

PIT services: the what and the why

Beth Plale, Jason Haga, Quan (Gabriel) Zhou

9/8/16



First: Building a concrete data fabric configuration

Beth Plale, Tobias Weigel
Taken from RDA PID Training, Garching, 2016/08/31

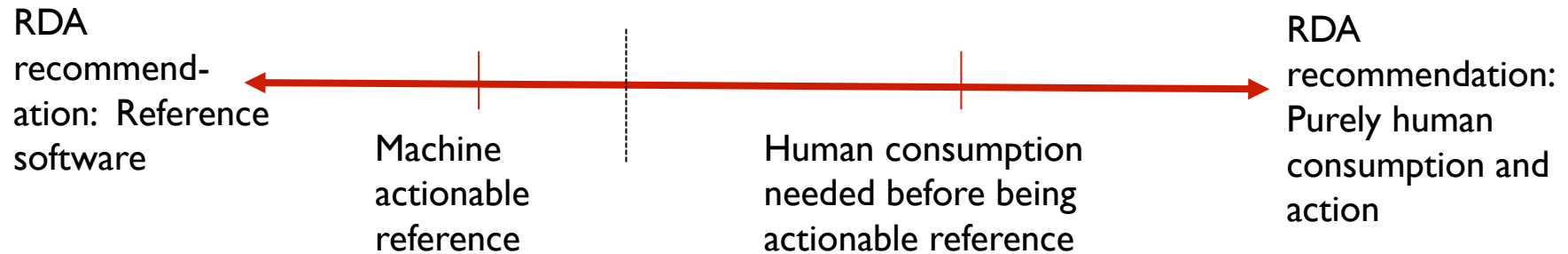
research data sharing without barriers
rd-alliance.org

-
- RDA Data Fabric has activity that examines fabric composition
 - Composing from RDA Recommendations (largely but not exclusively)
 - A couple Recommendations are around PIDs (session view not citation view)
 - Inductive (direct) approach to composition of component

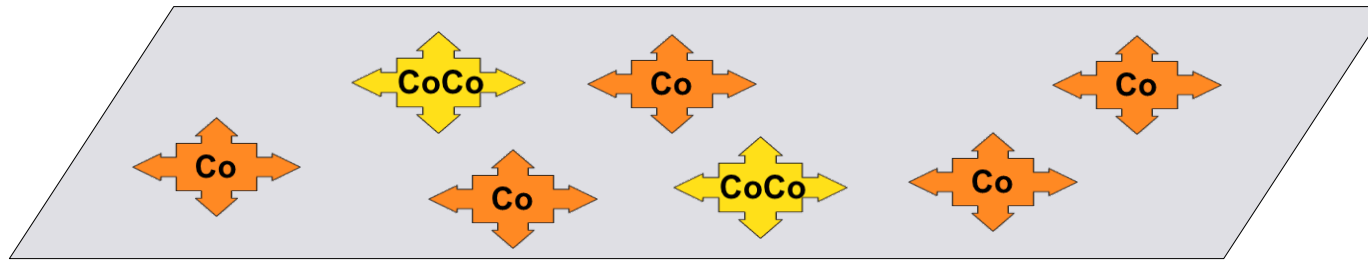


Dimensions of Testing: getting to Data Fabrics

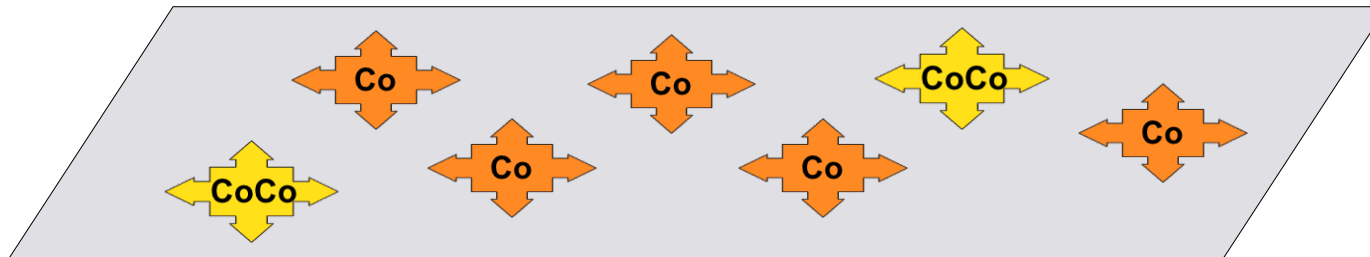
- RDA produces RDA Recommendations (i.e., outputs)
- Some technically-oriented RDA Recommendations have reference software with it, these are starting point
- Many technically oriented RDA recommendations do not have reference software, yet are machine actionable (a schema for instance). These also are of immediate interest to data fabric composition.
- Other recommendations play important but background roles in early composition



