

# **PDS Data Services Future Directions Discussion**

PDS MC F2F

Santa Barbara, California

Aug 13, 2019

# Outline

- Background, Objectives and Definitions
  - Approach: Address “The Roadmap”
- PDS Vision: Improved Support For:
  - Federated Architectures and  
International Standards for data discovery
  - API’s
  - Data Services
- Notional Plan
- Discussion

# Action Item

- March 2019 Face-to-Face: White Paper
  - “The decision was made by the group to proceed with writing a white paper that states the detailed path forward for PDS regarding our online presence. This will be of use for designing the CMU (architecture) study or will replace the CMU study as our (architectural) “roadmap” if the CMU study cannot be funded.”
  - Dan C. will take leadership role in writing the white paper.
  - Tim M. will lead Project Office management side with white paper.

# White Paper Status

- An *initial version* was circulated to the Node Leads in July 2019 for input, improvements/corrections and to drive discussion on the PDS forward path.
- Input received from all nodes (thank you!). Some nodes have indicated they have more and can provide additional details which we enthusiastically will accept.
- Intent is to update white paper and produce next version with MC and node inputs incorporated.

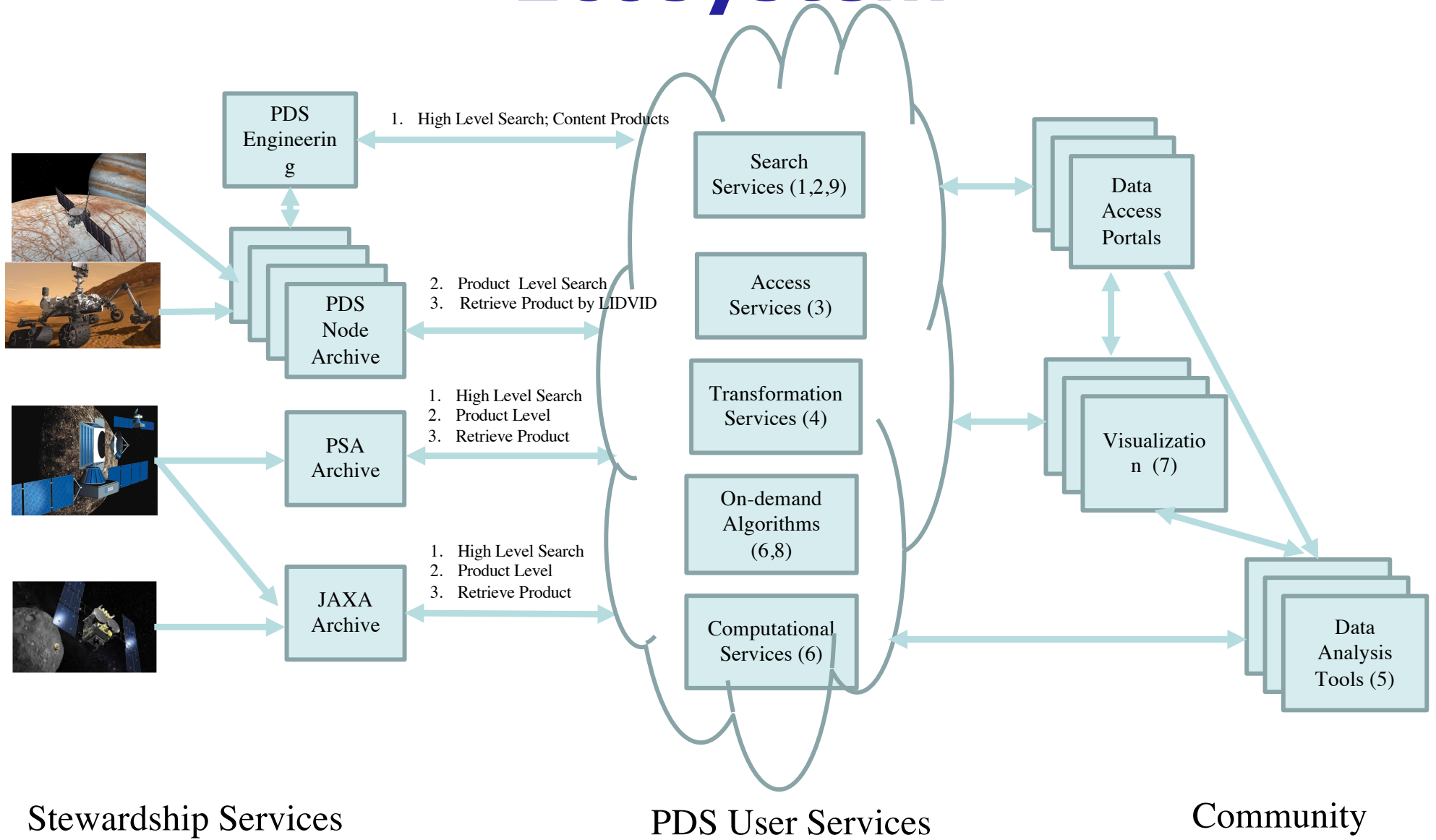
# White Paper Approach

- Focus on the overall PDS federated architecture with emphasis on data discovery.
- Generate a conceptual vision that can be used to drive discussions with PDS and community.
- Define the “As-Is” state of the PDS architecture
  - Capture detail in appendices (*This is an on-going validation process that needs node input – please forgive errors in this process!*)
  - Identify areas of improvement in the integration of the federation.
- Identify a “To-Be” state with an emphasis on improving integration and shared data services
  - Use 2017 Roadmap as a driver
  - Focus on data discovery (e.g., search and access)
  - Identify longer term opportunities and IT trends
- Benchmark to other data system efforts and plans
- Suggest a set of steps and phasing towards a “To Be” state as a basis for building a PDS-wide project plan for the federation.

