

Service infrastructure for cross-matching distributed datasets using OGSA-DAI and TAP

Mark Hollimanⁱ, Alastair C Humeⁱⁱ, Tilaye Alemuⁱⁱ, Robert G Mannⁱ, Keith Noddleⁱ

ⁱWide Field Astronomy Unit (WFAU), Institute for Astronomy, University of Edinburgh

ⁱⁱEdinburgh Parallel Computing Centre, University of Edinburgh

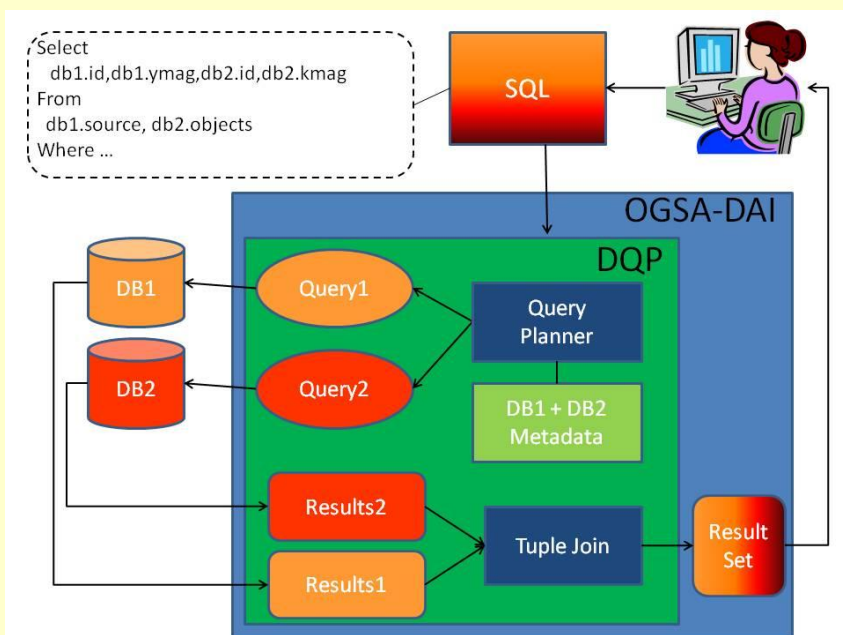
Objectives

One of the most powerful and important goals for VO developers has been to enable cross-match queries between disparate datasets for end users. This has only been achieved within the VO using the early SkyNode infrastructure and has not been reproduced using current IVOA standards.

To remedy that situation, the Wide Field Astronomy Unit (WFAU) has worked with the Edinburgh Parallel Computing Centre (EPCC) in leveraging the OGSA-DAI¹ grid middleware to enable cross catalogue queries on distributed Table Access Protocol² (TAP) services.

OGSA-DAI and DQP

OGSA-DAI is a grid middleware that operates as a web service which provides access to databases.



OGSA-DAI includes a component for performing join operations across distributed databases called Distributed Query Processing (DQP). DQP performs the tasks of parsing and optimizing the query, and builds and executes the query plan. It federates the metadata from the underlying databases into a single metadoc that acts as though all the tables are within a single database.

Possible JOIN types in OGSA-DAI

- In-memory join
One side stored in memory, other side streamed
- Partial in-memory join
Gets first results quickly, but all data stored to disk
- Ordered merge join
Both inputs ordered for a fully streamed join
- Parallel hash equi-join
- Batch joins using IN clauses
E.g. SELECT * FROM foo WHERE bar IN (x, y, z)

For More Details

Info on the different software referenced above can be found at these sites:

¹OGSA-DAI - <http://www.ogsadai.org.uk/>

²TAP - <http://www.ivoa.net/Documents/TAP/>

³VOExplorer - <http://www.astrogrid.org>

⁴TAPsh - <http://vo.ari.uni-heidelberg.de/soft/tapsh>

⁵CADC - <http://cadwww.dao.nrc.ca/caom/>

For questions or a demonstration contact Mark Holliman at msh@roe.ac.uk

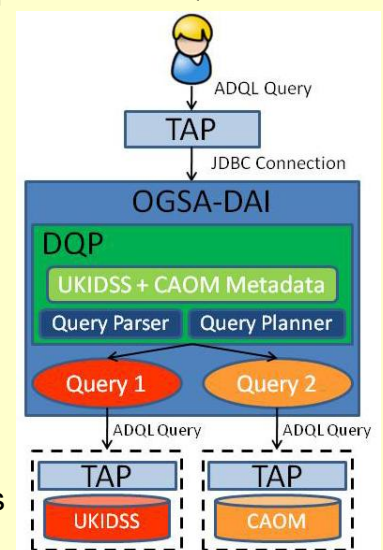
Table Access Protocol

TAP is a recently approved IVOA protocol for accessing tabular data in a standardized way. Current implementations allow users to submit queries in ADQL through standard HTTP GET and POST methods, and return results in VOTable format. Users can access TAP services using simple clients like web browsers or 'wget', or software like VODesktop³ and TAPsh⁴. There are a limited number of TAP services currently published on the VO, but more are expected in the near future.

Cross-matching Infrastructure

We have designed a 3 layered architecture that places the OGSA-DAI middleware above an arbitrary set of TAP services. OGSA-DAI combines the metadata from the underlying datasets, and this is then exposed through a single TAP interface at the top layer, enabling users to perform cross-match queries in ADQL on all of the federated datasets underneath.

We implemented a test infrastructure using an OGSA-DAI deployment that federates the UKIDSS DR3 TAP service and the Canadian Astronomy Data Centre's (CADC⁵) TAP service for CAOM data. These are independent TAP services, and no coordination or assistance was required from CADC in building the test system. We then submitted several scientifically relevant ADQL queries that cross-matched the different datasets to determine problems and performance for the system.



Overall the test system successfully demonstrated the viability of the service architecture, returning results for most queries. As such it merits further development to improve its capabilities. Several issues were discovered during tests relating to both the ODSA-DAI middleware and the TAP services which include language disparities between ADQL and DQP, poor performance on certain joins involving large datasets, and various minor software fixes. Others issues we are addressing include:

- No physical metadata is available for query planning
- Join algorithms are hardwired due to limits of cardinality estimation
- DQP does not use batch join algorithm
- Few available TAP services are deployed
- Existing TAP services vary in implementation & compliance

Future Development

Our ultimate goal is to deploy a registered TAP service utilizing the OGSA-DAI infrastructure to federate data from all the TAP services published on the VO. Such a service would be one of a kind in allowing users to cross-match all the databases and catalogues available.

In the meantime we are working on addressing the issues discovered through testing as well as implementing more join algorithms, improving join algorithm choice and execution, support for ADQL spatial functionality, and improved performance for joins of large (multi terabyte) databases.

