
Data Aggregation System Documentation

Release development

Valentin Kuznetsov

August 17, 2015

1	Contents:	3
1.1	Introduction	3
1.2	How DAS works?	3
1.3	How to use DAS?	4
1.4	Installation	13
1.5	Dependencies	15
1.6	DAS and KWS maps	16
1.7	Create DAS and KWS maps	16
1.8	Validate DAS and KWS maps	16
1.9	Import DAS and KWS maps	16
1.10	Fetch DAS and KWS maps	17
1.11	Setting up and customizing DAS installation	17
1.12	DAS architecture and internals	29
1.13	DAS configuration file	61
1.14	Release notes	63
1.15	References	92
2	Indices and tables	93
	Bibliography	95
	Python Module Index	97

- **Version:** development
- **Last modified:** August 17, 2015

Contents:

1.1 Introduction

DAS is an opensource virtual data service integration platform. The acronym DAS stands for *Data Aggregation System*. It provides a common layer above various data services, allowing end users to query one (or more) of them with a more natural, search-engine-like interface. It is being developed for the [CMS] High Energy Physics experiment on the [LHC], at CERN, Geneva, Switzerland. Although it is currently only used at the CMS experiment, it was designed to have a general-purpose architecture which would be applicable to other domains.

It provides several features:

1. a caching layer for underlying data-services
2. a common meta-data look up tool for these services
3. an aggregation layer for that meta-data

The main advantage of DAS is a uniform meta-data representation and consequent ability to perform look-ups via free text-based queries. The DAS internal data-format is JSON. All meta-data queried and stored by DAS are converted into this format, according to a mapping between the notation used by each data service and a set of common keys used by DAS.

1.2 How DAS works?

1.2.1 Resolving queries over data services

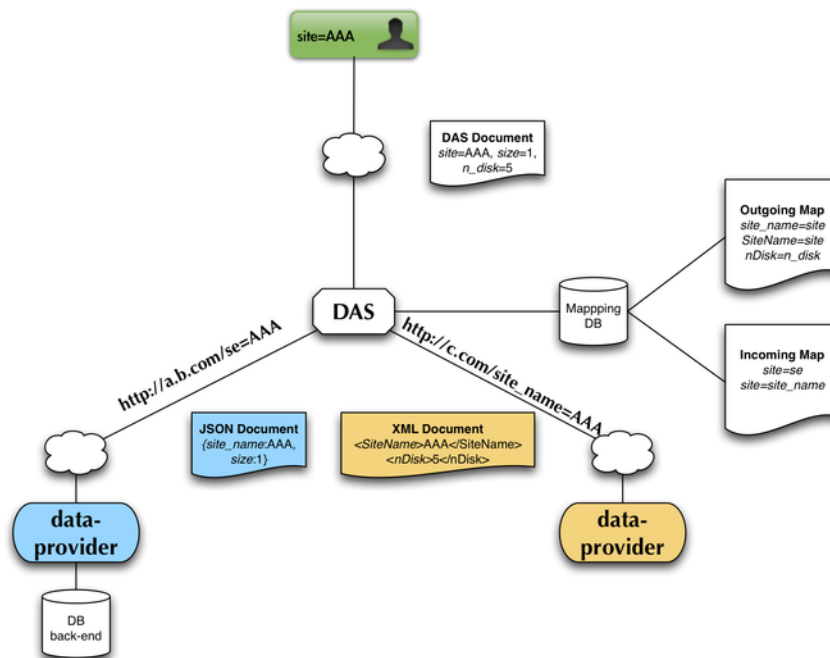
This diagram illustrates the request flow triggered by a user query.

The user makes the query *site=AAA*. DAS resolves it into several requests, by looking at the *daskeys* map for each available data service, which indicates that two services understand an input key *site*, which is transformed into queries for each of those services with parameter *se* and *site_name*.

DAS then makes the API calls (assuming the data isn't already available).

- *http://a.b.com/se=AAA*
- *http://c.com/site_name=AAA*

and retrieves the results, which are re-mapped into DAS records according to the *notation* map for each of those services.



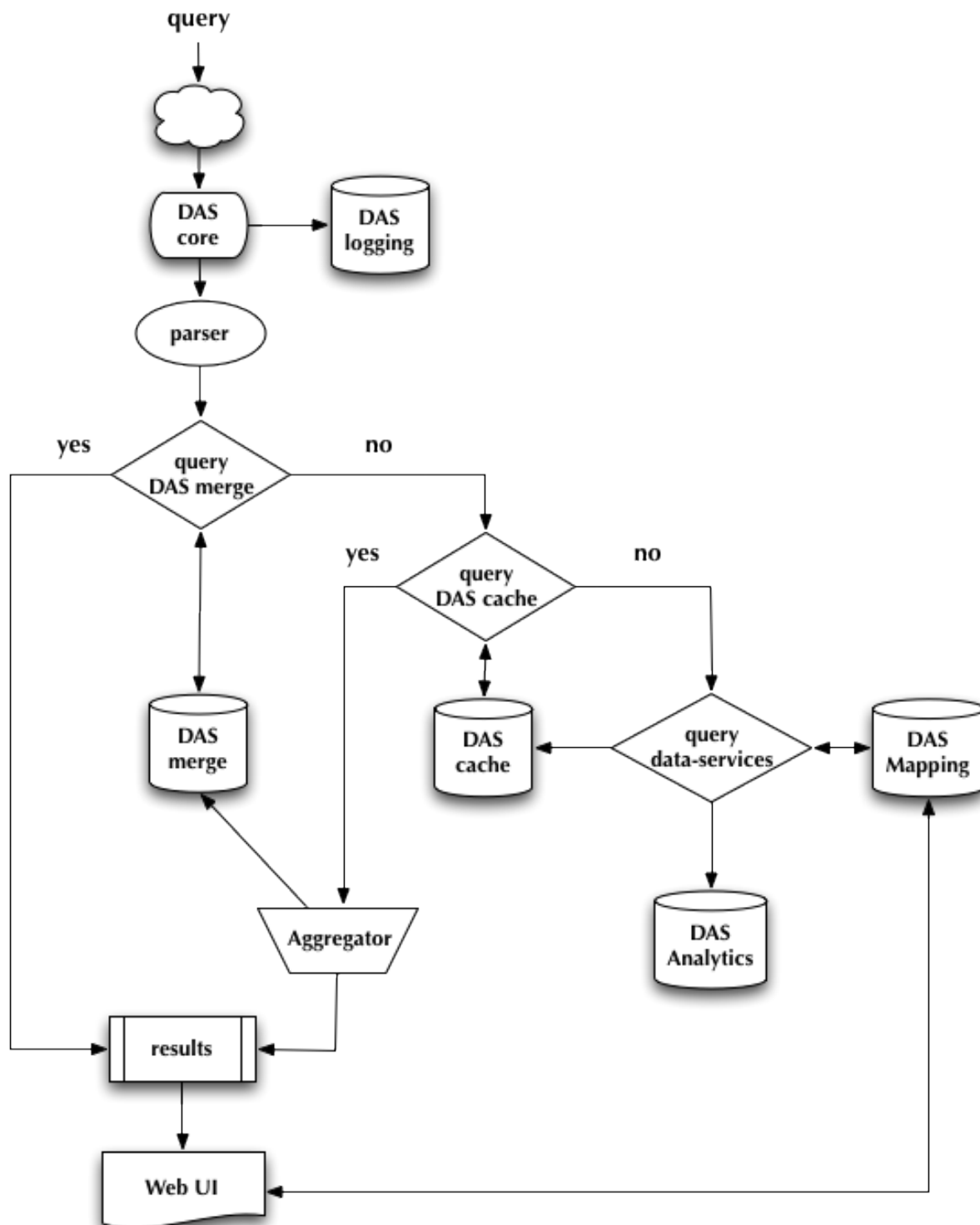
1.2.2 Query processing steps

DAS workflow consists of the following steps:

- parse input query
 - look-up for a superset of query conditions in DAS merge cache (ie, have we recently handled a query guaranteed to have returned the data being requested now)
 - * if unavailable
 - query raw cache
 - query services, to get the records that aren't available
 - write new records to raw cache
 - perform aggregation
 - write aggregation results to merged cache
 - * get results from merged cache
- present DAS records in Web UI (converting records if necessary)

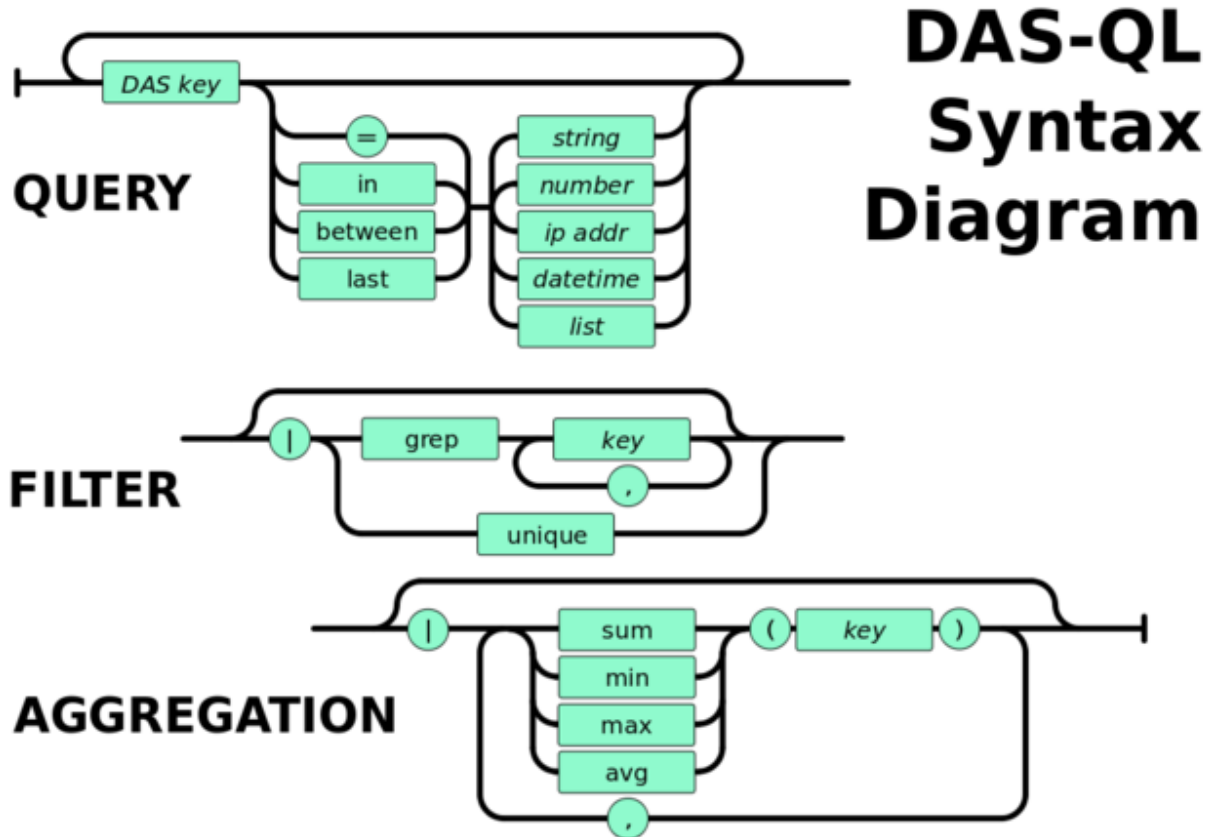
1.3 How to use DAS?

DAS can be used from web interface, through a command line interface, or accessed programically from other applications.



1.3.1 DAS Query Language

DAS Query Language (DAS-QL) provides intuitive, easy to use text based queries to underlying data-services. Its syntax is represnated in the following diagram



We can represent it via set of pipes, similar to UNIX pipes, the key look-up fields followed by filter and/or aggregator:

```
<key> <key> = <value> ... | <filter> <key.att> <op> <value>, ... | <aggregator>, ...
<key> date last <value h|m>
<key> date between [YYYYMMDD, YYYYMMDD]
```

Here the *<key>* is a valid DAS key, e.g. *dataset*, *run* and *<op>* is a valid DAS operator, e.g. *<*, *>*, *<=*, *>=*. The *<key.att>* represents a key attribute. They are deduced form the DAS records. For instance, if file record contain *size*, *id*, *md5* fields, all of them can be used as attributed of the *file* key, e.g. *file.md5*. But the attributes cannot appear in look-up part of the query.

The DAS query can be supplemented either by filters or aggregator functions. The filter consists of filter name and *key/key.attributes* value pair or *key/key.attributes* fields. The support filters are: *grep*, *unique*, *sort*. To *unique* filter does not require a *key*.

The supported aggregator functions are: *sum*, *count*, *min*, *max*, *avg*, *median*. A wild-card patterns are supported via asterisk character.

Here is just a few valid examples how to construct DAS query using file and dataset keys and file.size, dataset.name key attributes:

```

file dataset=/a/b/c | unique
file dataset=/a/b/c | grep file.name, file.size
file dataset=/a/b/c | sort file.name
file dataset=/a/b/c | grep file.name, file.size, | sort file.size
file dataset=/a/b/c | sum(file.size)
file dataset=/a/b/c | grep file.size>1 | sum(file.size)
file dataset=/a/b*

```

Special keywords

DAS has a several special keywords: *system*, *date*, *instance*, *records*.

- The *system* keyword is used to retrieve a records only from specified system (data-service), e.g. DBS.
- The *date* can be used in different queries and accepts values in YYYYMMDD format as well as can be specified as *last* value, e.g. *date last 24h*, *date last 60m*, where h, m are hours, minutes, respectively.
- The *records* keyword can be used to retrieve DAS records regardless from their content. For instance, if one user place a query *site=T1_CH_CERN**, the DAS requests data from several data-services (Phedex, SiteDB), while the output results will only show site related records. If user wants to see which other records exists in DAS cache for given parameter, he/she can use *records site=T1_CH_CERN** to do that. In that case user will get back all records (site, block records) associated with given condition.

1.3.2 Sample queries used at CMS

Here we provide concrete examples of DAS queries used in CMS.

Find primary dataset

```

primary_dataset=Cosmics
primary_dataset=Cosmics | grep primary_dataset.name
primary_dataset=Cosmics | count(primary_dataset.name)

```

Find dataset

```

dataset primary_dataset=Cosmics
dataset dataset=*Cosmic* site=T3_US_FSU
dataset release=CMSSW_2_0_8
dataset release=CMSSW_2_0_8 | grep dataset.name, dataset.nevents
dataset run=148126
dataset date=20101101

```

Find block

```

block dataset=/Wgamma/Winter09_IDEAL_V12_FastSim_v1/GEN-SIM-DIGI-RECO
block=/ExpressPhysics/Commissioning10-Express-v6/FEVT#f86bef6a-86c2-48bc-9f46-2e868c13d86e
block site=T3_US_Cornell*
block site=srm-cms.cern.ch | count(block.name), sum(block.replica.size), avg(block.replica.size), me

```

Find file

```

file dataset=/Wgamma/Winter09_IDEAL_V12_FastSim_v1/GEN-SIM-DIGI-RECO
file block=/ExpressPhysics/Commissioning10-Express-v6/FEVT#f86bef6a-86c2-48bc-9f46-2e868c13d86e
file dataset=/Wgamma/Winter09_IDEAL_V12_FastSim_v1/GEN-SIM-DIGI-RECO | grep file.name, file.size
file dataset=/Wgamma/Winter09_IDEAL_V12_FastSim_v1/GEN-SIM-DIGI-RECO | grep file.name, file.size>1500

```

```
file dataset=/Wgamma/Winter09_IDEAL_V12_FastSim_v1/GEN-SIM-DIGI-RECO | sum(file.size), count(file.na
file block=/ExpressPhysics/Commissioning10-Express-v6/FEVT* site=T2_CH_CAF
file run=148126 dataset=/ZeroBias/Run2010B-Dec4ReReco_v1/RECO
file dataset=/ExpressPhysics/Commissioning10-Express-v6/FEVT | grep file.size | max(file.size),min(f
```

Find lumi information

```
lumi file=/store/data/Run2010B/ZeroBias/RAW-RECO/v2/000/145/820/784478E3-52C2-DF11-A0CC-0018F3D0969A
```

Find parents/children of a given dataset/files

```
child dataset=/QCDpt30/Summer08_IDEAL_V9_v1/GEN-SIM-RAW
parent dataset=/QCDpt30/Summer08_IDEAL_V9_skim_hlt_v1/USER
child file=/store/mc/Summer08/QCDpt30/GEN-SIM-RAW/IDEAL_V9_v1/0000/1EAE7A08-187D-DD11-85B5-001CC47D03
parent file=/store/mc/Summer08/QCDpt30/USER/IDEAL_V9_skim_hlt_v1/0003/367E05A0-707E-DD11-B0B9-001CC47D03
```

Find information in local DBS instances

```
instance=cms_dbs_ph_analysis_02 dataset=/QCD_Pt*_TuneZ2_7TeV_pythia6/wteo-qcd_tunez2_pt*_pythia*
```

Find run information

```
run=148126
run in [148124,148126]
run date last 60m
run date between [20101010, 20101011]
run run_status=Complete
run reco_status=1
run dataset=/Monitor/Commissioning08-v1/RAW
```

Find site information

```
site=T1_CH_CERN
site=T1_CH_CERN | grep site.admin
```

Jobsummary information

```
jobsummary date last 24h
jobsummary site=T1_DE_KIT date last 24h
jobsummary user=ValentinKuznetsov
```

1.3.3 DAS Command Line Interface (CLI) tool

The DAS Command Line Interface (CLI) tool can be downloaded directly from DAS server. It is python-based tool and does not require any additional dependencies, although a python version of 2.6 and above is required. Its usage is very simple

```
Usage: das_client.py [options]
For more help please visit https://cmsweb.cern.ch/das/faq

Options:
  -h, --help                show this help message and exit
  -v VERBOSE, --verbose=VERBOSE
                           verbose output
```

```

--query=QUERY          specify query for your request
--host=HOST            host name of DAS cache server, default is
                      https://cmsweb.cern.ch
--idx=IDX              start index for returned result set, aka pagination,
                      use w/ limit (default is 0)
--limit=LIMIT          number of returned results (default is 10), use
                      --limit=0 to show all results
--format=FORMAT        specify return data format (json or plain), default
                      plain.
--threshold=THRESHOLD  query waiting threshold in sec, default is 5 minutes
--key=CKEY             specify private key file name
--cert=CERT            specify private certificate file name
--retry=RETRY          specify number of retries upon busy DAS server message
--das-headers          show DAS headers in JSON format
--base=BASE            specify power base for size_format, default is 10 (can
                      be 2)

```

The query parameter specifies an input *DAS query* <das_queries>, while the format parameter can be used to get results either in JSON or plain (suitable for cut and paste) data format. Here is an example of using das_client tool to retrieve information about dataset pattern

```

python das_client.py --query="dataset=/ZMM*/*/*"

Showing 1-10 out of 2 results, for more results use --idx/--limit options

/ZMM_14TeV/Summer12-DESIGN42_V17_SLHCTk-v1/GEN-SIM
/ZMM/Summer11-DESIGN42_V11_428_SLHC1-v1/GEN-SIM

```

And here is the same output using JSON data format, the auxiliary DAS headers are also requested:

```

python das_client.py --query="dataset=/ZMM*/*/*" --format=JSON --das-headers

{'apist': ['das_core', 'fakeDatasetPattern'],
 'ctime': 0.0015709400177,
 'data': [{'_id': '523dcd7f0ec3dc12198a44c5',
            'cache_id': ['523dcd7f0ec3dc12198a44c3'],
            'das': {'api': ['fakeDatasetPattern'],
                    'condition_keys': ['dataset.name'],
                    'expire': 1379782315.848377,
                    'instance': 'cms_dbs_prod_global',
                    'primary_key': 'dataset.name',
                    'record': 1,
                    'services': [{'dbs': ['dbs']}],
                    'system': ['dbs'],
                    'ts': 1379782015.863179},
            'das_id': ['523dcd7d0ec3dc12198a4498'],
            'dataset': [{'created_by': '/DC=ch/DC=cern/OU=computers/CN=wmagent/vocms216.cern.ch',
                        'creation_time': '2012-02-24 01:40:40',
                        'datatype': 'mc',
                        'modification_time': '2012-02-29 21:25:52',
                        'modified_by': '/DC=org/DC=doegrids/OU=People/CN=Alan Malta Rodrigues 4861',
                        'name': '/ZMM_14TeV/Summer12-DESIGN42_V17_SLHCTk-v1/GEN-SIM',
                        'status': 'VALID',
                        'tag': 'DESIGN42_V17::All'}],
            'qhash': 'e5ced95dd57a5cfel1a3126a22a85a301'},
            {'_id': '523dcd7f0ec3dc12198a44c6',

```

```

'cache_id': ['523dcd7f0ec3dc12198a44c4'],
'das': {'api': ['fakeDatasetPattern'],
        'condition_keys': ['dataset.name'],
        'expire': 1379782315.848377,
        'instance': 'cms_dbs_prod_global',
        'primary_key': 'dataset.name',
        'record': 1,
        'services': [{'dbs': ['dbs']}],
        'system': ['dbs'],
        'ts': 1379782015.863179},
'das_id': ['523dcd7d0ec3dc12198a4498'],
'dataset': [{'created_by': 'cmsprod@cmsprod01.hep.wisc.edu',
               'creation_time': '2011-12-29 17:47:25',
               'datatype': 'mc',
               'modification_time': '2012-01-05 17:40:17',
               'modified_by': '/DC=org/DC=doegrids/OU=People/CN=Ajit Kumar Mohapatra 867118',
               'name': '/ZMM/Summer11-DESIGN42_V11_428_SLHC1-v1/GEN-SIM',
               'status': 'VALID',
               'tag': 'DESIGN42_V11::All'}],
'qhash': 'e5ced95dd57a5cfe1a3126a22a85a301'}],
'incache': True,
'mongo_query': {'fields': ['dataset'],
                 'instance': 'cms_dbs_prod_global',
                 'spec': {'dataset.name': '/ZMM*/*/*'}},
'nresults': 2,
'status': 'ok',
'timestamp': 1379782017.68}

```

Using DAS CLI tool from other applications

It is possible to plug DAS CLI tool into other python applications. This can be done as following

```

from das_client import get_data

# invoke DAS CLI call for given host/query
# host: hostname of DAS server, e.g. https://cmsweb.cern.ch
# query: DAS query, e.g. dataset=/ZMM*/*/*
# idx: start index for pagination, e.g. 0
# limit: end index for pagination, e.g. 10, put 0 to get all results
# debug: True/False flag to get more debugging information
# threshold: 300 sec, is a default threshold to wait for DAS response
# ckey=None, cert=None are parameters which you can used to pass around
# your GRID credentials
# das_headers: True/False flag to get DAS headers, default is True

# please note that prior 1.9.X release the return type is str
# while from 1.9.X and on the return type is JSON

data = get_data(host, query, idx, limit, debug, threshold=300, ckey=None,
cert=None, das_headers=True)