# PDS Roadmap (2017): Highlights Around Data Services

- Discoverability - There is a need for PDS to both expand and deepen its search.
- Transformation – Increase support for translation services to deliver data in different formats.
- Integration with Other Archives – After international adoption of PDS4, the PDS is uniquely poised to lead efforts to make national and global archives interoperable.
- Information Technology – Ensure PDS continues to leverage modern technology approaches in search and data analytic support.
- Modernizing Metadata – Leveraging PDS4 metadata from stewardship to discovery.
- Transparency - Provide documented APIs for open access to data and services.
- Access Data – Build on PDS4 successes to continue to increase access to data.

# Towards a Planetary Data Ecosystem



# IPDA Statement of Progress

- “Since its founding in 2006, the International Planetary Data Alliance (IPDA) has made significant progress in the adoption of common standards, the development of compatible archives, and the establishment of open access policies. The efforts of IPDA, in working together, led to a common standard realized through the shared development of the Planetary Data System, version 4 (PDS4). **Based on PDS4 progress, the IPDA sees discoverability, seamless access to data holdings, and increased tool support for using high quality peer reviewed data across international archives as key foci for the future of the IPDA.**”

*Progress Report, International Planetary Data Alliance,  
January 29, 2018*



# AGU Emphasis

**Data & Emerging Technologies:** *Data is critical to scientific advancement and improving our understanding of how natural systems and phenomena operate and change. **Data should be openly accessible and archived for reuse into the future.** Emerging technologies are creating new instruments, sensor arrays, and platforms that enable the collection of new data types and/or improve the resolution, accuracy, and precision of data collection methodologies. **Frontier computational techniques and visualization tools are rapidly influencing the way we collect data and conduct science, thus forming a fertile breeding ground for new ideas and never-before-attempted science.***