Compositions of components



Composition (or Fabric) A



Composition B

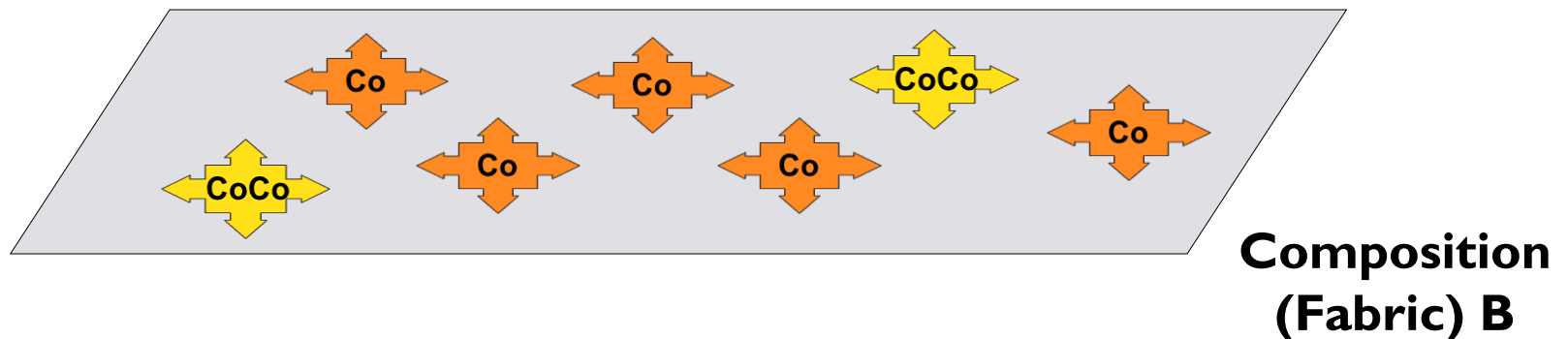


Compositions of components

Given nature of data (can't do much without understanding it), successful data fabric will likely:

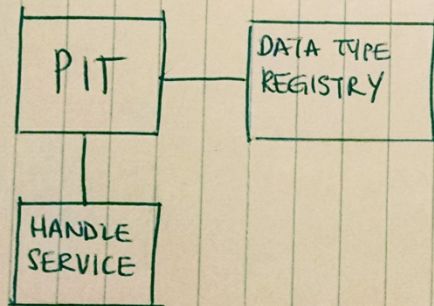
1. Run on possibly distributed e-infrastructure (EUDAT, NDS, ...)
2. Serve scholarly domain as domain infrastructure
3. Support multiple projects within that domain
4. And eventually result in cross-domain research

For 3 and 4 to be realized, composition of 2 must be shared across projects

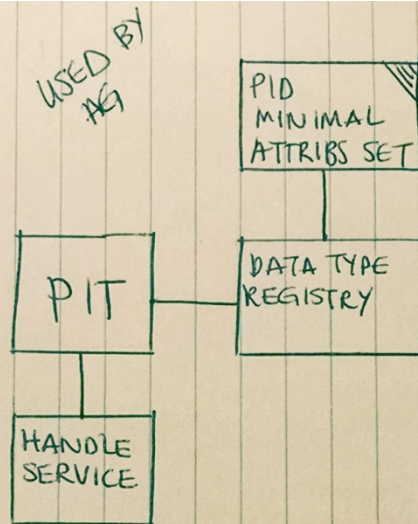


-
- RDA PIT WG Recommendation and RDA Data Type Registry Recommendation are starting point for composition, hence my interest in the topic