```

Please note, that aforementioned code snippet requires to load *das_client.py* which is distributed within CMSSW. Due to CMSSW install policies the version of *das_client.py* may be quite old. If you need up-to-date *das_client.py* functionality you can follow this recipe. The code below download *das_client.py* directly from cmsweb site, compile it and use it in your application:

```

import os
import json
import urllib2
import httplib
import tempfile

class HTTPSClientHdlr(urllib2.HTTPSHandler):
    """
    Simple HTTPS client authentication class based on provided
    key/ca information
    """
    def __init__(self, key=None, cert=None, level=0):
        if level:
            urllib2.HTTPSHandler.__init__(self, debuglevel=1)
        else:
            urllib2.HTTPSHandler.__init__(self)
        self.key = key
        self.cert = cert

    def https_open(self, req):
        """Open request method"""
        #Rather than pass in a reference to a connection class, we pass in
        # a reference to a function which, for all intents and purposes,
        # will behave as a constructor
        return self.do_open(self.get_connection, req)

    def get_connection(self, host, timeout=300):
        """Connection method"""
        if self.key:
            return httplib.HTTPSConnection(host, key_file=self.key,
                                           cert_file=self.cert)

        return httplib.HTTPSConnection(host)

class DASClient(object):
    """DASClient object"""
    def __init__(self, debug=0):
        super(DASClient, self).__init__()
        self.debug = debug
        self.get_data = self.load_das_client()

    def get_das_client(self, debug=0):
        "Download das_client code from cmsweb"
        url = 'https://cmsweb.cern.ch/das/cli'
        ckey = os.path.join(os.environ['HOME'], '.globus/userkey.pem')
        cert = os.path.join(os.environ['HOME'], '.globus/usercert.pem')
        req = urllib2.Request(url=url, headers={})
        if ckey and cert:
            hdlr = HTTPSClientHdlr(ckey, cert, debug)
        else:
            hdlr = urllib2.HTTPHandler(debuglevel=debug)
        opener = urllib2.build_opener(hdlr)
        fdesc = opener.open(req)
        cli = fdesc.read()
        fdesc.close()
        return cli

    def load_das_client(self):
        "Load DAS client module"

```

```
cli = self.get_das_client()
# compile python code as exec statement
obj = compile(cli, '<string>', 'exec')
# define execution namespace
namespace = {}
# execute compiled python code in given namespace
exec obj in namespace
# return get_data object from namespace
return namespace['get_data']

def call(self, query, idx=0, limit=0, debug=0):
    "Query DAS data-service"
    host = 'https://cmsweb.cern.ch'
    data = self.get_data(host, query, idx, limit, debug)
    if isinstance(data, basestring):
        return json.loads(data)
    return data

if __name__ == '__main__':
    das = DASClient()
    query = "/ZMM*/*/*"
    result = das.call(query)
    if result['status'] == 'ok':
        nres = result['nresults']
        data = result['data']
        print "Query=%s, #results=%s" % (query, nres)
        print data
```

Here we provide a simple example of how to use `das_client` to find dataset summary information.

```
# PLEASE NOTE: to use this example download das_client.py from
# cmsweb.cern.ch/das/cli

# system modules
import os
import sys
import json

from das_client import get_data

def drop_das_fields(row):
    "Drop DAS specific headers in given row"
    for key in ['das', 'das_id', 'cache_id', 'qhash']:
        if row.has_key(key):
            del row[key]

def get_info(query):
    "Helper function to get information for given query"
    host = 'https://cmsweb.cern.ch'
    idx = 0
    limit = 0
    debug = False
    data = get_data(host, query, idx, limit, debug)
    if isinstance(data, basestring):
        dasjson = json.loads(data)
    else:
        dasjson = data
    status = dasjson.get('status')
```



```

    if status == 'ok':
        data = dasjson.get('data')
        return data

def get_datasets(query):
    """Helper function to get list of datasets for given query pattern"""
    for row in get_info(query):
        for dataset in row['dataset']:
            yield dataset['name']

def get_summary(query):
    """
    Helper function to get dataset summary information either for a single
    dataset or dataset pattern
    """
    if query.find('*') == -1:
        print "\n### query", query
        data = get_info(query)
        for row in data:
            drop_das_fields(row)
            print row
    else:
        for dataset in get_datasets(query):
            query = "dataset=%s" % dataset
            data = get_info(query)
            print "\n### dataset", dataset
            for row in data:
                drop_das_fields(row)
                print row

if __name__ == '__main__':
    # query dataset pattern
    query = "dataset=/ZMM*/*/*"
    # query specific dataset in certain DBS instance
    query = "dataset=/8TeV_T2tt_2j_semilepts_200_75_FSim526_Summer12_minus_v2/alkalogue-MG154_START52
    get_summary(query)

```

1.4 Installation

1.4.1 Quick installation with pip and virtualenv

Basically on a Debian-like system you may just copy&paste all the commands below into a bash.

```

# get or install virtualenv, see http://www.virtualenv.org/en/latest/virtualenv.html#installation
# easiest is through apt-get or "pip install virtualenv" if available
sudo apt-get install python-virtualenv

# create virtual env
virtualenv dasenv # will use dir "dasenv"

# activate the virtualenv -- has to be run in every new shell
source dasenv/bin/activate

# get DAS source code
git clone git://github.com/dmwm/DAS.git

```

```
cd DAS

# install dependencies
export PYCURL_SSL_LIBRARY=gnutls # depending on your distro try: openssl or gnutls
pip install -r requirements_freezed.txt
python -m nltk.downloader -e words stopwords wordnet

# you also need to connect to or install a MongoDB server, e.g.
sudo apt-get install mongodb-server
# set MongoDB port to 8230 in /etc/mongodb.conf or change dasconfig
# sed -i -e 's/8230/27017/g' etc/das.cfg
# sed -i -e 's/8230/27017/g' bin/das_db_import

# install DAS
python setup.py install
source ./init_env.sh

# initialize database with (default) service mappings
das_create_json_maps src/python/DAS/services/maps
das_update_database src/python/DAS/services/maps/das_maps_dbs_prod.js

# download YUI (the location is set in init_env.sh)
(curl -o yui.zip -L http://yuilibrary.com/downloads/yui2/yui_2.9.0.zip && \
 unzip yui.zip && mkdir -p $YUI_ROOT && cp -R yui/* $YUI_ROOT/ )

# run the tests
source ./init_env.sh
touch /tmp/x509up_u$UID # or use grid-proxy-init if you require certs...
python setup.py test

# finally start das server
# P.S. don't forget "source dasenv/bin/activate" when restarting in a new shell
source ./init_env.sh
bash das_server start
# Finally you can access it at: http://localhost:8212/das/
```

If, while running tests or starting DAS, you experience problems like “pycurl: libcurl link-time ssl backend (gnutls) is different from compile-time ssl backend (openssl)”, try reinstalling (pycurl) using a different PYCURL_SSL_LIBRARY, e.g.:

```
PYCURL_SSL_LIBRARY=openssl
pip uninstall pycurl && pip install -r requirements_freezed.txt
```

For more details continue reading below.

Source code

DAS source code is freely available from [\[DAS\]](#) github repository. To install it on your system you need to use *git* which can be found from [\[GIT\]](#) web site.

1.5 Dependencies

DAS is written in python and relies on standard python modules. The design back-end is [MongoDB](#), which provides *schema-less* storage and a generic query language.

DAS depends on the following software:

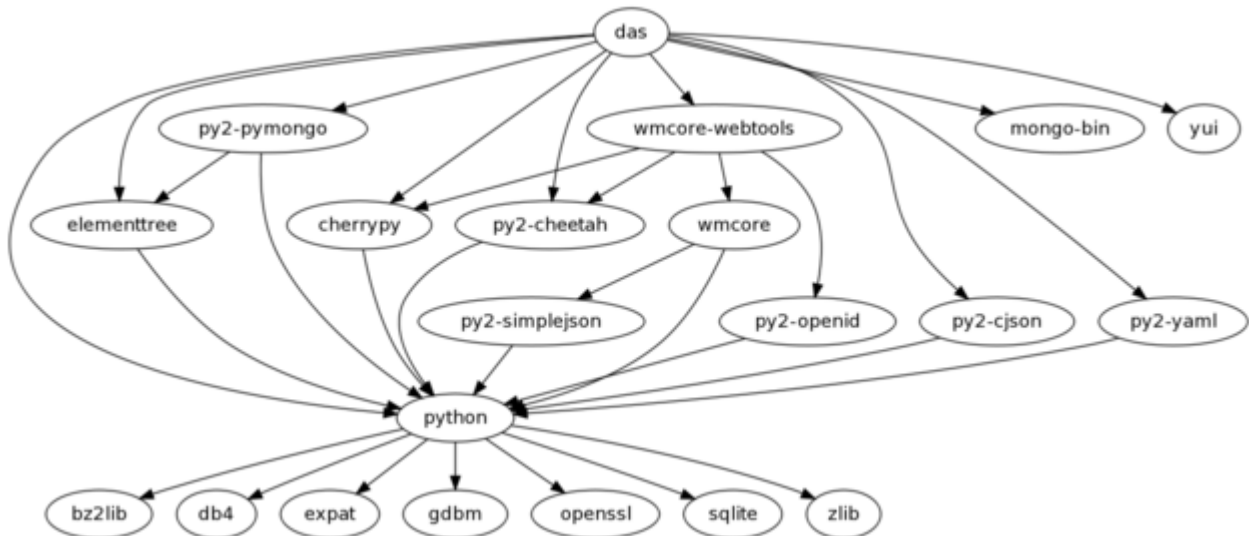
- MongoDB and pymongo module
- libcurl library
- YUI library (Yahoo UI, version 2)
- python modules:
 - yajl (Yet Another JSON library) or cJSON (C-JSON module)
 - CherryPy
 - Cheetah
 - PLY
 - PyYAML
 - pycurl

To install MongoDB visit [\[Mongodb\]](#) (make its bin directory available in your path). To install libcurl library visit [\[CURL\]](#). To install YUI library, visit Yahoo developer web site [\[YUI\]](#) and install *version 2* of their yui library.

To install python dependencies it is easier to use standard python installer *pip* (see above).

Dependencies in the CMS Environment

Below you can see current set of dependencies for DAS within the CMS environment:



The *wmcore* and *wmcore-webtools* modules are a CMS general purpose web framework based on CherryPy. It is possible to run DAS both within this framework and in an entirely standalone version using only CherryPy.

Below we list all dependencies clarifying their role for DAS

- *python*, DAS is written in python (2.6), see [\[Python\]](#);

- *cherrypy*, a generic python web framework, see [\[CPF\]](#);
- *yui* the Yahoo YUI Library for building richly interactive web applications, see [\[YUI\]](#);
- *elementtree* and its *cElementTree* counterpart are used as generic XML parser in DAS, both implementations are now part of python (2.5 and above);
- *MongoDB*, a document-oriented database, the DAS DB back-ends, see [\[Mongodb\]](#) and [\[MongodbOverview\]](#);
- *pymongo*, a MongoDB python driver, see [\[Pymongo\]](#);
- *yaml*, a human-readable data serialization format (and superset of JSON), a python YAML library is used for DAS maps and server configurations, see [\[YAML\]](#), [\[PyYAML\]](#);
- *Cheetah*, a python template framework, used for all DAS web templates, see [\[Cheetah\]](#);
- *sphinx*, a python documentation library servers all DAS documentation, see [\[Sphinx\]](#);
- *ipython*, an interactive python shell (optional, used in some admin tools), see [\[IPython\]](#);
- *cjson*, a C library providing a faster JSON decoder/encoder for python (optional), see [\[CJSON\]](#);
- *yajl*, a C library providing a faster JSON decoder/encoder for python (optional), see [\[YAJL\]](#);
- *pycurl* and *curl*, python module for libcurl library, see [\[PyCurl\]](#), [\[CURL\]](#);
- *PLY*, python Lexer and Yacc, see [\[PLY\]](#);

1.6 DAS and KWS maps

DAS relies on its maps to understand underlying data-services. The DAS tree contains set of YML files located at `services/{cms_maps,maps}` areas. These YML files contain information about underlying data-service APIs.

In order to generate DAS and KWS maps we provide series of tools which we'll describe here.

1.7 Create DAS and KWS maps

DAS and KWS maps can be generated via *das_create_json_maps* and *das_create_kws_maps*, respectively. KWS maps are generated via series of DAS queries (please adjust script itself if you need to change KWS coverage).

All maps will be stored in associated JS (JavaScript) files defined in aforementioned scripts. These files can be uploaded to github.com/dmwm/DASMaps repository and/or used to upload them into MongoDB.

1.8 Validate DAS and KWS maps

It is important to validate DAS and KWS maps since they have pre-defined structure. Moreover each map stores a hash of the map itself which is checked during validation. To perform validation please use *das_js_validate* script.

1.9 Import DAS and KWS maps

DAS and KWS maps need to be imported into MongoDB. This can be done via *das_js_import* script which accept location of DAS maps on local filesystem.

1.10 Fetch DAS and KWS maps

DAS maps can be fetched from github.com/dmwm/DASMaps repository via `das_js_fetch` script. These maps are good to go and uploaded to this repository by developers.

1.11 Setting up and customizing DAS installation

First set basic configuration on database to use and other server parameters.

To integrate DAS with your own custom services the biggest and most time consuming task is preparing the service mappings which define the services to be integrated.

1.11.1 DAS configuration file

DAS configuration consists of a single file, `$DAS_ROOT/etc/das.cfg`. Its structure is shown below:

```
[das]                                # DAS core configuration
verbose = 0                          # verbosity level, 0 means lowest
parserdir = /tmp                     # DAS PLY parser cache directory
multitask = True                     # enable multitasking for DAS core (threading)
core_workers = 10                    # number of DAS core workers who contact data-providers
api_workers = 2                      # number of API workers who run simultaneously
thread_weights = 'dbs:3','phedex:3' # thread weight for given services
error_expire = 300                   # expiration time for error records (in seconds)
emptyset_expire = 5                  # expiration time for empty records (in seconds)
services = dbs,phedex                # list of participated data-providers

[cacherequests]
Admin = 50                           # number of queries for admin user role
Unlimited = 10000                     # number of queries for unlimited user role
ProductionAccess = 5000              # number of user for production user role

[web_server]                         # DAS web server configuration parameters
thread_pool = 30                     # number of threads for CherryPy
socket_queue_size = 15               # queue size for requests while server is busy
host = 0.0.0.0                       # host IP, the 0.0.0.0 means visible everywhere
log_screen = True                    # print log to stdout
url_base = /das                      # DAS server url base
port = 8212                          # DAS server port
pid = /tmp/logs/dsw.pid              # DAS server pid file
status_update = 2500                 #
web_workers = 10                     # Number of DAS web server workers who handle user requests
queue_limit = 200                    # DAS server queue limit
adjust_input = True                  # Adjust user input (boolean)
dbs_daemon = True                    # Run DBSDaemon (boolean)
dbs_daemon_interval = 300            # interval for DBSDaemon update in sec
dbs_daemon_expire = 3600             # expiration timestamp for DBSDaemon records
hot_threshold = 100                  # a hot threshold for powerful users
onhold_daemon = True                 # Run onhold daemon for queries which put on hold after hot threshold

[dbs]                                # DBS server configuration
dbs_instances = prod,dev             # DBS instances
dbs_global_instance = prod           # name of the global DBS instance
dbs_global_url = http://a.b.c        # DBS data-provider URL
```

```
[mongodb]                # MongoDB configuration parameters
dburi = mongodb://localhost:8230 # MongoDB URI
bulkupdate_size = 5000    # size of bulk insert/update operations
dbname = das              # MongoDB database name
lifetime = 86400          # default lifetime (in seconds) for DAS records

[dasdb]                   # DAS DB cache parameters
dbname = das              # name of DAS cache database
cachecollection = cache  # name of cache collection
mergecollection = merge  # name of merge collection
mrcollection = mapreduce # name of mapreduce collection

[loggingdb]
capped_size = 104857600
collname = db
dbname = logging

[analyticsdb]             # AnalyticsDB configuration parameters
dbname = analytics        # name of analytics database
collname = db             # name of analytics collection
history = 5184000         # expiration time for records in an/ collection (in seconds)

[mappingdb]               # MappingDB configuration parameters
dbname = mapping          # name of mapping database
collname = db             # name of mapping collection

[parserdb]                # parserdb configuration parameters
dbname = parser           # parser database name
enable = True             # use it in DAS or not (boolean)
collname = db             # collection name
sizecap = 5242880         # size of capped collection

[dbs_phedex]              # dbs_phedex configuration parameters
expiration = 3600         # expiration time stamp
urls = http://dbs,https://phedex # DBS and Phedex URLs
```

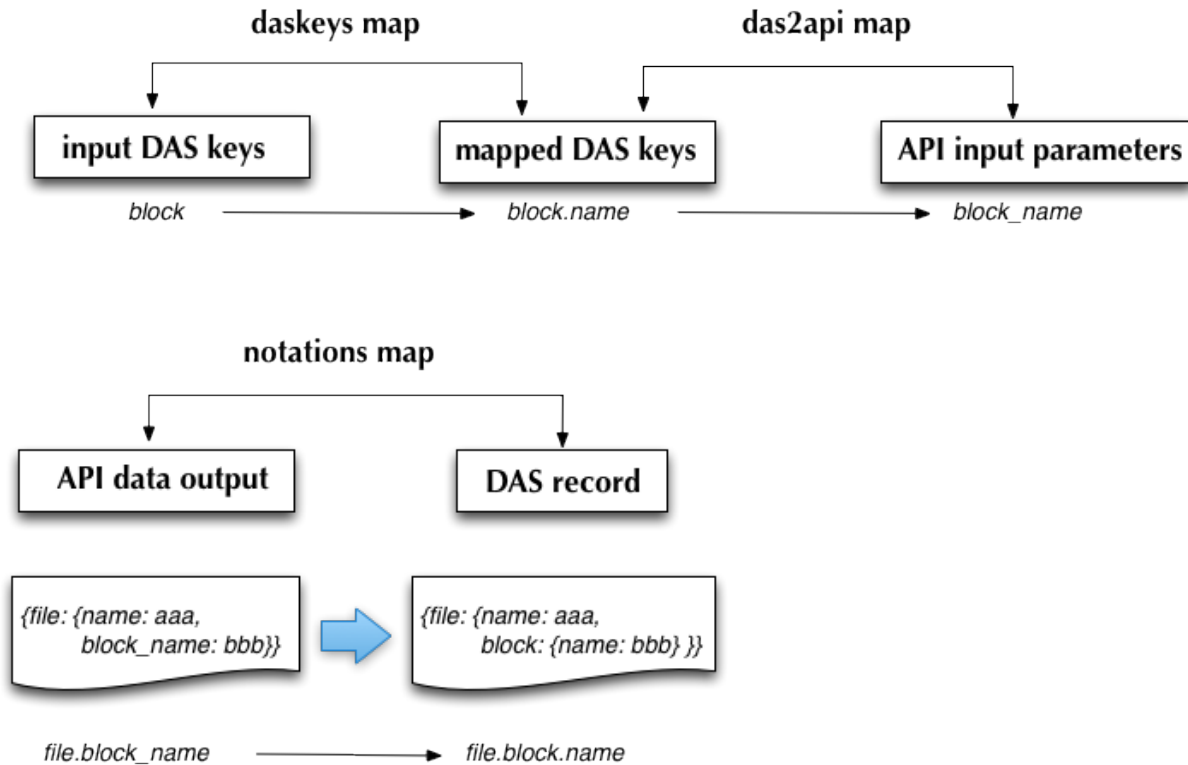
For up-to-date configuration parameters please see *utils/das_config.py*

1.11.2 DAS Mapping DB

DAS uses Query Language (QL) to look-up data from data-providers as well as its own cache. The data provided by various data services can come in variety of form and data formats, while DAS cache stores data records in **JSON** data format. Therefore we need to define certain mappings between DAS QL and data-provider API calls as well as DAS QL and data records in DAS cache. To serve this goal DAS relies on its Mapping DB which holds information about all the data-service APIs which are used by DAS, and the necessary mappings between DAS and API records

Each mapping file holds the following schema:

- **system**, the name of data-provider
- **format**, the data format used by data-provider, e.g. XML or JSON
- series of maps for APIs used in DAS workflow, where each API map has the following entries
 - **urn**, API alias name
 - **url**, the API URL
 - **params**, the set of input parameters for API in question



- **lookup**, the name of DAS look-up key the given API serves
- **das_map**, the list of maps which covers DAS QL keys, record names and API input argument, along with optional pattern; every map has the following parameters
 - * **das_key**, the name of DAS QL key, e.g. run
 - * **rec_key**, the DAS record key name, e.g. run.number
 - * **api_arg**, the corresponding API input argument name, e.g. run_number
 - * **pattern**, optional regex for input argument
- **wild_card**, the optional notation for wild-card usage in given API, e.g. * or %
- ckey and cert, the path to GRID credentials
- **notations** map which transforms API output into das record keys; this map consists of maps with the following structure
 - **api_output**, the key name returned by API record
 - **das_key**, the DAS record key
 - **api**, the name of API this mapping should be applied for

DAS also uses presentation map to translate DAS records into human readable form, e.g. to translate *file.nevents* into *Number of events*

The DAS maps use [YAML](#) data-format. Each file may contain several data-service API mappings, as well as auxiliary information about data-provider, e.g. data format, expiration timestamp, etc. For example here is a simple mapping file for google map APIs

```
system : google_maps
format : JSON
---
urn : google_geo_maps
url : "http://maps.google.com/maps/geo"
expire : 30
params : { "q" : "required", "output": "json" }
lookup : city
das_map : [
    { "das_key": "city", "rec_key": "city.name", "api_arg": "q" },
]
---
urn : google_geo_maps_zip
url : "http://maps.google.com/maps/geo"
expire : 30
params : { "q" : "required", "output": "json" }
lookup : zip
das_map : [
    { "das_key": "zip", "rec_key": "zip.code", "api_arg": "q" },
]
---
notations : [
    { "api_output": "zip.name", "rec_key": "zip.code", "api": "" },
    { "api_output": "name", "rec_key": "code", "api": "google_geo_maps_zip" },
]
```

As you can see it defines the data-provider name, `google_maps` (DAS call it system), the data format `JSON` used by this data-provider as well as three maps (for each API usage), separated by tripple dashes. The first one defines mapping for geo location for a given city key, the second defines geo location for a given zip key and mapping for notations used by DAS workflow. In particulat, DAS will map `zip.name` into `zip.code` for any api, and `name` into `code` for `google_geo_maps_zip` api (the meaning of these translation will become clear when we will discuss concrete example below).

As you may noticed, every mapping (the code between tripple dashes) has repeated strucute. It defines *urn*, *url*, *expire*, *params*, *lookup*, *das_map* values. The *urn* stands for uniform resource name, this alias is used by DAS to distinguish APIs and their usage pattern, the *url* is canonical URL for API in question, the *params* defines a dictionary of input parameters accepted by API, the *das_map* is mapping from DAS keys into DAS data records, and finally, *lookup* is the name of DAS key this map is designed for.

To accommodate different use cases of API usage the *params* structure may contain three types of parameter values: the **default**, **required** and **optional** values. The default value will be passed to API *as is*, the required value must be substituted by DAS workflow (it will be taken from the query provided by DAS user, if it will not be provide the API call will be discarded from DAS workflow) and the optional value which can be skipped by API call.

Example

In this section we show a concrete example of mappings used by DAS workflow for one of the data-services. Let's take the following DAS queries:

```
file file=X
file file=X status=VALID
```

These queries will correspond to the following DAS record structure:

```
{"file" : {"name": "X", "size":1, "nevents": 10, ...}}
```


The dots just indicate that structure can be more complex.

The file DAS key is mapped into `file.name` key.attribute value within DAS record, here the period divides key from attribute in aforementioned dictionary. Therefore `file.name` value is X, `file.size` value is 1, etc.

