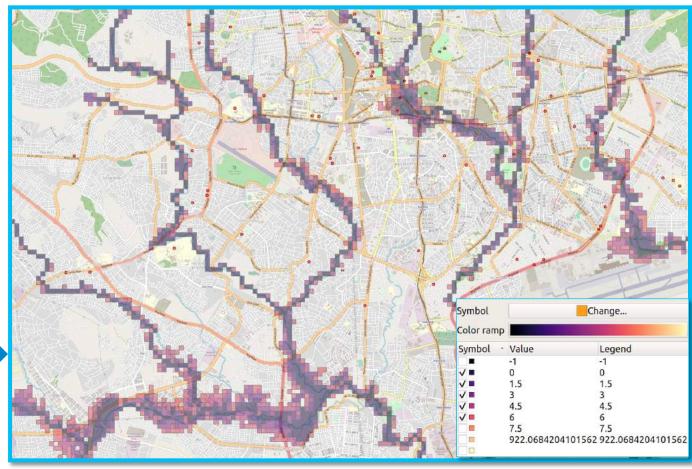


Results provide first pass vulnerability for flood risk

Pixels with colors on right of scale (brighter red) indicate higher change of inundation or flooding



Work with Suzanne Pierce, Daniel Hardesty Lewis, David Arctur Paola Passalacgua (UT), David Tarboton (Utah State); Misty Porter, Mary Hill (KU)

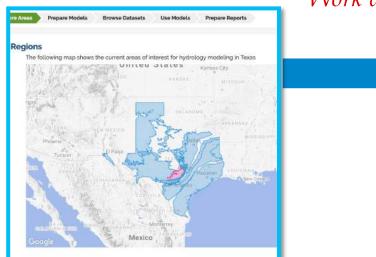


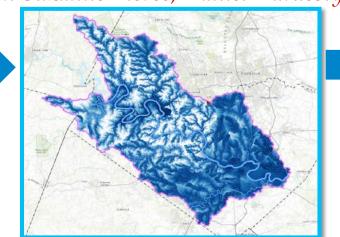
### Model Portability: Data Rich Regions



**Travis County** 

Work with Suzanne Pierce, Daniel Hardesty Lewis, David Arctur and Paola Passalacgua (UT)





To make the second of the seco

Combining

Terrain (10 m)

Population

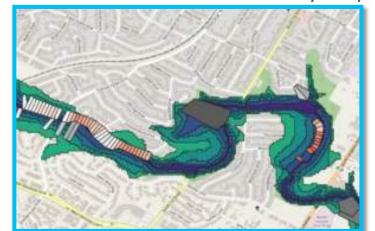
Urban Land Use

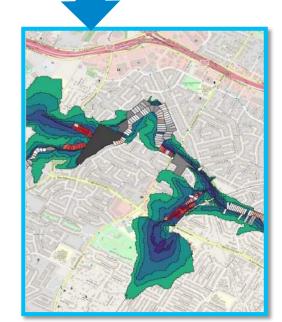
Comparison: Williamson County

Manually Constructed Buyout Map



MINT Generated Vulnerability Map





### A Perspective on the Future

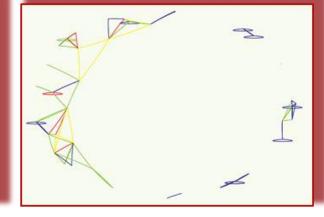
### Tackling Complex Scientific Phenomena

Single authorship



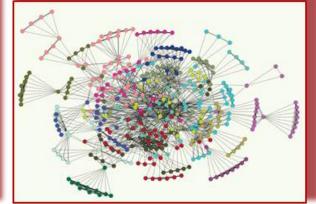
Co-authorship





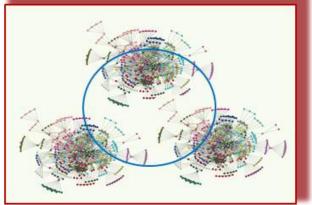
Large number of co-authors





Community as author







### Consider the Atlas Collaboration



#### Atlas collaboration

- Approximately 4,000 authors
- Worked in subgroups that coordinated with one another
- Collaboration lasted many years

Today, large scientific collaborations take significant time and effort and therefore are not very frequent

How can we change this?

