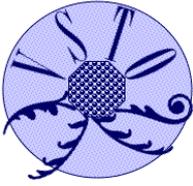


Semantically-Enabled Large-Scale Science Data Repositories



Peter Fox*

*HAO/ESSL/NCAR

**Deborah McGuinness^{\$\$}, Luca Cinquini^{%%}, Patrick West^{*},
Jose Garcia^{*}, Tony Darnell^{*}, James Benedict^{\$\$}, Don
Middleton^{%%}, Stan Solomon^{*}**

^{\$\$}McGuinness Associates

[#]Knowledge Systems and AI Lab, Stanford Univ.

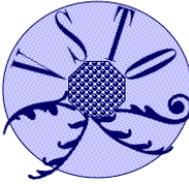
^{%%}SCD/CISL/NCAR

Also Rob Raskin, Krishna Sinha

Work funded by NSF and NASA

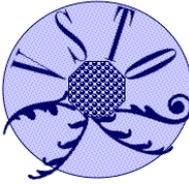


Outline



-
- Terminology and general introduction
 - Virtual Observatories and Data Integration
 - E.g. Virtual Solar-Terrestrial Observatory
 - Semantics and reasoning
 - Discussion

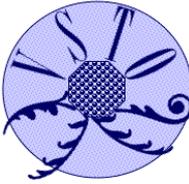
Terminology



- Workshop: A Virtual Observatory (VO) is a suite of software applications on a set of computers that allows users to uniformly find, access, and use resources (data, software, document, and image products and services using these) from a collection of distributed product repositories and service providers. A VO is a service that unites services and/or multiple repositories.
- VxOs - x is one discipline, domain, community, country
- VxyOs - x and y refer to different disciplines
- NB: VO also refers to Virtual Organization

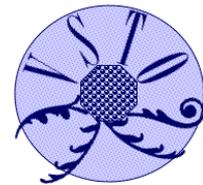


What should a VO do?