Here is an example of one of the DAS mapping records which can serve discussed DAS queries (please note that it may be several data-services which may provide the data for given DAS query).

```
urn: files
url : "https://cmsweb.cern.ch/dbs/prod/global/DBSReader/files/"
expire : 900
params : {
    "logical_file_name": "required",
    "detail": "True",
    "status": "optional",
}
lookup : file
das_map : [
    {"das_key": "file", "rec_key": "file.name", "api_arg": "logical_file_name",
    "pattern": "/.*.root"},
    {"das_key": "status", "rec_key": "status.name", "api_arg": "status",
    "pattern": "(VALID|INVALID)"},
]
```

This record defines files API with given URL and expire timestamp. It specifies the input parameters (params), in particular, `logical_file_name` is required by this API, the `detail` has default value `True` and `status` is an optional input parameter. The daskeys mapping defines mapping between DAS keys used by end-user and DAS record keys. For example

```
file file=X
```

will be mapped into the following API call:

```
https://cmsweb.cern.ch/dbs/prod/global/DBSReader/files?logical_file_name=X&detail=True
```

while:

```
file file=X status=VALID
```

will be mapped into:

```
https://cmsweb.cern.ch/dbs/prod/global/DBSReader/files?logical_file_name=X&detail=True&status=VALID
```

In both case, the data-provider will return back the following data-record, e.g.:

```
{"logical_file_name": "X", "size": 1, ...}
```

therefore we need another mapping from API data record into expected DAS record structure (as we discussed above):

```
{"file": {"name": "X", "size": 1, ...}}
```

To perform such translation DAS workflow consults `das2api` maps which defines them, e.g. `logical_file_name` maps into `file.name`, etc.

Sometimes, different data-services provides data records who have different notations, e.g. `fileName`, `file_name`, etc. To accommodate this differences DAS consults notation map to perform translation from one into another notation.

Finally, to translate DAS records into human readable form we need another mapping, the presentation one. It defines what should be presented at DAS UI level for a given DAS record. For example, we may want to display “File name” at DAS UI, instead of showing `file.name`. To perform this translation DAS uses presentation map.

1.11.3 How to add new data-service

DAS has pluggable architecture, so adding a new CMS data-service should be a relatively easy procedure. Here we discuss two different ways to add a new service into DAS.

Plug and play interface

This work is in progress.

A new data-service can register with DAS by providing a file describing the interface and available APIs. This configuration includes the data-service URL, the data format provided, an optional expiration timestamp for the data, the API name, necessary parameters and optional mapping onto DAS keys.

A new DAS interface will allow this information to be added via a simple configuration file. The data-service configuration files should be presented in [\[YAML\]](#) data-format.

An example configuration follows ¹:

```
# SiteDB API mapping to DAS
system : sitedb
format : JSON

# API record
---
# URI description
urn : CMSNametoAdmins
url : "https://a.b.com/sitedb/api"
params : {'name':''}
expire : 3600 # optional DAS uses internal default value

# DAS keys mapping defines mapping between query names, e.g. run,
# and its actual representation in DAS record, e.g. run.number
daskeys : [
  {'key':'site', 'map':'site.name', 'pattern':''},
  {'key':'admin', 'map':'email', 'pattern':''}
]

# DAS search keys to API input parameter mapping
das2api : [
  {'das_key':'site', 'api_param':'se', 'pattern':''}
]
---
# next API
---
# APIs notation mapping maps data-service output into
# DAS syntax, e.g
# {'site_name':'abc'} ==> {'site':{'name':'abc'}}
notation : [
  {'notation':'site_name', 'map':'site.name', 'api':''}
]
```

¹ This example demonstrates flexibility of YAML data-format and shows different representation styles.

The syntax consists of key:value pairs, where value can be in a form of string, list or dictionary. Hash sign (#) defines a comment, the three dashes (—) defines the record separator. Each record starts with definition of system and data format provided by data-service.

```
# comment
system: my_system_name
format: XML
```

Those definitions will be applied to each API defined later in a map file. The API section followed after the record separator and should define: *urn*, *url*, *expire*, *params* and *daskeys*.

```
# API section
---
urn: api_alias
url: "http://a.b.com/method"
expire: 3600 # in seconds
params: {} # dictionary of data-service input parameters
daskeys: [{}, {}] # list of dictionaries for DAS key maps
```

- the *urn* is the API name or identifier (any name different from the API name itself) and used solely inside of DAS
- the *url* defines the data-service URL
- the *params* are data-service input parameters
- the *daskeys* is a list of maps between data-service input parameters and DAS internal key representation. For instance when we say *site* we might mean site CMS name or site SE/CE name. So the DAS key will be *site* while DAS internal key representation may be *site.name* or *site.sename*. So, each entry in *daskeys* list is defined as the following dictionary: {'key':value, 'map':value, 'pattern':''}, where pattern is a regular expression which can be used to differentiate between different arguments where they have different structures.
- the (optional) *das2api* map defines mapping between DAS internal key and data-service input parameter. For instance, *site.name* DAS key can be mapping into *_name_* data-service input parameter.

The next API record can be followed by the next record separator, e.g.

```
---
# API record 1
urn: api_alias1
url: "http://a.b.com/method1"
expire: 3600 # in seconds
params: {} # dictionary of data-service input parameters
daskeys: [{}, {}] # list of dictionaries for DAS key maps
---
# API record 2
urn: api_alias2
url: "http://a.b.com/method2"
expire: 1800 # in seconds
params: {} # dictionary of data-service input parameters
daskeys: [{}, {}] # list of dictionaries for DAS key maps
```

At the end of DAS map there is an optional *notation* mapping, which defines data-service output mapping back into DAS internal key representation (including converting from flat to hierarchical structures if necessary).

```
---
# APIs notation mapping maps data-service output into
# DAS syntax, e.g
# {'site_name':'abc'} ==> {'site':{'name':'abc'}}
```

```
notation : [
    {'notation': 'site_name', 'map': 'site.name', 'api': ''}
]
```

For instance, if your data service returns `runNumber` and in DAS we use `run_number` you'll define this mapping in *notation* section.

To summarize, the YAML map file provides

- system name
- underlying data format used by this service for its meta-data
- the list of apis records, each record contains the following:
 - urn name, DAS will use it as API name
 - url of data-service
 - expiration timestamp (how long its data can live in DAS)
 - input parameters, provide a dictionary
 - list of daskeys, where each key contains its name *key*, the mapping within a DAS record, *map*, and appropriate pattern
 - list of API to DAS notations (if any); different API can yield data in different notations, for instance, `siteName` and `site_name`. To accommodate these syntactic differences we use this mapping.
- notation mapping between data-service provider output and DAS

Add new service via API

You can manually add new service by extending `DAS.services.abstract_service.DASAbstractService` and overriding the *api* method.

To do so we need to create a new class inherited from `DAS.services.abstract_service.DASAbstractService`.

```
class MyDataService(DASAbstractService):
    """
    Helper class to provide access to MyData service
    """
    def __init__(self, config):
        DASAbstractService.__init__(self, 'mydata', config)
        self.map = self.dasmapping.servicemap(self.name)
        map_validator(self.map)
```

optionally the class can override `.. function:: def api(self, query)` method of `DAS.services.abstract_service.DASAbstractService` Here is an example of such an implementation

```
def api(self, query):
    """My API implementation"""
    api      = self.map.keys()[0] # get API from internal map
    url      = self.map[api]['url']
    expire   = self.map[api]['expire']
    args     = dict(self.map[api]['params']) # get args from internal map
    time0    = time.time()
    dasrows = function(url, args) # get data and convert to DAS records
```

```
ctime    = time.time() - time0
self.write_to_cache(query, expire, url, api, args, dasrows, ctime)
```

The hypothetical function call should contact the data-service and fetch, parse and yield data. Please note that we encourage the use of python generators [\[Gen\]](#) in function implementations.

1.11.4 DAS CMS Operations

Here we outline CMS specific operations to build/install and maintain DAS system. Please note, all instructions below refer to \$DAS_VER as DAS version.

Building RPM

For generic CMS build procedure please refer to this [page](#). We build SLC5 RPMs on vocms82 build node using 64-bit architecture. Here is a list of build steps:

```
# clean the area
rm -rf bin bootstraptmp build BUILD cmsset_default.* \
    common RPMs slc5_amd64_gcc434 SOURCES SPECS SRPMS tmp var

# get PKGTOOLS and CMSDIST
cvs co -r $DAS_VER PKGTOOLS
cvs co -r dg20091203b-comp-base CMSDIST
cvs update -r 1.59 CMSDIST/python.spec

# update/modify appropriate spec's, e.g.
# cvs up -A CMSDIST/das.spec

# perform the build
export SCRAM_ARCH=slc5_amd64_gcc434
PKGTOOLS/cmsBuild --architecture=$SCRAM_ARCH --cfg=./build.cfg
```

The build.cfg file has the following content:

```
[globals]
assumeYes: True
onlyOnce: True
testTag: False
trace: True
tag: cmp
repository: comp
[bootstrap]
priority: -30
doNotBootstrap: True
repositoryDir: comp
[build das]
compilingProcesses: 6
workersPoolSize: 2
priority: -20
#[upload das]
#priority: 100
#syncBack: True
```

Build logs and packages are located in BUILD area. Once build is complete all CMS RPMs are installed under slc5_amd64_gcc434 area. Please verify that DAS RPMs has been installed. Once everything is ok, request from CERN operator to upload DAS RPMs into CMS COMP repository. It can be done using the following set of commands:

```
eval `ssh-agent -s`  
ssh-add -t 36000  
export SCRAM_ARCH=slc5_amd64_gcc434  
PKGTOOLS/cmsBuild --architecture=$SCRAM_ARCH --cfg=./upload.cfg
```

where upload.cfg is similar to build.cfg with last three lines commented out.

Installing RPMs

DAS follows CMS build/install generic procedure. Here we outline all necessary steps to install CMS DAS RPMs into your area

```
export PROJ_DIR=/data/projects/das # modify accordingly  
export SCRAM_ARCH=slc5_amd64_gcc434  
export APT_VERSION=0.5.15lorg3.2-cmp  
export V=$DAS_VER  
  
mkdir -p $PROJ_DIR  
  
wget -O$PROJ_DIR/bootstrap.sh http://cmsrep.cern.ch/cmssw/cms/bootstrap.sh  
chmod +x $PROJ_DIR/bootstrap.sh  
# perform this step only once  
$PROJ_DIR/bootstrap.sh -repository comp -arch $SCRAM_ARCH -path $PROJ_DIR setup  
cd $PROJ_DIR  
source $SCRAM_ARCH/external/apt/$APT_VERSION/etc/profile.d/init.sh  
  
apt-get update  
apt-get install cms+das+$V
```

Updating RPM

Please following these steps to update DAS RPM in your area:

```
export PROJ_DIR=/data/projects/das  
export SCRAM_ARCH=slc5_amd64_gcc434  
export APT_VERSION=0.5.15lorg3.2-cmp  
export V=$DAS_VER  
  
cd $PROJ_DIR  
source $SCRAM_ARCH/external/apt/$APT_VERSION/etc/profile.d/init.sh  
  
apt-get update  
apt-get install cms+das+$V
```

setup.sh

In order to run DAS installed via CMS RPMs you need to setup your environment. Here we provide a simple steps which you can follow to create a single setup.sh file and use it afterwards:

```

echo "ver=$V" > $PROJ_DIR/setup.sh
echo "export PROJ_DIR=$PROJ_DIR" >> $PROJ_DIR/setup.sh
echo "export SCRAM_ARCH=$SCRAM_ARCH" >> $PROJ_DIR/setup.sh
echo -e "export APT_VERSION=$APT_VERSION\n" >> $PROJ_DIR/setup.sh
echo 'source $SCRAM_ARCH/external/apt/$APT_VERSION/etc/profile.d/init.sh' >> $PROJ_DIR/setup.sh
echo 'source $PROJ_DIR/slc4_ia32_gcc345/cms/das/$ver/etc/profile.d/init.sh' >> $PROJ_DIR/setup.sh

```

1.11.5 DAS operations

Running DAS services

DAS consists of multi-threaded DAS web and keyword search KWS servers. Please refer to [DAS CMS operations](#) section for deployment instructions and [DAS configuration](#) section for description of configuration parameters.

By CMS conventions DAS uses the following ports

- 8212 for DAS web server
- 8214 for DAS KWS server

In order to start each of server you need to setup your environment, see [setup.sh](#) file. For that use the following steps:

```

cd /data/projects/das/ # change accordingly
source setup.sh

```

Setting up DAS maps

Data-services are registered via the service mappings (DAS Maps) which define the relationships between the services and how their inputs or outputs shall be transformed.

To initialize or update the DAS maps (along with other metadata) call the bin/das_update_database with location of maps. The CMS maps are located in \$DAS_ROOT/src/python/DAS/services/cms_maps directory.

Apache redirect rules

Here we outline Apache redirect rules which can be used to serve DAS services. Please note we used localhost IP, 127.0.0.1 for reference, which should be substituted with actual hostname of the node where DAS services will run.

Rewrite rules for apache configuration file, e.g. httpd.conf

```

# Rewrite rules
# uncomment this line if you compile apache with dynamic modules
#LoadModule rewrite_module modules/mod_rewrite.so
# uncomment this line if you compile your modules within apache
RewriteEngine on

# DAS rewrite rules
RewriteRule ^(/das(/.*)?)$ https://127.0.0.1$1 [R=301,L]

Include conf/extra/httpd-ssl.conf

```

Rules for SSL rewrites:

```
RewriteRule ^/das(/.*)?$ http://127.0.0.1:8212/das/$1 [P,L]
```

MongoDB server

DAS uses [MongoDB](#) as its back-end for Mapping, Analytics, Logging, Cache and Merge DBs.

RPM installation will supply its proper configuration file. You must start cache back-end before starting the DAS cache server.

MongoDB server operates via standard UNIX init script:

```
$MONGO_ROOT/etc/profile.d/mongo_init.sh start|stop|status
```

The MongoDB database is located in `$MONGO_ROOT/db`, while logs are in `$MONGO_ROOT/logs`.

DAS web server

DAS web server is multi-threaded server. It can be started using the following command

```
das_server start|stop|status|restart
```

The `das_server` should be in your path once you setup your CMS DAS environment, see [setup.sh](#), otherwise please locate it under `$DAS_ROOT/bin/` area.

It consists of N threads used by DAS web server, plus M threads used by DAS core to run data retrieval processes concurrently. DAS configuration has the following parameters

```
config.web_server.thread_pool = 30
config.web_server.socket_queue_size = 15
config.web_server.number_of_workers = 8
config.web_server.queue_limit = 20
```

The first two (`thread_pool` and `socket_queue_size`) are used by underlying CherryPy server to control its internal thread pool and queue, while second pair (`number_of_workers` and `queue_limit`) are used by DAS core to control number of worker threads used internally. The DAS core multitasking can be turned off by using

```
config.das.multitask = False
```

parameter.

DAS administration

DAS RPMs provide a set of tools for administration tasks. They are located at `$DAS_ROOT/bin`.

- `das_server` is a DAS server init script;
- `das_cli` is DAS stand-alone CLI tool, it doesn't require neither cache or web DAS servers;
- `das_code_quality.sh` is a bash script to check DAS code quality. It is based on `pylint` tool, see [\[PYLINT\]](#).
- `das_config` is a tool to create DAS configuration file;
- `das_js_fetch` script fetches DAS maps from remote github repository
- `das_js_validate` script validates DAS maps

- `das_js_import` script imports DAS and KWS maps into MongoDB
- `das_js_update` script fetches, validate and import DAS and KWS maps into MongoDB
- `das_maps_yaml2json` is a tool that generates json DB dumps from YML dasmaps.
- `das_mapreduce` is a tool to create map/reduce function for DAS;
- `das_stress_test` script provides ability to perform stress tests against DAS server

1.12 DAS architecture and internals

This section describes DAS internal modules and datastructures (including records stored in the database).

1.12.1 DAS architecture

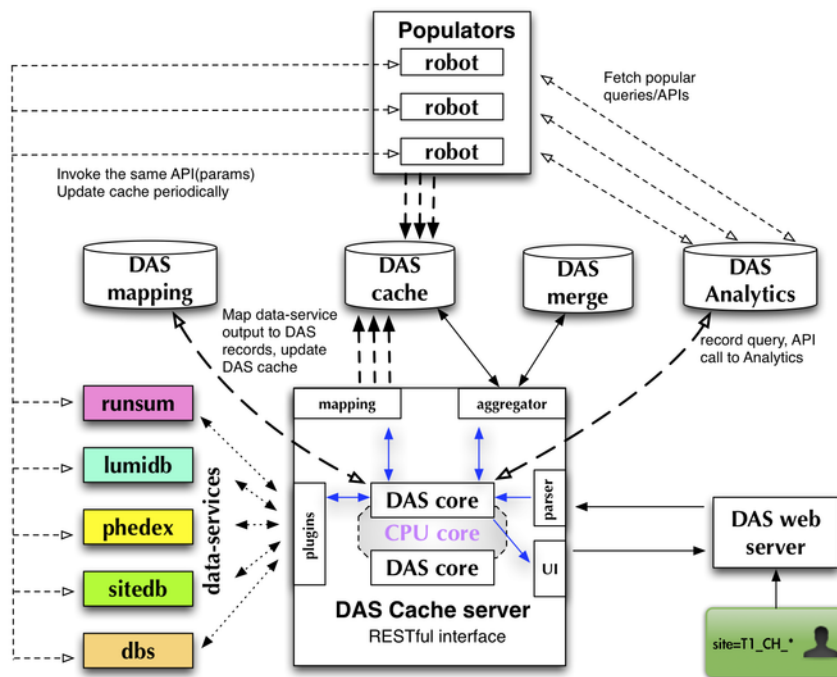
DAS architecture is based on several components:

- common core library
 - analytics DB
 - mapping DB
 - caching DB
 - merging DB
 - logging DB
- data-service plugins, each plugin contains
 - data-service handler class
 - API map
 - notation map
- cache server
- client web server

The last two components, cache and client servers, are optional. The code itself can work without cache/client servers with the CLI tool which uses core libraries. But their existence allows the introduction of DAS pre-fetch strategies, DAS robots, which can significantly improve responsiveness of the system and add multi-user support into DAS. The following picture represents current DAS architecture:

It consists of DAS web server with RESTful interface, DAS cache server, DAS Analytics/Mapping/Cache DBs and DAS robots (for pre-fetching queries). The DAS cache server uses multithreading to consume and work on several user requests at the same time. All queries are written to the DAS Analytics DB. A mapping between data-service and DAS notations is stored in the DAS Mapping DB. Communication with end-users is done via set of REST calls. User can make GET/POST/DELETE requests to fetch or delete data in DAS, respectively. The DAS workflow can be summarised as:

- DAS cache-server receives a query from the client (either DAS web server or DAS CLI)
- The input query is parsed and the selection key(s) and condition(s) identified and mapped to the appropriate data-services
- The query is added to the DAS Analytics DB
- The DAS cache is checked for existing data



- if available
 - * Data is retrieved from the cache
- otherwise:
 - * The necessary data services are queried
 - * The results are parsed, transformed and inserted into the DAS cache
 - * The user receives a message that the data is being located
- The data is formatted and returned to the user

For more information please see the [DAS workflow](#) page. The DAS DBs use the MongoDB document-oriented database (see [\[MongoDB\]](#)), although during design/evaluation process we considered a number of other technologies, such as different RDMS flavors, memcached [\[Memcached\]](#) and other key-value based data stores (eg CouchDB [\[Couchdb\]](#)), etc.

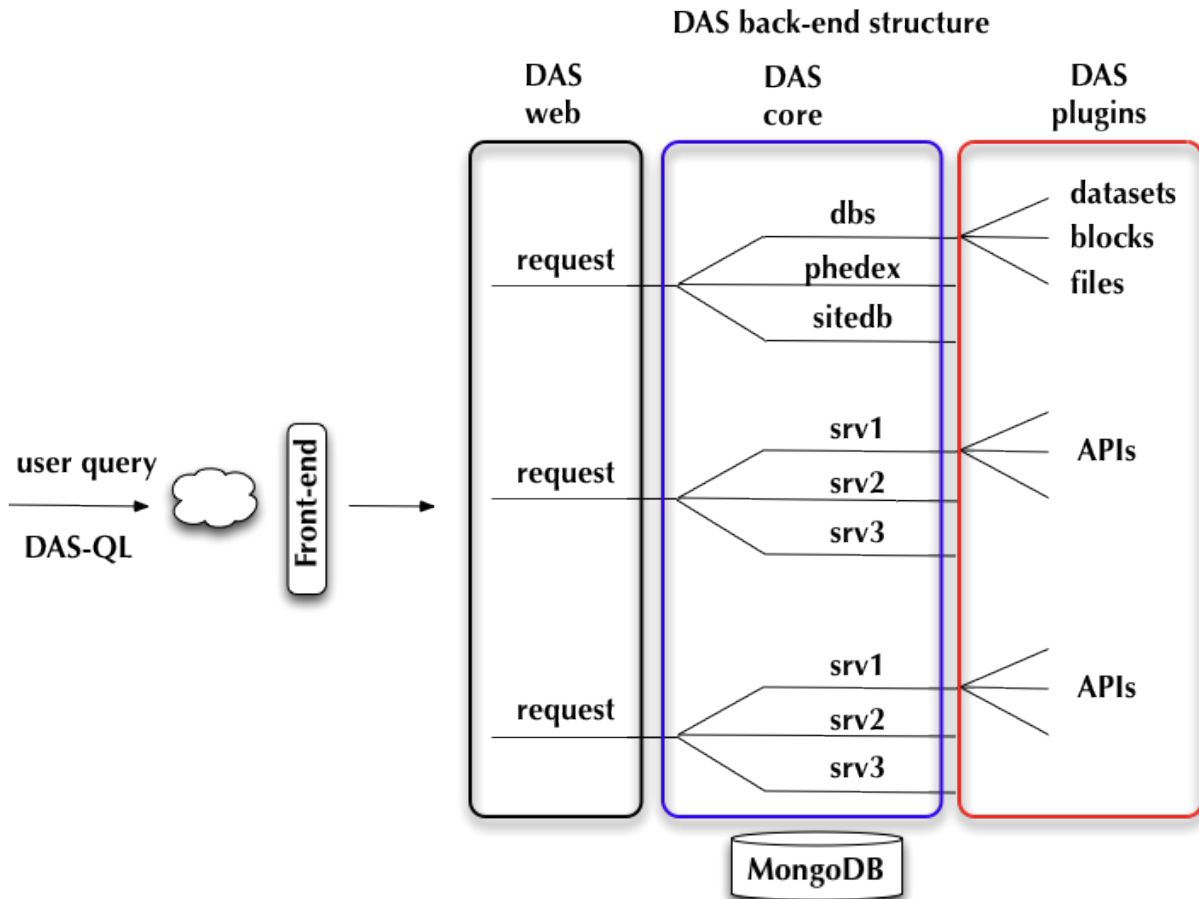
1.12.2 DAS server

DAS server is multi-threaded application which runs within CherryPy web framework. It consists of the following main componens:

- DAS web server
- DAS core engine, which by itself consist of
 - DAS abstract data-service
 - Data-provider plugins (data-service implementations)
- DAS analytics daemons

- DAS command-line client
- Data-provider daemons, e.g. DBS daemon, etc.
- Various monitors, e.g. MongoDB connection monitor, etc.
- MongoDB

Upon user request, the front-end validates user input and pass it to DAS web server. It decompose user query into series of requests to underlying core engine, who by itself invokes multiple APIs to fetch data from data-provider and place them into MongoDB. Later this data are serverd back to the user and stay in cache for period of time determined by data-providers, see Figure:



Thread structure of DAS server

Below we outline a typical layout of DAS server threads:

- 1 main thread
- 1 HTTP server thread (CherryPy)
- 1 TimeoutMonitor thread (CherryPy)
- 1 dbs_phedex_monitor thread (dbs_phedex combined service)
- 1 dbs_phedex worker thread (dbs_phedex combined service)

- 1 lumi_service thread (lumi service)
- N_{CP} CherryPy threads
- N_{DAS} worker threads, they are allocated as following:
 - N_{web} threads for web workers
 - N_{core} threads for DAS core workers
 - N_{api} threads for DAS service APIs

In addition DAS configuration uses `das.thread_weights` parameter to weight certain API threads. It is defined as a list of `srv:weight` pairs where each service gets N_{api} number of threads.