# PDS Federated Architecture Today

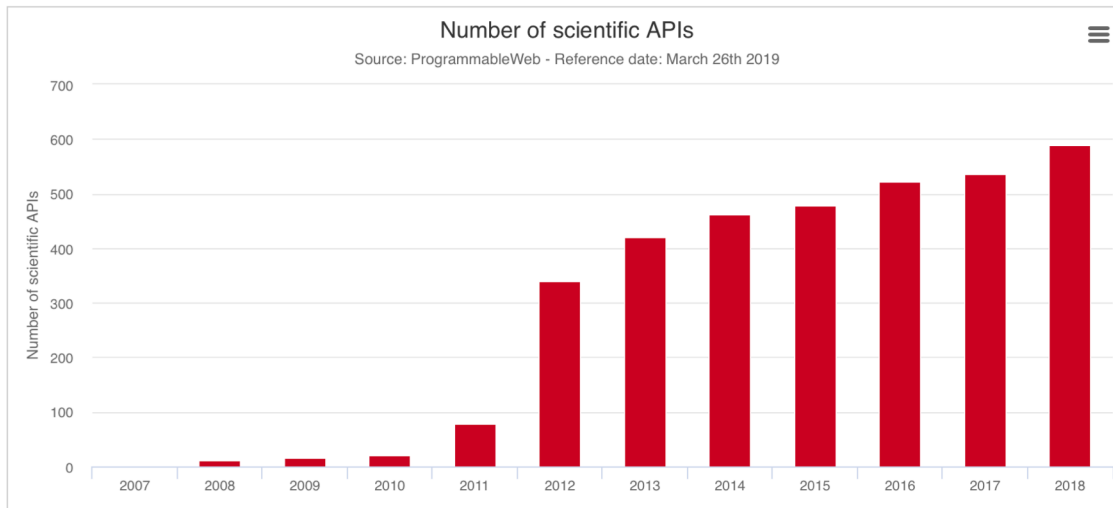
- A robust information model and process for archiving and stewarding data.
  - *This is a critical foundation for international discovery and use!*
  - *It is the de facto standard world-wide which is a huge opportunity*
- A multi-level search implementation following a system-of-systems approach
  - Users can get to the data
  - We have some great node services
  - We can link internationally

# PDS Architecture Today: Key Challenges

- Inconsistency in common protocols/standards around data discovery across PDS and the archive community systems: APIs; linking between search engines
- Inconsistent metadata (PDS3); differing use of parameters (PDS3, PDS4, other) for search/query
- Adoption of differing search methodologies and technologies (*this is a fast moving target*)
- End-to-end “search chain” and its impact on user interface/user experience (UI/UX)
- Lack of explicit, consistent, and sharable web services for sharing PDS data and common functions (“data services”) across the federation
- Differing approaches to navigating mission support and data access

# Why are APIs important?

An API approach is an architectural approach that revolves around providing a program interface to a set of services to different applications serving different types of consumers. (Wikipedia)



## NETWORKWORLD

### SOFTWARE QUALITY

By Ole Lensmar, Network World  
MAY 28, 2013 03:40 PM PT

### How open data and APIs fuel innovation

Relating Legos to open data, APIs, and quality software.

Just like many fellow developers and technology geeks, I was an avid Lego-builder during my youth. Those small plastic pieces gave me the possibility to create anything my imagination had in stock for me – and the cool thing was (and still is), the more basic the blocks were, the more freedom they gave me. Getting older I changed, and so did Lego - I can still sit with my friends complaining about the "dark years" of Lego, when those basic building blocks turned into pre-molded pieces of wings, buildings, or animals, putting a definitive stop to the creative outlet that provided so much joy in our youth. As stated by

# API Examples @ NASA

The screenshot shows the Earthdata website with the following content:

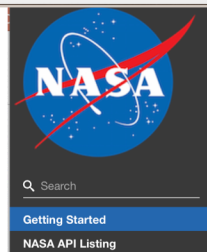
- Header: EARTHDATA Powered by EOSDIS. Navigation: ABOUT, DATA, COLLABORATE, LEARN. Search bar: Search datasets, news, articles, and information.
- Breadcrumbs: Collaborate > Open Data, Services and Software Policies > Application Program Interfaces (APIs)
- Section: **Application Program Interfaces (APIs)**
- Overview section with links:
  - Common Metadata Repository (CMR) APIs
  - Earthdata Login APIs
  - Earthdata Search APIs
  - Earthdata Tools APIs
  - Global Imagery Browse Services (GIBS) APIs
  - Science Data Processing Software (SDPS) APIs
  - Distributed Active Archive Center (DAAC) APIs
  - Open-source Project for a Network Data Access Protocol (OPeNDAP)
- Section: **Earthdata Developer Resource**
- Text: The API pages are the central location for all publicly accessible developer documentation related to EOSDIS enterprise services and applications, including:
  - Common Metadata Repository (CMR) APIs**
    - CMR is a spatial and temporal metadata registry that stores metadata from a variety of science disciplines and domains. CMR is intended to enable broader use of NASA's EOS data by providing a more uniform view of NASA's substantial and diverse data holdings. CMR interfaces with clients and users through various APIs; CMR is an open system.
  - Distributed Active Archive Center (DAAC) APIs**
    - NASA's Distributed Active Archive Centers (DAACs), located throughout the United States, are custodians of EOS mission data and ensure that data will be easily accessible to users. A number of APIs are available for direct access to these data holdings.
  - Earthdata Login APIs**
    - Earthdata Login provides user profile management and authentication services, freeing up your application from the problems associated with managing user databases. Earthdata Login also provides an application programming interface (API) that can be used to query the user database and retrieve user information.