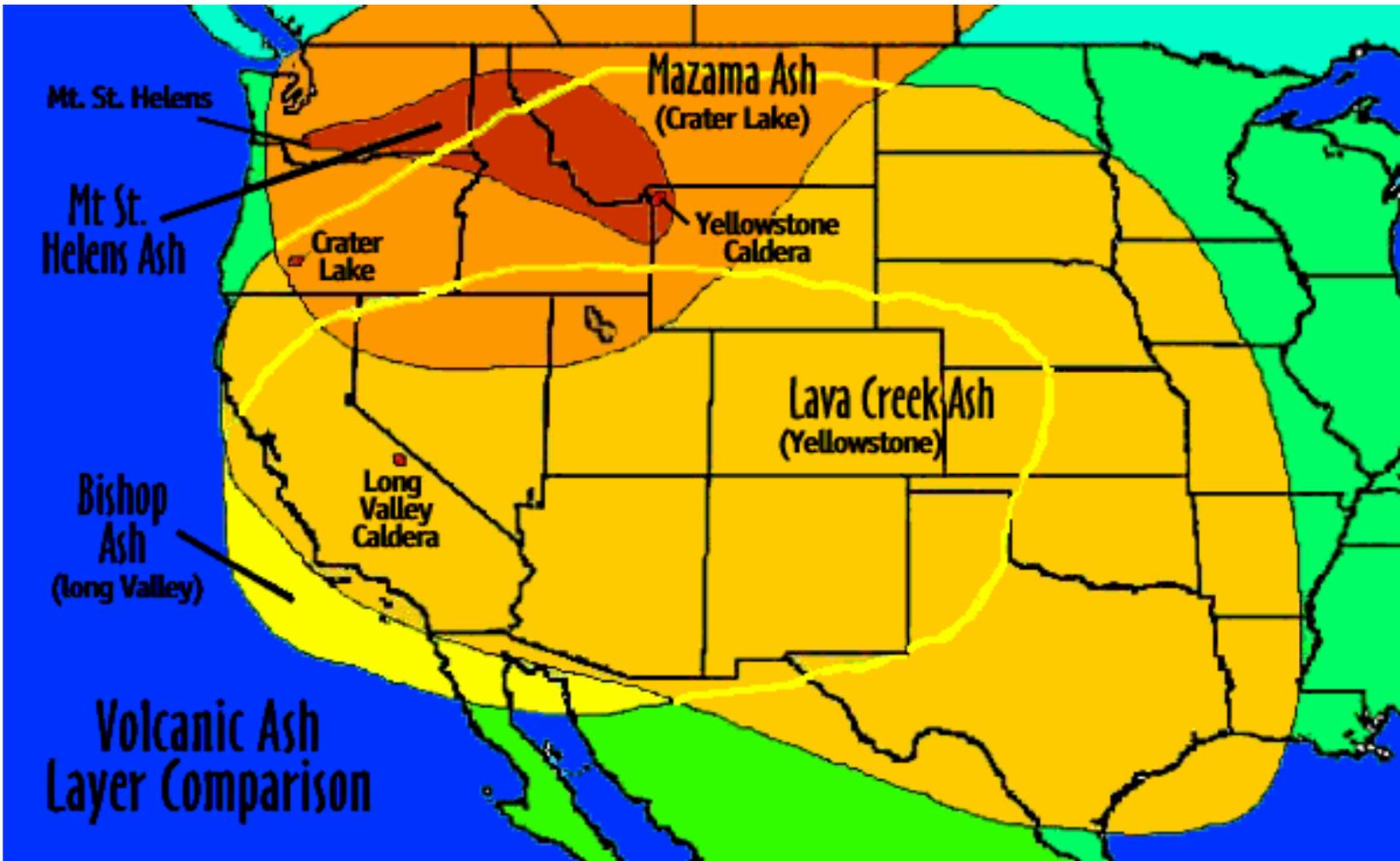


- Make “standard” scientific research much more efficient.
 - Even the principal investigator (PI) teams should want to use them.
 - Must improve on existing services (mission and PI sites, etc.). VOs will not replace these, but will use them in new ways.
 - Access for young researchers, non-experts, other disciplines,
- Enable new, global problems to be solved.
 - Rapidly gain integrated views, e.g. from the solar origin to the terrestrial effects of an event.
 - Find meaningful data related to any particular observation or model.
 - (Ultimately) answer “higher-order” queries such as “Show me the data from cases where a large coronal mass ejection observed by the Solar-Orbiting Heliospheric Observatory was also observed *in situ*.” (science-speak) or “What happens when the Sun disrupts the Earth’s environment” (general public)”

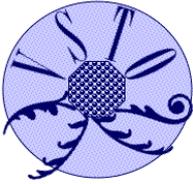
Compilation of distribution of volcanic ash associated with large eruptions. Note the continental scale ash fall associated with Yellowstone eruption ~600,000 years ago. Geologic databases provide the information about the magnitude of the eruption, and its impact on atmospheric chemistry and reflectance associated with particulate matter requires integration of concepts that bridge terrestrial and atmospheric ontologies.



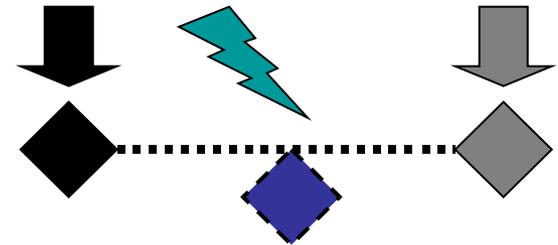
1



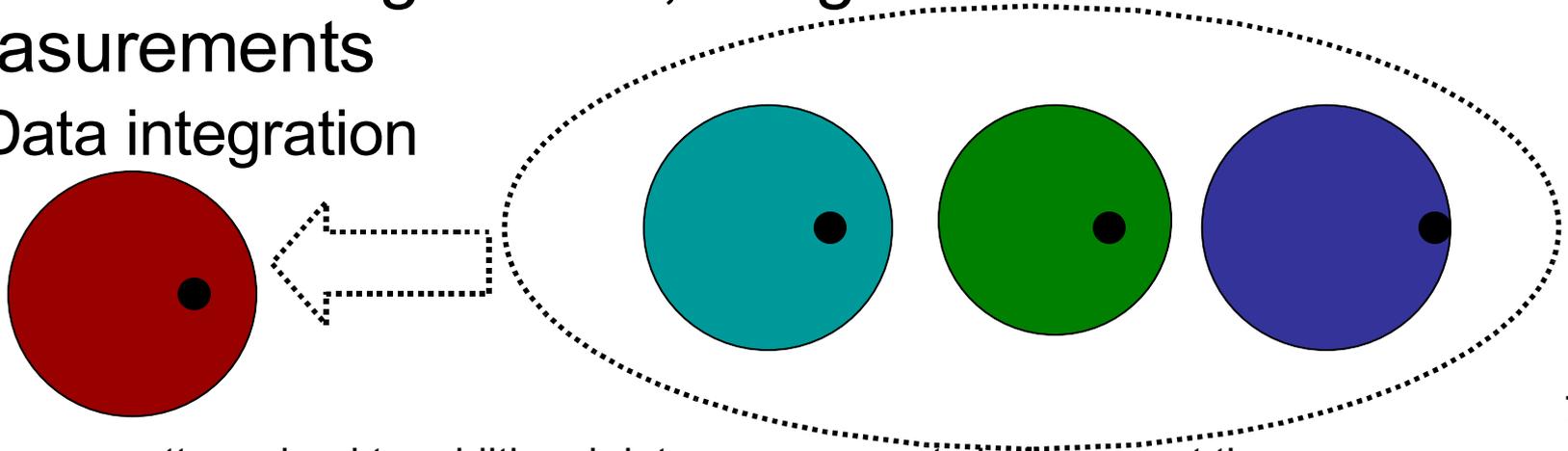
Virtual Observatory schematic



- Conceptual examples:
- In-situ: Virtual measurements
 - Sensors, etc. everywhere
 - Related measurements

