

Hypothesis-Driven Data Analysis of Science Repositories

Yolanda Gil



Information Sciences Institute
and Department of Computer Science

University of Southern California

<http://www.isi.edu/~gil>

@yolandagil
gil@isi.edu



Acknowledgments



- *WINGS contributors:* Varun Ratnakar, Daniel Garijo, Rajiv Mayani, Ricky Sethi, Hyunjoon Jo, Jihie Kim, Yan Liu, Dave Kale (USC), Ralph Bergmann (U Trier), William Cheung (HKBU), Oscar Corcho (UPM), Pedro Gonzalez & Gonzalo Castro (UCM), Paul Groth (VUA)
- *WINGS collaborators:* Chris Mattmann & Paul Ramirez & Dan Crichton & Rishi Verma (JPL), Ewa Deelman & Gaurang Mehta & Karan Vahi (USC), Sofus Macskassy (ISI), Natalia Villanueva & Ari Kassin (UTEP)
- *Organic Data Science contributors:* Felix Michel and Matheus Hauder (TUM); Varun Ratnakar (ISI); Chris Duffy (PSU); Paul Hanson, Hilary Dugan, Craig Snortheim (U Wisconsin); Jordan Read (USGS); Neda Jahanshad (USC), Julien Emile-Geay (USC), Nick McKay (NAU)
- *DISK contributors:* Parag Mallick, Raval Adusumilli, Hunter Boyce (Stanford); Suzanne Pierce, John Gentle (UT Austin)
- *Biomedicine:* Parag Mallick & Raval Adusumilli & Hunter Boyce (Stanford U.), Phil Bourne & Sarah Kinnings (UCSD), Chris Mason (Cornell); Joel Saltz & Tahsin Kurk (Emory U.); Jill Mesirov & Michael Reich (Broad); Randall Wetzel (CHLA); Shannon McWeeney & Christina Zhang (OHSU)
- *Geosciences:* Suzanne Pierce & John Gentle (UT Austin), Chris Duffy (PSU); Paul Hanson (U Wisconsin); Tom Harmon & Sandra Villamizar (U Merced); Tom Jordan & Phil Maechlin (USC); Kim Olsen (SDSU)
- *And many others!*



Community Workshops



NSF Workshop on Discovery Informatics

February 2-3, 2012

Arlington, VA

Final Workshop Report

August 31, 2012



Intelligent Systems for Geosciences



<http://www.IS-GEO.org>

Outline

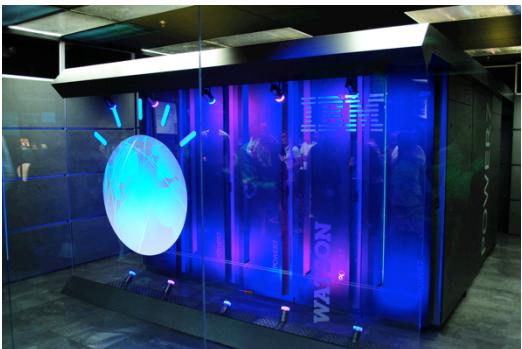
- Artificial intelligence and scientific discovery
 - The knowledge systems tradition
- Our recent work on capturing knowledge about data analysis strategies
 - Hypothesis-driven data analysis
- Representing and capturing data analysis knowledge
 - About data, software, methods, meta-analysis
- Summary of AI challenges

Outline

- Artificial intelligence and scientific discovery
 - The knowledge systems tradition
- Our recent work on capturing knowledge about data analysis strategies
 - Hypothesis-driven data analysis
- Representing and capturing data analysis knowledge
 - About data, software, methods, meta-analysis
- Summary of AI challenges

AI's Coming of Age

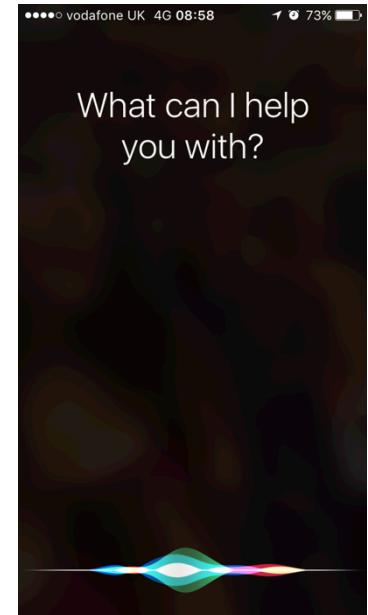
IBM Watson



Google Knowledge Graph



Apple Siri



Netflix Recommenders



Tesla AutoPilot



RoboCup Soccer



[https://en.wikipedia.org/wiki/Watson_\(computer\)](https://en.wikipedia.org/wiki/Watson_(computer))

<https://en.wikipedia.org/wiki/Siri>

https://commons.wikimedia.org/wiki/File:Google_Knowledge_Panel.png

<https://commons.wikimedia.org/wiki/File:13-06-28-robocup-eindhoven-005.jpg>

http://www.greencarreports.com/news/1100482_tesla-autopilot-the-10-most-important-things-you-need-to-know

<https://en.wikipedia.org/wiki/Netflix>

AI's Coming of Age

IBM Watson



Google Knowledge Graph



Apple Siri



Problem Solving Tradition

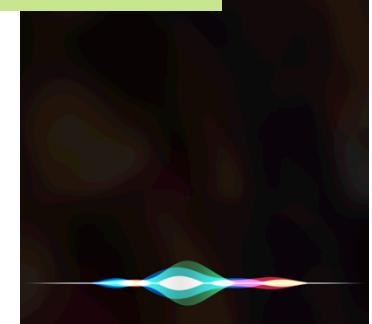
Birth: July 4, 1743, Charlottesville, VA
Presidential term: March 4, 1801 – March 4, 1809
Spouse: Martha Jefferson (m. 1772–1782)
Party: Democratic-Republican Party
Awards: AIA Gold Medal

Get updates about Thomas Jefferson

People also search for

John Adams George Washington Benjamin Franklin James Madison Alexander Hamilton

Feedback



Recommenders



Autonomous Drivers

Robot Teams

Data Systems Tradition



[https://en.wikipedia.org/wiki/Watson_\(computer\)](https://en.wikipedia.org/wiki/Watson_(computer))#/media/File:IBM_Watson.PNG

<https://en.wikipedia.org/wiki/Siri>#/media/File:SirioniOS9.png

https://commons.wikimedia.org/wiki/File:Google_Knowledge_Panel.png

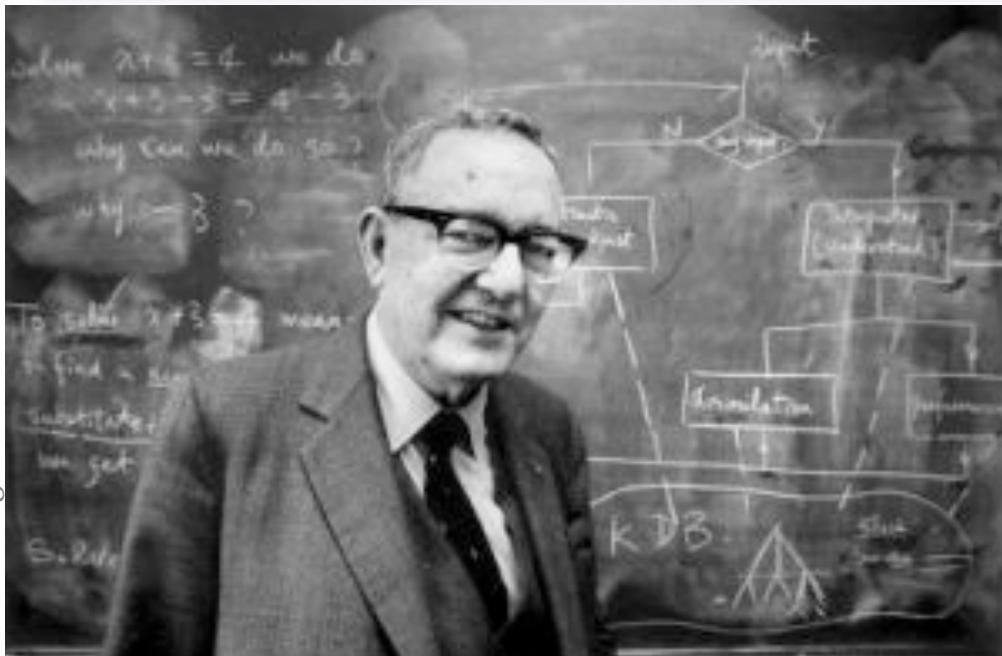
<https://commons.wikimedia.org/wiki/File:13-06-28-robocup-eindhoven-005.jpg>

http://www.greencarreports.com/news/1100482_tesla-autopilot-the-10-most-important-things-you-need-to-know