The screenshot shows the NASA Open APIs website with the following content:

- Header: NASA Open APIs
- Section: **Exoplanet Archive**
- Introduction:

The Exoplanet Archive API allows programmatic access to NASA's Exoplanet Archive database. This API contains a ton of options so to get started please visit [this page](#) for introductory materials. To see [what data](#) is available in this API [visit here](#) and also be sure to check out [best-practices and troubleshooting](#) in case you get stuck. Happy planet hunting!
- Figure: A diagram showing a yellow star with a planet orbiting it. The planet's brightness is plotted against time, showing a characteristic dip during a transit.
- Example Queries table:

Example API	URL
<a href="#">Confirmed planets in the Kepler field</a>	<a href="https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=exoplanets&amp;format=ipack&amp;where=pl_kepf1ag=1">https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=exoplanets&amp;format=ipack&amp;where=pl_kepf1ag=1</a>
<a href="#">Confirmed planets that transit their host stars</a>	<a href="https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=exoplanets&amp;format=ipack&amp;where=pl_transit1ag=1">https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=exoplanets&amp;format=ipack&amp;where=pl_transit1ag=1</a>
<a href="#">All planetary candidates smaller than 2Re with equilibrium temperatures between 180-303K</a>	<a href="https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=cumulative&amp;where=koi_prac&lt;2 and koi_teq-180 and koi_teq&lt;303 and koi_disposition like 'CANDIDATE'">https://exoplanetarchive.ipac.caltech.edu/cgi-bin/nstEDAPI/ngh-nstEDAPI?table=cumulative&amp;where=koi_prac&lt;2 and koi_teq-180 and koi_teq&lt;303 and koi_disposition like 'CANDIDATE'</a>



{NASA APIs}

# Three Pronged Strategy: Open Data, Standards Models, Shared APIs

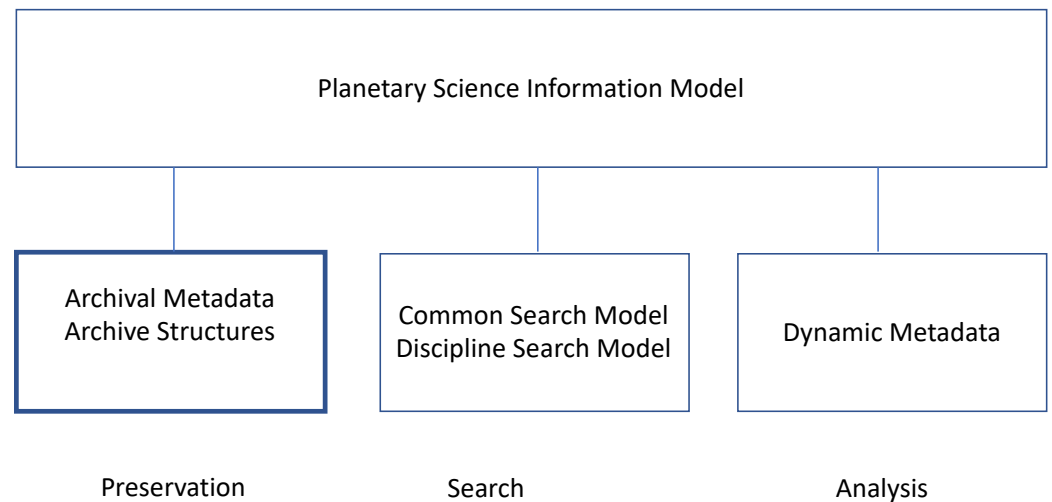
- Open Data – Planetary data online and open for use (*in the DNA of PDS*)
- Standard Data Models – Planetary data described by a consistent model (*PDS4*)
- Shared APIs – Well documented interfaces to get to data and services for use across PDS and by the community (*we need to do more of this!*)
  - *APIs directly driven by the PDS4 information model (e.g., discipline query models)*