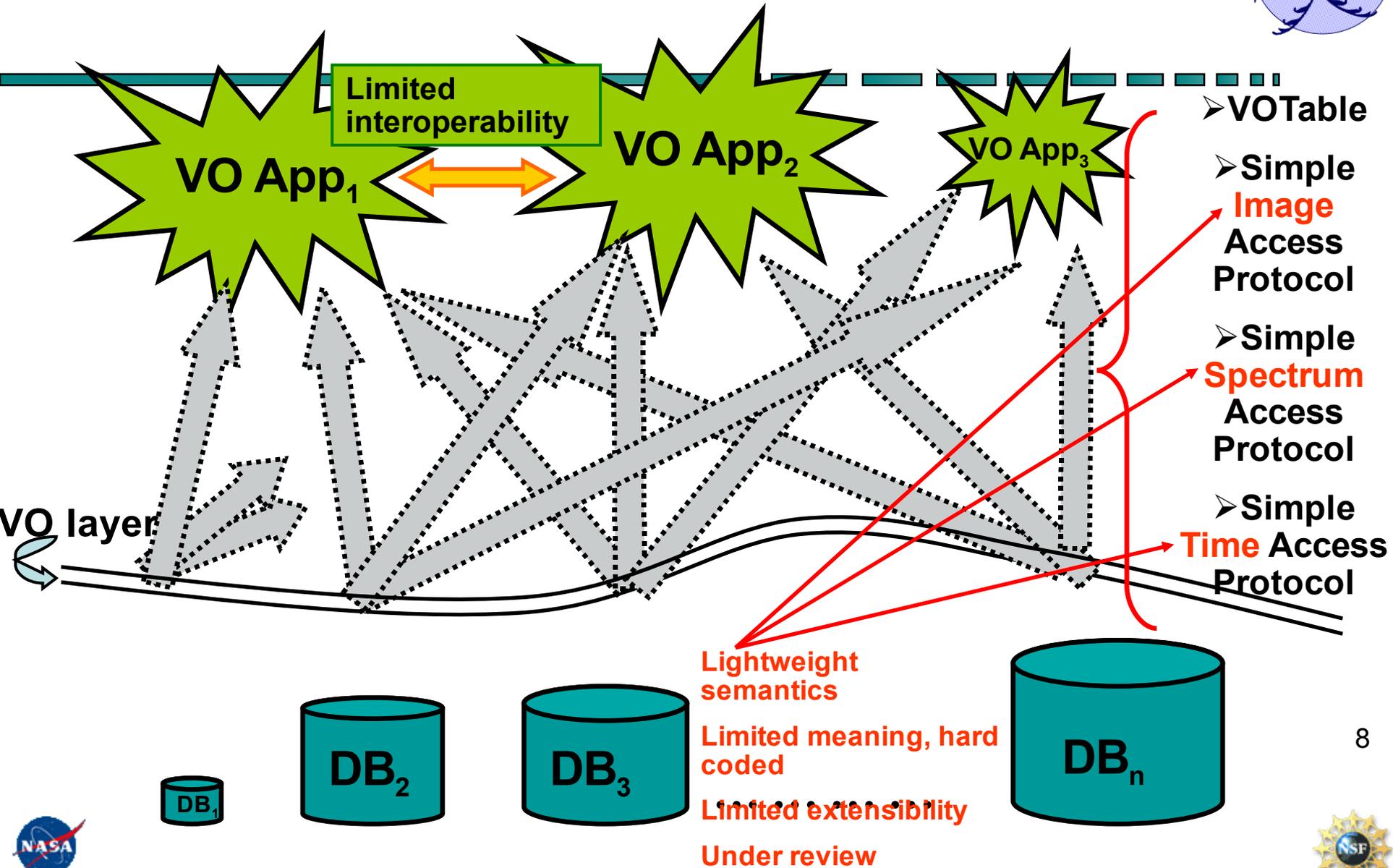
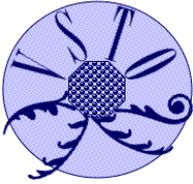


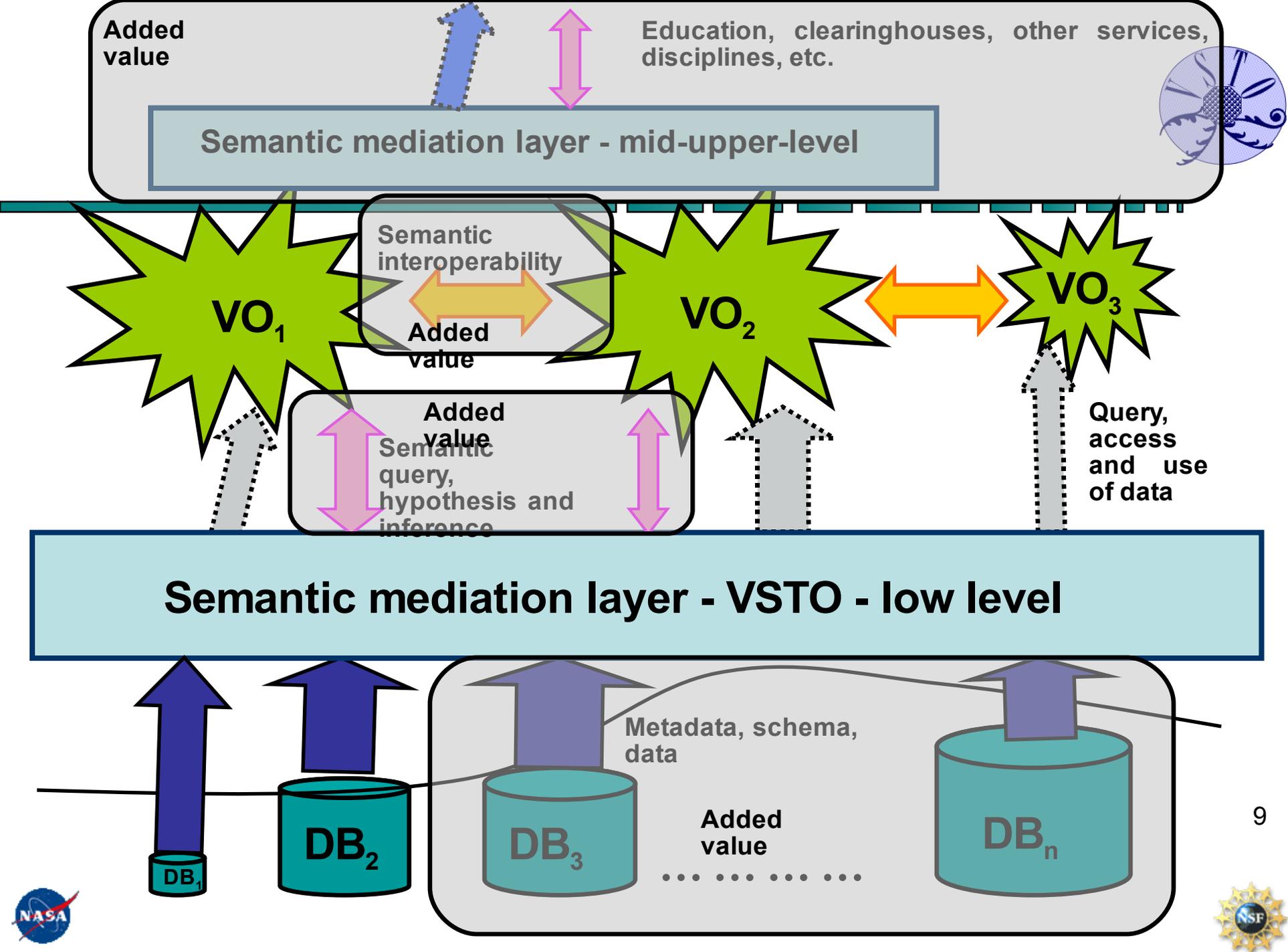
- Remote sensing: Virtual, integrative measurements
 - Data integration