<https://en.wikipedia.org/wiki/Netflix#/media/File:NetflixDVD.jpg>

Artificial Intelligence and Scientific Discovery: The Problem Solving Tradition

Pittsburgh Post Gazette Archives



Outline

- Artificial intelligence and scientific discovery
 - The knowledge systems tradition
- Our recent work on capturing knowledge about data analysis strategies
 - Hypothesis-driven data analysis
- Representing and capturing data analysis knowledge
 - About data, software, methods, meta-analysis
- Summary of AI challenges

DISK: Automated Discovery of Scientific Knowledge

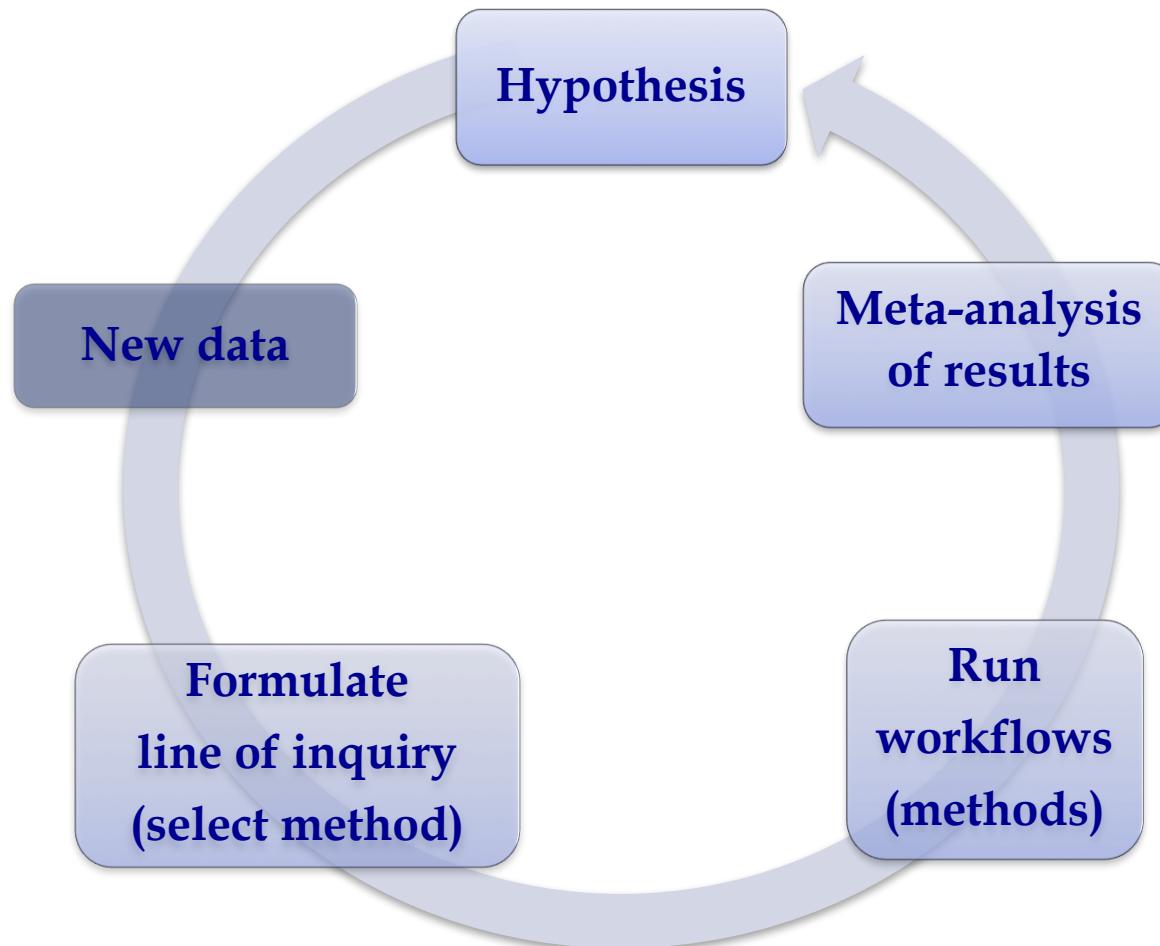
Work with P. Mallick, R. Adusumilli, H Boyce (Stanford U.)

- Long-Term Goal: Human directs automated intelligent system to explore hypotheses of interest
 - Hypothesis-driven data analysis and discovery
 - Systematic and reproducible analyses
 - Report of findings with explanations (“Friday” meeting)

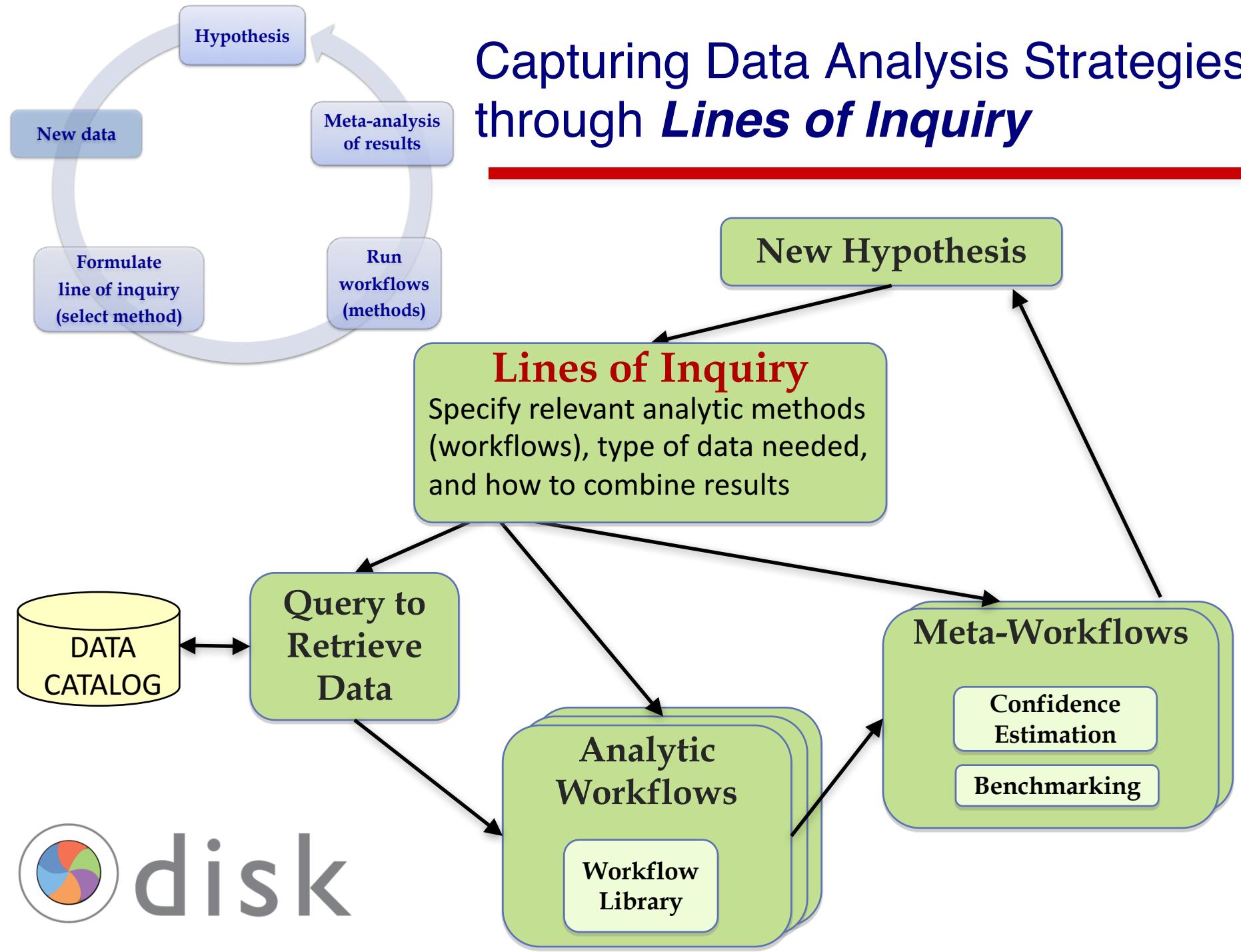
- Approach: Intelligent system that captures common data analysis strategies used by scientists in a domain
 - Build on WINGS intelligent workflow system that can adapt data analysis given the constraints of algorithmic steps