### The Importance of Process



Freestyle Chess Champion Anson Williams

"The winner was revealed to be not a grandmaster with a state-of-the-art PC but a pair of amateur American chess players using three computers at the same time. Their skill at manipulating and "coaching" their computers to look very deeply into positions effectively counteracted the superior chess understanding of their grandmaster opponents and the greater computational power of other participants.

Weak human + machine + better process was superior to a strong computer alone and, more remarkably, superior to a strong human + machine + inferior process."

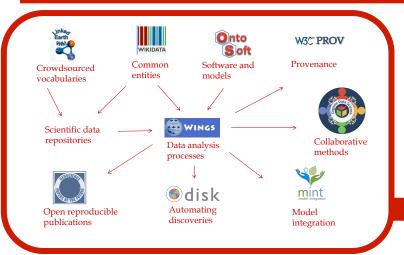
- Garry Kasparov, 2010

### Future Al Systems as Partners for Discovery [Gil DSJ 2017]

### Thoughtful AI: Principles for Partnership

Rationality	Behavior is governed by explicit knowledge structures
Context	Seek to understand the purpose and scope of tasks
Initiative	Proactively learn new knowledge relevant to their task
Networking	Access external sources of knowledge and capabilities
Articulation	Respond with persuasive justifications and arguments
Systems	Facilitate integration & collaboration with humans/systems
Ethics	Behavior that conveys scope and uncertainty

These are important research challenges for Al



#### Thoughtful AI: Principles for Partnership

Rationality	Behavior is governed by explicit knowledge structures
Context	Seek to understand the purpose and scope of tasks
Initiative	Proactively learn new knowledge relevant to their task
Networking	Access external sources of knowledge and capabilities
Articulation	Respond with persuasive justifications and arguments
Systems	Facilitate integration & collaboration with humans/systems
Ethics	Behavior that conveys scope and uncertainty

#### Scientific Paper of the Future

#### Modern Paper

#### Text:

Narrative of the method, some data is in tables, figures/plots, and the software used is mentioned

#### Data:

Include data as supplementary materials and pointers to data repositories

#### Reproducible Publication

#### Softwar

For data preparation, data analysis, and visualization

Provenance and methods: Workflow/scripts specifying dataflow, codes, configuration files, parameter settings, and runtime dependencies

#### **Open Science**

#### Sharing:

Deposit data and software (and provenance/workflow) in publicly shared repositories

#### Open licenses:

Open source licenses for data and software (and provenance/workflow)

#### Metadata:

Structured descriptions of the characteristics of data and software (and provenance/workflow)

#### **Digital Scholarship**

#### Persistent identifiers:

For data, software, and authors (and provenance/workflow)

#### Citations:

Citations for data and software (and provenance/workflow)

**YOUR IDEAS HERE!** 

### Formulating Grand Challenges

#### COMMUNICATIONS

ACM

January 2019

**REVIEW ARTICLES** 

#### Intelligent Systems for Geosciences: An Essential Research Agenda

By Yolanda Gil, Suzanne A. Pierce, Hassan Babaie, Arindam Banerjee, Kirk Borne, Gary Bust, Michelle Cheatham, Imme Ebert-phoff, Carla Gomes, Mary Hill, John Horel, Leslie Hsu, Jim Kinter, Craig Knoblock, David Krum, Vipin Kumar, Pierre Lermusiaux, Yan Liu, Chris North, Victor Pankratius, Shanan Peters, Beth Plale, Allen Pope, Sai Ravela, Juan Restrepo, Aaron Ridley, Hanan Samet, Shashi Shekhar



Many aspects of geosciences pose novel problems for intelligent systems research. Geoscience data is challenging because it tends to be uncertain, intermittent, sparse, multiresolution, and multiscale. Geosciences processes and objects often have amorphous spatiotemporal boundaries. The lack of ground truth makes model evaluation, testing, and comparison difficult. Overcoming these challenges requires breakthroughs that would significantly transform intelligent systems, while greatly benefitting the geosciences in turn. Although there have been significant and