Therefore the total number of threads is quite high (range in first hundred) and it is determined by the following formula

$$N_{threads} = N_{main} + N_{CP} + N_{DAS}$$

$$N_{DAS} = N_{web} + N_{core} + N_{api}$$

N_{main} equals to sum of main, timeout, http, dbs_phedex and lumi threads N_{CP} is defined in DAS configuration file, typical value is 30 N_{web} is defined in DAS config file, see `web_server.web_workers` N_{core} is defined in DAS config file, see `das.core_workers` N_{api} is defined in DAS config file, see `das.api_workers`

For example, use the following configuration parameters

```
[das]
multitask = True           # enable multitasking for DAS core (threading)
core_workers = 10         # number of DAS core workers who contact data-providers
api_workers = 2           # number of API workers who run simultaneously
thread_weights = 'dbs:3','phedex:3' # thread weight for given services
das.services = dbs,phedex,dashboard,monitor,runregistry,sitedb2,tier0,conddb,google_maps,postalcode,

[web_server]              # DAS web server configuration parameters
thread_pool = 30          # number of threads for CherryPy
web_workers = 10          # Number of DAS web server workers who handle user requests
dbs_daemon = True         # Run DBSDaemon (boolean)
onhold_daemon = True      # Run onhold daemon for queries which put on hold after hot threshold
```

we get the DAS server running with 151 threads

$$N_{main} = 7, N_{CP} = 30, N_{DAS} = 114$$

where N_{DAS} has the following breakdown

$$N_{web} = 20, N_{core} = 60, N_{api} = 34$$

here we calculated N_{api} as following: we have 13 services, each of them uses 2 API workers (as specified in das configuration), but dbs and phedex data-services are weighed with weight 3, therefore the total number of dbs and phedex workers is 6, respectively. To sum up the numbers we have: 11 services with 2 API workers plus 6 workers for dbs and 6 workers for phedex.

Debugging DAS server

There is nice way to get a snapshot of current activity of DAS server by sending SIGUSR1 signal to DAS server, e.g. upon executing `kill -SIGUSR1 <PID>` command you'll get the following output in DAS log

```
# Thread: DASAbstractService:dbs:PluginTaskManager(4706848768)
File: "/opt/local/Library/Frameworks/Python.framework/Versions/2.7/lib/python2.7/threading.py", line
    self.__bootstrap_inner()
File: "/opt/local/Library/Frameworks/Python.framework/Versions/2.7/lib/python2.7/threading.py", line
    self.run()
File: "/Users/vk/CMS/GIT/github/DAS/src/python/DAS/utils/task_manager.py", line 39, in run
    task = self._tasks.get()
File: "/opt/local/Library/Frameworks/Python.framework/Versions/2.7/lib/python2.7/Queue.py", line 168,
    self.not_empty.wait()
File: "/opt/local/Library/Frameworks/Python.framework/Versions/2.7/lib/python2.7/threading.py", line
    waiter.acquire()
....
.... and similar output for all other DAS threads
....
```

1.12.3 DAS records

DAS data objects

DAS needs to deal with a variety of different data object representations. Data from providers may have both different formats (eg XML, JSON), and different ways of storing hierarchical information. The structure of response data is not known to DAS a-priori. Therefore it needs to treat them as data objects. Here we define what it means for DAS and provide examples of DAS data objects or DAS records.

There are basically two types of data objects: flat and hierarchical ones.

- Flat

```
{"dataset": "abc", "size": 1, "location": "CERN"}
```

- Hierarchical

```
{"dataset": {"name": "abc", "size": 1, "location": "CERN"}}
```

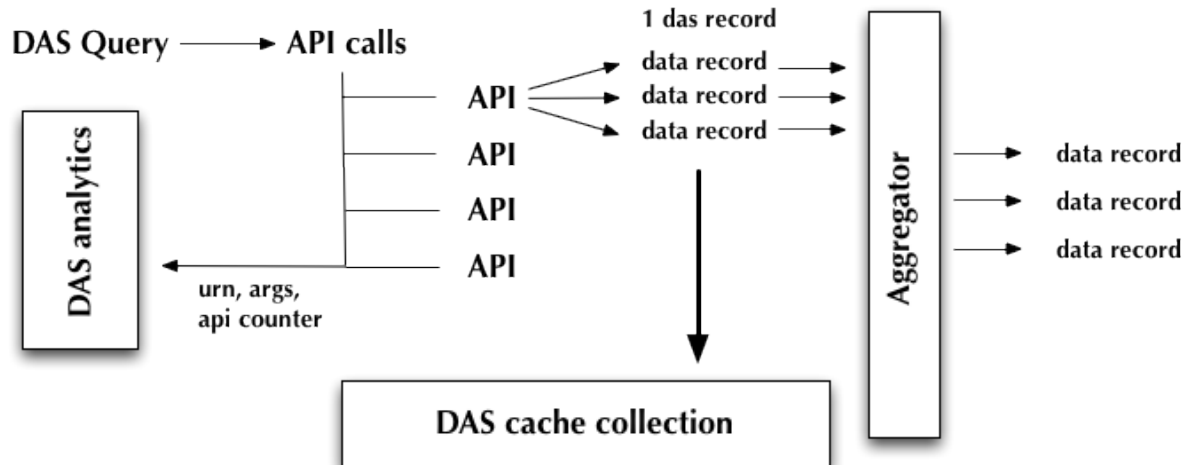
The first has the disadvantage of not being able to easily tell what this object represents, whereas this is not the case for the latter. It is clear in that case that the object principally represents a *dataset*. However, all good things come with a cost; the hierarchical structures are much more expensive to parse, both in python and MongoDB. In DAS we store objects in hierarchical structures, but try to minimize the nesting depth. This allows us to talk about the key/attributes of the object in a more natural way, and simplifies their aggregation. For instance, if two different data-providers serve information about files and file objects containing the *name* attribute, it will be trivial to merge the objects based on the *name* value.

DAS records represent meta-data in *[JSON]* data format. This is a lightweight, near universal format which can represent complex, nested structures via values, dictionaries, lists. JSON is a native JavaScript data format (JavaScript Object Notation), and is represented by the dictionary type in python. A DAS record is just a collection of data supplied by data-services which participate in DAS. DAS wraps each data-service record with auxiliary meta-data such as internal ID, references to other DAS records, and a DAS header. The DAS header contains information about underlying API calls made to data-provider. For example:

```
{"query": "{\"fields\": null, \"spec\": {\"block.name\": \"/abc\"}}",
  "_id": "4b6c8919e2194e1669000002",
  "qhash": "aa8bcf183d916ea3befbdfbcbf40940a",
  "das": {"status": "ok", "qhash": "aa8bcf183d916ea3befbdfbcbf40940a",
    "ctime": [0.55365610122680664, 0.54806804656982422],
    "url": ["http://a.v.com/api", "http://c.d.com/api"],
```

```
"timestamp": 1265404185.2611251,
"lookup_keys": ["block.name", "file.name", "block.name"],
"system": ["combined"],
"services": [{"combined": ["dbs", "phedex"]}],
"record": 0,
"api": ["blockReplicas", "listBlocks"],
"expire": 1265407785.2611251}
}
```

DAS workflow produce the following set of records



- data records, which contains data coming out from data-services, every data record contains a pointer to das record
- das records, which contains information how we retrieve data
- analytics records, which contains information about API calls

1.12.4 DAS raw cache

The DAS raw cache holds all *raw* data-service output, converted into uniform DAS records.

DAS cache records

The DAS cache keeps metadata for the data-services in a form of DAS record. Each data-service output is converted into a DAS record according to *DAS mapping*.

data records

Each data record contains data-service meta-data. The structure of the data service response is a-priori unknown to DAS. Since DAS operates with *JSON* almost any data structure can be stored into DAS, e.g. dictionaries, lists, strings, numerals, etc. The only fields DAS appends to the record are:

- das, contains DAS expiration timestamp
- das_id, refers to query2apis data record

- `primary_key`, provide a primary key to the stored data

For example, here is a data record from [SiteDB](#)

```
{u'_id': ObjectId('4b4f4cb0e2194e72b2000002'),
 u'das': {u'expire': 1263531376.062233},
 u'das_id': u'4b4f4cb0e2194e72b2000001',
 u'primary_key': u'site.name',
 u'site': {u'name': u'Tl_CH_CERN', u'sitename': u'CERN'}}
```

query2apis records

This type of DAS record contains information about underlying API calls made by DAS upon provided user query. It contains the following keys

- `das`, a dictionary of DAS operations
 - `api`, a list of API calls made by DAS upon provided user query
 - `ctime`, a call time spent for every API call
 - `expire`, a shortest expiration time stamp among all API calls
 - `lookup_keys`, a DAS look-up key for provided user query
 - `qhash`, a md5 hash of input query (in MongoDB syntax)
 - `status`, a status request field
 - `system`, a corresponding list of data-service names
 - `timestamp`, a timestamp of last API call
 - `url`, a list of URLs for API calls
 - `version`, reserved for future use
- `query`, an input user query in a form of MongoDB syntax

Here is an example `query2apis` record for the following user input `site=Tl_CH_CERN`

```
{u'_id': ObjectId('4b4f4cb0e2194e72b2000001'),
 u'das': {u'api': [u'CMStoSiteName',
                  u'CMStoSiteName',
                  u'CMStoSAMName',
                  u'CMSNametoAdmins',
                  u'CMSNametoSE',
                  u'SEtoCMSName',
                  u'CMSNametoCE',
                  u'nodes',
                  u'blockReplicas'],
          u'ctime': [1.190140962600708,
                    1.190140962600708,
                    0.71966314315795898,
                    0.72777295112609863,
                    0.7784569263458252,
                    0.75019693374633789,
                    0.74393796920776367,
                    0.28762698173522949,
                    0.30852007865905762],
          u'expire': 1263489980.1307981,
```



```

    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001',
    u'4b4f4cb0e2194e72b2000001'],
  u'site': [{u'ce': u'cel26.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce201.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'cel31.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'cel03.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'cel28.cern.ch', u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'Josh.Bendavid@cern.ch', u'forename': u'Josh',
      u'surname': u'Bendavid', u'title': u'Data Manager'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'gowdy@cern.ch', u'forename': u'Stephen',
      u'surname': u'Gowdy', u'title': u'T0 Operator'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'gowdy@cern.ch', u'forename': u'Stephen',
      u'surname': u'Gowdy', u'title': u'Site Executive'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'gowdy@cern.ch', u'forename': u'Stephen',
      u'surname': u'Gowdy', u'title': u'Data Manager'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'gowdy@cern.ch', u'forename': u'Stephen',
      u'surname': u'Gowdy', u'title': u'Site Admin'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'dmason@fnal.gov', u'forename': u'David',
      u'surname': u'Mason', u'title': u'Data Manager'},
      u'name': u'T1_CH_CERN'},
    {u'ce': u'cel32.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'cel30.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'cel27.cern.ch', u'name': u'T1_CH_CERN'},
    {u'name': u'T1_CH_CERN, u'samname': u'CERN-PROD'},
    {u'name': u'T1_CH_CERN, u'sitename': u'CERN'},
    {u'admin': {u'email': u'Victor.Zhiltsov@cern.ch', u'forename': u'Victor',
      u'surname': u'Zhiltsov', u'title': u'Data Manager'},
      u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'Peter.Kreuzer@cern.ch', u'forename': u'Peter',
      u'surname': u'Kreuzer', u'title': u'Site Admin'},

```

```

    u'name': u'T1_CH_CERN'},
    {u'ce': u'ce125.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce112.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce129.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce133.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce202.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce106.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce105.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce111.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce104.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce113.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce107.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce114.cern.ch', u'name': u'T1_CH_CERN'},
    {u'ce': u'ce124.cern.ch', u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'Peter.Kreuzer@cern.ch', u'forename': u'Peter',
                u'surname': u'Kreuzer', u'title': u'Site Executive'},
               u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'Christoph.Paus@cern.ch', u'forename': u'Christoph',
                u'surname': u'Paus', u'title': u'Data Manager'},
               u'name': u'T1_CH_CERN'},
    {u'admin': {u'email': u'ceballos@cern.ch', u'forename': u'Guillermo',
                u'surname': u'Gomez-Ceballos', u'title': u'Data Manager'},
               u'name': u'T1_CH_CERN'}}}]

```

1.12.6 Query suggestions based on Keyword Search

Virtual Data Integration (such as this system, DAS) provides a coherent interface for querying heterogeneous data sources (e.g., web services, proprietary systems) with minimum upfront effort. Still, this requires its users to learn a new query language and to get acquainted with data organization which may pose problems even to proficient users.

The keyword search plugin, proposes a ranked list of structured queries along with their explanations.

Below is an example of keyword search suggestions for an unstructured query:

The screenshot shows the DAS keyword search interface. At the top, there are controls for results format (list), results per page (10), database instance (int/global), and autocompletion (enable). A search bar contains the query "size of relval datasets number events>100". Below the search bar, a link "Show DAS keys description" is visible. The interface then displays a section titled "Did you mean any of the queries below?" with a filter by entity set to "any". Several query suggestions are listed, each with a colored bar indicating the entity type: dataset (green), file (blue), group (orange), and site (red). The suggestions are:

- dataset group=RelVal | grep dataset.size dataset.nevents>100
- dataset dataset=*relval* | grep dataset.size, dataset.name, datas
- dataset group=RelVal | grep dataset.nevents>100
- file dataset=*relval* | grep file.size, dataset.name, file.nevents>100

 A tooltip explains the first suggestion: "Explanation: find Dataset size (i.e. dataset.size) for each dataset where group=RelVal AND Number of events (i.e. dataset.nevents) > 100". To the right, a section titled "Coloring of query suggestions:" explains the color coding: green for "entity to be retrieved", blue for "filter (an input to service(s))", and red for "expensive filter (applied only after retrieving all data)".

The results can be even more accurate if the query is at least partially structured.

Implementation details

The keyword search operates in “offline” mode - matching the keywords to the available metadata, such as the constraints on inputs accepted by services, or the list of values for some fields.

The code is still rather messy and needs some cleanup and generalization, however it provides a working open source implementation of the keyword search based query suggestions for data service integration (and it could quite easily be generalized for some other domain than the one of the CMS Experiment).

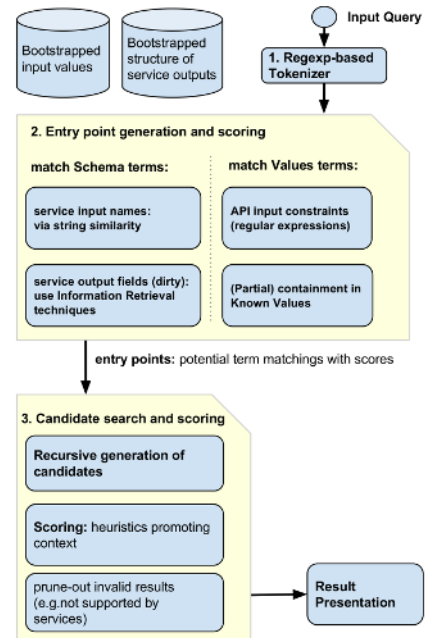


Fig. 1.1: The major components and steps in generating the suggestions

For more details on the implementation see the [code documentation](#) and the code itself.

More information and lessons learned may be found on in the following publications and whitepapers:

DAS KWS Poster

DAS KWS Paper

<https://github.com/vidma/das-paper/blob/master/paper.pdf>

1.12.7 DAS logging DB

The DAS logging DB holds information about all requests made to DAS. All records are stored in a [capped collection](#) of MongoDB. Capped collections are fixed sized collections that have a very high performance auto-LRU age-out feature (age out is based on insertion order) and maintain insertion order for the objects in the collection.

Logging DB records

```
{
  "args": {
    "query": "site=T1_CH_CERN",
    "qhash": "7c8ba62a07ff2820c217ae3b51686383",
    "ip": "127.0.0.1",
    "hostname": "",
    "port": 65238,
    "headers": {
      "Remote-Addr": "127.0.0.1",
      "Accept-Encoding": "identity",
      "Host": "localhost:8211",
      "Accept": "application/json",
      "User-Agent": "Python-urllib/2.6",
      "Connection": "close"
    },
    "timestamp": 1263488174.364929,
    "path": "/status",
    "_id": "4b4f4caee2194e72ae000003",
    "method": "GET"
  }
}
```

Keyword Search over Data Service Integration for Accurate Results

University of Amsterdam
 Vrije Universiteit Amsterdam
 Amsterdam University of Applied Sciences
 Amsterdam University of Professional Education

Value for Research on
 Current, Incomplete, & No
 clear integration

Summary

Recent developments in data integration have been largely in the area of heterogeneous data integration, with various data integration systems (e.g., data integration, data integration, data integration) being developed. However, there is still a need for a system that can integrate data from different sources and provide a unified view of the data. This is the goal of the proposed system. The system is designed to be a data integration system that can integrate data from different sources and provide a unified view of the data. The system is designed to be a data integration system that can integrate data from different sources and provide a unified view of the data. The system is designed to be a data integration system that can integrate data from different sources and provide a unified view of the data.

Challenges

- A general approach to integrate data from different sources is needed. This is a challenging task because of the heterogeneity of the data sources and the complexity of the integration process.
- A general approach to integrate data from different sources is needed. This is a challenging task because of the heterogeneity of the data sources and the complexity of the integration process.
- A general approach to integrate data from different sources is needed. This is a challenging task because of the heterogeneity of the data sources and the complexity of the integration process.

Context: a system for Virtual Data Integration

Virtual Data Integration System (VDIS)

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

• A system for integrating data from different sources

The workflow

Data integration workflow

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

The workflow consists of several steps: 1. Data integration: The data from different sources is integrated into a single view. 2. Query processing: The query is processed against the integrated view. 3. Result generation: The results of the query are generated and returned to the user.

1.12.8 Code documentation

DAS Core classes

DAS aggregators

DAS core module

DAS key learning module

DAS mapping module

DAS mongo cache module

DAS parser

DAS parser cache module

DAS QL module

DAS Query module

DAS robot module

DAS SON manipulator

DAS services

DAS operates using a provided list of data-services and their definitions. The nature of these services is unimportant provided that they provide a some form of API which DAS can call and aggregate the data returned. Each service needs to be registered in DAS by providing an appropriate configuration file and (optionally) a service handler class..

CMS services

Each CMS data-service is represented by a mapping and, optionally, by a plugin class. The data-service map contains description of the data-service, e.g. URL, URN, expiry timeout as well as API and notations maps.

- the API map relates DAS keys and API input parameters. It contains the following items:
 - *api*, name of the API
 - *params*, a list of API input parameters together with regex patterns accepted to check the format of or identify ambiguous values.
 - *record* represents DAS record. Each record has
 - * *daskeys*, a list of DAS maps; each map relates keys in the user query to the appropriate DAS representation
 - *key*, a DAS key used in DAS queries, e.g. *block*
 - *map*, a DAS record representation of the *key*, e.g. *block.name*
 - *pattern*, a regex pattern for DAS key

- * *das2api*, is a map between DAS key representations and API input parameters
 - *api_param*, an API input parameter, e.g. *se*
 - *das_key*, a DAS key it represents, e.g. *site.se*
 - *pattern*, a regex pattern for *api_param*

- Notation map represents a mapping between data-service output and DAS records. It is optional.

Please use these links [API map](#) and [API notation](#) for concrete examples.

DAS abstract service

DAS generic service

DAS map reader module

DAS web modules

DAS autocomplete module

DAS cms representation module

Set of DAS codes

DAS web status codes

`DAS.web.das_codes.decode_code (code)`
Return human readable string for provided code ID

`DAS.web.das_codes.web_code (error)`
Return DAS WEB code for provided error string

DAS das representation module

DAS server

DAS test data-service server

DAS web server

DAS web manager

DAS dbs daemon module

DAS web tools

DAS web utils

Keyword Search Modules

DAS Keyword Search modules

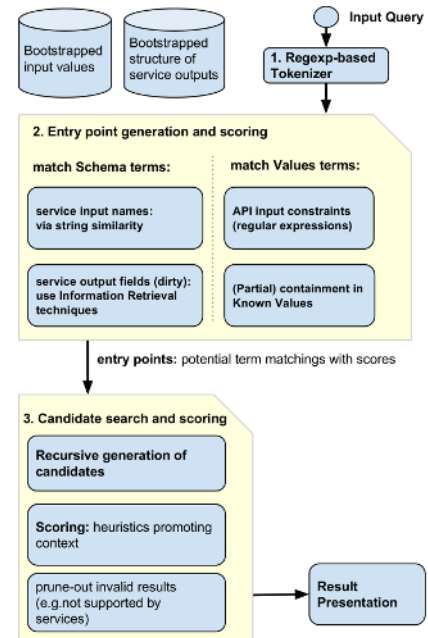
Module author: Vidmantas Zemleris <vidmantas.zemleris@gmail.com>

- *Overview*
- *Adapters to access Metadata*
 - *The schema adapter*
 - *Input Values Tracker*
 - *DAS Query Language definitions*
- *Tokenizing and parsing the query*
- *Entity matchers*
 - *Value matching*
 - *Name matching*
 - *Name matching: multi-term chunks representing field names*
- *Generating and Ranking the Query suggestion*
 - *A ranker implemented in Cython and built into a C extension*
- *Presenting the Results to the user*

Overview

The basic keyword search workflow:

- The query is tokenized (including shallow parsing for key=val patterns and quotes)
- Then by a number of entity matchers the entry points are generated (matches of each keyword or the few nearby keywords into some of the schema or value terms which make part of a structured query)
- Exploring the various combinations of the entry points the candidates for query suggestions are evaluated and ranked
- The top results are presented to the users (generating a valid DAS query and a providing a readable query description)



Adapters to access Metadata

Contain meta-data related functions:

- accessing integration schema: fields, values, constraints on inputs/queries
- tracking fields available
- tracking known (input field) values

The schema adapter Provide a layer of abstraction between Keyword Search and the Data Integration System.

```
DAS.keywordsearch.metadata.schema_adapter2.ApiDef
    alias of ApiInputParamsEntry
```

```
class DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter
    provides an adapter between keyword search and Data integration system.
```

```
classmethod are_wildcards_allowed(entity, wildcards, params)
```

Whether wildcards are allowed for given inputs

currently only these simple rules allowed: site=W* dataset=W* dataset site=T1_CH_* dataset site=T1_CH_* dataset=/A/B/C file dataset=/DoubleMu/Run2012A-Zmmg-13Jul2012-v1/RAW-RECO site=T1_* file block=/A/B/C#D file file=W* dataset=FULL file file=W* block=FULL

these are supported (probably) because of DAS-wrappers file dataset=*DoubleMuParked25ns* file dataset=*DoubleMuParked25ns* site=T2_RU_JINR

```
check_result_field_match(fieldname)
```

checks for complete match to a result field

```
entities_for_inputs(params)
```

lists entities that could be retrieved with given input params

```
get_api_param_definitions()
```

returns a list of API input requirements

get_result_field_title (*result_entity, field, technical=False, html=True*)
 returns name (and optionally title) of output field

init (*dascore=None*)
 initialization or re-initialization

list_result_fields (*entity=None, inputs=None*)
 lists attributes available in all service outputs (aggregated)

validate_input_params (*params, entity=None, final_step=False, wildcards=None*)
 checks if DIS can answer query with given params.

validate_input_params_lookupbased (*params, entity=None, final_step=False, wildcards=None*)
 checks if DIS can answer query with given params.

Gathers the list of fields available in service outputs

`DAS.keywordsearch.metadata.das_output_fields_adapter.flatten` (*list_of_lists*)
 Flatten one level of nesting

`DAS.keywordsearch.metadata.das_output_fields_adapter.get_outputs_field_list` (*dascore*)
 makes a list of output fields available in each DAS entity this is taken from keylearning collection.