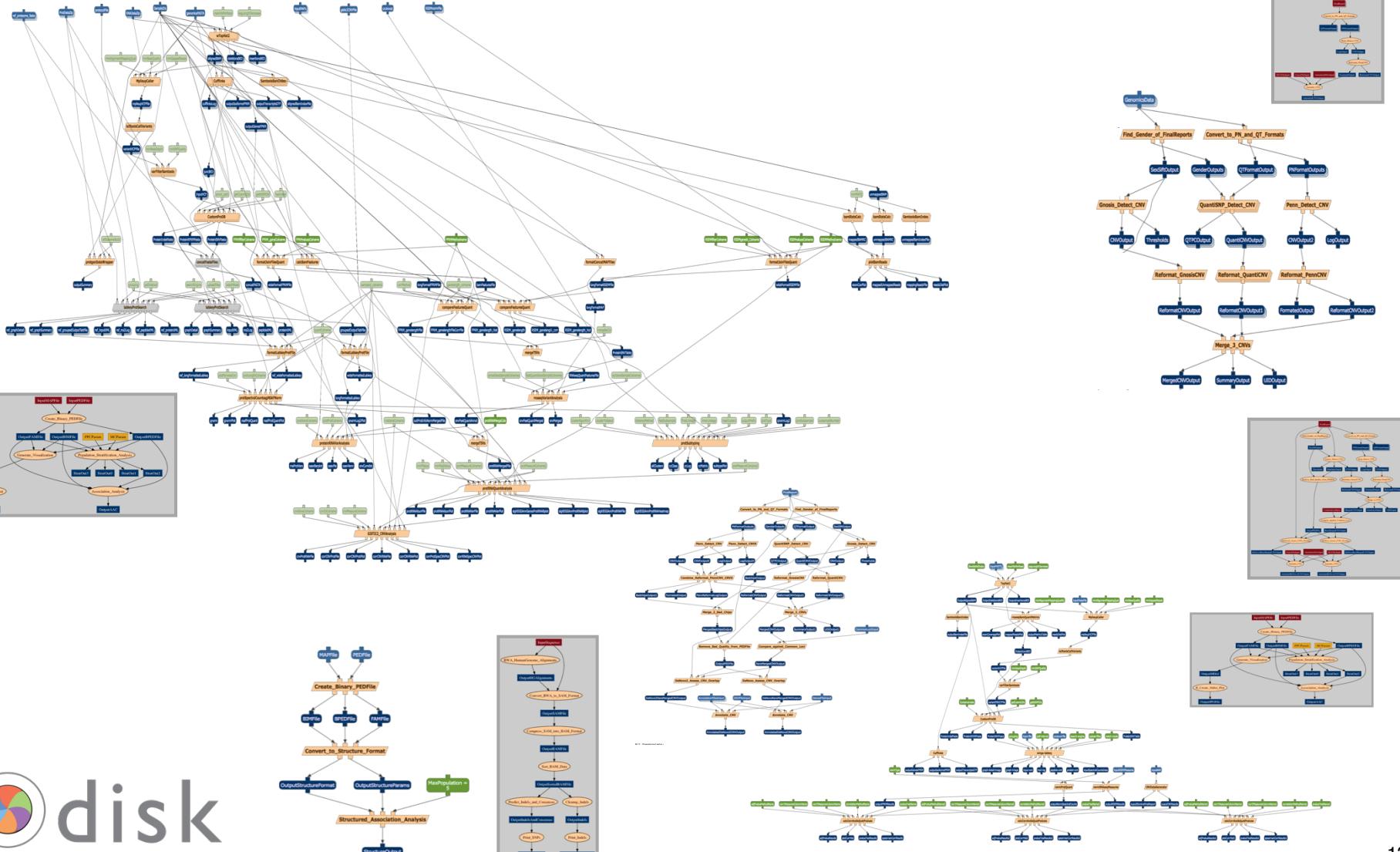
Scientific Data Analysis Today: Inefficient, Incomplete, Irreproducible, One-Shot



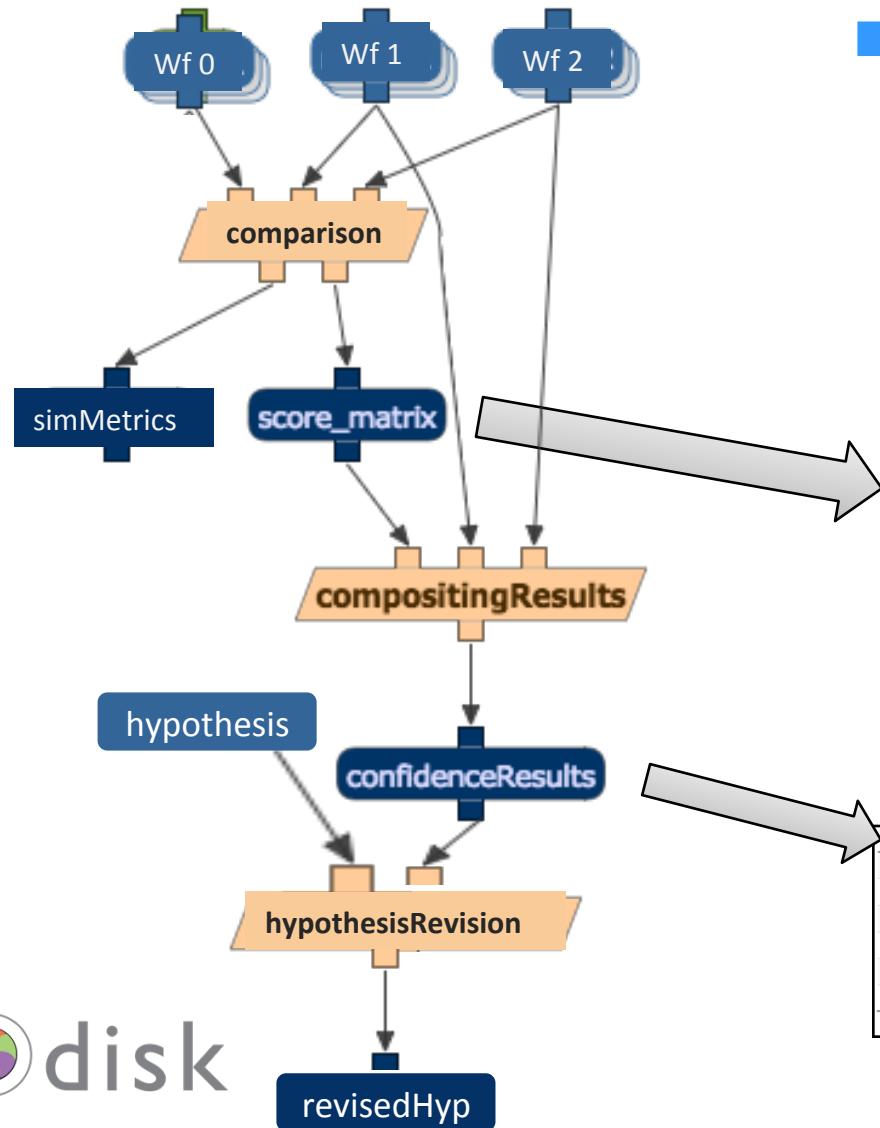
Capturing Data Analysis Strategies through *Lines of Inquiry*



