# Building on Existing Communities: the Virtual Astronomical Observatory (and NIST)

Robert Hanisch

Space Telescope Science Institute

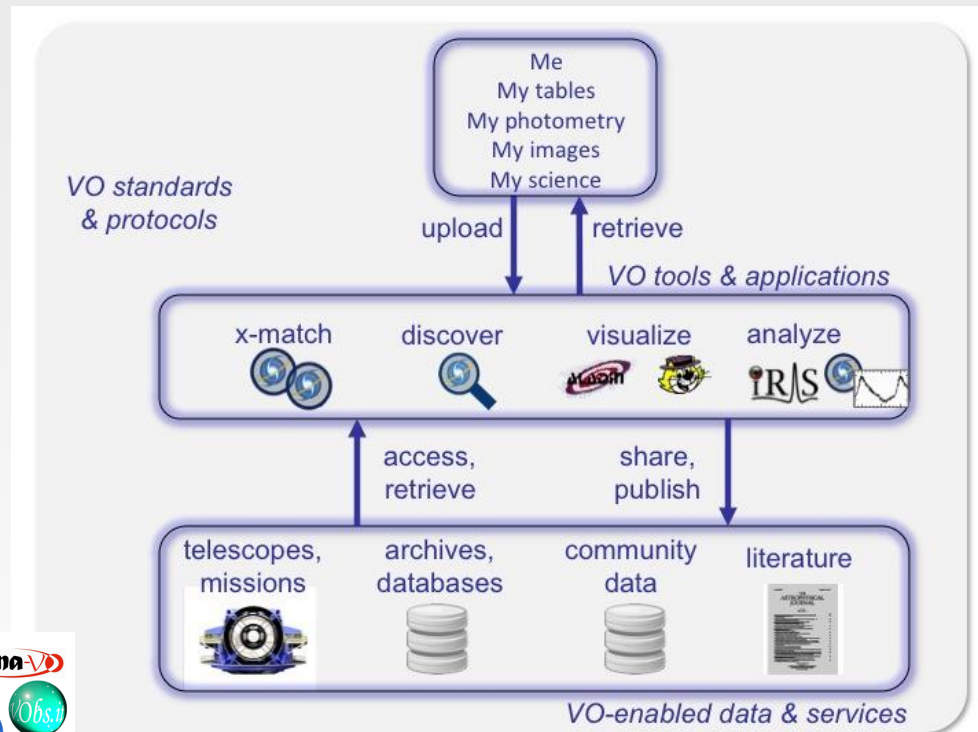Director, Virtual Astronomical Observatory

# Data in astronomy

- ~70 major data centers and observatories with substantial on-line data holdings
- ~10,000 data "resources" (catalogs, surveys, archives)
- Data centers host from a few to ~100s TB each, currently at least 2 PB total
- Current growth rate ~0.5 PB/yr, increasing
- Current request rate ~1 PB/yr
- Future surveys will increase data rates to PB/day
  - "For LSST, the telescope is a peripheral to the data system" (T. Tyson)