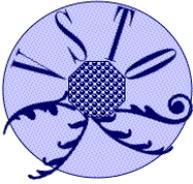
- Both usage patterns lead to additional data management challenges at the source **and** for users; now managing virtual 'datasets'

The Astronomy approach



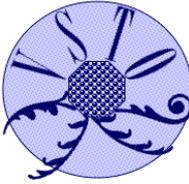


Semantic Mediation Layer



- Ontology - capturing concepts of Parameters, Instruments, Date/Time, Data Product (and associated classes, properties) and Service Classes
- Maps queries to underlying data
- Generates access requests for metadata, data
- Allows queries, reasoning, analysis, new hypothesis generation, testing, explanation, etc.

Integrative use-cases:



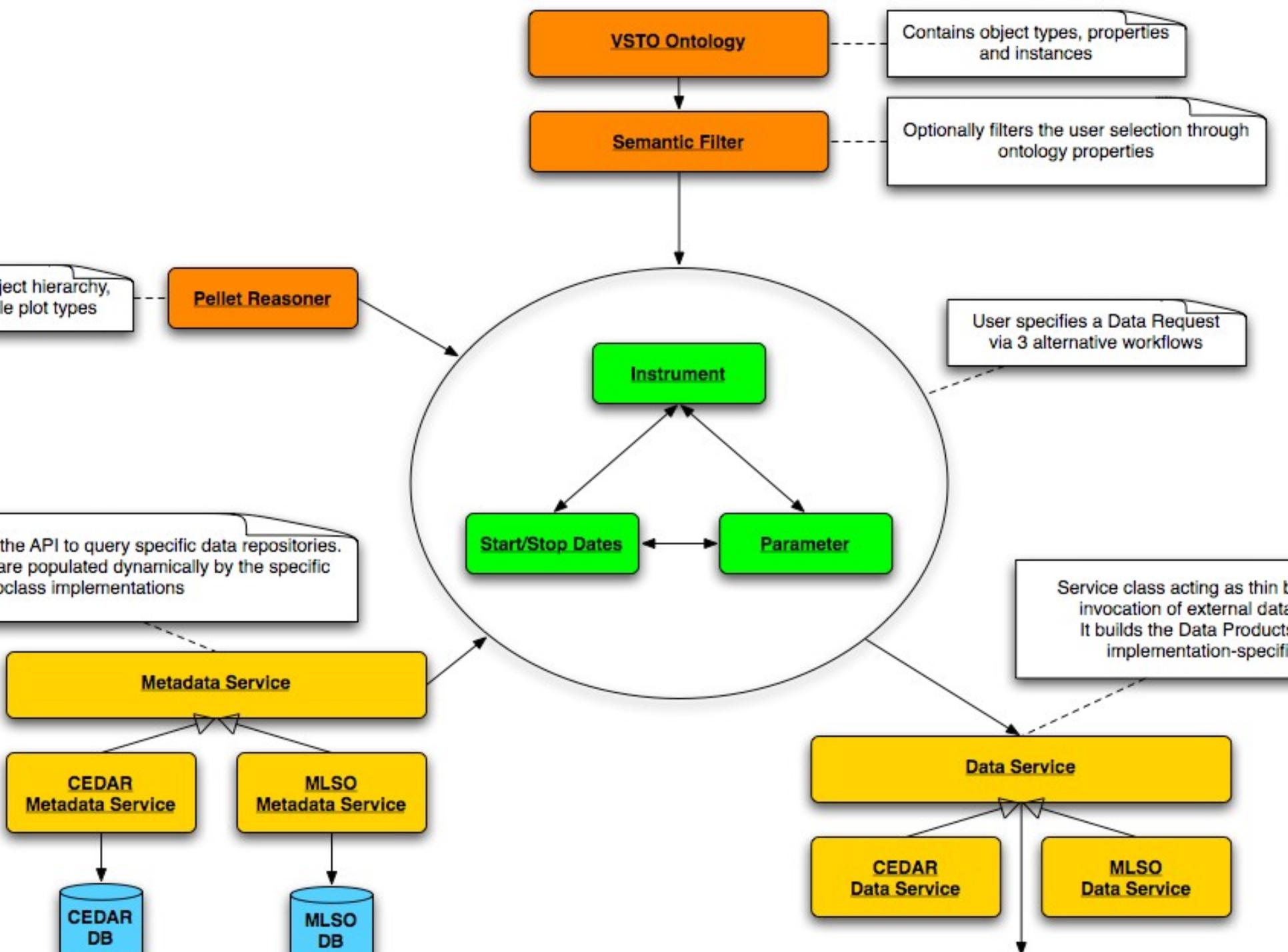
Find data which represents the state of the neutral atmosphere anywhere above 100km and toward the arctic circle (above 45N) at any time of **high geomagnetic activity.**

**Translate this into a complete query for data.
Was all the needed information recorded?**

Information needs to be inferred (and integrated) from the use-case

What is returned: Data from instruments, indices and models.