Knowledge about Multi-Omics Data Analysis Captured in Workflows



Knowledge about Evidence Aggregation Captured in Meta-Workflows



- After running the workflows, meta-workflows analyze their results and generate a combined confidence value

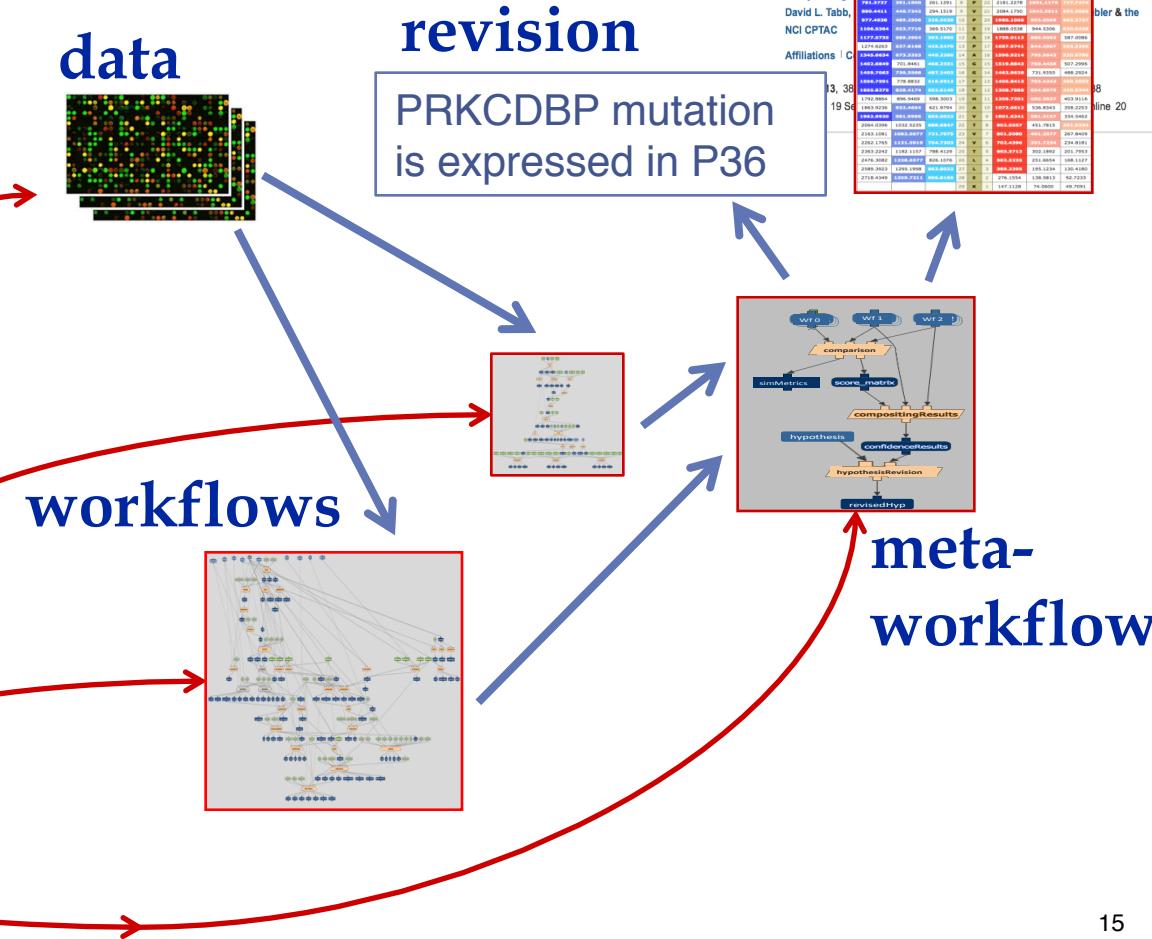
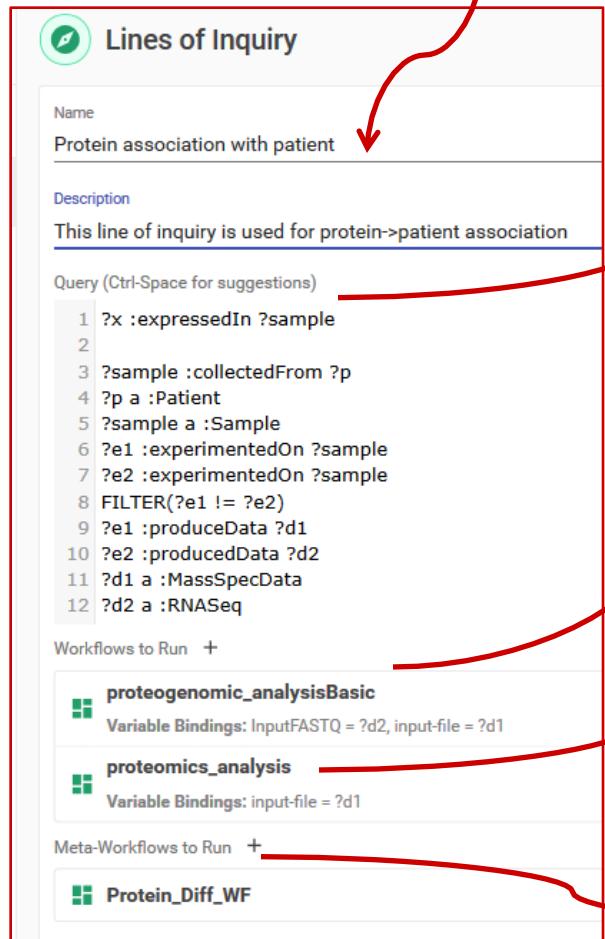
| | Workflow 0 | Workflow 1 | Workflow 2 |
|------------|------------|------------|------------|
| Workflow 0 | 1.00 | 0.96 | 0.72 |
| Workflow 1 | 0.96 | 1.00 | 0.84 |
| Workflow 2 | 0.72 | 0.84 | 1.00 |

| Scan | Peptide | Peptide Prophet | Protein | Confidence |
|------|---|-----------------|--|------------|
| 9226 | R.E.SA.LE.PG.PV.PEA.PAG.GPV.HAV.TV.VILLE.KL | 0.9958 | NP_659477_r2682123-R&P2s1051992.L158P | 0.9589 |
| 9401 | K.A.EY.LAS.IF.GT.EK.D | 0.9797 | NP_659424_W95G | 0.9998 |
| 7556 | K.I.SNP.WQSPSGTLPAL.R.T | 0.9975 | NP_002446_677:AGG | 0.972 |
| 8699 | R.S.DP.YTL.NV.LY.GP.DV.PT.SPSK.A | 0.9975 | NP_002474_T100L.G239W | 0.972 |
| 9614 | K.L.FM.VD.SIP.K.V | 0.874 | NP_659403_r7041710.1863V2s10961700.V1502M.S1716L.A1784G.r10961689.Q2143P | 0.468 |
| 9282 | M.A.EY.LAS.IF.GT.EK.D.K.V | 0.802 | NP_659424_W95G | 0.9998 |