*How do astronomers navigate through all of this data?*

# The Virtual Observatory

*The VO is a data discovery, access, and integration facility*

- Images, spectra, time series
- Catalogs, databases
- Transient event notices
- Software and services
- Application inter-communication
- Distributed computing
  - authentication, authorization, process management
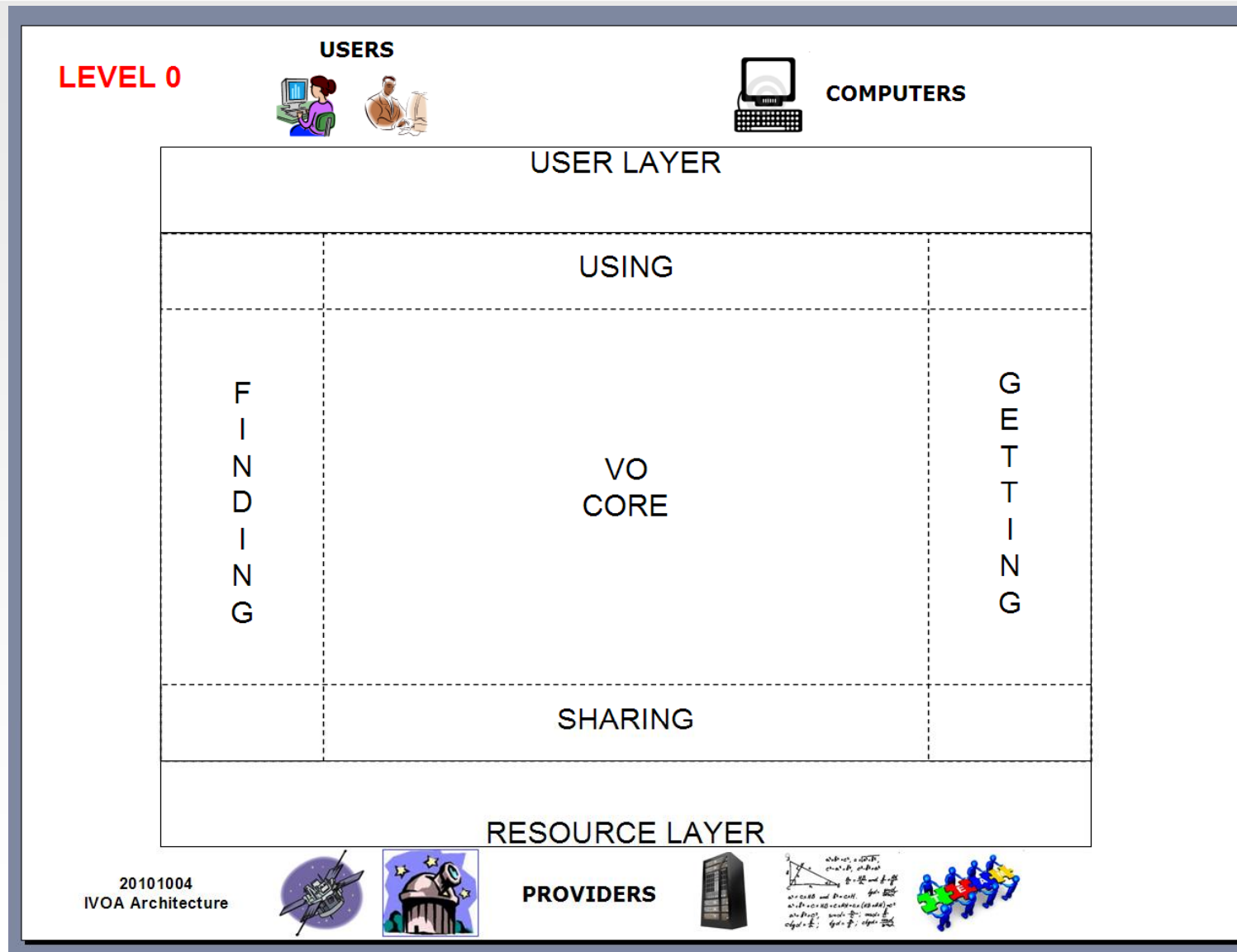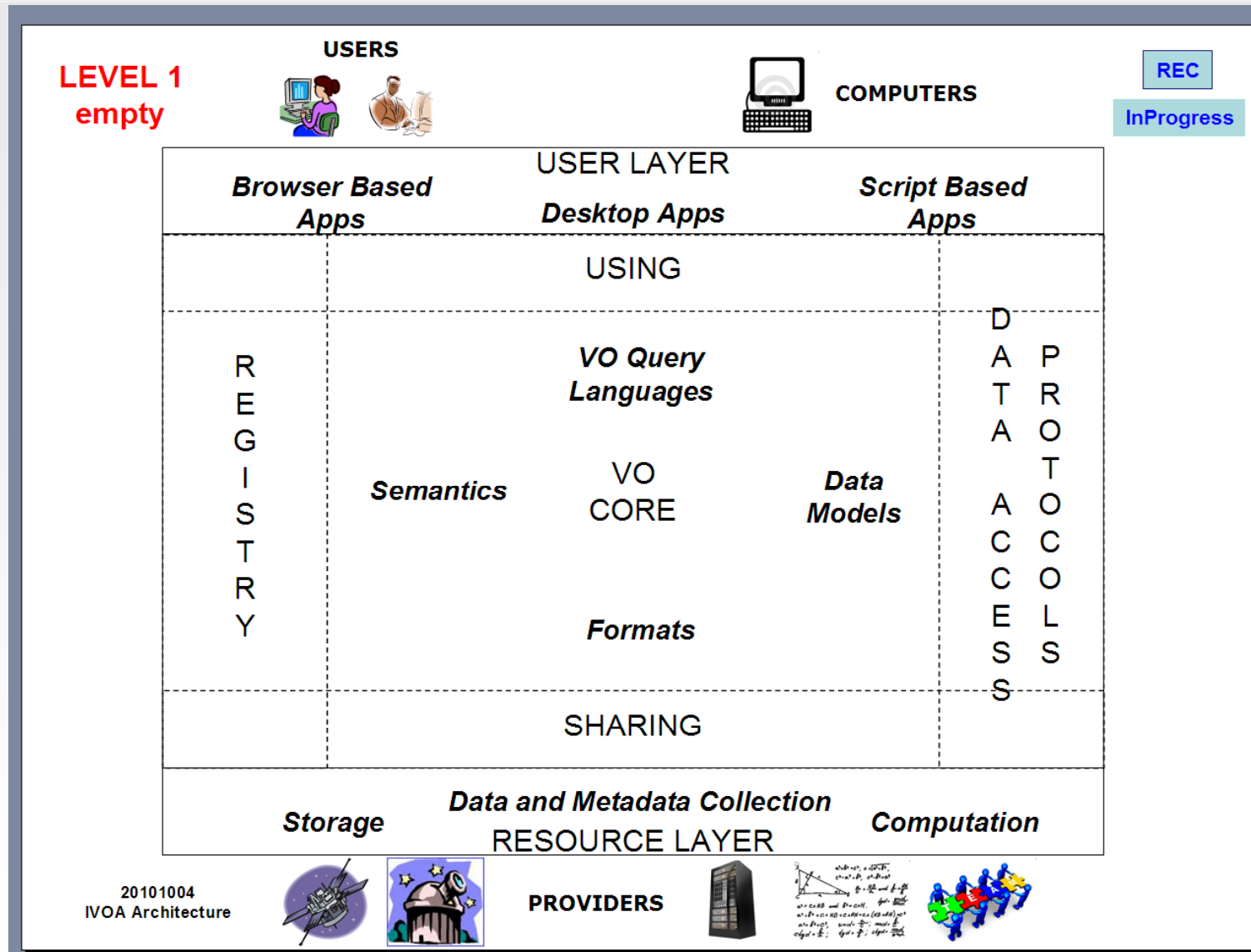- International coordination collaboration IVOA W3C)