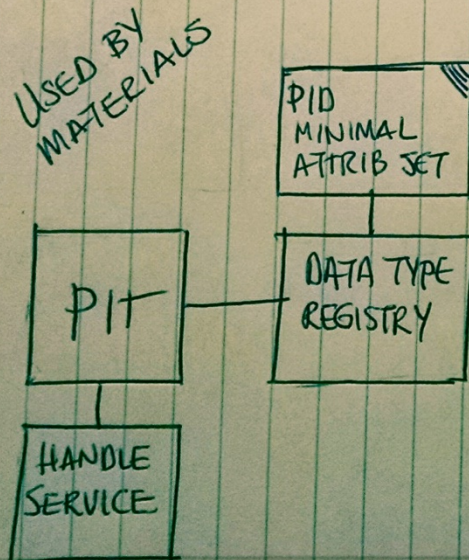




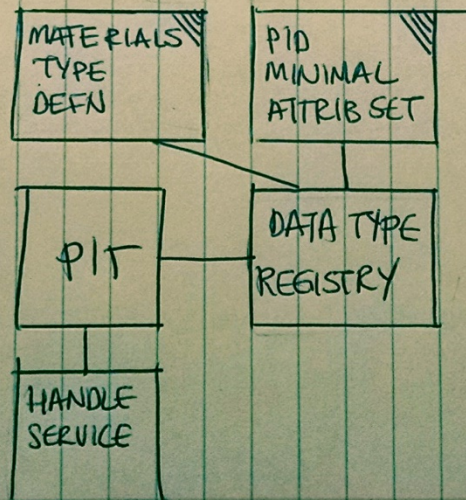
COMPOSITION (FABRIC) 1




COMPOSITION 2

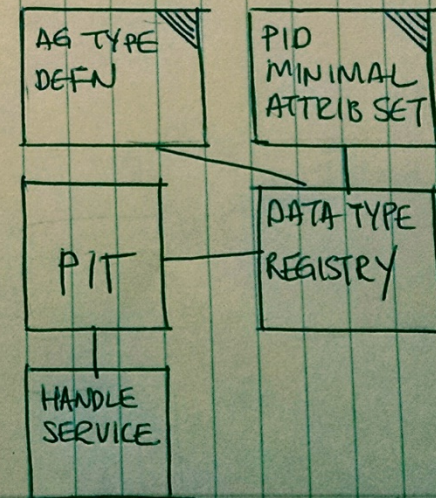


COMPOSITION 3

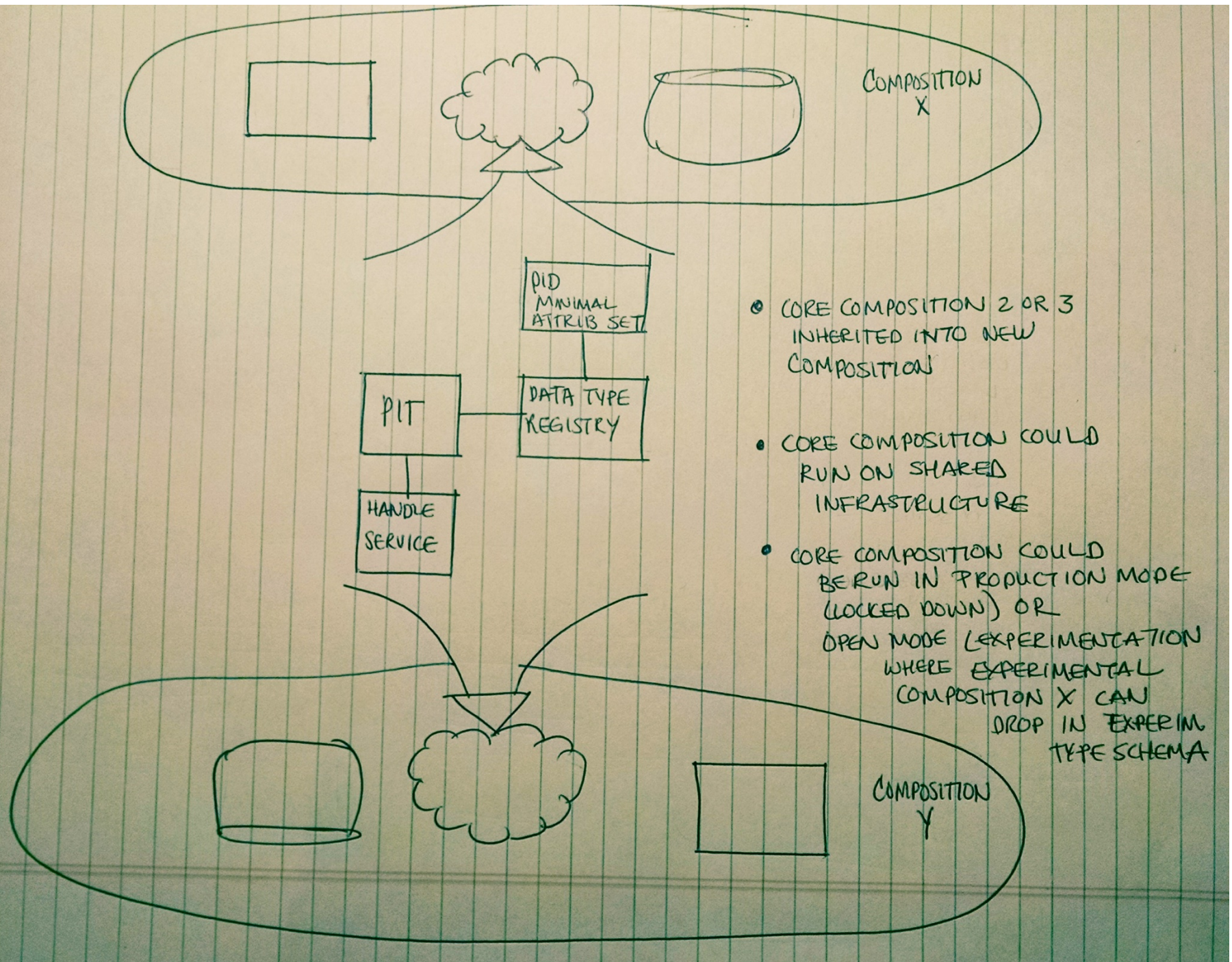


COMPOSITION 4

- COMPOSITION 1 IS NOT USEABLE AS IS
- COMPOSITION 2 \equiv COMPOSITION 3 SO SUBJECT TO "CORE" DESIGNATION
- COMPOSITION 4 AND COMPOSITION 5 START TO GET AT UNIQUE COMMUNITY INSTANCES OF A FABRIC
-  : DOESN'T EXIST TODAY
- COMPOSITIONS 4 AND 5 HAVE $N=5$