Virtual Solar Terrestrial Observatory

Home Data Communities About Us Login

Start by Instrument | Start by Dates | Start by Parameter

Semantic filtering by domain or instrument hierarchy

Data Request Summary

1. Instrument:

2. Start Date:
Stop Date:

3. Parameters:

Input Step 1 of 3: Choose Instrument

Please select an instrument.

You may filter the instruments selection by one of the following criteria:

Filter by Physical Domain:

Filter by Instrument Type:

Show Instrument Code

[?] Instrument:

- OpticalInstrument > Interferometer > FabryPerot > Arecibo P.R. Fabry-Perot [?]
- OpticalInstrument > Interferometer > FabryPerot > Millstone Hill Fabry-Perot [?]
- OpticalInstrument > Interferometer > FabryPerot > Peach Mountain Fabry-Perot [?]
- OpticalInstrument > Photometer > Chromospheric Helium Imaging Photometer [?]
- OpticalInstrument > Photometer > MK3-K Coronameter [?]
- OpticalInstrument > Photometer > MK4-K Coronameter [?]
- OpticalInstrument > Photometer > H-alpha prominence and solar disk monitor [?]
- Radar > IncoherentScatterRadar > Irkutsk Russia I.S. Radar [?]

Partial exposure of Instrument class hierarchy - users seem to LIKE THIS



Virtual Solar Terrestrial Observatory

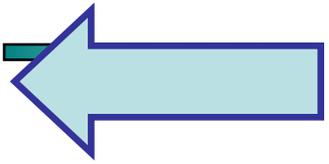
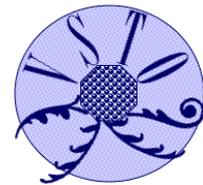
NCAR

Home | Data | Communities | About Us | Logout

Start by Instrument | Start by Dates | Start by Parameter

Data Workflow #1c

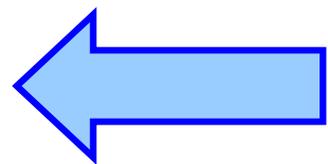
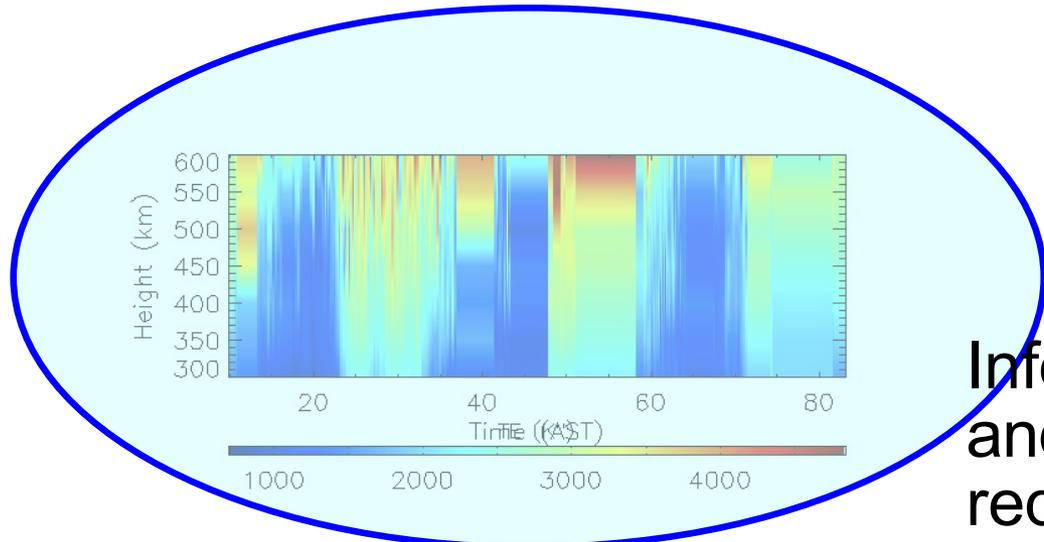
Data Request Summary		Available Output	
1. Parameter: NeutralTemperature		Data Files: ▶ STREAM [?] ▶ DAS [?] ▶ INFO [?] ▶ TAB [?] ▶ OPeNDAP [?] ▶ IDL [?] ▶ FLAT [?]	
2. Start Date: 2000/05/01 Stop Date: 2000/05/11		Data Plots: ▶ Time Series [?]	
3. Instrument: Millstone Hill Fabry-Perot		Change Input Click on the Back button to change your data selection, or Cancel to end the workflow <input type="button" value="Back"/> <input type="button" value="Cancel"/>	



Inferred plot type and return formats for data products

Instrument: 53 - Irkutsk Russia I.S. Radar
 Operating Modes:
 53/9801 - Nu Te Ti Vi
 Parameters:
 560 - te - Electron temperature
 Starting: February 09, 1999
 Ending: February 13, 1999

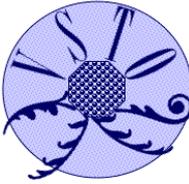
These plots are produced for visual browsing of the data and should not be used in publications without citing the data provider and CEDARWEB.