# Overall Federated Architectural Concept

- The architecture must handle *variety* and *distributed services* as a core architectural principles.
- With PDS4 now in place, increase focus on integration of the federation and improving open access to data and services so they can be shared across PDS, IPDA and the community.
  - Improve integration and sharing of distributed search services (both EN and DN)
  - Drive consistent protocols for major functions (access, search, transform) as a basis for PDS and the international community to adopt
  - Improve web presence (UI/UX, mission support pages, etc) and search chain by leveraging of shared services and APIs
  - Explicitly publish common services for use by the community 15

# Leveraging PDS4 Information Model to Drive User Services

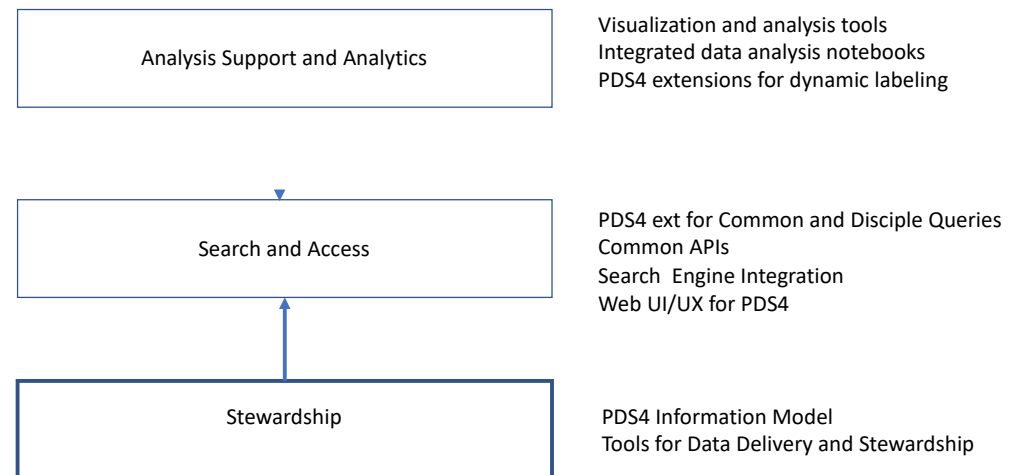
- Primary focus has been on stewardship for planetary science data
  - Given data variety, a model-driven architecture is central to the future of a data-driven environment
- Three areas can be supported Model Views
  - Stewardship and Management
  - Discovery and Search
  - Analysis
- Future capability needs
  - Explicit discipline search models
  - Non-archive metadata to label data for machine learning and other discovery approaches