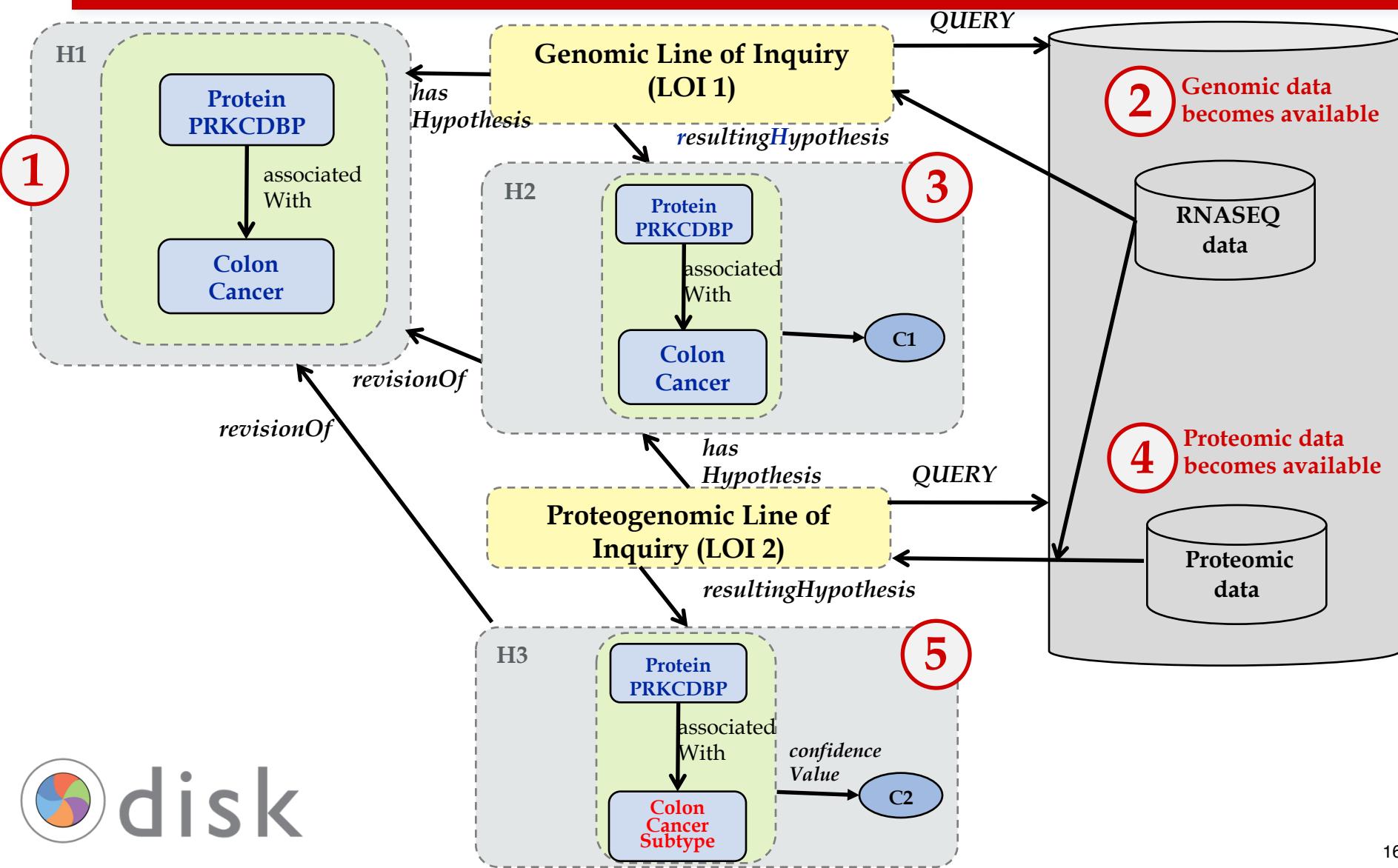
Knowledge about Data Analytic Strategies Captured in Lines of Inquiry

hypothesis

Protein PRKCDBP is expressed in samples of patient P36



Representing Hypothesis Evolution as New Data is Available and Analyzed



Addressing Problems with Current Science Data Analysis Practice



- Integrative studies of multi-source data are rare
 - Requires collaborations to cover all the specialized expertise needed, so they take years

“Automated Hypothesis Testing with Large Scientific Data Repositories.”

Y. Gil, D. Garijo, V. Ratnakar, R. Mayani, R. Adusumilli, H. Boyce, and P. Mallick.
Proceedings of the 4th Annual Conference on Advances in Cognitive Systems (ACS), 2016.

- Data resources are constantly growing
 - Scientists rarely re-evaluate evidence for hypotheses as new data becomes available

“Towards Continuous Scientific Data Analysis and Hypothesis Evolution.”

Y. Gil, D. Garijo, V. Ratnakar, R. Mayani, R. Adusumilli, H. Boyce, A. Srivastava and P. Mallick. *To appear in the Proceedings of the 31st Annual Conference of the Association for the Advancement of Artificial Intelligence (AAAI)*, 2017.

Representation Challenges in Hypothesis-Driven Discovery from Scientific Data Repositories

- A domain-independent framework that works across scientific disciplines
- Hypothesis representation and evolution
- Meta-reasoning strategies to select and prioritize analyses
- Aggregation of evidence from multiple types of data/observations
- Generation of explanations from provenance records
- Learning to improve from user feedback
- Identifying “interestingness” in results
- Designing new data analysis methods
- Incorporating science knowledge and theories to guide hypothesis formation

Outline

- Artificial intelligence and scientific discovery
 - The knowledge systems tradition
- Our recent work on capturing knowledge about data analysis strategies
 - Hypothesis-driven data analysis
- Representing and capturing data analysis knowledge
 - About data, software, methods, meta-analysis
- Summary of AI challenges

Representing and Capturing Scientific Knowledge

Data



Software



Workflows



Meta-Workflows



Provenance

W3C® PROV

OPMW Workflow repository

Knowledge about Data: Linked Earth

Work with Julien-Emile Geay of USC and Nick McKay of NAU

Palmyra Atoll

Structured Properties

| | | |
|-------------------------------|--------------------------------------|--------|
| main type (GND) | geographical feature | [edit] |
| is in the administrative unit | United States Minor Outlying Islands | [edit] |

Porites

Structured Properties

[add fact]

[x] Property:Name Topic:Finger Coral [hide]

- [\[x\] http://dbpedia.org/resource/Porites](#)
- [\[add source\]](#)

Wikipedia Entry  [go to original Wikipedia article](#)

Porites is a genus of stony coral; they are SPS (Small Polyp Stony) corals. They are characterised by a finger-like morphology. Members of this genus have widely spaced



Geochemistry datasets

| | Archive | Interpretation | MeasurementMaterial | MeasurementStandard | MeasurementUnits |
|-------------------|---------------|----------------|---------------------|---------------------|------------------|
| Lake Bosumtwi | LakeSediments | Lake Level | Authigenic Calcite | VPDB | Permil |
| Quelccaya | IceCore | | Ice | VSMOW | Permil |
| Palmyra coral 20C | Coral | SST,SSS | Skeletal aragonite | VPDB | Permil |



Palmyra coral 20C

Data

- **DOWNLOAD**

From: <http://www.ncdc.noaa.gov/paleo/metadata/noaa-coral-1865.html>

Structured Properties

| | | |
|--------------------------------|----------------------|-------------|
| [x] SiteName | Palmyra | (By Julien) |
| [x] Archive | Coral | (By Julien) |
| [x] Domain(s) | Climate,geochemistry | (By Julien) |
| [x] Forward model | 10.1029/2011GL048224 | (By Julien) |
| [x] Genus | Porites | (By Julien) |
| [x] Interpretation | SST,SSS | (By Nick) |
| [x] Measurement | | |
| [x] MeasurementMaterial | | |
| [x] MeasurementStandard | | |
| [x] MeasurementUnits | | |
| [x] Reference | | |
| [x] Species | | |

AI opportunities:

- collection
- normalization
- organization

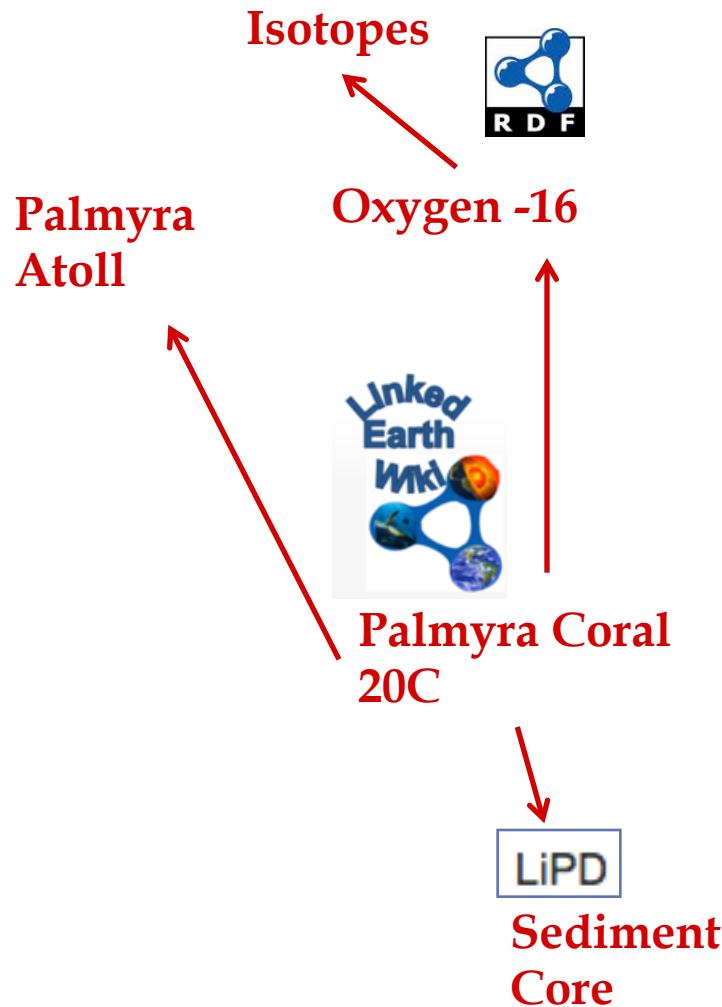
Credits

Users who have contributed to this Page:

- [Julien \(43 Edits\)](#)
- [Nick \(34 Edits\)](#)



Linked Data and Entities: Semantic Web Objects as RDF + URIs



Knowledge about Software: OntoSoft

*Work with C. Duffy (PSU), C. Mattmann (JPL),
S. Peckham (CU), and E. Robinson (ESIP)*

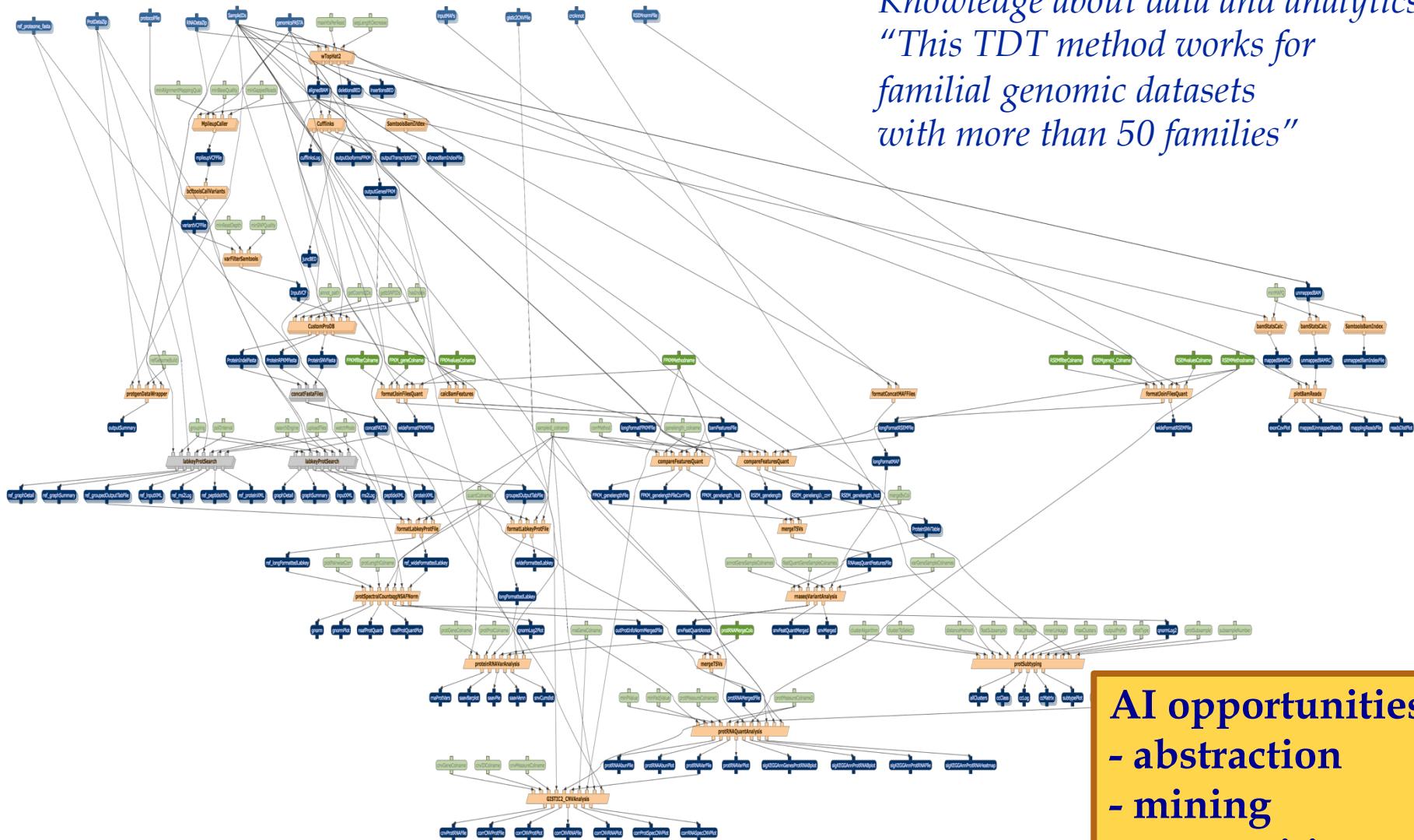
Compare Software

DrEICH algorithm, PIHM, PIHMGIS, TauDEM, WBMsed

| PIHM | PIHMGIS | DrEICH | TauDEM | WBMsed |
|---|---|---|---|---|
|  |  |  |  |  |
| What are domain specific keywords for this software ? (eg: hydrology, climate) | | | | |
| Geomorphology, Hydrological, Bedrock channel ero- | Basins, Continental | Basins, GIS | Hydrologically corrected DEM, Watershed | Sediment flux, Global land |
| What Operating Systems can the software run on ? | | | | |
| Unix Linux | Unix Windows Linux Mac OS | Unix Windows Linux Mac OS | Unix Windows Linux Mac OS | |
| Is there any test data available for the software ? | | | | |
| Test Data Location: http://onlinelibrary.wiley.com/doi/10.1002/2013WR015167/full | Test Data Location: http://source-forg.net/projects/pihmmodel/ | | Test Data Location: http://csdms.colorado.edu/wiki/Model:TauDEM#Testing | Test Data Location: http://csdms.colorado.edu/wiki/Model:WBMsed#Testing |
| Test Data Description: Two test DEMs are included in the repository, | Test Data Description: Upper Juniata River 875 km^2: see: http://source-forg.net/projects/pihmmodel/ | | Test Data Description: The Logan River DEM is a small test dataset useful | Test Data Description: Extensive input dataset is available on the CSDMS |