Spring 2016

# Artificial Intelligence to Win the Nobel Prize and Beyond: Creating the Engine for Scientific Discovery

Hiroaki Kitano

Al reproducing articles

Al as research assistant

Al as co-author

### The Next Two Decades

2025: AI can generate automatically new complex scientific analyses using open data

2025: Al detects when it is missing knowledge and can seek and read new scientific papers on target topics

2030: Al can generate and test sophisticated hypotheses about complex physical phenomena

2030: AI can reproduce the results in 80% of the articles in a scientific journal

2035: Al can describe a scientific experiment and discuss sophisticated aspects of it

2035: AI can compare scientific experiments and papers and contrast their merits

2040: Al can teach advanced theories in some scientific domain effectively to students

2040: Al can formulate research questions and generate novel contributions in some scientific domain

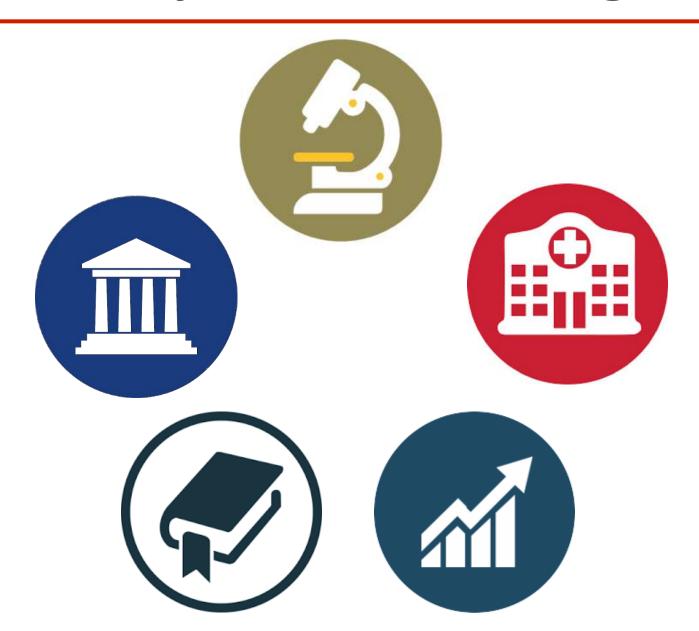
Al reproducing articles

Al as research assistant

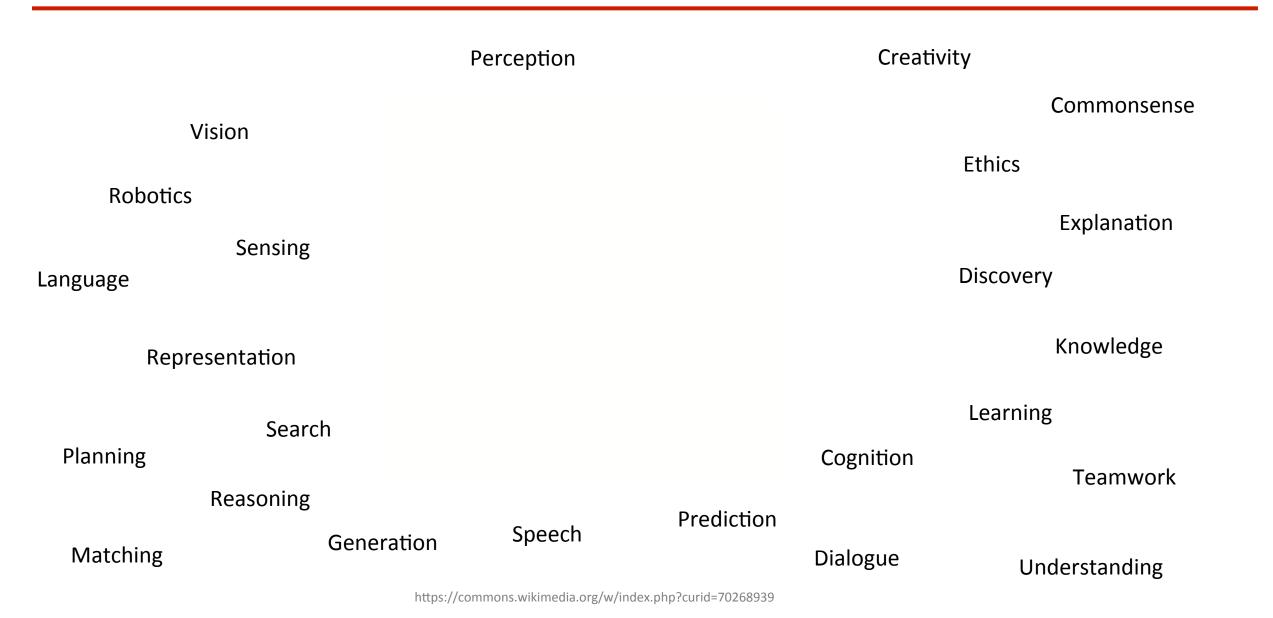
Al as co-author

2030 2035 2040

### Al to Address Major Future Challenges



### Diversity and Breadth of Advances in Al



- The AI community has always been
  - Visionary
  - Broad
  - Inclusive
  - Interdisciplinary
  - Determined
- And dare I say
  - Successful

So... the answer may be yes?

- Al researchers have been:
  - Visionary
  - Broad
  - Inclusive
  - Interdisciplinary
  - Determined
- And dare I say
  - Successful

So... the answer may be yes?

- Humans:
  - Not systematic
  - Errors
  - Biases
  - Poor reporting

So... the answer is definitely yes?

- The AI community has always been
  - Visionary
  - Broad
  - Inclusive