`DAS.keywordsearch.metadata.das_output_fields_adapter.get_titles_by_field` (*dascore*)
 returns a dict of titles taken from presentation cache

`DAS.keywordsearch.metadata.das_output_fields_adapter.is_reserved_field` (*field, re-sult_type*)
 returns whether the field is reserved, e.g. **.error, *.reason, qhash*

`DAS.keywordsearch.metadata.das_output_fields_adapter.print_debug` (*dascore, fields_by_entity, re-sults_by_entity*)
 verbose output for `get_outputs_field_list`

`DAS.keywordsearch.metadata.das_output_fields_adapter.result_contained_errors` (*rec*)
 decide whether keylearning record contain errors (i.e. as responses from services contained errors) and whether the record shall be excluded

Input Values Tracker

DAS Query Language definitions defines DASQL and keyword search features, e.g. what shall be considered as:
 * word * simple operators * aggregation operators (not implemented)

`DAS.keywordsearch.metadata.das_ql.flatten` (*list_of_lists*)
 Flatten one level of nesting

`DAS.keywordsearch.metadata.das_ql.get_operator_synonyms` ()
 return synonyms for das aggregation operators (not used yet)

Tokenizing and parsing the query

Module description:

- first clean up input keyword query (rm extra spaces, standardize notation)
- then it tokenizes the query into:

- individual query terms
- compound query terms in brackets (e.g. “number of events”)
- phrases: “terms operator value” (e.g. nevent > 1, “number of events”=100)

DAS.keywordsearch.tokenizer.**cleanup_query**(*query*)

Returns cleaned query by applying a number of transformation patterns that removes spaces and simplifies the conditions

```
>>> cleanup_query('number of events = 33')
'number of events=33'
```

```
>>> cleanup_query('number of events > 33')
'number of events>33'
```

```
>>> cleanup_query('more than 33 events')
'>33 events'
```

```
>>> cleanup_query('X more than 33 events')
'X>33 events'
```

```
>>> cleanup_query('find datasets where X more than 33 events')
'datasets where X>33 events'
```

```
>>> cleanup_query('=2012-02-01')
' = 20120201'
```

```
>>> cleanup_query('>= 2012-02-01')
'>= 20120201'
```

DAS.keywordsearch.tokenizer.**get_keyword_without_operator**(*keyword*)
splits keyword on operator

```
>>> get_keyword_without_operator('number of events >= 10')
'number of events'
```

```
>>> get_keyword_without_operator('dataset')
'dataset'
```

```
>>> get_keyword_without_operator('dataset=Zmm')
'dataset'
```

DAS.keywordsearch.tokenizer.**get_operator_and_param**(*keyword*)
splits keyword on operator

```
>>> get_operator_and_param('number of events >= 10')
{'type': 'filter', 'param': '10', 'op': '>='}
```

```
>>> get_operator_and_param('dataset')
```

```
>>> get_operator_and_param('dataset=Zmm')
{'type': 'filter', 'param': 'Zmm', 'op': '='}
```

DAS.keywordsearch.tokenizer.**test_operator_containment** (*keyword*)
returns whether a keyword token contains an operator (this is useful then processing a list of tokens, as only the last token may have an operator)

```
>>> test_operator_containment('number of events >= 10')
True
```

```
>>> test_operator_containment('number')
False
```

DAS.keywordsearch.tokenizer.**tokenize** (*query*)
tokenizes the query retaining the phrases in brackets together it also tries to group “word operator word” sequences together, such as

```
"number of events">10 or dataset=/Zmm/*/raw-reco
```

so it could be used for further processing.

special characters currently allowed in data values include: _*/-

For example:

```
>>> tokenize('file dataset=/Zmm/*/raw-reco lumi=20853 nevents>10'
['file', 'dataset=/Zmm/*/raw-reco', 'lumi=20853', 'nevents>10', 'number of events>10', '/Zmm']

>>> tokenize('file dataset=/Zmm/*/raw-reco lumi=20853 dataset.nevents>10'
['file', 'dataset=/Zmm/*/raw-reco', 'lumi=20853', 'dataset.nevents>10', 'number of events>10']

>>> tokenize("file dataset=/Zmm/*/raw-reco lumi=20853 dataset.nevents>10"
['file', 'dataset=/Zmm/*/raw-reco', 'lumi=20853', 'dataset.nevents>10', 'number of events>10']

>>> tokenize('user=vidmasze@cern.ch')
['user=vidmasze@cern.ch']
```

Entity matchers

Contain entity matching related functions:

- Name matching / custom String distance
- Chunk matching (multi-word terms into names of service output fields)
- Value matching
- And CMS specific dataset matching

Value matching module provide custom Levenshtein distance function

```
DAS.keywordsearch.entity_matchers.string_dist_levenstein.levenshtein(string1,  
                                                                    string2,  
                                                                    sub-  
                                                                    cost=3,  
                                                                    mod-  
                                                                    ifica-  
                                                                    tion_middle_cost=2)
```

string-edit distance returning min cost of edits needed

```
DAS.keywordsearch.entity_matchers.string_dist_levenstein.levenshtein_normalized(string1,  
                                                                    string2,  
                                                                    sub-  
                                                                    cost=2,  
                                                                    max-  
                                                                    cost=3)
```

return a levenshtein distance normalized between [0-1]

Name matching

Name matching: multi-term chunks representing field names Modules for matching chunks of keywords into attributes of service outputs. This is obtained by using information retrieval techniques.

Generating and Ranking the Query suggestion

A ranker implemented in Cython and built into a C extension A ranker combine scores of individual keywords to make up the final score. It evaluates the possible combinations and provides a ranked list of results (i.e. query suggestions).

the source code is in `DAS.keywordsearch.rankers.fast_recursive_ranker` which is compiled into `DAS.extensions.fast_recursive_ranker` (with help of `DAS/keywordsearch/rankers/build_cython.py`)

Presenting the Results to the user

Presentation of query suggestions The module contain functions for presenting the results as DASQL and formatting/coloring them in HTML.

```
DAS.keywordsearch.presentation.result_presentation.dasql_to_nl(dasql_tuple)
```

Returns natural language representation of a generated DAS query so to explain users what does it mean.

```
DAS.keywordsearch.presentation.result_presentation.fescape(value)
```

escape a value to be included in html

```
DAS.keywordsearch.presentation.result_presentation.result_to_dasql(result,  
                                                                    frmt='text',  
                                                                    shorten_html=True,  
                                                                    max_value_len=26)
```

returns proposed query as DASQL in there formats:

- text, standard DASQL
- html, colorified DASQL with long values shortened down (if `shorten_html` is specified)

`DAS.keywordsearch.presentation.result_presentation.shorten_value` (*value*,
max_value_len)
provide a shorter version of a very long value for displaying in (html) results. Examples include long dataset or block names.

DAS tools

DAS provides a useful set of tools.

DAS config tools

DAS config generator

class `DAS.tools.create_das_config.ConfigOptionParser`
option parser

get_opt ()
Returns parse list of options

`DAS.tools.create_das_config.main` ()
Main function

DAS admin tools

DAS stress test module

DAS bench tools

DAS benchmark tool

class `DAS.tools.das_bench.NClientsOptionParser`
client option parser

get_opt ()
Returns parse list of options

class `DAS.tools.das_bench.UrlRequest` (*method*, **args*, ***kwargs*)
URL requestor class which supports all RESTful request methods. It is based on `urllib2.Request` class and overwrite request method. Usage: `UrlRequest(method, url=url, data=data)`, where `method` is GET, POST, PUT, DELETE.

get_method ()
Return request method

`DAS.tools.das_bench.avg_std` (*input_file*)
Calculate average and standard deviation

`DAS.tools.das_bench.gen_passwd` (*length*=8, *chars*=`'ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz0123456789'`)
Random string generator, code based on <http://code.activestate.com/recipes/59873-random-password-generation/>

`DAS.tools.das_bench.main` ()
Main routine

`DAS.tools.das_bench.make_plot` (*xxx*, *yyy*, *std*=None, *name*=`'das_cache.pdf'`, *xlabel*=`'Number of clients'`, *ylabel*=`'Time/request (sec)'`, *yscale*=None, *title*=`''`)
Make standard plot time vs nclients using matplotlib

`DAS.tools.das_bench.natcasecmp` (*aaa, bbb*)
Natural string comparison, ignores case.

`DAS.tools.das_bench.natcmp` (*aaa, bbb*)
Natural string comparison, case sensitive.

`DAS.tools.das_bench.natsort` (*seq, icmp=<function natcmp>*)
In-place natural string sort.

`DAS.tools.das_bench.natsort_key` (*sss*)
Used internally to get a tuple by which s is sorted.

`DAS.tools.das_bench.natsorted` (*seq, icmp=<function natcmp>*)
Returns a copy of seq, sorted by natural string sort.

`DAS.tools.das_bench.random_index` (*bound*)
Generate random number for given upper bound

`DAS.tools.das_bench.runjob` (*nclients, host, method, params, headers, idx, limit, debug=0, logname='spammer', dasquery=None*)
Run spammer for provided number of parallel clients, host name and method (API). The output data are stored into lognameN.log, where logname is an optional parameter with default as spammer.

`DAS.tools.das_bench.spammer` (*stream, host, method, params, headers, write_lock, debug=0*)
start new process for each request

`DAS.tools.das_bench.try_int` (*sss*)
Convert to integer if possible.

`DAS.tools.das_bench.urlrequest` (*stream, url, headers, write_lock, debug=0*)
URL request function

DAS CLI tool

DAS client

DAS command line tool

```
class DAS.tools.das_client.DASOptionParser
    DAS cache client option parser

    get_opt ()
        Returns parse list of options

class DAS.tools.das_client.HTTPSClientAuthHandler (key=None, cert=None, level=0)
    Simple HTTPS client authentication class based on provided key/ca information

    get_connection (host, timeout=300)
        Connection method

    https_open (req)
        Open request method

DAS.tools.das_client.check_auth (key)
    Check if user runs das_client with key/cert and warn users to switch

DAS.tools.das_client.check_glidein ()
    Check glidein environment and exit if it is set

DAS.tools.das_client.convert_time (val)
    Convert given timestamp into human readable format
```

`DAS.tools.das_client.extract_value (row, key)`
Generator which extracts row[key] value

`DAS.tools.das_client.fullpath (path)`
Expand path to full path

`DAS.tools.das_client.get_data (host, query, idx, limit, debug, threshold=300, ckey=None, cert=None, das_headers=True)`
Contact DAS server and retrieve data for given DAS query

`DAS.tools.das_client.get_value (data, filters, base=10)`
Filter data from a row for given list of filters

`DAS.tools.das_client.keys_attrs (lkey, oformat, host, ckey, cert, debug=0)`
Contact host for list of key/attributes pairs

`DAS.tools.das_client.main ()`
Main function

`DAS.tools.das_client.prim_value (row)`
Extract primary key value from DAS record

`DAS.tools.das_client.print_from_cache (cache, query)`
print the list of files reading it from cache

`DAS.tools.das_client.print_summary (rec)`
Print summary record information on stdout

`DAS.tools.das_client.size_format (uinput, ibase=0)`
Format file size utility, it converts file size into KB, MB, GB, TB, PB units

`DAS.tools.das_client.unique_filter (rows)`
Unique filter drop duplicate rows.

`DAS.tools.das_client.x509 ()`
Helper function to get x509 either from env or tmp file

DAS map reduce tools

DAS mapping db tools

DAS robot tools

DAS utilities

DAS provides a useful set of utilities.

CERN SSO authentication utils

CERN SSO toolkit. Provides `get_data` method which allow to get data behind CERN SSO protected site.

`class DAS.utils.cern_sso_auth.HTTPSClientAuthHandler (key=None, cert=None, level=0)`
Simple HTTPS client authentication class based on provided key/ca information

`get_connection (host, timeout=300)`
Connection method

`https_open (req)`
Open request method

`DAS.utils.cern_sso_auth.get_data(url, key, cert, debug=0)`
Main routine to get data from data service behind CERN SSO. Return file-like descriptor object (similar to open).

`DAS.utils.cern_sso_auth.timestamp()`
Construct timestamp used by Shibboleth

DAS config utils

Config utilities

`DAS.utils.das_config.das_configfile(fname=None)`
Return DAS configuration file name \$DAS_ROOT/etc/das.cfg

`DAS.utils.das_config.das_readconfig(fname=None, debug=False)`
Return DAS configuration

`DAS.utils.das_config.das_readconfig_helper(fname=None, debug=False)`
Read DAS configuration file and store DAS parameters into returning dictionary.

`DAS.utils.das_config.read_configparser(dasconfig)`
Read DAS configuration

`DAS.utils.das_config.read_wmcore(filename)`
Read DAS python configuration file and store DAS parameters into returning dictionary.

`DAS.utils.das_config.wmcore_config(filename)`
Return WMCore config object for given file name

`DAS.utils.das_config.write_configparser(dasconfig, use_default)`
Write DAS configuration file

DAS DB utils

DAS option module

DAS option's class

class `DAS.utils.das_option.DASOption(section, name, itype='string', default=None, validator=None, destination=None, description='')`

Class representing a single DAS option, independent of storage mechanism. Must define at least a section and name.

In configparser

[section] name = value

In WMCORE config

config.section.name = value

The type parameter forces conversion as appropriate. The default parameter allows the option to not be specified in the config file. If default is not set then not providing this key will throw an exception. The validator parameter should be a single-argument function which returns true if the value supplied is appropriate. The destination argument is for values which shouldn't end up in `config[section][name] = value` but rather `config[destination] = value` Description is provided as a future option for some sort of automatic documentation.

get_from_configparser (config)
Extract a value from a configparser object.

get_from_wmcore (*config*)

As per get_from_configparser.

write_to_configparser (*config, use_default=False*)

Write the current value to a configparser option. This assumes it has already been read or the values in self.value somehow changed, otherwise if use_default is set only those values with defaults are set.

write_to_wmcore (*config, use_default=False*)

As per write_to_configparser.

DAS timer module

DAS json wrapper

JSON wrapper around different JSON python implementations. We use simplejson (json), cJSON and yaJL JSON implementation.

NOTE: different JSON implementation handle floats in different way Here are few examples

..doctest:

```
r1={"ts":time.time()}
print r1
{'ts': 1374255843.891289}
```

Python json:

..doctest:

```
print json.dumps(r1), json.loads(json.dumps(r1))
{"ts": 1374255843.891289} {'ts': 1374255843.891289}
```

CJSON:

..doctest:: print cJSON.encode(r1), cJSON.decode(cJSON.encode(r1)) {"ts": 1374255843.89} {'ts': 1374255843.89}

YAJL:

..doctest:

```
print yaJL.dumps(r1), yaJL.loads(yaJL.dumps(r1))
{"ts":1.37426e+09} {'ts': 1374260000.0}
```

Therefore when records contains timestamp it is ADVISED to round it to integer. Then json/cjson implementations will agree on input/output, while yaJL will still differ (for that reason we can't use yaJL).

class DAS.utils.jsonwrapper.**JSONDecoder** (***kwargs*)
JSONDecoder wrapper

decode (*istring*)

Decode JSON method

raw_decode (*istring*)

Decode given string

class DAS.utils.jsonwrapper.**JSONEncoder** (***kwargs*)
JSONEncoder wrapper

encode (*idict*)

Decode JSON method

iterencode (*idict*)
Encode input dict

`DAS.utils.jsonwrapper.dump` (*doc, source*)
Use `json.dump` for back-ward compatibility, since `cjson` doesn't provide this method. The `dump` method works on file-descriptor objects.

`DAS.utils.jsonwrapper.dumps` (*idict, **kwargs*)
Based on default `MODULE` invoke appropriate JSON encoding API call

`DAS.utils.jsonwrapper.load` (*source*)
Use `json.load` for back-ward compatibility, since `cjson` doesn't provide this method. The `load` method works on file-descriptor objects.

`DAS.utils.jsonwrapper.loads` (*idict, **kwargs*)
Based on default `MODULE` invoke appropriate JSON decoding API call

DAS logger

General purpose DAS logger class. `PrintManager` class is based on the following work <http://stackoverflow.com/questions/245304/how-do-i-get-the-name-of-a-function-or-method-from-within-a-python-function-or-m> <http://stackoverflow.com/questions/251464/how-to-get-the-function-name-as-string-in-python>

class `DAS.utils.logger.NullHandler` (*level=0*)
Do nothing logger

emit (*record*)
This method does nothing.

handle (*record*)
This method does nothing.

class `DAS.utils.logger.PrintManager` (*name='function', verbose=0*)
`PrintManager` class

debug (*msg*)
print debug messages

error (*msg*)
print warning messages

info (*msg*)
print info messages

warning (*msg*)
print warning messages

`DAS.utils.logger.funcname` ()
Extract caller name from a stack

`DAS.utils.logger.print_msg` (*msg, cls, prefix=''*)
Print message in a form `cls::caller msg`, suitable for class usage

`DAS.utils.logger.set_cherrypy_logger` (*hdlr, level*)
set up logging for `CherryPy`

DAS pycurl manager

DAS query utils

DAS query utils.

`DAS.utils.query_utils.compare_dicts(input_dict, exist_dict)`

Helper function for `compare_specs`. It compares key/val pairs of Mongo dict conditions, e.g. `{ '$gt': 10 }`. Return true if `exist_dict` is superset of `input_dict`

`DAS.utils.query_utils.compare_specs(input_query, exist_query)`

Function to compare set of fields and specs of two input mongo queries. Return True if results of `exist_query` are superset of result for `input_query`.

`DAS.utils.query_utils.compare_str(query1, query2)`

Function to compare string from specs of query. Return True if `query2` is superset of `query1`. `query1 & query2` is the string in the pattern:

```
([a-zA-Z0-9_-\#/*\.] ) *
```

* is the sign indicates that a sub string of *:

```
([a-zA-Z0-9_-\#/*\.] ) *
```

case 1. if `query2` is flat query(w/out *) then `query1` must be the same flat one

case 2. if `query1` is start/end w/ * then `query2` must start/end with *

case 3. if `query2` is start/end w/out * then `query1` must start/end with `query2[0]/query[-1]`

case 4. `query1 & query2` both include *

Way to perform a comparison is splitting:

`query1` into `X0*X1*X2*X3` `query2` into `Xa*Xb*Xc`

foreach X in (`Xa`, `Xb`, `Xc`):

case 5. X is '':

continue

special case:

when `X0 & Xa` are '' or when `X3 & Xc` are '' we already cover it in case 2

case 6. X not in `query1` then return False

case 7. X in `query1` begin at index:

case 7-1. X is the first X not '' we looked up in query1.(`Xa`) `last_idx = index` ; continue

case 7-2. X is not the first:

try to find the smallest `Xb > Xa` if and Only if we could find a sequence:

satisfy `Xc > Xb > Xa`, otherwise return False '=' will happen when `X0 = Xa` then we could return True

`DAS.utils.query_utils.convert2pattern(query)`

In MongoDB patterns are specified via regular expression. Convert input query condition into regular expression patterns. Return new MongoDB compiled w/ regex query and query w/ debug info.

`DAS.utils.query_utils.decode_mongo_query(query)`

Decode query from storage format into mongo format.

`DAS.utils.query_utils.encode_mongo_query(query)`

Encode mongo query into storage format. MongoDB does not allow storage of dict with keys containing "." or MongoDB operators, e.g. \$lt. So we convert input mongo query spec into list of dicts whose "key"/"value" are mongo query spec key/values. For example

```
spec:{"block.name":"aaa"}

converted into

spec:[{"key":"block.name", "value":"'aaa'"}]
```

Conversion is done using JSON dumps method.

DAS regex expressions

Regular expression patterns

`DAS.utils.regex.word_chars(word, equal=True)`

Creates a pattern of given word as series of its characters, e.g. for given word dataset I'll get `^d$|^da$|^dat$|^data$|^datas$|^datase$|^dataset$` which can be used later in regular expressions

DAS task manager

DAS threadpool utils

URL utils

Generic utils

1.12.9 DAS performance benchmarks

Performance benchmarks

CMS data used for Benchmarks

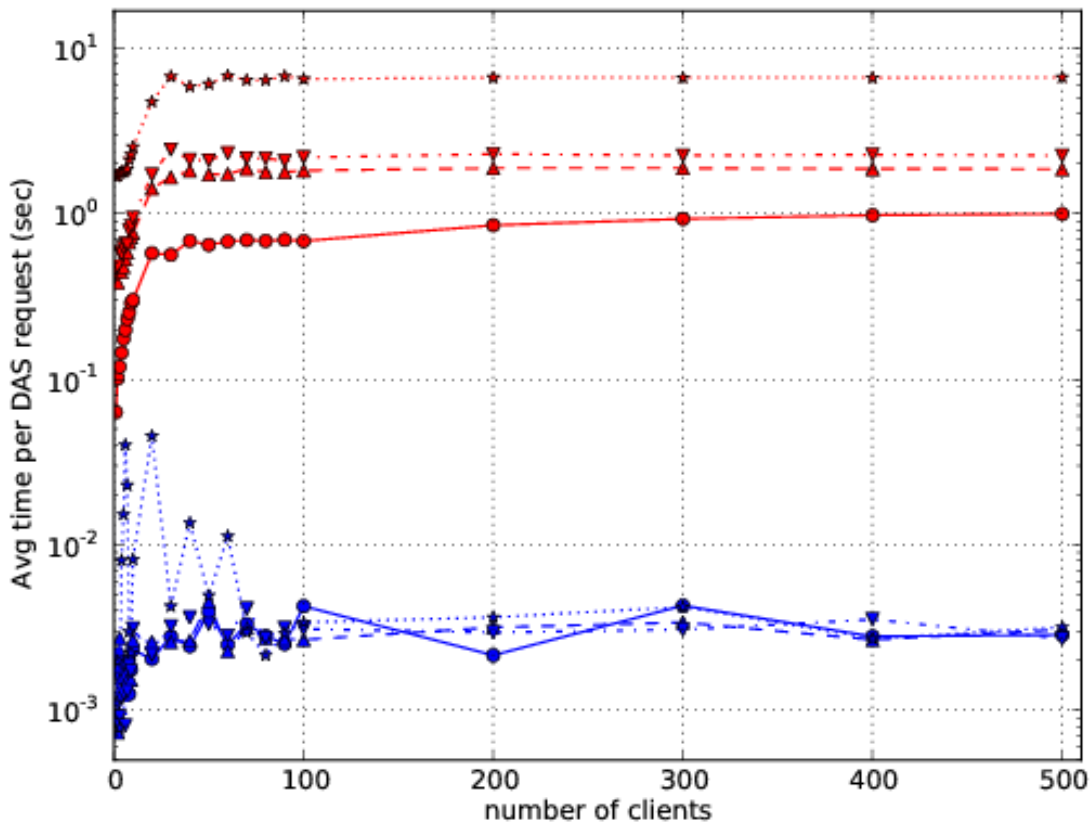
We used ~50K datasets from CMS DBS system and ~450K block records from both DBS and Phedex CMS systems. All of them were populated into DAS cache up-front, since I was only interested in read tests (DAS have an ability to populate the cache). The tests consist of three different types of queries:

- all clients use fixed value for DAS queries, e.g. `dataset=/a/b/c` or `block=/a/b/c#123`
- all clients use fixed pattern for DAS queries, e.g. `dataset=/a*` or `block=/a*`
- all clients use random patterns, e.g. `dataset=/Z*` or `block=/a*`, `block=/Z*`

Once the query has been placed into DAS cache server we retrieve only first record out of the result set and ship it back to the client. The respond time is measured as the total time DAS server spends for a particular query.