AI opportunities:
- functional descriptions
- organization
- linking to data

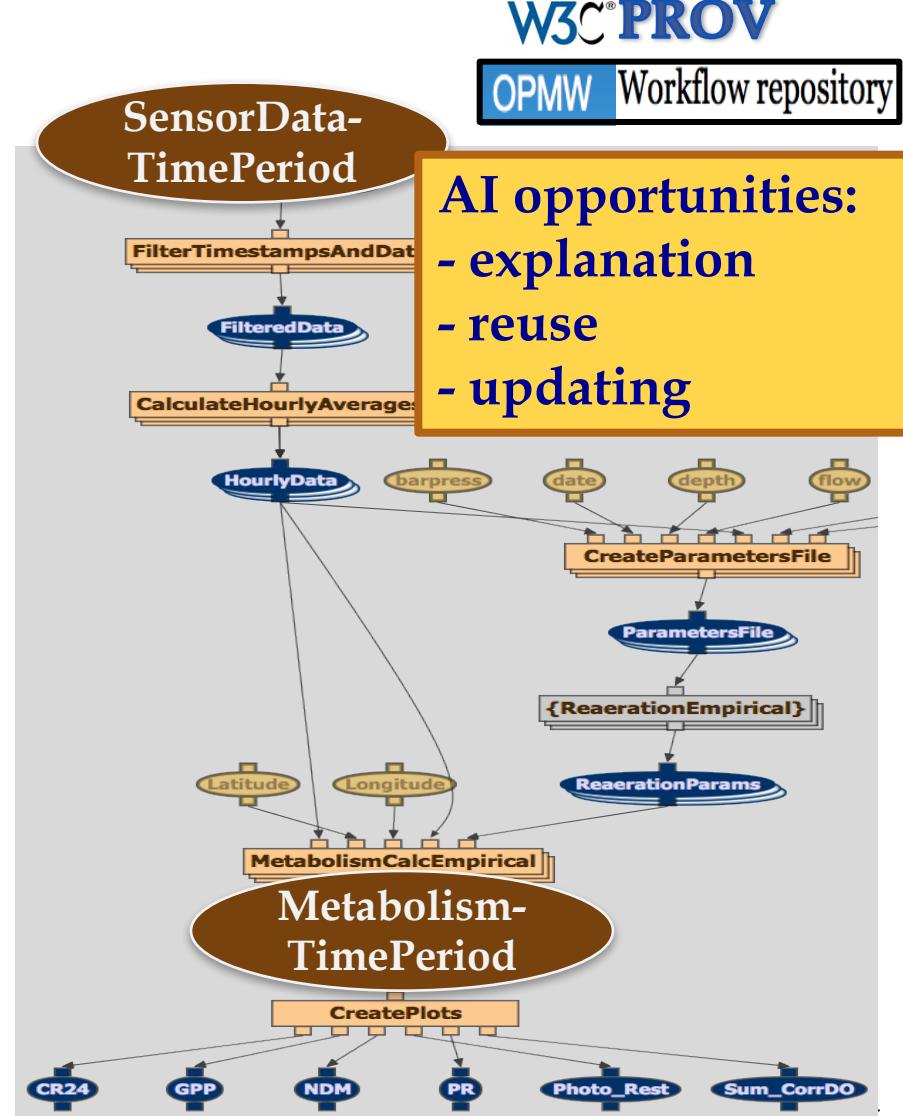
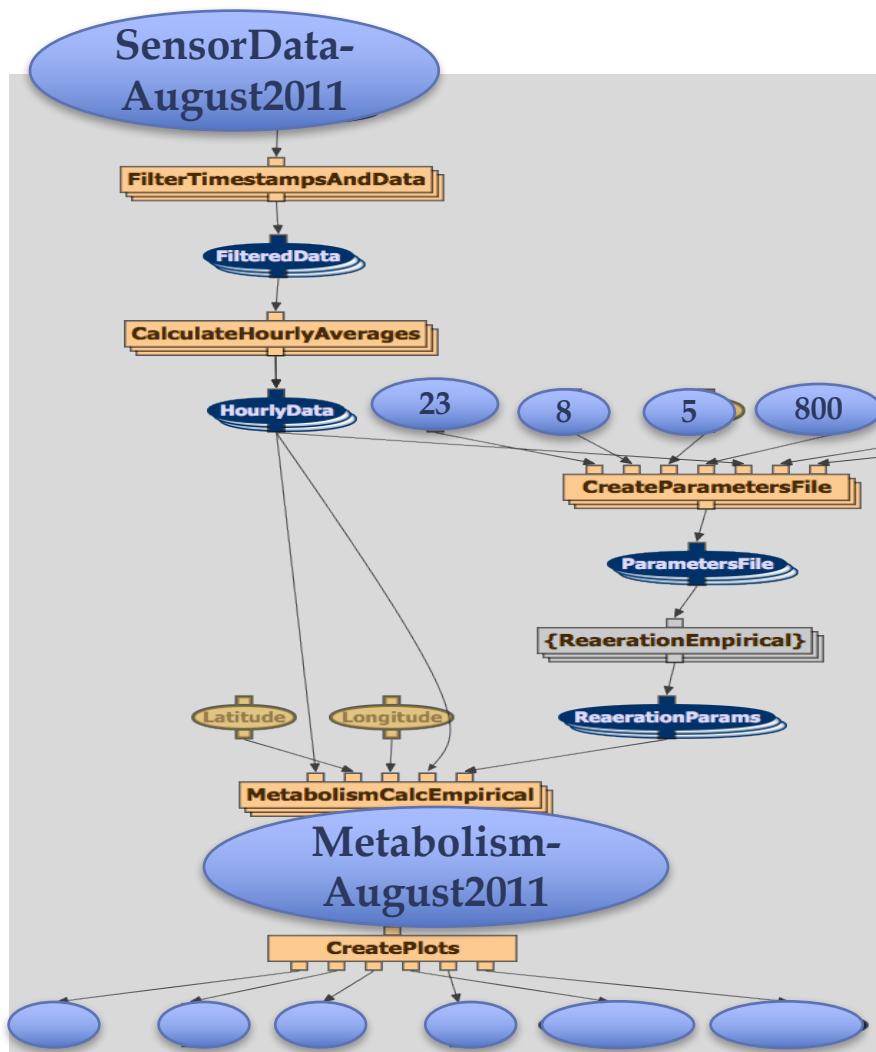
Knowledge About Analytic Methods: WINGS



*Knowledge about data and analytics:
“This TDT method works for
familial genomic datasets
with more than 50 families”*