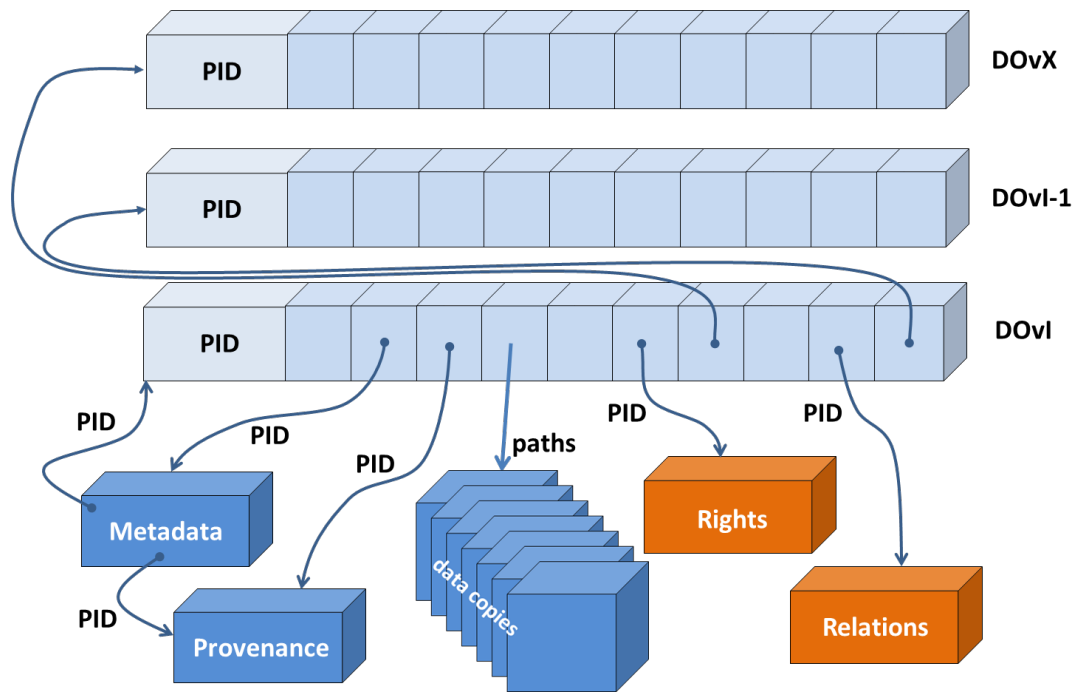
COMPOSITION 5



-
- We have seen that RDA has $n=2$ shared services around PIDs, can it go to $n=3$
 - Can it get to agreement around minimal PID attributes?