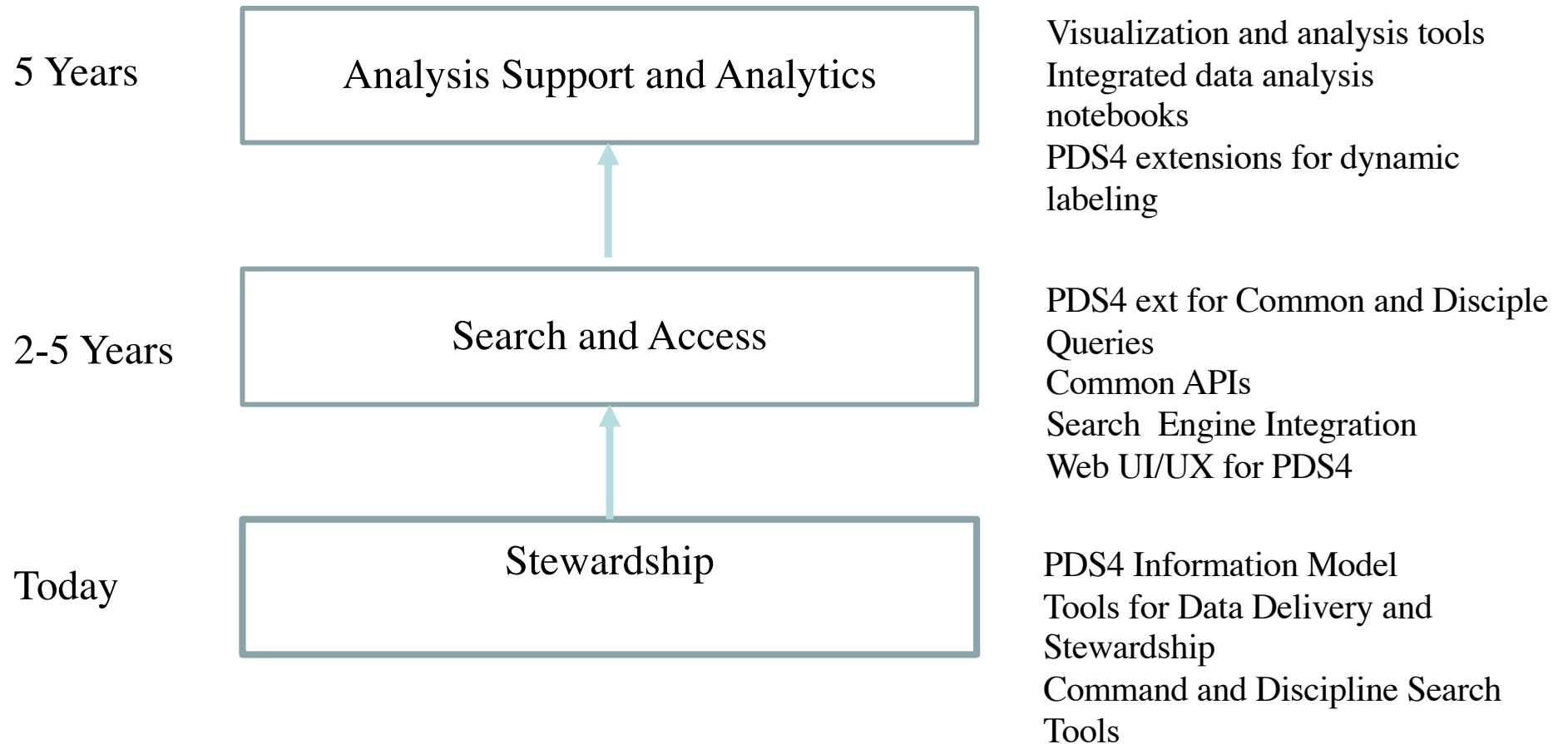


# A Process for Unifying the Architecture of the PDS Federation

- Expand the information model to define common and discipline specific queries.
- Support dynamic metadata for enhanced product level searching.
- Support cross-code, cross-product federated searching using modern search engines.
- Develop and use consistent APIs for searching, access, and retrieving data across systems.
- Build and register consistent archive support pages to help users boot strap into archives.
- Develop a next generation web design leveraging PDS4 metadata, data services, and integration of search across PDS to improve user experience.



# Notional Data Services Evolution



*This needs to be driven by a rigorous project plan based on available resources  
(funds and staff)*

# Example Services

- Search
- Access
- Transform
- Computation
- Tool Integration
- Dynamic Indexing

# Other Planning Considerations

- Community Collaborations
- Driving UI/UX capabilities forward
- The role of cloud computing to enable data services
- PDS role in integrating data analytic capabilities
- Outreach on the vision

# Next Steps

- Tie off federated architecture white paper
  - Update with node inputs; validate detail
- Achieve MC agreement on executive summary
  - Consider whether to also publish a brief/vision on future directions of PDS
- Formulate a more detailed project plan with resources for discussion with MC
- Consider whether to stand up a WG at some point

# Discussion

# Backup

# State of Data Services Mapping Table

Capability	IPDA an PDS Support
1a. High Level Search within an Archive (PSA, PDS, etc)	PDAP; EPN-TAP; PDS Search Service API; PDS and PSA supporting services interfaces.
1b. High Level Search across archives (PSA, PDS, etc.)	PDS Search Service API; Search integration (e.g., passing search parameters) inconsistently implemented.
2a. Product Level Search within a Site	Differing support for REST-based API access
2b. Product Level Search across Sites	No PDS or IPDA-wide product-level search;
3. Retrieve Product	Most archives/node provide HTTP access for download. Inconsistent services for download (label, data).
4. Transform	PDS library exists but no service
5. Tool Integration	Integration with tools such as ISIS; no support for Jupyter notebooks; no support for languages such as R
6. Computational Support	Limited support for on-the-fly processing and running analytical results.
7. Visualization Support	Solar System Treks; VESPA
8. Metadata Extraction	Local node activities
9. Indexing on dynamic metadata	Local node activities