Benchmark results

First, we tested our CherryPy server and verified that it can sustain a load of 500 parallel clients at the level of 10^{-5} sec. Then we populated MongoDB with 50k dataset and 500k block records from DBS and Phedex CMS systems. We performed the read test of MongoDB and DAS using 1-500 parallel clients with current set of CMS datasets and block meta-data, 50K and 450K, respectively. Then we populated MongoDB with 10x and 100x of statistics and repeat the tests. The plot showing below represents comparison of DAS (red lines) versus MongoDB (blue lines) read tests for 50k (circles), 500k (down triangles), 5M (up triangles) and 50M (stars):

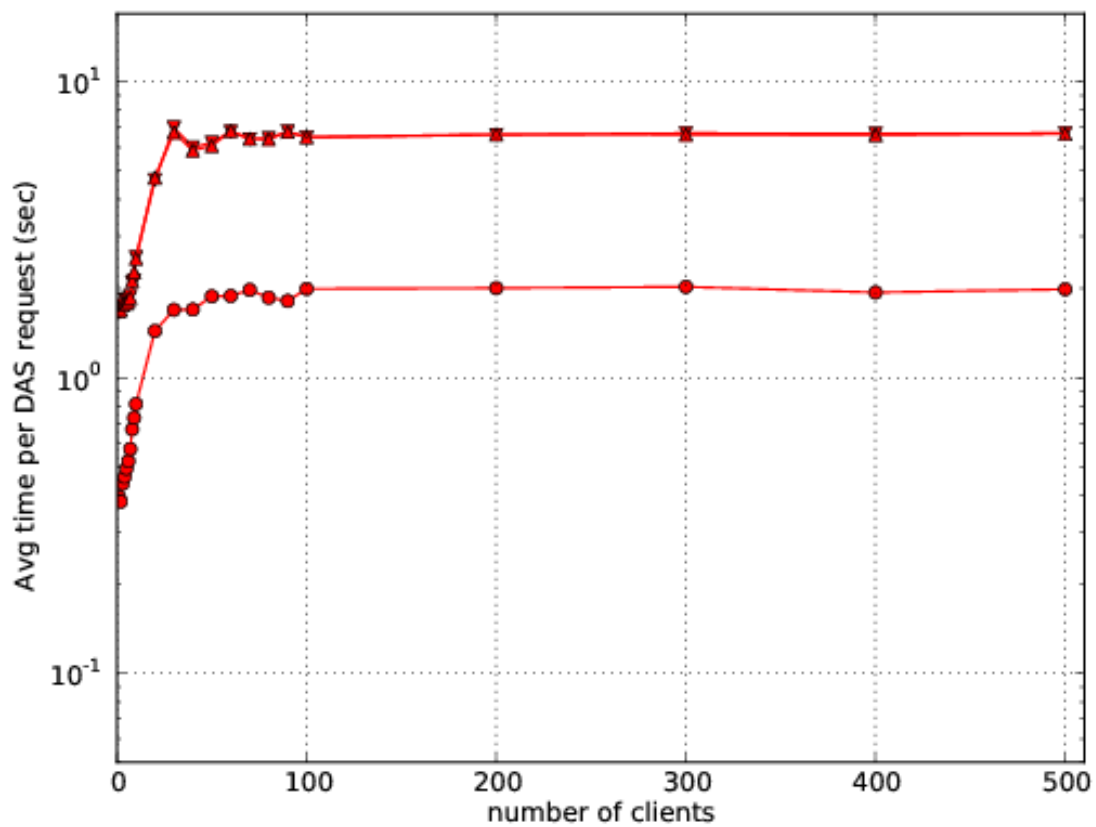


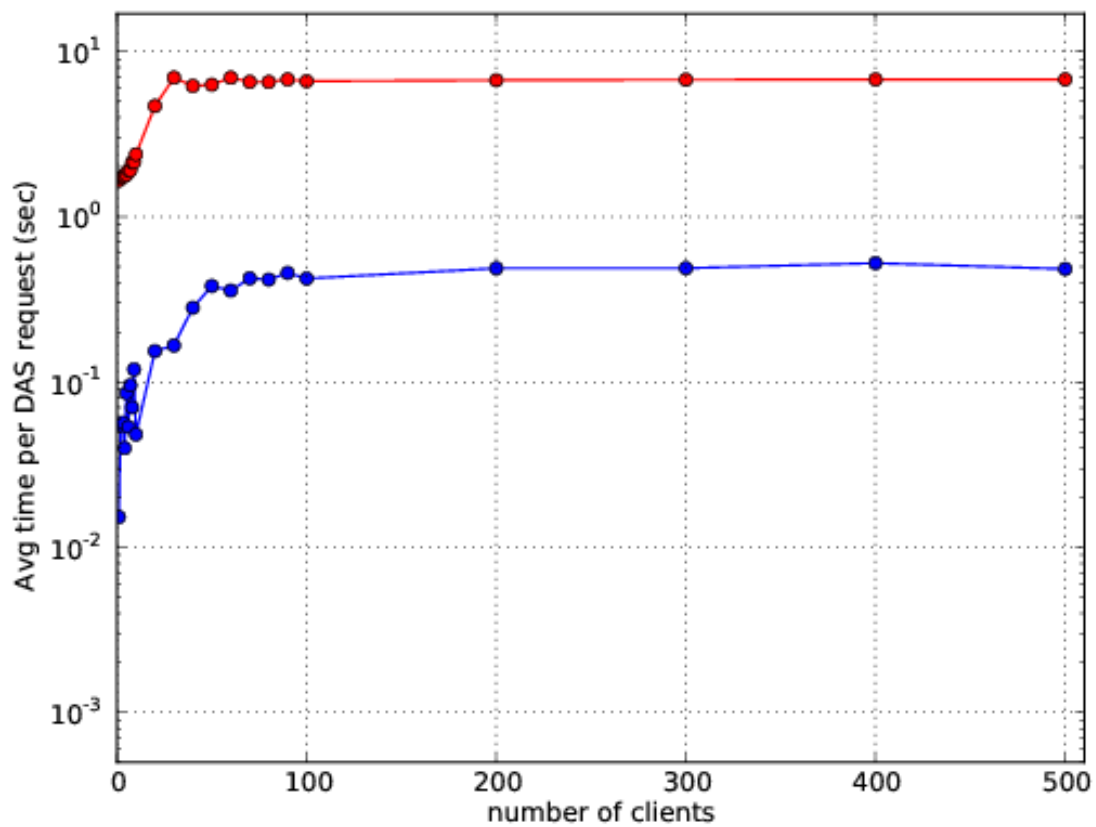
We found these results very satisfactory. As was expected MongoDB can easily sustain such load at the level of few milli-seconds. The DAS numbers also seems reasonable since DAS workflow is much more complicated. It includes DAS parsing, query analysis, analytics, etc. The most important, the DAS performance seems to be driven by MongoDB back-end and has constant scale factor which can be tuned later.

Next we performed three tests discussed above with 10x of block meta-data statistics.

The curve with circles points represents test #1, i.e. fixed key-value, while top/down triangles represents pattern value and random pattern value, tests #2 and #3, respectively. As can be seen pattern tests are differ by the order of magnitude from fixed key-value, but almost identical among each other.

Finally, we tested DAS/MongoDB with random queries and random data access, by asking to return me a single record from entire collection (not only the first one as shown above). For that purpose we generated a random index and used idx/limit for MongoDB queries. Here is the results





The blue line shows MongoDB performance, while red shows the DAS. This time the difference between DAS and MongoDB is only one order of magnitude differ with respect to first shown plot and driven by DAS workflow.

Benchmarks tool

The DAS provides a benchmarking tool, `das_bench`

```
das_bench --help

Usage: das_bench [options]

Examples:
  Benchmark Keyword Search:
  das_bench --url=https://cmsweb-testbed.cern.ch/das/kws_async --nclients=20 --dasquery="/DoubleMu/A

  Benchmark DAS Homepage (repeating each request 10 times):
  src/python/DAS/tools/das_bench.py --url=https://das-test-dbs2.cern.ch/das --nclients=30 --dasquery=

Options:
  -h, --help                show this help message and exit
  -v DEBUG, --verbose=DEBUG
                             verbose output
  --url=URL                  specify URL to test, e.g.
                             http://localhost:8211/rest/test
  --accept=ACCEPT            specify URL Accept header, default application/json
  --idx-bound=IDX            specify index bound, by default it is 0
  --logname=LOGNAME          specify log name prefix where results of N client
                             test will be stored
  --nclients=NCLIENTS       specify max number of clients
  --minclients=NCLIENTS     specify min number of clients, default 1
  --repeat=REPEAT            repeat each benchmark multiple times
  --dasquery=DASQUERY        specify DAS query to test, e.g. dataset
  --output=FILENAME          specify name of output file for matplotlib output,
                             default is results.pdf, can also be file.png etc
```

which can be used to benchmark DAS.

Performance of individual components (profiling)

DAS has been profiled on vocms67:

- 8 core, Intel Xeon CPU @ 2.33GHz, cache size 6144 KB
- 16 GB of RAM
- Kernel 2.6.18-164.11.1.el5 #1 SMP Wed Jan 20 12:36:24 CET 2010 x86_64 x86_64 x86_64 GNU/Linux

The DAS benchmarking is performed using the following query

```
das_cli --query="block" --no-output
```

Latest results are shown below:

```
...
INFO    0x9e91ea8> DASMongocache::update_cache, ['dbs'] yield 387137 rows
...
```



```

INFO    0x9e91ea8> DASMongocache::update_cache, ['phedex'] yield 189901 rows
...
INFO    0x9e91ea8> DASMongocache::merge_records, merging 577038 records

DAS execution time (phedex) 106.446726799 sec
DAS execution time (dbs) 72.2084658146 sec
DAS execution time (merge) 62.8879590034 sec
DAS execution time 241.767010927 sec, Wed, 20 Jan 2010 15:54:33 GMT

```

The largest contributors to execution time are

```

das_cli --query="block" --verbose=1 --profile --no-output

Mon Jan 18 22:43:27 2010    profile.dat

54420138 function calls (54256630 primitive calls) in 247.423 CPU seconds

Ordered by: internal time

ncalls  tottime  percall  cumtime  percall filename:lineno(function)
560649   78.018    0.000   78.018    0.000 {method 'recv' of '_socket.socket' objects}
  992    23.301    0.023   23.301    0.023 {pymongo._cbson._insert_message}
1627587   19.484    0.000   19.484    0.000 {DAS.extensions.das_speed_utils._dict_handler}
  467    17.295    0.037   21.042    0.045 {pymongo._cbson._to_dicts}
 23773   16.945    0.001   16.945    0.001 {built-in method feed}
576626   12.101    0.000   77.974    0.000 /data/projects/das/COMP/DAS/src/python/DAS/utils/utils.py
1627587    9.383    0.000   28.867    0.000 /data/projects/das/COMP/DAS/src/python/DAS/utils/utils.py
969267    7.716    0.000   10.656    0.000 /data/projects/das/slc5_amd64_gcc434/external/py2-pymongo
392636    4.499    0.000   47.851    0.000 /data/projects/das/COMP/DAS/src/python/DAS/utils/utils.py
  1      3.877   3.877   72.942   72.942 /data/projects/das/COMP/DAS/src/python/DAS/core/das_mongo
1153242    3.644    0.000    5.443    0.000 /data/projects/das/COMP/DAS/src/python/DAS/utils/utils.py
576626    3.042    0.000   89.345    0.000 /data/projects/das/COMP/DAS/src/python/DAS/core/das_mongo
576715    3.002    0.000    8.917    0.000 /data/projects/das/slc5_amd64_gcc434/external/py2-pymongo
969267    2.798    0.000   17.809    0.000 /data/projects/das/slc5_amd64_gcc434/external/py2-pymongo
.....

```

1.13 DAS configuration file

DAS configuration consists of a single file, \$DAS_ROOT/etc/das.cfg. Its structure is shown below:

```

[das]                                # DAS core configuration
verbose = 0                          # verbosity level, 0 means lowest
parserdir = /tmp                     # DAS PLY parser cache directory
multitask = True                     # enable multitasking for DAS core (threading)
core_workers = 10                    # number of DAS core workers who contact data-providers
api_workers = 2                      # number of API workers who run simultaneously
thread_weights = 'dbs:3','phedex:3' # thread weight for given services
error_expire = 300                   # expiration time for error records (in seconds)
emptyset_expire = 5                  # expiration time for empty records (in seconds)
services = dbs,phedex                # list of participated data-providers

[cacherequests]
Admin = 50                           # number of queries for admin user role
Unlimited = 10000                     # number of queries for unlimited user role
ProductionAccess = 5000              # number of user for production user role

```

```
[web_server]                # DAS web server configuration parameters
thread_pool = 30             # number of threads for CherryPy
socket_queue_size = 15       # queue size for requests while server is busy
host = 0.0.0.0               # host IP, the 0.0.0.0 means visible everywhere
log_screen = True            # print log to stdout
url_base = /das              # DAS server url base
port = 8212                  # DAS server port
pid = /tmp/logs/dsw.pid      # DAS server pid file
status_update = 2500         #
web_workers = 10             # Number of DAS web server workers who handle user requests
queue_limit = 200            # DAS server queue limit
adjust_input = True          # Adjust user input (boolean)
dbs_daemon = True            # Run DBSDaemon (boolean)
dbs_daemon_interval = 300    # interval for DBSDaemon update in sec
dbs_daemon_expire = 3600     # expiration timestamp for DBSDaemon records
hot_threshold = 100          # a hot threshold for powerful users
onhold_daemon = True         # Run onhold daemon for queries which put on hold after hot threshold

[dbs]                        # DBS server configuration
dbs_instances = prod,dev     # DBS instances
dbs_global_instance = prod   # name of the global DBS instance
dbs_global_url = http://a.b.c # DBS data-provider URL

[mongodb]                   # MongoDB configuration parameters
dburi = mongodb://localhost:8230 # MongoDB URI
bulkupdate_size = 5000       # size of bulk insert/update operations
dbname = das                 # MongoDB database name
lifetime = 86400             # default lifetime (in seconds) for DAS records

[dasdb]                     # DAS DB cache parameters
dbname = das                 # name of DAS cache database
cachecollection = cache      # name of cache collection
mergecollection = merge      # name of merge collection
mrcollection = mapreduce     # name of mapreduce collection

[loggingdb]
capped_size = 104857600
collname = db
dbname = logging

[analyticsdb]               # AnalyticsDB configuration parameters
dbname = analytics           # name of analytics database
collname = db                # name of analytics collection
history = 5184000            # expiration time for records in an/ collection (in seconds)

[mappingdb]                 # MappingDB configuration parameters
dbname = mapping             # name of mapping database
collname = db                # name of mapping collection

[parserdb]                  # parserdb configuration parameters
dbname = parser              # parser database name
enable = True                # use it in DAS or not (boolean)
collname = db                # collection name
sizecap = 5242880            # size of capped collection

[dbs_phedex]                # dbs_phedex configuration parameters
expiration = 3600            # expiration time stamp
urls = http://dbs,https://phedex # DBS and Phedex URLs
```

For up-to-date configuration parameters please see *utils/das_config.py*

1.14 Release notes

Release 2.X.Y series

This release series is targeted to DAS production stability and quality.

- 2.15.X
 - replace DASPLY/DASParserDB with DASQueryParser. It parses input DAS query into MongoDB spec and provides thread-safe code.
 - enable custom `das_ql_parser` module instead of `PLY`
 - add security module
 - modify `das_client` to check glidein/auth settings, add warning for end-users w/o auth, change version
- 2.14.X
 - work on IB issue, 4234 and DAS server timeouts
 - Add Singleton class
 - Add `MongoConnection` singleton class and use it in `DBConnection`
 - adjust code to work with `pymongo2/pymongo3` drivers
- 2.13.X
 - futurize stage1 changes in place
 - adjust code to use both `pymongo` version branches, version 2.X and version 3.X. Developers of `pymongo` driver changed significantly naming convention for both `MongoClient` and `find` APIs, therefore we wrap our code into independent `das_pymongo` module which handle the differences
 - Move `pycurl` options from the code and put them into DAS configuration Increase threshold for `CONNECTTIMEOUT` option.
- 2.12.X
 - add `nlTK_legacy` module to keep functions from `NLTK-2` release
 - create `db.collection` with empty storage (required by `mongo3/pymongo2.8`)
- 2.11.X
 - fix issues: 4212
- 2.10.X
 - fix issues: 4207, 4209, 4210
- 2.9.X
 - replaced hints ajax calls with template parsing of DAS records which carry hints parts
 - Add support for dataset `prepid=XXX` queries (done via `Reqmgr`, `MCM`, `DBS3` calls)
 - Add `prepid` support via `DBS3` datasets call
 - fixed: 4201, 4200, 4198, 4197, 4196, 4190, 4124, 4188
- 2.8.X

- show type of config link, ticket 4189
- adjust `das_client` to show DAS keys/attributes via `–list-attributes` option
- add Keys menu which shows list of DAS keys/attributes
- add functionality to accumulate information about DAS keys/attributes within DAS keylearning DB (ticket #4190)
- fixed: 4188, 4178
- 2.7.X
 - remove DBS2 code base
 - fix 4186 (it requires DBS3 fix for their issue dmwm/DBS#381)
 - add pattern for data tier and data types
- 2.6.X
 - Fix issue with 3 slashes and hints on web UI
 - Perform DAS/MongoDB separation
 - Update DAS docs accordingly with recent changes
 - Remove analytics
 - Re-factor DASMaps management scripts, now we'll use `das_js_fetch`, `das_js_validate`, `das_js_import`, `das_js_update`
 - Re-factor DASMaps creation scripts, now we'll use `das_create_json_maps`, `das_create_kws_maps`
 - * drop from KWS maps `_id/oid` feilds
 - Simplify hints templates
 - Remove `db_monitor` thread due to DAS/MongoDB separation
 - Re-factor `bin/das_server` script to support DAS/MongoDB separation
 - Add support for `child dataset=X release=Y site=Z` query
 - Add support for `dataset parent=X release=Y site=Z` query
- 2.5.X
 - Fixed 4166, 4165, 4164, 4162, 4160, 4158, 4156, 4155, 4153, 4148, 4142
 - Re-factor ReqMgr config queries to fetch configs from ReqMgr and WMStats
 - Add dataset suggestion hints when no results are found in default DBS
 - Fixed query disappearance under high load
 - Set `DAS_SERVER` in `DAS/__init__.py` module which will be used in HTTP User-Agent header. The `DAS_SERVER` uses version which will be updated upon deployment to specific release name. Therefore the User-Agent string will change once we deploy new DAS release. It has the following pattern: `das-server/<DAS release>::python/<python version>`, e.g. `das-server/2.4.8::python/2.7`
 - Change DASMapping to be singleton
 - Switch to `<dbns_namespace>/<dbns_instance>` schema for DBS3 maps, e.g. `prod/global`, `int/global`, `prod/phys01`, `int/phys01`
 - Reorganize DAS map upload (new scripts: `das_create_json_maps`, `das_update_database`)
 - Add support for Travis CI

- Improve error handling, add DASHTMLParser to parse HTML content coming from data-providers (e.g. when service is down)
- Change behavior of DAS CLI to show bytes as is and upon user request its representation in certain base, e.g. `-base=2` will show MiB notations, etc.
- Support DBS3 instances with slashes, e.g. `prod/global`, `int/global`
- Extend list of supported DBS3 instances, in addition to `global` and `phys0[1-3]`, I added instances with slashes, e.g. `int/global`. If instances does not have slash it is considered to come from prod DBS URL.
- 2.4.X
 - Re-factor RequestManager code to use internal store instead of MongoDB one
 - Re-evaluate racing conditions:
 - * the clean-up worker should use lock to safely wipe out expired records
 - * the `remove_expired` method should only care about given dasquery records
 - * add additional delta when check expiration timestamp, it is required to protect code from situation when records can be wiped out during request operation
 - Fixed issues: 4098, 4097, 4090, 4095, 4093, 4089, 4085, 4082, 4081, 4079, 4077
 - Wrapped `dasply` parser call into `spawn` function, this fix intermittent problem with `dasply` under high load
 - Added `spawn_manager` along with its unit test
 - Re-factor code to use individual DB connections instead of sharing one in various places. This work address DAS instability issue reported in #4024.
 - Added cookie support in DAS client
 - Added support of DAS client version, it is configured via `config.web_server.check_clients` option
 - * check DAS client version in DAS server
 - * issue warning status and associative message for the client about version mismatch
 - Added code to re-initiate DAS request in `check_pid`, issue #4060
 - Prepared server/client for distributed deployment on cmsweb
- 2.3.X
 - Fix issue 4077 (don't yield record from empty MCM response)
 - Work on DAS deployment procedure
 - * two set of DAS maps (one for production URLs and another for testbed)
 - Separate DAS web server with KWS one
 - * run KWS on dedicated port (8214) with single-threaded DASCore
 - Assign type for records in mapping db; DASMapping code re-factoring to reduce latency of the queries
 - Fix query default assignment (issue #4055)
 - Provide ability to look-up config files used by ReqMgr for dataset creation (issue #4045)
 - * Add config look-up for given dataset which bypass MCM data-service
 - Add clean-up worker for DAS caches; fixed #4050
 - Additional work on stability of DAS server under high-load (issue #4024)
 - * add exhaust option to all MongoDB find calls

- * add index for mapping db
- * add DAS clean-up daemon, issue #4050
- * reduce pagination usage
- * increase MongoDB pool size, turn on `auto_start_request` and `safe` option for MongoClient
- * step away from db connection singleton
- Add outputdataset API for ReqMgr data-service (issue #4043)
- Fix python3 warnings
- 2.2.X
 - Fixed MCM prepid issue, switch to produces rest API
 - Merge and integrated KWS search
- 2.1.X
 - Make `das-headers` mandatory in `das_client.py` and keep `–das-headers` option and `das_headers` input parameter in `get_data` API for backward compatibility
 - Replaced all `has_key` dict calls with `key in dict` statement (work towards python3 standard)
 - Add `check_services` call to DAS core to check status of services
 - Pass `write_concern` flag to MongoClient, by default it is off
 - Fixed #4032
 - Re-factor core/web code to propagate error record back to end-user and setup error status code in this case
 - Throw error records when `urlfetch_getdata` fails
 - Move `set_misses` into `write_to_cache`
 - Made adjustments to DBS3 data-service based on recent changes of DBS3 APIs
- 2.0.X
 - Add `services` attribute to `das` part of data record, it shows which DAS services were used, while `system` attribute used to show which CMS systems were used to produce data record(s)
 - Turn off `dbs_phedex`, it producing too much load, instead use individual services
 - Re-evaluate lifetime of records in DAS cache: the clean-up should be done either for `qhash/das.expire` pair (less then current `tstamp`) or for records which live in cache long enough, via `das.exire<tstamp-rec_ttl`
 - Introduce `dasdb.record_ttl` configuration parameter int `das config`
 - Fix issue4023
 - Changes to allow DAS run with DBS2/DBS3 in a mix mode
 - Extend download LFN link to download web page, issue 4022
 - Add Status link to DAS header and let users to see status of DAS queue
 - Re-factor `DASMapping` code, see ticket 4021
 - Add support for `mcm dataset=/a/b/c` query; first it looks-up information from ReqMgr to get its info for given dataset, then it parse ReqMgr info and extracts Prepid and passes it to MCM data-service.
 - Add MCM links on dataset summary page when information is provided by reqmgr data-service (MC datasets)

- Add code to support MCM (PREP) data-service (issue 3449), user can look-up mcm info by using the following query: mcm prepId=<PREP-ID>
- Remove timestamp attribute from passed dict to md5hash function, it is required due to dynamic nature of timestamp which leads to modification of the hash of the record
- Add new stress tool, see bin/das_stress_tool
- Round timestamp for map records as well as for dasheader due to inconsistent behavior of json parsers, see note in jsonwrapper module
- Fix issue4017: add hash to all DAS map records; add verification of hash into DASMapping check_maps method
- Fix issue4016: add aux-record called arecord; arecord contains count of corresponding map record, map record type and a system. Adjust DASMapping check_maps method to perform full check of DAS maps by comparing count field from aux-record with actual number of maps in DAS mapping DB
- Apply common set of indexes for both cache/merge collection to properly get/merge records
- Allow runs DBS3 API to yield individual records
- Support block tier=GEN-SIM date between [20120223, 20120224] query via blocksummaries DBS3 API
- Switch from block_names to block_name as input parameter for blocksummaries DBS3 API; handle correctly incorrect values for dates in DBS3 blocksummaries API
- Fix issues 4014, 4013, 4009
- Add lumi4block_run and dataset4block DBS3 APIs
- fix run input parameter for all DBS3 APIs
- Add runsummaries API