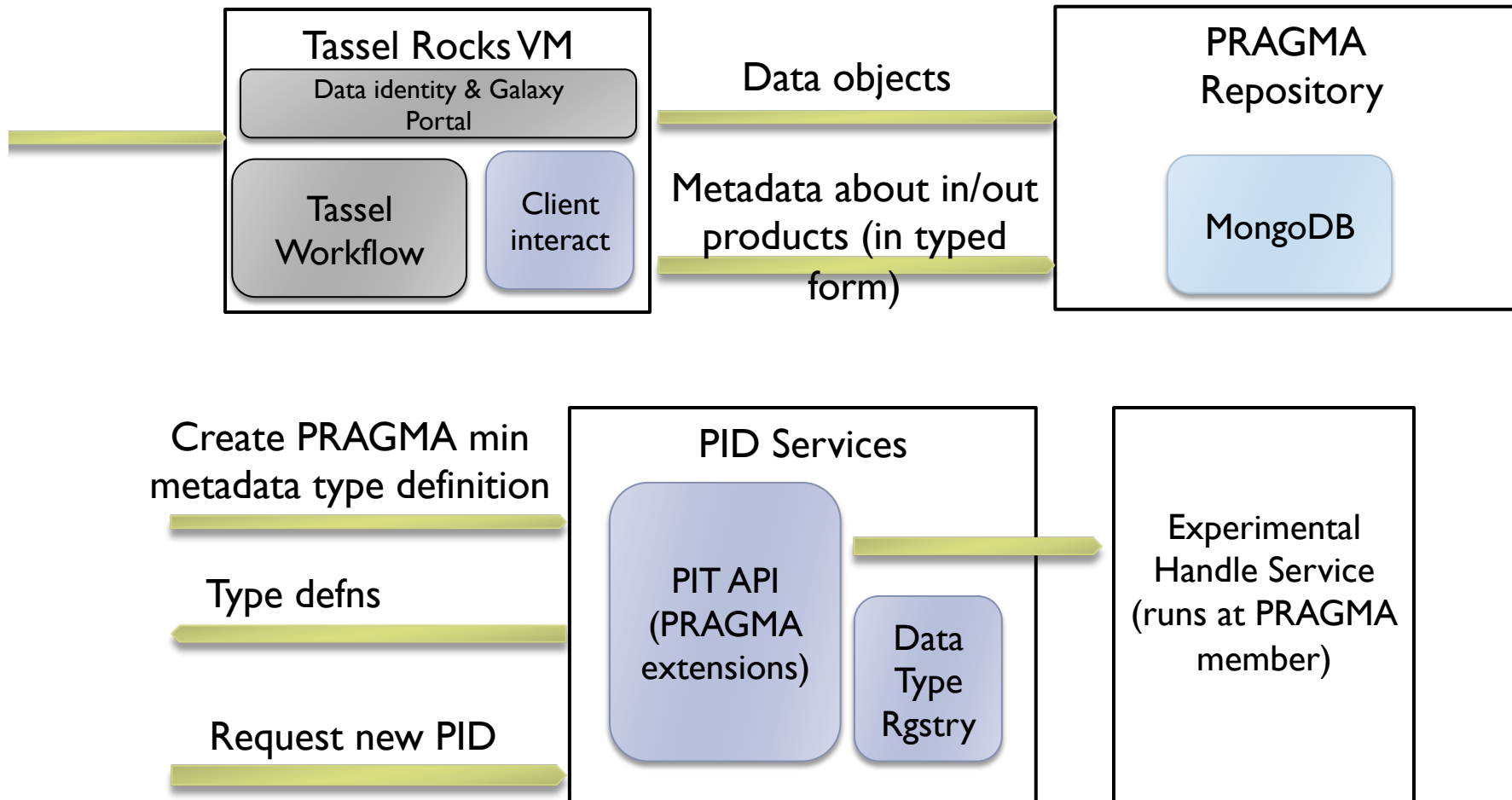
Example PID minimal metadata



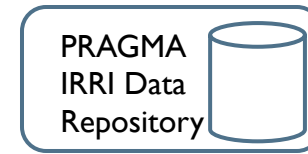
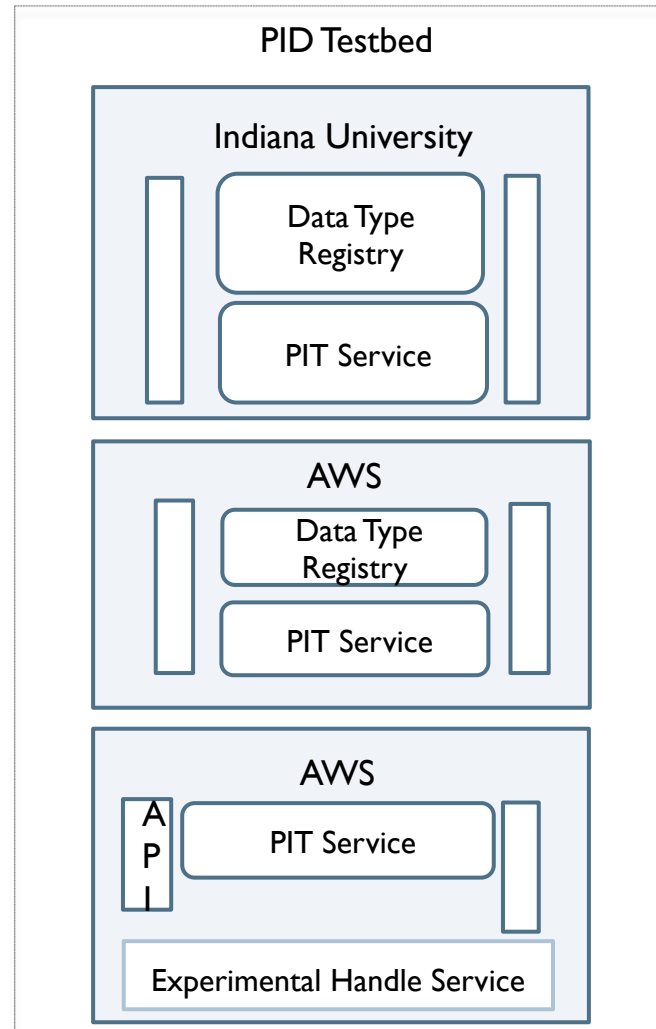
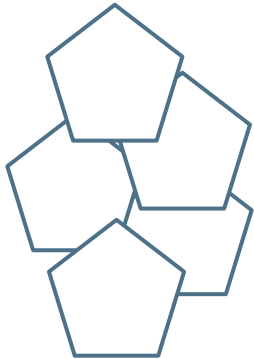
- Define PRAGMA defined definition (for experimental purposes ... works for rice genomics, Airbox, weather data repository ...)