VO standards & protocols

Me
My tables
My photometry
My images
My science

upload    retrieve

VO tools & applications

x-match    discover    visualize    analyze

access, retrieve    share, publish

telescopes, missions    archives, databases    community data    literature

VO-enabled data & services

# Virtual Observatory capabilities

- Data exchange / interoperability / multi-λ (co-observing)

  *Data Access Layer (SIAP, SSAP / time series)*

- Query and cross-match across distributed databases

  *Cone Search, Table Access Protocol*

- Remote (but managed) access to centralized computing and data storage resources

  *VOSpace, Single-Sign-On (OpenID), SciDrive*

- Transient event notification, scalable to $10^6$ messages/night

  *VOEvent*

- Data mining, characterization, classification, statistical analysis

  *VOStat, Data Mining and Exploration toolkit*

# VO architecture



LEVEL 0

USERS

COMPUTERS

USER LAYER

USING

FINDING

VO CORE

GETTING

SHARING

RESOURCE LAYER

20101004
IVOA Architecture

PROVIDERS

# VO architecture



LEVEL 1
empty

USERS

COMPUTERS

REC

InProgress

USER LAYER

Browser Based
Apps

Desktop Apps

Script Based
Apps

USING

R E G I S T R Y

Semantics

VO Query
Languages

VO
CORE

Formats

Data
Models

D A T A A C C E S S

P R O T O C O L S

SHARING

Storage

Data and Metadata Collection

RESOURCE LAYER

Computation

20101004
IVOA Architecture

PROVIDERS

# VO architecture

# Key to discovery:  Registry

- Used to discover and locate *resources*—data and services—that can be used in a VO application

- Resource:  anything that is describable and identifiable.
  – Besides data and services:  organizations, projects, software, standards

- Registry: a list of resource descriptions
  – Expressed as structured metadata in XML

    to enable automated processing and searching

    Metadata based on Dublin Core

# Registry framework



harvest

(pull)