Release 1.X.Y series

- 1.12.X
 - Fix wildcards to provide more informative messages in text mode
 - Fix issues: 3997, 3975
 - Replace phedex_tier_pattern with phedex_node_pattern
 - Get rid of empty_record, query. Instead, introduce das.record with different codes. Codes are defined in utils/utils.py record_codes function. Add mongodb index on codes; modified queries to look-up das/data-records using new das.record field
 - Fix issue with ply_query parameter
 - Add extra slash to avoid one round trip
 - Work on support new run parameter w/ DBS3 APIs, now DAS is capable to use run-range/run-list queries into DBS3
 - Use json.dumps to printout JSON dict to stdout
- 1.11.X
 - Add support for block,run,lumi dataset=/a/b/c queries
 - Add plistlib python module w/ None modifications to handle DAS XML output
 - Add list of attributes for config output

- Add summary4block_run API
 - Highlight unknown global tags in web UI
 - Re-factor the code: add insert_query_records which scan input DAS query and insert query records into DAS cache, then it yields list of acknowledged data-services which used by call API for data retrieval
 - Extend incache API to work with query or data records by providing query_record flag with default value of False (check data records)
 - Take care of potential failure of PLY parser. Use few trials on given input and then give-up
 - Fix bug in task manager when I mix-up return type of spawn function which cause task fails under race conditions
 - Add support for summary dataset=/a/b/c query without run conditions
 - Add support for run range in DBS2 summary dataset/run query
 - Add expand_lumis helper function into das aggregators which flatten lumi lists, e.g. [[1,3], [5,7]] into [1,2,3,5,6,7]. This allows correctly count number of lumis in DAS records
 - Implement support for comp-ops queries, e.g. find run, lumi for given dataset and optional run range find file, lumi for given dataset and optional run range find file, lumi, run for given dataset and optional run range this work is done via new urlfetch_getdata module
- 1.10.X
 - Add urlfetch_pycurl module to fetch content from multiple urls
 - Use custom db_monitor which check MongoDB connection as well as periodically reload DAS maps
 - Add preliminary support for file block=/a/b/c#123 runs site query (must have urlfetch proxy)
 - Allow user to get DBS file into regardless of its status, ticket 3992
 - Add indexes for file.name,dataset.name.block.name and run.run_number in DAS cache collection to prevent error on sorting entities
 - Add support for block dataset run in/between [1,2] query, ticket 3974
 - Apply file.name index to allow MongoDB to sort the files, ticket 3988 this is required in rare case when number of files is very large and MongoDB give up on sorting without the index. I may apply similar index on block as well since their number in dataset can be large as well.
 - Add constrain on block name for lumi block=/a/b/c#123 queries, ticket 3977
 - Add pyurlfetch client
 - Add proxy_getdata to request data from external urlproxy server, ticket 3986; should be used to fetch data concurrently
 - Add support for file dataset=/a/b/c run in [1,2,3] site=T2_CH_CERN, ticket 3982 (requires external url-proxy server, see 3986)
 - Split fakeDatasetSummary into fakeDatasetPattern and fakeDatasetSummary to support look-up of valid datasets for given pattern and any dataset info for givan dataset path; ticket 3990
 - Add draft code to accommodate file dataset=/a/b/c run in [1,2,3] site=X query (still under development)
 - Add url_proxy module which can work with pyurlfecth or Go proxy server
 - Add get_proxy, proxy_getdata and implementation (still experimental) of proxy usage within DBS3 module
 - Re-wrote update_query_record API; update ctime for query records

- Separate insertion of query and data records
- Remove analytics calls from abstract service, current analytics implementation require full re-design, it does not make any good so far
- Add distinguishing message in ticket issue title for no apis/no results errors
- Add fakeFiles4BlockRun API to cover file block=/a/b/c#123 run in [1,2,3] queries required by CMSSW Integration Builds (IB).
- Fix file block=/a/b/c#123 query (DBS should contribute to it)
- Add dataset pattern constraints for all DBS/DBS3 queries
- Remove listLFNs since listFiles cover the use case to look-up file for a given dataset
- Add filelumis4dataset API to support file,lumi dataset=/a/b/c queries
- Add support for run IN [1,2,3] queries, this will be allowed in DBS/DBS3, CondDB, RunRegistry data-services
- Upgrade to Prototype.js 1.7
- Remove lumi API from CondDB mapping; add lumi API to RunRegistry mapping; clean-up RunRegistry code and remove v2 APIs, the v3 is default now
- Re-factor Vidmantas code: move wild-card errors into separate template; sanitize template parameters; clean-up code
- Add das_exceptions module, move all Wild-card exception into this module
- Improve web UI links with box_attention for submitting DAS tickets, ticket #3969
- 1.9.X
 - Fix ticket #3967 (preserve DAS records order while removing duplicates)
 - Fix ticket #3966 (strip-off zero in das filters)
 - Add JS function to handle Event (hide DAS keys window) via ESC
 - Resolve double counting issue, ticket #3965
 - Add Show DAS keys description to web UI
 - Wrap combined_site4dataset API call into try/except block and show exception on web UI. This will help to catch transient missing values from combined data-service for site dataset=/a/b/c queries.
 - Add DASKEY EQUAL VALUE VALUE error condition to DAS PLY parser to cover the case when user cut-and-paste some value and it has empty space, e.g. dataset=/a/b/c om
 - Always use upper() for DBS status since it is stored in upper-case in DBS DB
 - Add function to print DAS summary records
 - Add DAS SERVER BUSY message to web server, ticket #3945
 - Read prim_key from mapping DB rather than lookup_keys in das_mongocache module (with fallback to lookup_keys)
 - Fix verbose printout for pycurl_manager module
 - Add support for summary dataset=/a/b/c run=123, ticket #3960
 - Re-factor das_client to be used in other python application; change return type from str to json in get_data API; add das-headers flag to explicitly ask for DAS headers, by default drop DAS headers
 - Re-factor dasmongocache code to support multiple APIs responses for single DAS key

- Add `api=das_core` to `dasheader` when we first register query record
- Extend DAS aggregator utility to support multiple APIs response for single DAS key
- Add `db_monitor` threads to `DASMapping/DASMongocache` classes
- Switch from explicit show/hide links to dynamic show/hide which switch via `ToggleTag` JS function
- Adjust web UI with Eric's suggestions to show service names in color boxes; remove DAS color map line in result output
- Revert to base 10 in `size_format`
- Add `update_filters` method to `DASQuery` class to allow upgrade its filters with spec keys; this is useful on web UI, when end-user specifies a filter and we need to show primary key of the record
- Wrote `check_filters` function to test applied filters in a given query and invoke it within `nresults` method, ticket #3958
- Collapse lumi list from DBS3, ticket #3954
- Remove `db` url/instances from DAS configuration and read this information directly from DAS maps; fixed #3955
- 1.8.X
 - Add support of lumi `block=/a/b/c#123` and `block file=/path/f.root` queries both in DBS and DBS3
 - Do not check field keys in a query, e.g. allow to get partial results
 - Fix plain web view when using DAS filters
 - Extend DAS support for file dataset `=/a/b/c` run between [1,2] queries
 - Keep links around even if data service reports the error
 - Catch error in combined data-service and report them to UI
 - Protect `qxml_parser` from stream errors
 - Convert regex strings into raw strings
 - Separate curl cache into get/post instances to avoid racing condition for cached curl objects
 - Convert das timestamp into presentation datetime format
 - Queue type can be specified via `qtype` parameter in web section of DAS configuration file
 - Extend `task_manager` to support `PriorityQueue`
 - Revert default to `cjson` instead of `yajl` module, since later contains a bug which incorrectly rounds off large numbers; there is also an outstanding issue with potential memory leak
 - Remove dataset summary look-up information for dataset pattern queries to match DBS2 behavior and reduce DAS/DBS latency, see 9254ae2..86138bd
 - Replace `range` with `xrange` since later returns generator rather than list
 - Add capability to dump DAS status stack by sending `SIGQUIT` signal to DAS server, e.g. upon the following call `kill -3 <PID>` DAS server will dump into its logs the current snapshot of all its threads
 - Apply Vidmantas wildcard patch to improve usage of dataset patterns on web UI
 - Fix Phedex checksum parsing
 - Switch to new PyMongo driver, version 2.4
 - * change Connection to MongoClient

- * remove safe=True for all insert/update/remove operation on mongo db collection, since it is default with MongoClient
- DAS CLI changes:
 - * Add exit codes
 - * Add -retry option which allows user to decide if s/he wants to proceed with request when DAS server is busy; retry follows log⁵ function
 - * Set init waiting time to 2 sec and max to 20 sec; use cycle for sleep time, e.g. when we reach the max drop to init waiting time and start cycle again. This behavior reduce overall waiting time for end-users
- Fix issue with DBS3 global instance look-up
- Switch to HTML5 doctype
- New schema for DAS maps
 - * re-factor code to handle new schema
 - * change all maps/cms_maps according to new schema
 - * add new documentation for new schame, see mappings.rst
- Add support to look-up INVALID files in DBS2/DBS3
- Enable dbs_phedex combined engine
- Add new thread module to deal with threads in DAS
- Switch from low-level thread.start_new_thread to new DAS thread module, assign each thread a name
- Properly handle MongoDB connection errors and print out nice output about their failure (thread name, time stamps, etc.)
- 1.7.X
 - Switch from PRODUCTION to VALID dataset access type in DBS3
 - Adjust das_core and das_mongocache to optionally use dasquery.hashes
 - * hashes can be assigned at run-time for pattern queries, e.g. dataset=/abc
 - * hashes can be used to look-up data once this field is filled up
 - Let DBSDaemon optionally write dataset hashes, this can be used to enhance dataset pattern look-up in DAS cache, see ticket #3932
 - Add hashes data member and property to DASQuery class
 - Work on DBS3 APIs
 - Fix issue with forward/backward calls in a browser which cause existing page to use ajaxCheckPid. I added reload call which enforces browser to load page content with actual data
 - * revisit ajaxCheckPid and check_pid functions. Removed ahash, simplify check_pid, use reload at the end of the request/check_pid handshake
 - Add fakeDataset4Site DBS2 API to look-up datasets for a given site, ticket #3084
 - * DBS3 will provide new API for that
 - Change DAS configuration to accept web_service.services who lists local DAS service, e.g. dbs_phedex, dbs_lumi
 - Modify dbs_phedex service to initialize via DAS maps

- Add lumi_service into combined module
- Introduced services mapping key
- Adjust combined map file to use services mapping key
- Switch to pycurl HTTP manager, which shows significant performance boost
- Work on pycurl_manager to make it complaint with httplib counterpart
- 1.6.X
 - Add new logging flag to enable/disable logging DAS DB requests into logging db (new flag is dasdb.logging and its values either True or False)
 - Change pymongo.objectid to bson.objectid, pymongo.code to bson.code since pymongo structure has been changed (since 2.2.1 pymongo version)
 - Introduce new dataset populator tool which should fetch all DBS datasets and keep them alive in DAS cache (not yet enabled)
 - Move DAS into github.com/dmwm organization
 - Extend das_dateformat to accept full timestamp (isoformat); provide set of unit tests for das_dateformat; fix web UI to accept date in full isoformat (user will need to provide quotes around timestamp, e.g. '20120101 01:01:01'); fixes #3931
 - Set verbose mode only when parserdb option is enabled
- 1.5.X
 - Add SERVICES into global scope to allow cross service usage, e.g. site look-up for DBS dataset records
 - Add site look-up for user based datasets, ticket #3432
 - Revisit onhold daemon and cache requests flaw
 - * Start onhold daemon within init call (ensure MongoDB connection)
 - * Check DAS cache first for CLI requests regardless if pid presence in a request
 - * Put requests on hold only if user exceeds its threshold and server is busy, otherwise pass it through
 - Set DAS times, ticket #3758
 - Convert RR times into DAS date format (isoformat)
 - Fix ticket #3796
- 1.4.X
 - Move code to github
 - Fix bug in testing for numbers, SiteDB now contains unicode entries
 - Add HTTP links into record UI representation
 - Call clean-up method upon request/cache web methods.
 - Add htlKeyDescription, gtKey into RunRegistry, ticket #3735
 - Improve no result message, ticket #3724
 - Update error message with HTTPError thrown by data-provider, ticket #3718
 - Fix das_client to proper handle DAS filters, ticket #3706
 - Change Error to External service error message, see ticket #3697
 - Skip reqmgr API call if user provide dataset pattern, ticket #3691

- Enable cache threshold reading via SiteDB group authorization
- Add support for block dataset=/bla run=123 query, ticket #3688
- Fix tickets #3636, #3639
- 1.3.X
 - Add new method for SiteDB2 which returns api data from DAS cache
 - Add parse_dn function to get user info from user DN
 - Add new threshold function which parse user DN and return threshold (it consults sitedb and look-up user role, if role is DASSuperUser it assigns new threshold)
 - Add suport_hot_threshold config parameter to specify hot threshold for super users
 - Extend check_pid to use argument hash (resolve issue with compeing queries who can use different filters)
 - Do not rely on Referrer settings, ticket #3563
 - Fix tickets #3555, #3556
 - Fix plain view, ticket #3509
 - Fix xml/json/plain requests via direct URL call
 - Clean-up web server and checkargs
 - Add sort filer to web UI
 - Add sort filter, users will be able to use it as following file dataset=/a/b/c | sort file.size, file dataset=/a/b/c | sort file.size- The default order is ascending. To reverse it, user will need to add minus sign at the end of the sort key, e.g. file.size-
 - Re-factor code to support multiple filters. They now part of DASQuery object. All filters are stored as a dict, e.g. {'grep': <filter list>, 'unique': 1, 'sort': 'file.size'}
 - Add sitedb links for site/user DAS queries
 - Re-factor code which serves JS/CSS/YUI files; reduce number of client/server round-trips to load those files on a page
 - fix ddict internal loop bug
 - add representation of dict/list values for given key attributes, e.g. user will be able to select block.replica and see list of dicts on web page
- 1.2.X
 - Pass instance parameter into das_duplicates template, ticket #3338
 - Add qhash into data records (simplify their look-up in mongocache manager)
 - Simplify query submission for web interface (removed obsolete code from web server)
 - Fix issue with sum coroutines (handle None values)
 - Avoid unnecessary updates for DAS meta-records
 - Made das core status code more explicit
 - Remove ensure_index from parser.db since it's capped collection
 - Made QLManager being a singleton
 - Add safe=True for all inserts into das.cache/merge collection to avoid late records arrival in busy multi-threaded environment

- Add trailing slash for condDB URL (to avoid redirection)
- Show data-service name in error message
- Show dataset status field
- Add support to pass array of values into DAS filter, ticket #3350 but so far array needs to consist of single element (still need to fix PLY)
- Update TFC API rules (fix its regex in phedex mapping)
- Init site.name with node.name when appropriate
- Fill admin info in new SiteDB when user look-up the site
- Switch to new SiteDB
- Switch to new REST RunRegistry API
- Remove dbs instance from phedex subscription URL and only allow DBS global link, ticket #3284
- Fix issue with invalid query while doing sort in tableview (ticket #3281) discard qhash from the tableview presentation layer
- Implement onhold request queue. This will be used to slow down users who sequentially abuse DAS server. See ticket #3145 for details.
- Add qhash into DASquery __str__
- Fix issue with downloading config from gridfs, ticket 3245
- Fix DBS run in query with wide run range, use gte/lte operators instead
- Fix issue with recursive calls while retrieve dict keys
- Eliminate duplicates in plain view, ticket 3222
- Fix fakeFiles4DatasetRunLumis API call and check its required parameters
- Fix plain view with filter usage, ticket #3216
- Add support for dataset group=X site=T3_XX_XXXX or dataset group=X site=a.b.com queries via block-replicas Phedex API, ticket #3209
- Fix IP look-up for das_stats, ticket #3208
- Provide match between various SiteDB2 APIs in order to build combined record
- Remove ts field and its index from das.cache collection, it is only needed for das.merge
- Work on integration with new SiteDB, ticket #2514
- Switch to qhash look-up procedure, ticket #3153
- Fix DBS summary info, ticket #3146
- Do not reflect request headers, ticket #3147
- Fix DBSDaemon to work with https for DBS3
- Add ability to DAS CLI to show duplicates in records, ticket #3120
- Parse Phedex checksum and split its value into Adler32/checksum, ticket #3119, 3120
- Remove from dataset look-up for a given file constrain to look-up only VALID datasets, when user provide a file I need to look-up dataset and provide its status, ticket #3123
- Resolved issue with duplicates of competing, but similar queries at web UI.
- Changed task manager to accept given pid for tasks.

- Generated pid at web layer; check status of input query in a cache and find similar one (if found check status of similar request and generate results upon its completion); moved check_pid code from web server into its one template; adjusted ajaxCheckPid call to accept external method parameter (such that I can use different methods, e.g. check_pid and check_similar_pid)
- Fixed several issues with handling StringIO (delivered by pycurl)
- 1.1.X
 - Extend not equal filter to support patterns, ticket #3078
 - Reduce number of DAS threads by half (the default values for workers was too high)
 - Name all TaskManagers to simplify their debugging
 - Configure number of TaskManager for DASCore/DASAbstractService via das configuration file
 - Fix issue with data look-up from different DBS instances (introduce instance in das part of the record), ticket #3058
 - Switch to generic DASQuery interface. A new class is used as a placeholder for all DAS queries. Code has been refactored to accept new DASQuery interface
 - Revisited analytics code based on Gordon submission: code-refactoring; new tasks (QueryMaitainer, QueryRunner, AnalyticsClenup, etc); code alignment with DAS core reorganization, ticket #1974
 - Fix issue with XML parser when data stream does not come from data-service, e.g. data-service through HTTP error and DAS data layer creates HTTP JSON record
 - Fix bug in db_monitor who should check if DB connection is alive and reset DB cursor, ticket #2986
 - Changes for new analytics (das_singleton, etc.)
 - Add new tool, das_stats.py, which dumps DAS statistics from DAS logdb
 - Add tooltip template and tooltips for block/dataset/replica presence; ticket #2946
 - Move creation of logdb from web server into mongocache (mongodb layer); created new DASLogdb class which will responsible for logdb; add insert/deletion records into logdb; change record in logdb to carry type (e.g. web, cache, merge) and date (in a form of yyyyymmdd) for better querying
 - add gen_counter function to count number of records in generator and yield back records themselves
 - add support for != operator in DAS filters and precise match of value in array, via filter=[X] syntax, ticket #2884
 - match nresults with get_from_cache method, i.e. apply similar techniques for different types of DAS queries, w/ filters, aggregators, etc.
 - properly encode/decode DAS queries with value patterns
 - fix issue with system keyword
 - allow usage of combined dbs_phedex service regardless of DBS, now works with both DBS2 and DBS3
 - Fix unique filter usage in das client, add additions to convert timestamp/size into human readable format, ticket #2792
 - Retire DASLogger in favor of new PrintManager
 - code re-factoring to address duplicates issue; ticket #2848
 - add dataset/block/replica presence, according to ticket #2858; made changes to maps
- 1.0.X
 - add support for release file=lfm query, ticket #2837

- add creation_time/modification_time/created_by/modified_by into DBS maps, ticket #2843
- fix duplicates when applying filters/aggregators to the query, tickets #2802, #2803
- fix issue with MongoDB 2.x index lookup (error: cannot index parallel arrays).
- test DAS with MongoDB 2.0.1
- remove IP lookup in phedex plugin, ticket #2788
- require 3 slashes for dataset/block pattern while using fileReplicas API, ticket #2789
- switch DBS3 URL to official one on cmsweb; add dbs3 map into cms_maps
- migrate from http to https for all Phedex URLs; ticket 2755
- switch default format for DAS CLI; ticket 2734
- add support for 'file dataset=/a/b/c run=1 lumi=80' queries both in DBS2/DBS3, ticket #2602
- prohibit queries with ambiguous value for certain key, ticket #2657
- protect filter look-up when DAS cache is filled with error record, ticket #2655
- fix makepy to accept DBS instance; ticket #2646
- fix data type conversion in C-extension, ticket #2594
- fix duplicates shown in using DAS CLI, ticket #2593
- add Phedex subscription link, fixes #2588
- initial support for new SiteDB implementation
- change the behavior of compare_spec to only compare specs with the same key content, otherwise it leads to wrong results when one query followed by another with additional key, e.g. file dataset=abc followed by file dataset=abc site=X. This lead compare_spec to identify later query as subset of former one, but cache has not had site in records, ticket #2497
- add new data retrieval manager based on pycurl library; partial resolution for ticket #2480
- fix plain format for das CLI while using aggregators, ticket 2447
- add dataset name to block queries
- add DAS timestamp to all records; add link to TC; fixes #2429, #2392
- re-factor das web server, and put DAS records representation on web UI into separate layer. Create abstract representation class and current CMS representation. See ticket 1975.

Release 0.9.X series

- 0.9.X
 - change RunRegistry URL
 - fix issue with showing DAS error records when data-service is down, see ticket #2230
 - add DBS prod local instances, ticket 2200
 - fix issue with empty record set, see tickets #2174, 2183, 2184
 - upon user request highlight in bold search values; dim off other links; adjust CSS and das_row template, ticket #2080
 - add support for key/cert in DAS map records, fixes #2068

- move DotDict into stand-alone module, fixes #2047
- fix block child/parent relationship, tickets 2066, 2067
- integrate DAS with FileMover, add Download links to FM for file records, ticket #2060
- add filter/aggregator builder, fixes #978
- remove several run attributes from DBS2 output, since this information belong to CondDB and is not present in DBS3 output
- add das_diff utility to check merged records for inconsistencies. This is done during merge step. The keys to compare are configurable via presentation map. So far I enable block/file/run keys and check for inconsistencies in size/nfiles/nevents in them
- replace ajax XHR recursive calls with pattern matching and onSuccess/onException in ajaxCheck-Pid/check_pid bundle
- walk through every exception in a code and use print_exc as a default method to print out exception message. Adjust all exception to PEP 3110 syntax
- code clean-up
- replace traceback with custom print_exc function which prints all exceptions in the following format: msg, timestamp, exp_type, exc_msg, file_location
- remove extra cherryypy logging, clean-up DAS server logs

Release 0.8.X series

- 0.8.X
 - resolve double requests issue, ticket #1881, see discussion on HN <https://hypernews.cern.ch/HyperNews/CMS/get/webInterfaces/708.html>
 - Adjust RequestManager to store timestamp and handle stale requests
 - Make DBSDaemon be aware of different DBS instances, ticket #1857
 - fix getdata to assign proper timestamp in case of mis-behaved data-services ticket #1841
 - add dbs_daemon configuration into DAS config, which handles DBS parameters for DBSDaemon (useful for testing DBS2/DBS3)
 - add TFC Phedex API
 - add HTTP Expires handling into getdata
 - made a new module utils/url_utils.py to keep url related functions in one place; remove duplicate getdata implementation in combined/dbs_phedex module
 - add dbs_daemon whose task to fetch all DBS dataset; this info is stored into separate collection and can be used for autocompletion mode
 - improve autocompletion
 - work on scalability of DAS web server, ticket #1791

Release 0.7.X series

This release series is targeted to DAS usability. We collected users requests in terms of DAS functionality and usability. All changes made towards making DAS easy to use for end-users.