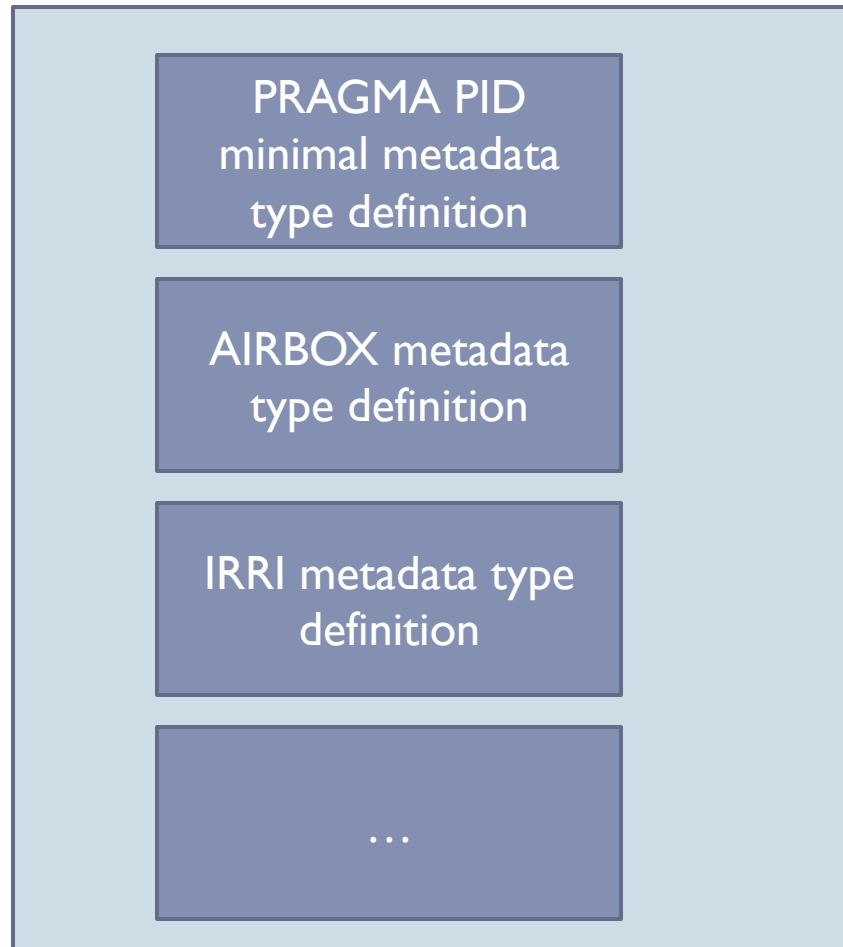
PRAGMA Rice Genomics: major components



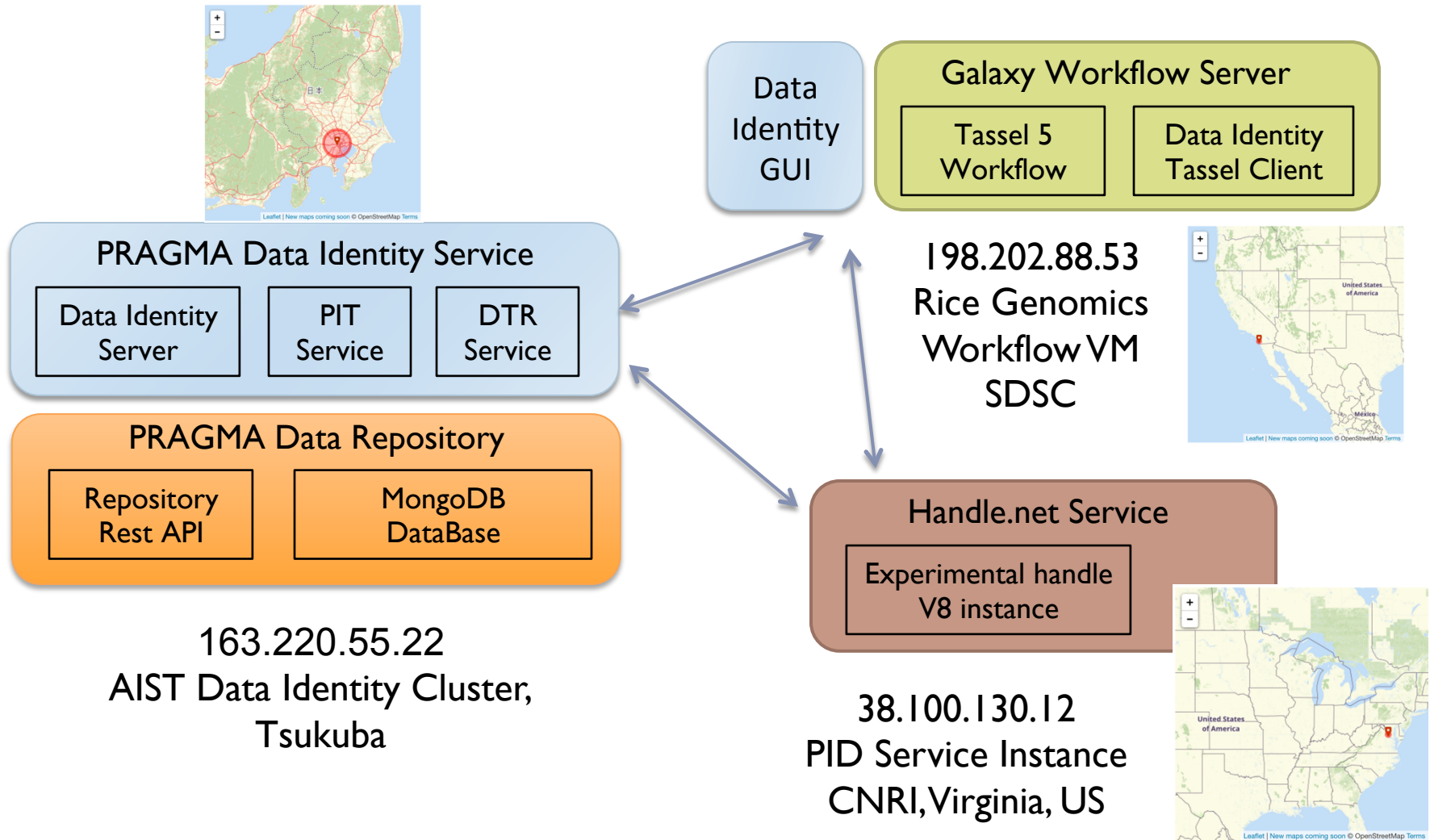
Discovery services :
community evaluation
clients