Inferred plot type and required axes data

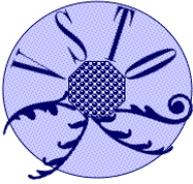


VSTO Notable progress



- Conceptual model and architecture developed by combined team; KR experts, domain experts, and software engineers
- Semantic framework developed and built with a small, cohesive, carefully chosen team in a relatively short time (deployments in 1st year)
- **Production** portal released, includes security, etc. with community migration (and so far endorsement)
- VSTO ontology version 0.4, (vsto.owl)
- Web Services encapsulation of semantic interfaces being documented
- More Solar Terrestrial use-cases to drive the completion of the ontologies - filling out the instrument ontology
- Evaluating ontologies for broader use (volcanoes, climate, ...)

What has KR done for us?



- In addition to valued added noted previously - some of which is transparent
- Reduced the need for 8 steps to query to 3 and reduced choices at each stage
- Allowed scientists to get data from instruments they never knew of before
- Allowed augmentation and validation of data
- Useful and related data provided without having to be an expert to ask for it
- Integration and use (e.g. plotting) based on inference

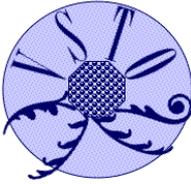
16



Ask and answer questions not possible before

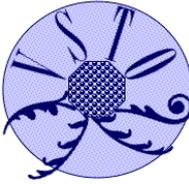


Knowledge engineering, etc.



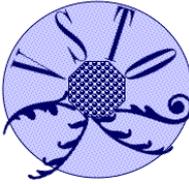
-
- Many controlled vocabularies as starting points.
 - Strive for compatibility with “best practice” controlled vocabularies, taxonomies, and ontologies.
 - Designed our own ontologies as dictated by use-case needs constantly with the goal of reusability and extensibility. (Provided VSTO modules back to at least one ontology suite with a much broader scope.)
 - Early design HIGHLY collaborative in design and implementation – critical and continued contributions from domain scientists and knowledge representation.
 - Result is fairly extensible by entire team.

Issues for Virtual Observatories



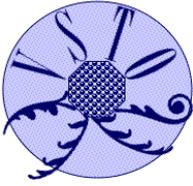
-
- Scaling to large numbers of data providers
 - Crossing disciplines
 - Security, access to resources, policies
 - Branding and attribution (where did this data come from and who gets the credit, is it the correct version, is this an authoritative source?)
 - Provenance/derivation (propagating key information as it passes through a variety of services, copies of processing algorithms, ...)
 - Data quality, preservation, stewardship, rescue
 - Interoperability at a variety of levels (~3)

Final remarks

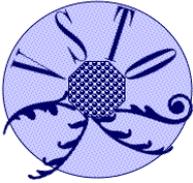


- Many geoscience VOs are in production
- Informatics efforts in Geosciences are exploding
 - GeoInformatics Town Hall at Fall AGU meeting Dec. 11 2006 in San Francisco, many cyberinfrastructure sessions
 - VO conference - April 2007 in Denver, CO
 - e-monograph to document state of VOs
 - NEW Journal of Earth Science Informatics
 - Special issue of Computers and Geosciences: “Knowledge Representation in Earth and Space Science Cyberinfrastructure”
- Ongoing activities for VOs through 2008 under the auspices of the Electronic Geophysical Year (eGY; www.egy.org)
- Contact pfox@ucar.edu, d1m@cs.stanford.edu