Full
Searchable
Registry

replicate

Local
Publishing
Registry

OAI/PMH

Full
Searchable
Registry

Data
Centers

Local
Publishing
Registry

search
queries

Users,
applications

# Data discovery

# Data discovery

# SciDrive: astro-centric cloud storage



**Automatic Metadata Extraction**

Extract tabular data from:
- CSV
- FITS
- TIFF
- Excel

Extract metadata from:
- FITS
- Image files (TIFF, JPG)

Automatically upload tables into relational databases:
- CasJobs/MyDB
- SQLShare

*Controlled data sharing*
*Single sign-on*
*Deployable as virtual machine*

# The VO concept elsewhere

- Space Science
  - Virtual Heliophysics Observatory (HELIO)
  - Virtual Radiation Belt Observatory (ViRBO)
  - Virtual Space Physics Observatory (VSPO)
  - Virtual Magnetospheric Observatory (VMO)
  - Virtual Ionosphere Thermosphere Mesosphere Observatory (VITMO)
  - Virtual Solar-Terrestrial Observatory (VSTO)
  - Virtual Sun/Earth Observatory (VSEO)
- Virtual Solar Observatory
- Planetary Science Virtual Observatory
- Deep Carbon Virtual Observatory
- Virtual Brain Observatory

# Data management at

- I move to NIST 7/28/2014 as Director, Office of Data and Informatics, Material Measurement Laboratory
  - Materials science, chemistry, biology
  - Materials Genome Initiative
- Foster a culture of data management, curation, re-use in a bench-scientist / PI-dominated organization having a strong record of providing "gold standard" data
- Inward-looking challenges
  - Tools, support, advice, common platforms, solution broker
  - Big data, lots of small/medium data
- Outward-looking challenges
  - Service directory
  - Modern web interfaces, APIs, better service integration
  - Get better sense of what communities want from NIST
- Define standards, standard practices
- Collaboration: other government agencies, universities, domain repositories

# NDS and domain repositories

- Domain repositories are discipline-specific
- Various business models in use; long-term sustainability is a major challenge*
- Potential NDS roles
  - Customizable data management and curation tools built on a common substrate
  - Access to cloud-like storage but at non-commercial rates
  - A directory of ontology-building and metadata management tools
  - A directory of domain repositories
  - Accreditation services
  - Advice, referral services, "genius bar"

* "Sustaining Domain Repositories for Digital Data: A White Paper," C. Ember & R. Hanisch, eds. http://datacommunity.icpsr.umich.edu/sites/default/files/WhitePaper_ICPSR_SDRDD_121113.pdf

# Technologies/standards to build on

- Just use the VO standards!
  - OK, seriously…  NIH syndrome
  - Much could be re-used in terms of architecture
  - Generic, collection-level metadata
- Cross-talk with Research Data Alliance (ANDS, EUDAT)
  - Data Citation WG
  - Data Description Registry Interoperability WG
  - Data Type Registries WG
  - Domain Repositories IG
  - Long Tail of Research Data IG
  - Metadata IG
  - Metadata Standards Directory WG
  - Preservation e-Infrastructure WG
  - and others…

*Dataverse, Dryad, iRODS, DSpace, etc.*

# Lessons learned re/ federation

- It takes more time than you think
  - Community consensus requires buy-in early and throughout
- Top-down imposition of standards likely to fail
- Balance requirements coming from a research-oriented community with innovation in IT
- Marketing is very important
  - Managing expectations
  - Build it, and they might come
- Coordination at the international level is essential
  - But takes time and effort
- Data models – sometimes seem obvious, more often not
- Metadata collection and curation are eternal but essential tasks

# Lessons learned re/ federation

- For example, the Cancer Biomedical Informatics Grid (caBIG) [$350M]
  - "…goal was to provide shared computing for biomedical research and to develop software tools and standard formats for information exchange."
  - "The program grew too rapidly without careful prioritization or a cost-effective business model."
  - "…software is overdesigned, difficult to use, or lacking support and documentation."
  - "*The failure to link the mission objectives to the technology shows how important user acceptance and buy-in can be.*"  (M. Biddick, Fusion PPT)

  J. Foley, *InformationWeek*, 4/8/2011

  http://www.informationweek.com/architecture/report-blasts-problem-plagued-cancer-research-grid/d/d-id/1097068?