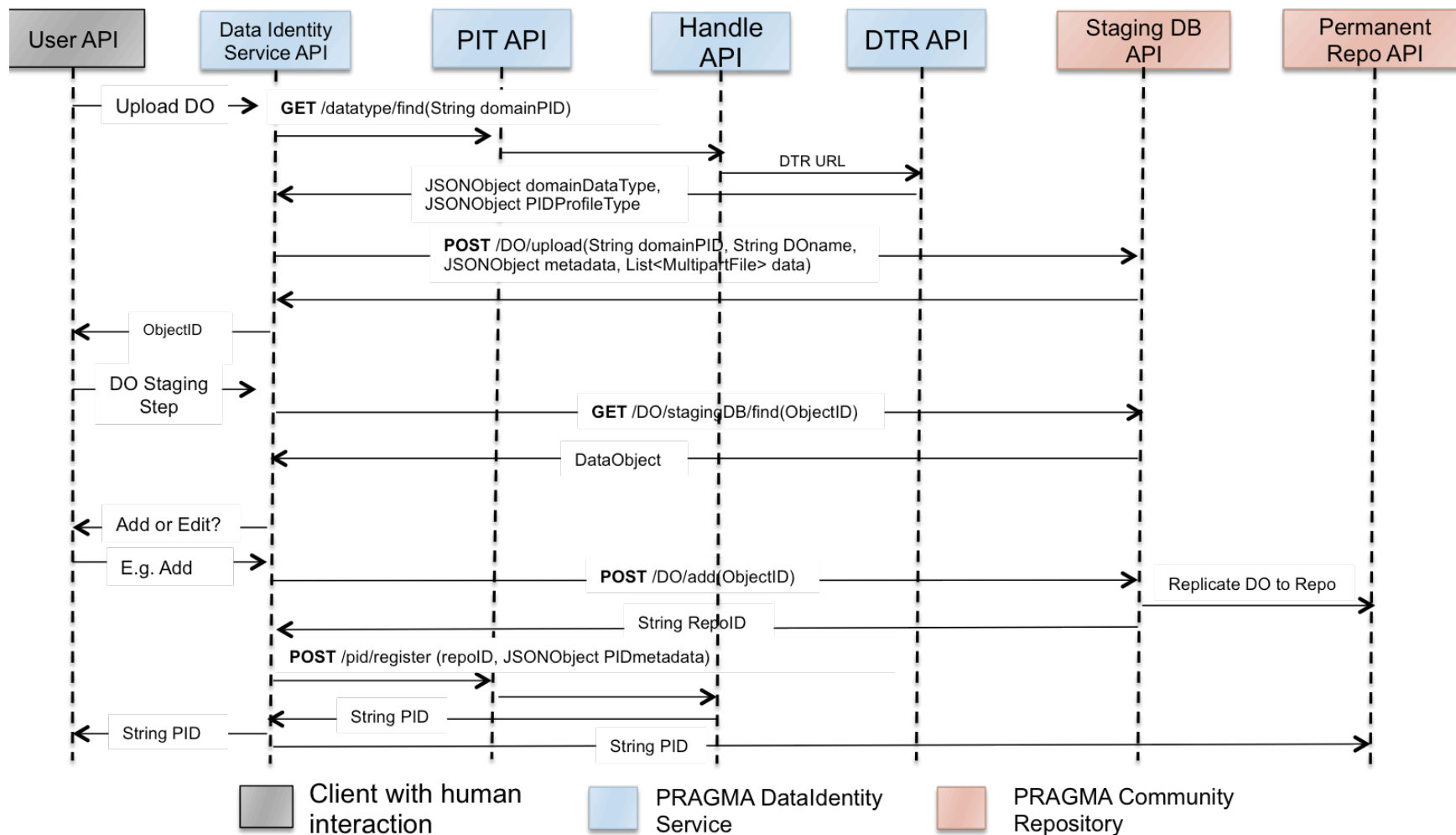
Data Type Registry holds type definitions (aka, schemas for both the PID minimal metadata defn of a community plus metadata)