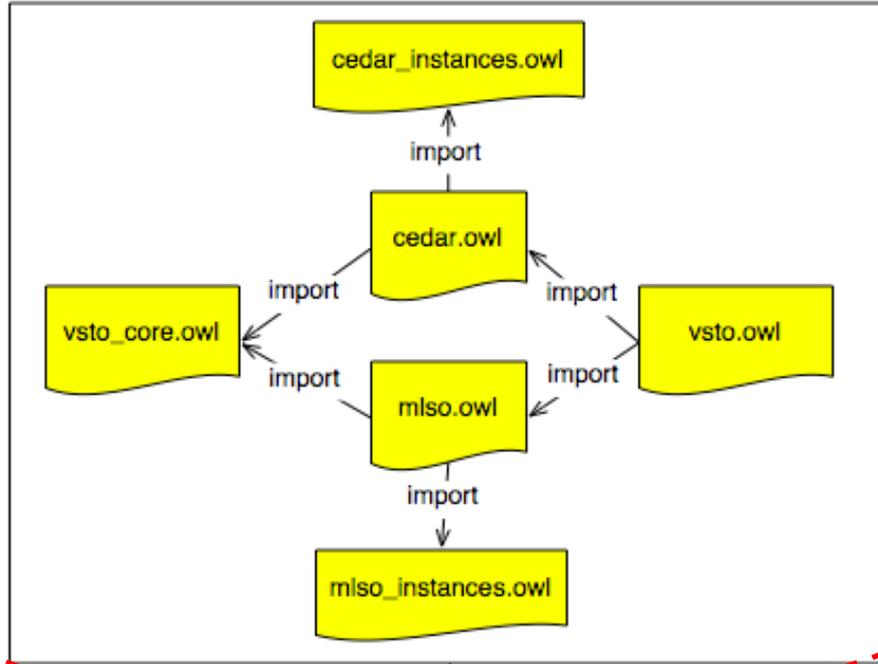
Garage



VSTO_SOFTWARE-ARCHITECTURE

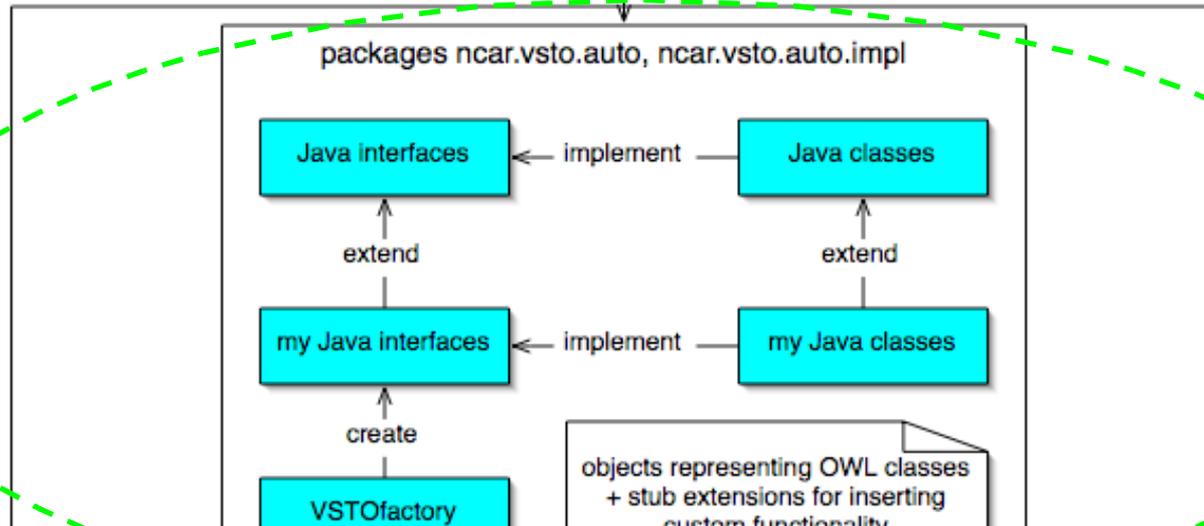


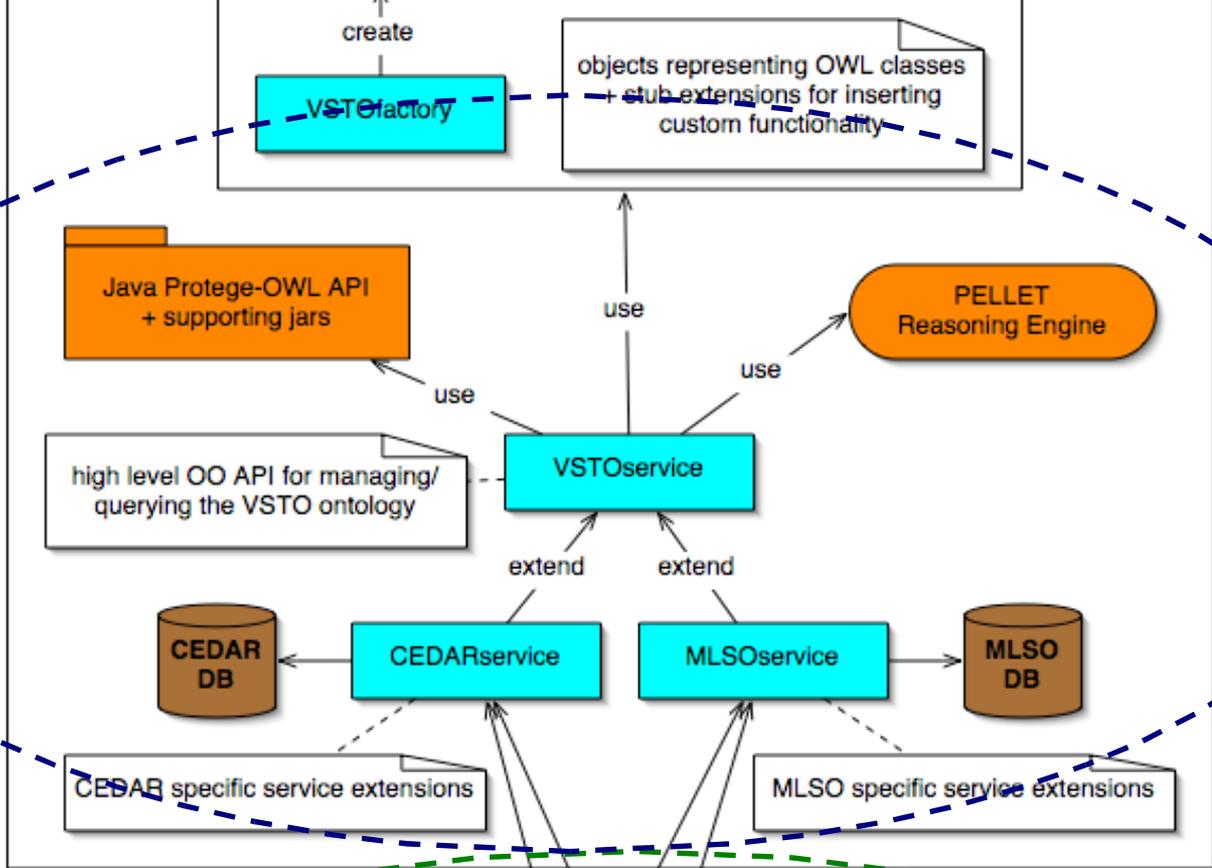
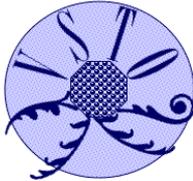
OWL ONTOLOGIES