# Search

- Goal: *Integrated* search of PSD4-compliant federated archives
  - There is no single search that will ever provide *comprehensive* discovery across all planetary types using a single search string
- Web services PDS can offer:
  1. A consistent search protocol at both the high level and product level
  2. A set of parameters for forwarding search strings
  3. A centralized index for product level search of common attributes (e.g., LIDVID)

# Access

- Goal: Return every product, including both labels and data, from any PDS4 compliant archive
- Web services PDS can offer:
  1. Retrieve base product product based on LIDVID
  2. Retrieve label for a product based on LIDVID
  3. Retrieve related ancillary information

# Transform

- Goal: Provide standard library of transformations
- Web services PDS can offer:
  1. Common transformations across all archives for PDS4 data offered as web service

# Compute

- Goal: Provide standard processing services (e.g., subsetting, coordinate translation, etc)
- Web services PDS can offer
  1. Geometry services
  2. Subsetting services
  3. Other routine processing as data increases

# Tool Integration

- Goal: Support integration with common tools and frameworks
- Web services PDS can offer
  - Jupyter Notebooks for planetary science data
  - Python, R, and other language bindings to search, retrieve, transform, and compute on data
  - Shared visualization tools

# Dynamic Indexing

- Goal: Further enhance discovery and use of data through auto labeling using machine learning feature detection and classification methods based on PDS4 model
- Web services PDS can offer
  - Shared ML libraries for labeling including image and other data
- Note: this implies the ability to use dynamic metadata for search as well as user supplied metadata for analysis

# Community Collaborations

- IPDA – Providing consistent APIs and data services will increase interoperability opportunities with international archives.
- Community – Providing consistent APIs and data services will enable the community to build analysis tools on top of PDS.

# UI/UX Considerations

- Improving User Interface/User Experience (UI/UX) and leveraging PDS4 is critical to the future of improving the PDS.
  - This needs to be part of the plan!
- It is important to have a *consistent architectural foundation* on which to build an improved UI/UX web presence



# Cloud Computing to Enable Data Services

- Part of the forward plan needs to determine the role of cloud computing for
  - Shared data services (*this is where cloud could provide a big advantage*)
  - Storage (primary, secondary/backup)
- Cloud services provide
  - Standards APIs for access, search tool support, data management, computation, machine learning, etc.

*Cost and other trades need to be made as part of the implementation plan*

# PDS Role in Evolving Data Analytics

- Increasing technical capabilities are opening new opportunities for AI, ML, and visualization.
  - Increasing support in many different science disciplines to apply these techniques to improve how data is used and analyzed.
- Part of enabling an international planetary data ecosystem
  - PDS is leading the way in exploring its role through the Planetary Science Data Analytics and Informatics Conference and Planetary Data Workshop.
  - PDS can enable this but it needs to define how far it should go with its resources and its role.

# Outreach

- Socialize PDS vision and plans at community meetings
- Consider writing a joint PDS paper outlining the future

# Enabling Data Services Functional Capabilities

1. High Level Search across the federation and IPDA
2. Product-level Search
  - a. Within an archive (comprehensive)
  - b. Across federated archives for different scenarios\*
3. Retrieve products from an archive
4. Transform a PDS4 product
5. Integrate modern data analysis tools
6. Compute on PDS4 products
7. Interactive visualization for PDS4 data products
8. Metadata extraction using ML
9. Indexing on dynamic metadata

*It's understood that little common metadata exists for every PDS product (e.g, LIDVID, LID, Time, Instrument, Target, etc)*

# Definitions

- *Architecture (systems and software)* is “the fundamental organization of a system embodied in its components, their **relationships** to each other, and to the environment, and the principles guiding its design and evolution.” (ISO/IEC/IEEE 42010:2011, IEEE-1471)
- We have often used three critical views to describe an architecture: data lifecycle, data architecture (models, dictionaries, etc), technology (services, tools, etc).
- What is also important are the protocols and standards that support the integration of software services and tools for PDS, IPDA, and users.