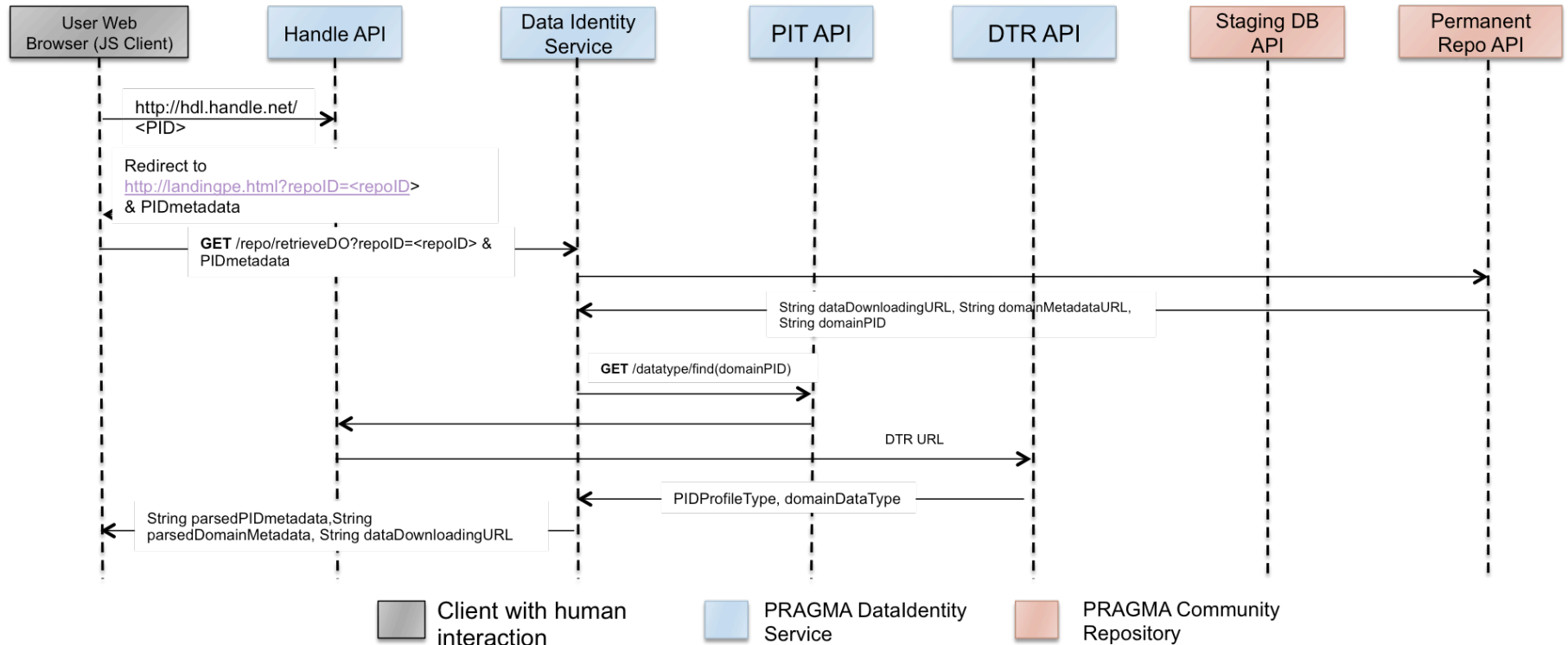
Deployment Diagram



DO Upload Timeline Diagram



DO Retrieval Timeline Diagram



Middleware Service Timeline Diagram