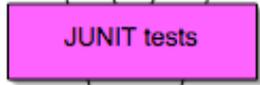
automatic generation

JAVA OBJECT MODEL

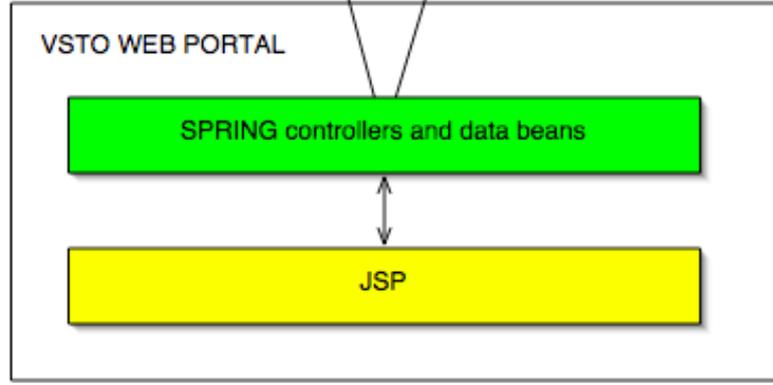




USE CASES WORKFLOW
SIMULATION PROGRAMS

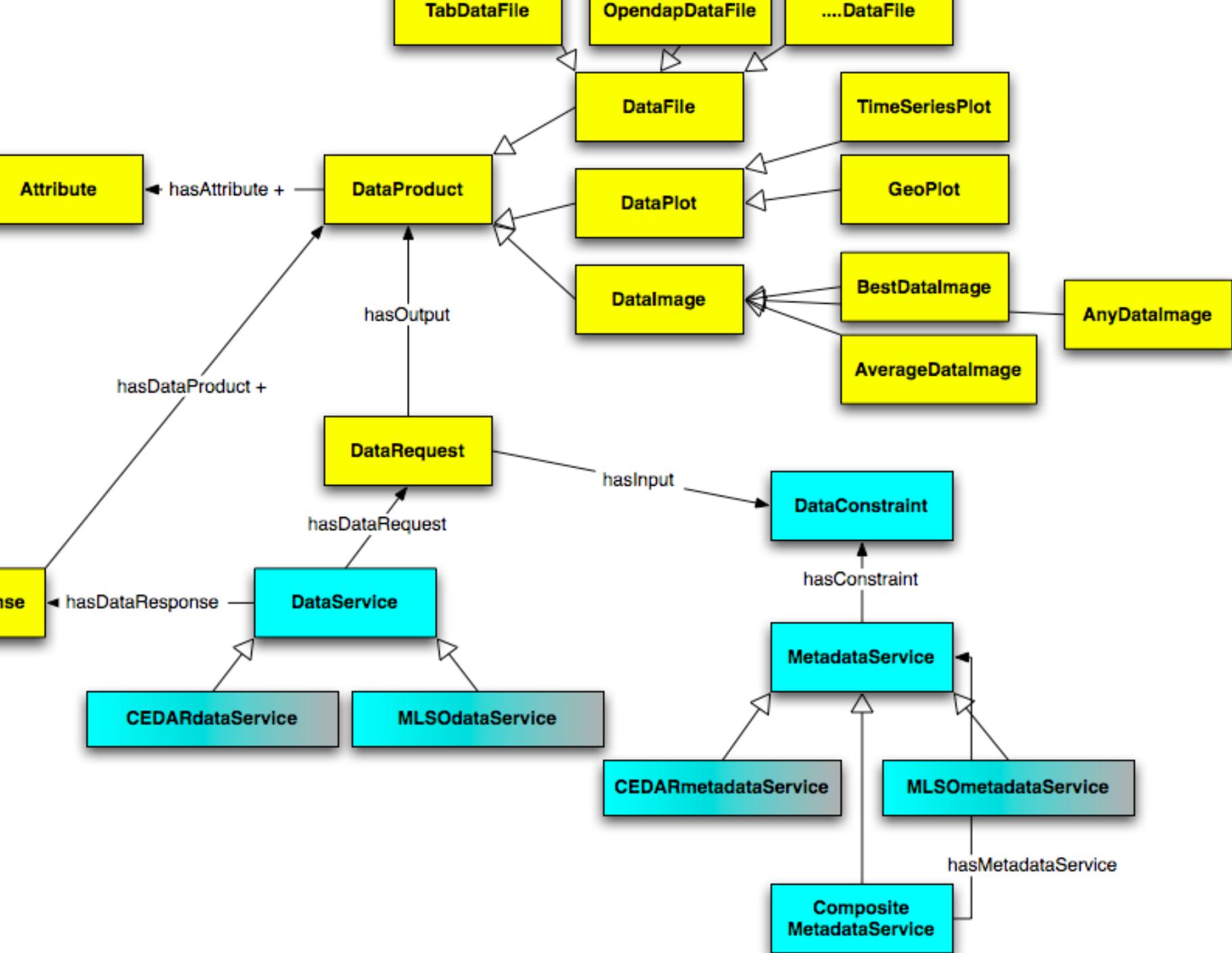


USER INTERFACE
CONTROL COMPONENTS

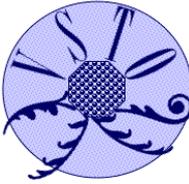


USER INTERFACE VIEWS





Why we were led to semantics



- When we integrate, we integrate concepts, terms
- In the past we would ask, guess, research a lot, or give up
- It's pretty much about **meaning**
- Semantics can really help find, access, **integrate, use, explain, trust...**
- What if you...
 - could not only use your data and tools but remote colleague's data and tools?
 - understood their assumptions, constraints, etc and could evaluate applicability?
 - knew whose research currently (or in the future) would benefit from your results?
 - knew whose results were consistent (or inconsistent) with yours?...