  - Interdisciplinary
  - Determined
- And dare I say
  - Successful

Maybe the answer is no:



Pennicillin discovery resulted from human error...

- Not systematic
- Errors
- Biases
- Poor reporting

So... the answer may be yes?

So... the answer is definitely yes?

- The AI community has always been
  - Visionary
  - Broad
  - Inclusive
  - Interdisciplinary
  - Determined
- And dare I say
  - Successful

So... the answer may be yes?

Maybe the answer is no:



Pennicillin discovery resulted from human error...



...and humans make unique contributions...

- Not systematic
- Errors
- Biases
- Poor reporting

So... the answer is definitely yes?

### Thank you!





































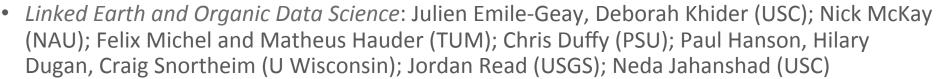








- Varun Ratnakar, Daniel Garijo, Deborah Khider, Maximiliano Osorio, Hernan Vargas (USC)
- Workflows: Jihie Kim, Ewa Deelman, Karan Vahi; Rafael Ferreira, Rajiv Mayani, Hyunjoon Jo, Yan Liu, Dave Kale (USC); Ralph Bergmann (U Trier); William Cheung (HKBU); Oscar Corcho (UPM); Pedro Gonzalez, Gonzalo Castro (UCM); Paul Groth (UA); Ricky Sethi (FSU); Carole Goble (UM); Chris Mattmann, Paul Ramirez, Dan Crichton, Rishi Verma (JPL); Natalia Villanueva (UTEP)



- Biomedical workflows: Phil Bourne, Sarah Kinnings (UCSD); Chris Mason (Cornell); Joel Saltz, Tahsin Kurk (Emory U.); Jill Mesirov, Michael Reich (Broad); Shannon McWeeney, Christina Zhang (OHSU); Parag Mallick, Ravali Adusumilli, Hunter Boyce (Stanford U.)
- Geosciences workflows: Paul Hanson (U Wisconsin), Tom Harmon & Sandra Villamizar (U Merced), Tom Jordan & Phil Maechlin (USC), Kim Olsen (SDSU); Suzanne Pierce (UT); Chris Duffy & Armen Kemanian (PSU); Scott Peckham & Maria Stoica (CU)
- And many others!