- AI opportunities:
 - abstraction
 - mining
 - composition

Knowledge About Provenance: W3C PROV



Knowledge about Meta-Processes: Organic Data Science



Organic Data Science

1

2

3 All Tasks | My Tasks (36)

4 computer science search

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

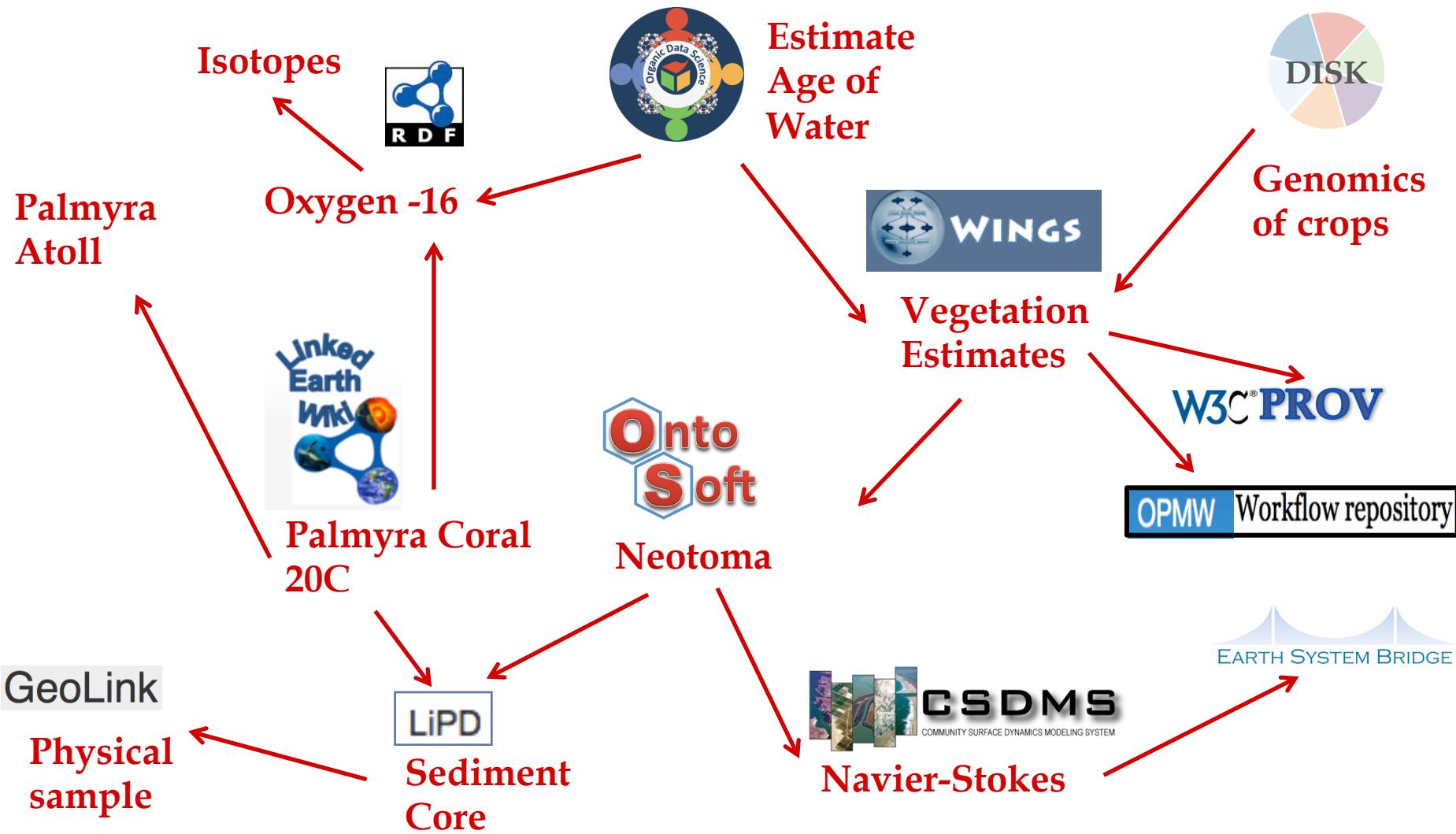
110

111

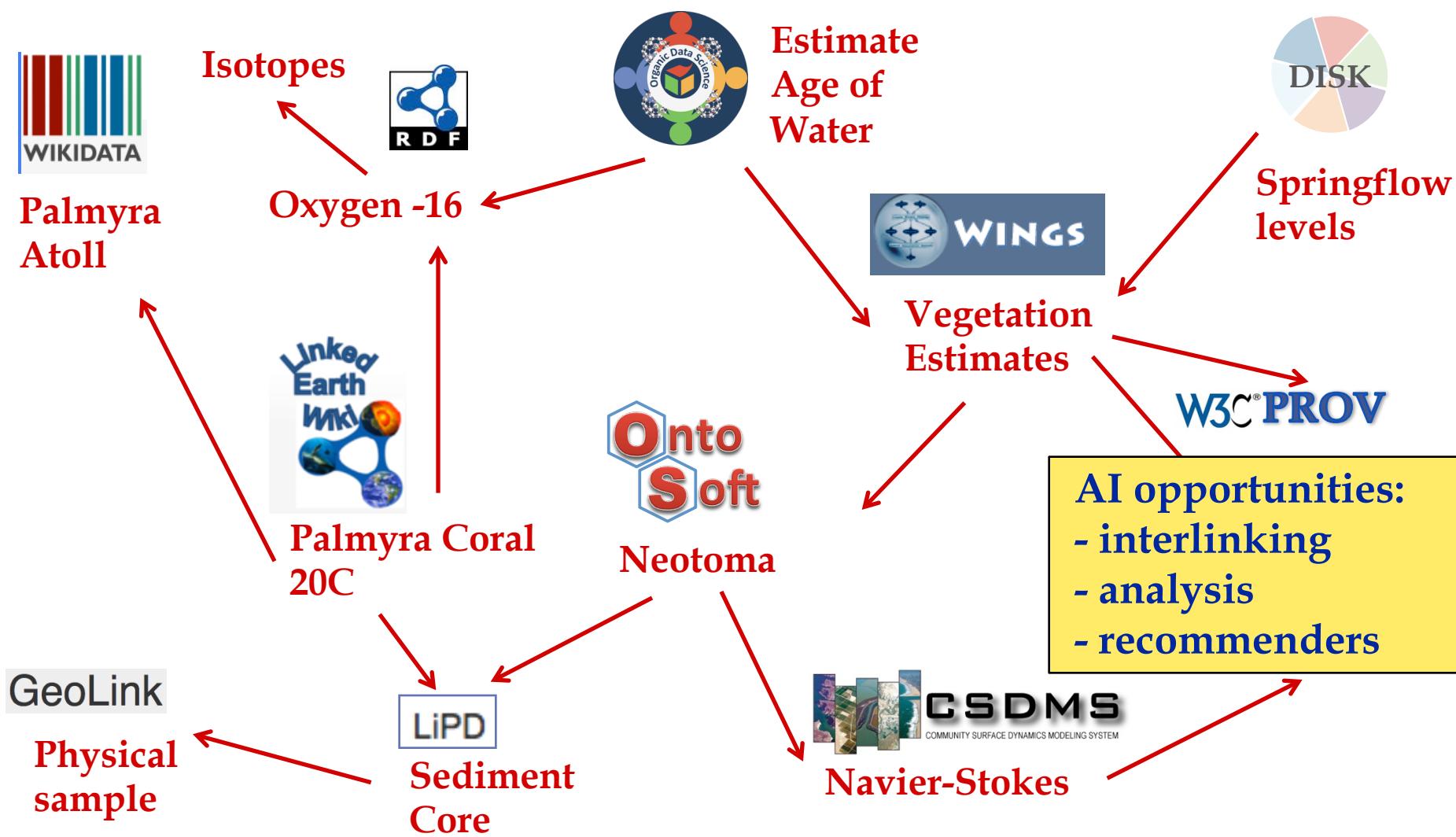
112

113

Capture and Interlink Scientific Knowledge



Linked Data and Linked Knowledge



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)

Outline

- Artificial intelligence and scientific discovery
 - The knowledge systems tradition
- Our recent work on capturing knowledge about data analysis strategies
 - Hypothesis-driven data analysis
- Representing and capturing data analysis knowledge
 - About data, software, methods, meta-analysis
- Summary of AI challenges

Artificial Intelligence and Scientific Discovery: Challenges for Knowledge Systems

Knowledge Representation Challenges

Metadata:

- collection
- normalization
- organization

Software:

- functional descrs.
- organization
- linking to data

Methods:

- abstraction
- mining
- composition

Provenance:

- explanation
- reuse
- updating

Meta-analysis:

- collaboration
- group formation
- community health

Problem Solving Challenges

Hypothesis-driven discovery:

- A domain-independent framework across scientific disciplines
- Hypothesis representation and evolution
- Capturing analytic knowledge
- Meta-reasoning strategies to select and prioritize analyses
- Aggregation of evidence from multiple types of data/observations
- Generation of explanations from provenance records
- Learning to improve from user feedback
- Identifying “interestingness” in results
- Designing new data analysis methods
- Incorporating science models to guide hypothesis formation