- 0.7.X
 - ticket #1727, issue with index/sort while getting records from the cache
 - revisit how to retrieve unique records from DAS cache
 - add DAS query builder into autocomplete
 - extend reflex to support free-text based queries
 - add DBS status keyword to allow to select dataset with different statuses in DBS, the default status is VALID, ticket #1608
 - add datatype to select different type of data, e.g. MC, data, calib, etc.
 - if possible get IP address of SE and create appropriate link to ip service
 - calculate run duration from RR output
 - add conddb map into cms_maps
 - add initial support for search without DAS keywords
 - apply unique filter permanently for output results
 - add help cards to front web page to help users get use with DAS syntax
 - work on CondDB APIs
 - fix issue with IE
 - turn off multitask for analytics services
 - add query examples into front-page
 - get file present fraction for site view (users want to know if dataset is completed on a site or not)
 - fix PLY to accept yln as a value, can be used to check openness of the block
 - add create_indexes into das_db module to allow consistently create/ensure indexes in DAS code

Release 0.6.X series

This release series is targeted towards DAS production version. We switched from implicit to explicit data retrieval model; removed DAS cache server and re-design DAS web server; add multitasking support.

- 0.6.5
 - handle auto-connection recovery for DBSPhedexService
 - fix site/se hyperlinks
- 0.6.4
 - create new DBSPhedexService to answer the dataset/site questions. it uses internal MongoDB to collect info from DBS3/Phedex data-services and map-reduce operation to extract desired info.
- 0.6.3
 - support system parameter in DAS queries, e.g. block block=/a/b/c#123 system=phedex
 - add condition_keys into DAS records, this will assure that look-up conditions will be applied properly. For instance, user1 requested dataset site=abc release=1 and user2 requested dataset site=abc. The results of user1 should not be shown in user2 queries since it is superset of previous query. Therefore each cache look-up is supplemented by condition_keys

- add support for the following queries: dataset release=CMSSW_4_2_0 site=cmssrm.fnal.gov dataset release=CMSSW_4_2_0 site=T1_US_FNAL
- add new combined DAS plugin to allow combined queries across different data services. For instance, user can request to find all datasets at given Tier site for a given release. To accomplish this request I need to query both DBS/Phedex. Provided plugin just do that.
- add new method/tempalte to get file py snippets
- re-factor code which provide table view for DAS web UI
- add new phedex URN to lookup files for a given dataset/site
- put instance as separate key into mongo query (it's ignored everywhere except DBS)
- work on web UI (remove view code/yaml), put dbs instances, remember user settings for view/instance on a page
- add physics group to DBS2 queries
- add support to look-up of sites for a given dataset/block
- allow to use pattern in filters, e.g. block.replica.site=*T1*
- add filters values into short record view
- add links to Release, Children, Parents, Configs into dataset record info
- add support to look-up release for a given dataset
- add support to look-up cofiguration files for given dataset
- add fakeConfig, fakeRelease4Dataset APIs in DBS2
- add support for CondDB
- add hyperlinks to DAS record content (support only name, se, run_number), ticket #1313
- adjust das configuration to use single server (remove cache_server bits)
- switch to single server, ticket #1125
 - * remove web/das_web.py, web/das_cache.py
- switch to MongoDB 1.8.0
- 0.6.2
 - das config supports new parameters queue_limit, number_of_workers)
 - add server busy feature (check queue size vs nworkers, reject requests above threshold), ticket #1315
 - show results of agg. functions for key.size in human readable format, e.g. GB
 - simplify DASCACHEMgr
 - fix unique filter #1290
 - add missing fakeRun4File API to allow look-up run for a given file, fixes #1285
 - remove 'in' from supported list of operator, users advised to use 'between' operator
 - DBS3 support added, ticket #949
 - fix #1278
 - fix #1032; re-structure the code to create individual per data-srv query records instead of a single one. Now, each request creates 1 das query record plus one query record per data-srv. This allows to assign different expire timestamp for data-srv's and achieve desired scalability for data-service API calls.

- re-wrote task_manager using threads, due to problems with multiprocessing modules
- re-wrote cache method for DAS web servers to use new task_manager
- adjust das_client to use new type of PID returned by task_manager upon request. The PID is a hash of passed args plus time stamp
- bump to new version to easy distinguish code evolution
- 0.6.1
 - replace gevent with multiprocessing module
 - add task_manager which uses multiprocessing module and provides the same API as gevent
- 0.6.0
 - code refactoring to move from implicit data look-up to explicit one. The 0.5.X series retrieved all data from multiple sources based on query constrains, e.g. dataset=/a/b/c query cause to get datasets, files, block which match the constraint. While new code makes precise matching between query and API and retrieve only selected data, in a case above it will retrieve only dataset, but not files. To get files users must explicitly specify it in a query, e.g. file dataset=/a/b/c
 - constrain PLY to reject ambiguous queries with more then one condition, without specifying selection key, e.g. dataset=/a/b/c site=T1 is not allowed anymore and proper exception will be thrown. User must specify what they want to select, dataset, block, site.
 - protect aggregator functions from NULL results
 - new multiprocessing pool class
 - use gevent (if present, see <http://www.gevent.org/>) to handle data retrieval concurrently
 - switch to YAJL JSON parser
 - add error_expire to control how long expire records live in cache, fixes #1240
 - fix monitor plugin to handle connection errors

Release 0.5.X series

This release series is targeted to DAS stability. We redesigned DAS-QL parser to be based on PLY framework; re-write DAS analytics; add benchmarking tools; performed stress tests and code audit DAS servers.

- 0.5.11
 - change RunRegistry API
 - fix showing result string in web UI when using aggregators
 - bug fix for das_client with sparse records
 - add new das_web_srv, a single DAS web server (not enabled though)
 - fix das_top template to use TRACE rather then savannah
- 0.5.10
 - add DAS cache server time into the web page, fixes #941
 - remove obsolete yuijson code from DAS web server
 - use DASLogger in workers (instead of DummyLogger) when verbosity level is on. This allows to get proper printouts in debug mode.

- fix bug in compare_specs, where it was not capable to identify that str value can be equal to unicode value (add unittest for that).
- classified logger messages, move a lot of info into debug
- change adjust_params in abstract interface to accept API as well
- adjust DBS2 plugin to use adjust_params for specific APIs, e.g. listPrimaryDatasets, to accept other parameters, fix #934
- add new DAS keyword, parent, and allow parent look-up for dataset/file via appropriate DBS2 APIs
- extend usage of records DAS keyword to the following cases
 - * look-up all records in DAS cache and apply conditions, e.g. records | grep file.size>1, file.size<10
 - * look-up all records in DAS cache regardless of their content (good/bad records), do not apply das.empty_record condition to passed empty spec
- Fix filter->spec overwrite, ticket #958
- Add cache_cleaner into cache server, its task is periodically clean-up expired records in das.cache, das.merge, analytics.db
- Fix bug in expire_timestamp
- Remove loose query condition which leads to pattern look-up (ticket #960)
- Fix but in das_ply to handle correctly date
 - * add new date regex
 - * split t_DATE into t_DATE, t_DATE_STR
- add support for fake queries in DBS plugin to fake non-existing DBS API via DBS-QL
- remove details from DSB listFiles
- add adjust_params to phedex plugin
- adjust parameters in phedex map, blockReplicas can be invoked with passed dataset
- update cms_maps with fake DBS2 APIs
- add DAS_DB_KEYWORDS (records, queries, popular)
- add abstract support to query DAS (popular) queries, a concrete implementation will be added later
- fix #998
- fix SiteDB maps
- fix host parameter in das_cache_client
- remove sys.exit in das_admin to allow combination of multiple options together
- fix compare_specs to address a bug when query with value A is considered as similar to next query with value A*
- fix get_status to wait for completion of DAS core workflow
- fix merge insert problem when records exceed MongoDB BSON limit (4MB), put those records into GridFS
- fix nresults to return correct number of found results when applying a filter, e.g. monitor | grep monitor.node=T3_US_UCLA
- replace listProcessedDatasets with fakeDatasetSummary, since it's better suits dataset queries. DBS3 will provide proper API to look-up dataset out of provided dataset path, release, tier, primary_dataset.

- fix listLFNs to supply file as primary key
- comment out pass_api call to prevent from non-merge situation, must revisit the code
 - * fix issue with missing merge step when das record disapper from cache
- bug fix to prevent from null string in number of events
- increase expire time stamp for dashboard, due to problem described in 1032 ticket. I need to revisit code and make das record/service rather than combined one to utilize cache better. Meanwhile align expire timestamp wrt to DBS/Phedex
- add DBS support to look-up file via provided run (so far using fake API)
- use fakseDataset4Run instead of fakeFile4Run, since it's much faster. Users will be able to find dataset for a given run and then find files for a given dataset
- fix issue with JSON'ifying HTTP error dict
- replace DAS error placement from savannah to TRAC
- add new special keyword, instance, to allow query results from local DBS instances. The keyword itself it neutral and can be applied to any system. Add new abstract method url_instance in abstract_service which can be used by sub-systems to add actual logic how to adjust sub-system URL to specific instance needs.
- replace connection_monitor with dascore_monitor to better handle connection/DASCore absense due to loosing connection to MongoDB
- propagate parser error to user, adjust both DAS cache/web servers
- fix queries with date clause, ticket #1112
- add filter view to show filtered data in plain/text, ticket #959
- add first implementation of tabular representation, ticket #979, based on YUI DataSource table with dynamic JSON/AJAX table feeder
- add jsonstreamer
- add cache method to web server (part of future merge between cache/web servers)
- add das_client which talks to web server; on a web server side made usage of multiprocessing module to handle client requests. Each request spawns a new process.
- visualize record's system by colors on web UI, ticket #977
- add child/parent look-up for dataset/files
- work on DAS PLY/web UI to make errors messages more clear, especially adjust to handle DBS-QL queries
- added dbsql_vs_dasql template which guides how to construct DAS QL expressions for given DBS QL ones
- fix concurrency problem/query race conditions in DAS core
- remove fakeListFile4Site from DBS maps since DBS3 does not cover this use case
- modified das_client to allow other tools use it as API
- fix DBS/phedex maps to match dashes/underscores in SE patterns
- add adjust_params into SiteDB to allow to use patterns in a way SiteDB does it (no asterisks)
- disable expert interface
- update analytics in DAS core when we found a match

- 0.5.9
 - fix issue with <, > operators and numeric values in filters
 - add tier into DBS listProcessedDatasets API as input parameter, so user can query as “dataset primary_dataset=ZJetToEE_Pt* tier=*GEN*”
 - DBS2 API provides typos in their output, e.g. primary_datatset, processed_datatset, add those typos into DAS map to make those attributes queryable.
 - Add lumi into DBS map, as well as its presentation UI keys
- 0.5.8
 - Finish work to make presentation layer more interactive, ticket #880
 - * create hyperlinks for primary DAS keys
 - * round numbers for number of events, etc.
 - * present file/block size in GB notations
 - add new “link” key into presentation to indicate that given key should be used for hyperlinks
 - add reverse look-up from presentation key into DAS key
 - add cache for presentation keys in DAS mapping class
 - update DAS chep paper, it is accepted as CMS Note CR-2010/230
 - fix issue with similar queries, e.g. dataset=/a/b/c is the same as dataset dataset=/a/b/c
 - improve presentation layer and add links
 - * replace link from boolean to a list of record in presentation YAML file
 - * the link key in presentation now refers to list of records, where each record is a dict of name/query. The name is shown on a web UI under the Links:, whiel query represents DAS query to get this value, for example {“name”:“Files”, “query”:“file dataset=%s”}
 - fix issue with counting results in a cache
 - make dataset query look-up close to DD view, fixes #821
 - add YAJL (Yet Another JSON Library) as experimental JSON module, see <http://lloyd.github.com/yajl/> and its python binding.
 - add keylearning and autocompletion, ticket #50
 - add parse_filter, parse_filters functions to parse input list of filters, they used by core/mongocache to yield/count results when filters are passed DAS-QL. This addresses several Oli use cases when multiple filters will be passed to DAS query, e.g. file dataset=/a/b/c | grep file.size>1, file.size<100
 - add special DAS key records, which can be used to look-up records regardless of condition/filter content, e.g. the DAS query site=T1_CH_CERN only shows site records, while other info can be pulled to DAS. So to look-up all records for given condition user can use records site=T1_CH_CERN
 - remove obsolete code from das_parser.py
- 0.5.7
 - Fix dbport/dbhost vs uri bug for das expert interface
 - Created new self-contained unit test framework to test CMS data-services
 - * add new DASTestDataService class which represents DAS test integration web server
 - * provide unit test against DAS test data web service

- * add new configuration for DAS testDataService server
 - * perform queries against local DAS test data service, all queries can be persistent and adjusted in unittest
 - * add fake dbs/phedex/sitedb/ip/zip services into DAS testDataService
 - remove all handlers before initialization of DASLogger
 - add NullHandler
 - add collection parameter to DAS core get_from_cache method
 - add unit test for web.utils
 - add delete_db_collection to mapping/analytics classes
 - remove obsolete templates, e.g. das_admin, mapreduce.
 - sanitize DAS templates, #545
 - Fix issues with showing records while applying DAS filters, #853
 - Move opensearch into das_opensearch.tmpl
 - Fix dbs/presentation maps
 - Add size_format function
 - Updated performance plot
 - make presentation layer more friendly, fixes #848, #879, #880
 - add new configuration parameter status_update, which allow to tune up DAS web server AJAX status update message (in msec)
 - re-factor DAS web server code (eliminate unnecessary AJAX calls; implement new pagination via server calls, rather JS; make form and all view methods to be internal; added check_data method; redesign AJAX status method)
 - Make admin tool be transparent to Ipython
 - Add new functions/unit tests for date conversion, e.g. to_seconds, next_day, prev_day
 - fix date issue with dashboard/runregistry services, fixes #888. Now user will be able to retrieve information for a certain date
- 0.5.6
 - add usable analytics system; this consists of a daemon (analytics_controller) which schedules tasks (which might spawn other tasks), several worker processes which actually perform these tasks and a cherrypy server which provides some information and control of the analytics tasks
 - the initial set of tasks are
 - * Test - prints spam and spawns more copies of itself, as might be expected
 - * QueryRunner - duplicates DAS Robot, issues a fixed query at regular intervals
 - * QueryMaintainer - given a query, looks up expiry times for all associated records and reschedules itself shortly before expiry to force an update
 - * ValueHotspot - identifies the most used values for a given key, and spawns QueryMaintainers to keep them in the cache until the next analysis
 - * KeyHotspot - identifies the most used query keys, and spawns ValueHotspot instances to keep their most popular values maintained in the cache

- provides a cli utility, `das_analytics_task` allowing one-off tasks to be run without starting the analytics server
 - fix apicall records in `analytics_db` so that for a given set of all parameters except expiry, there is only one record
 - fix `genkey` function to properly compare dictionaries with different insert histories but identical content
 - alter `analyticsdb` query records to store an array of call times rather than one record per query, with a configurable history time
 - append “/” to `$base` to avoid `/das?query` patterns
 - Updates for analytics server, add JSON methods, add help section to web page
 - Analytics CLI
 - Add ability to learn data-service output keys, fixes #424
 - Add new class `DASQuery`
 - Add analytics server pid into analytics configuration
 - Prepend python to all shell scripts to avoid permission problem
 - fix `db`s blockpath map
 - add new presentation layouts for various services
 - increase `ajaxStatus` lookup time
 - fix issue with date, in the case when date was specified as a range, e.g. date last 24h, the merge records incorrectly record the date value
- 0.5.5
 - fix map-reduce parsing using `DAS PLY`
 - introduce `das_mapreduces()` function which look-up MR functions in `das.mapreduce` collection
 - fixes for `Tier0,DBS3` services
 - fix core when no services is available, it returns an empty result set
 - fix `DAS` parser cache to properly store MongoDB queries. By default MongoDB does not allow usage of \$ sign in dictionary keys, since it is used in `MongoQL`. To fix the issue we encode the query as dict of key/value/operator and decode it back upon retrieval.
 - fix `DAS PLY` to support value assignment in filters, e.g. `block | grep site=T1`
 - Fixes for `Dashboard`, `RunRegistry` services
 - Eliminate `DAS_PYTHONPATH`, automatically detect `DAS` code location
 - Drop off `ez_setup` in favor python `distutils`, re-wrote `setup.py` to use only `distutils`
 - add `opensearch` plugin
 - fix issue with `DAS PLY` shift/reduce conflict (issue with `COMMA/list_for_filter`)
 - add to `DAS PLY` special keys, `date` and `system`, to allow queries like `run date last 24h`, `jobsummary date last 24h`. Prevent queries like `run last 24h` since it leads to ambiguous conditions.
 - add support for `GridFS`; `parse2gridfs` generator pass docs whose size less then MongoDB limit (4MB) or store doc into `GridFS`. In later case the doc in `DAS` workflow is replaced with `gridfs` pointer (issue #611)
 - add new method to `DAS` cache server to get data from `GridFS` for provided file id
 - fix `DAS` son manipulator to support `gridfs_id`

- fix das_config to explicitly use DAS_CONFIG environment
- fix bug with expire timestamp update from analytics
- add support for “test” and “clean” action in setup.py; remove das_test in favor standard python setup.py test
- add weighted producer into querryspammer toolkit; this allows to mimic real time behavior of most popular queries and ability to invoke DAS robots for them (up-coming)
- fix #52, now both min and max das aggregators return _id of the record
- return None as db instances when MongoDB is down
- add avg/median functions to result object; modified result object to hold result and rec counter; add helper das function to associate with each aggregators, e.g. das_min
- drop dbhost/dbport in favor of dburi, which can be a list of MongoDB uris (to be used for connection with MongoDB replica sets)
- replace host/port to URI for MongoDB specs, this will allow to specify replication sets in DAS config
- use bson.son import SON to be compatible with newer version of pymongo
- use col.count() vs col.find().count(), since former is O(1) operation wrt O(N) in later case
- 0.5.3 - 0.5.4 series
 - Clean-up %post and do not package docs over there
 - All names in bin are adjusted to one schema: das_<task>.
 - All scripts in bin are changed to use /bin/sh or /bin/bash and use \${1+}\$@” instead of “\$@”
 - bin area has been clean-up, e.g. das_doc, dassh is removed, etc.
 - Remove runsum_keys in runsum_service.py since it is obsolete code
 - Fix issue w/ root.close() for runsum_service.py (parser function)
 - Remove session from plotfairy
 - Remove encode4admin
 - Add urllib.quote(param) for das_services.tmpl and das_tables.tmpl
 - fix #446
 - das_jsonstable.tmpl is removed since it's obsolete and no one is using it.
 - Remove das_help.tmpl and /das/help since it is obsolete
 - Remove das_admin.py since it is obsolete
 - Reviewed decorator in web/tools.py and commented out unused decorators, exposexml, exposeplist. I want to keep them around upon they become relevant for DAS long terms.
 - Fix issue with wrap2das methods and made them internal.
 - Add checkargs decorator to validate input parameters for das_web
 - Change socket_queue_size to 100
 - Set engine.autoreload_on=False, request.show_tracebacks=False. Verified that server runs in production mode by default.
 - Add parameters validation for das_web/das_expert.
 - fix #493, allow relocation of PLY parsertab.py

- fix #494, allow usage of HTTP Expires if data-services provide that
 - change eval(x) into eval(x, { "__builtins__": None }, {}) for those cases when fail to use json.load(x). Some data-service are not fully compliant and the issue with them need to be resolved at their end.
 - Use singleton class for Connection to reduce number of ESTABLISHED connections seeing on server. For details see http://groups.google.com/group/mongodb-user/browse_thread/thread/67d77a62059568d7# <https://svnweb.cern.ch/trac/CMSDMWM/ticket/529>
 - use isinstance instead of types.typeXXX
 - make generic cern_sso_auth.py to authenticate with CERN SSO system
 - make das_map to accept external map dir parameter which specify locations of DAS maps
 - fix queryspammer to handle generators; add weights
 - unify DAS configuration via das_option
 - Remove das docs from RPM, will run it stand-alone elsewhere
 - Move checkargs into DAS.web.utils; reuse this decorator for all DAS servers to sanitize input arguments; added new unit test for it
 - Introduce DAS server codes, they resides in DAS.web.das_codes
 - Change DAS server behavior to return HTTPError. The passed message contains DAS server error code.
 - fix #525, #542.
 - fix issue with counting of empty records, #455
 - Handle the case when MongoDB is down. Both DAS servers can handle now outage of MongoDB either at start-up or during their operations. Adjust code to use a single mongodb host/port across all databases, fix #566
 - Remove from all unit test hardcoded value for mongodb host/port, instead use those from DAS configuration file
 - Use calendar.timegm instead of time.mktime to correctly convert timestamp into sec since epoch; protect expire timestamp overwrite if expires timestamp is less then local time
 - Add empty_record=0 into DAS records, to allow consistent look-up
 - Added DAS_PYTHONROOT, DAS_TMPLROOT, DAS_IMAGESROOT, DAS_CSSROOT, DAS_JSROOT to allow DAS code relocation
- 0.5.0 till 0.5.2
 - based on Gordon series of patches the following changes has been implemented
 - * new analytics package, which keeps track of all input queries
 - * new DAS PLY parser/lexer to confirm DAS QL
 - * added new queryspammer tool
 - added spammer into DAS cache client, to perform benchmarking of DAS cache server
 - added a few method to DAS cache server for performance measurements of bare CherryPy, CherryPy+MongoDB, CherryPy+MongoDB+DAS
 - remove white/back list in favor of explicit configuration of DAS services via DAS configuration systems (both das.cfg and das_cms.py)
 - added index on das.expire
 - fixed issue with SON manipulator (conversion to str for das_id, cache_id)

- enable checks for DAS key value patterns
- added URN's to query record
- added empty records into DAS merge to prevent cases when no results aggregated for user request
 - * empty records are filtered by web interface
 - * values for empty records are adjusted to avoid presence of special \$ key, e.g. we cannot store to MongoDB records with { '\$in': [1,2]}
- new das_bench tool
- fixed regex expression for DAS QL pattern, see http://groups.google.com/group/mongodb-user/browse_thread/thread/8507223a70de7d51
- various speed-up enhancements (missing indexes, empty records, regex bug, etc.)
- added new RunRegistry CMS data-service
- updated DAS documentation (proof-reading, DAS QL section, etc.)
- remove src/python/ply to avoid overlap with system default ply and added src/python/parser to keep parsertab.py around

Release 0.4.X series

The most significant part of this release is new plug-and-play mechanism to add new data-services. This is done via data-service map creation. Each map is represented data-service URI (URL, input parameters, API, etc.).

- 0.4.13 till 0.4.18
 - adjustment to CMS environment and SLA requirements
 - ability to read both cfg and CMS python configuration files
 - replacement of Admin to Expert interface and new authentication scheme via DN (user certificates) passed by front-end
 - new mongodb admin.dns collection
 - add PID to cherrypy das_server configuration
- 0.4.12
 - added unique filter
 - change value of verbose/debug options in all cli tools to be 0, instead of None, since it's type suppose to be int
 - add new example section to web FAQ
 - re-define logger/logformat in debug mode; the logger is used StreamHandler in this mode, while logformat doesn't use time stamp. This is usefull for DAS CLI mode, when -verbose=1 flag is used.
 - add "word1 word2" pattern to t_WORD for das_lexer, it's going to be used by searching keywords in cmsswconfig service and can be potentially used elsewhere to support multiple keywords per single DAS key
 - fix bug with apicall which should preceed update_cache
 - add simple enc/dec schema for DAS admin authentication
 - add logger configuration into das.cfg
 - separate logger streams into das.log, das_web.log and das_cache.log

- das_lexer supports floats
- Add ability for filter to select specific values, e.g. run=123 | grep PD=MinBias right now only equal condition is working, in future may extend into support of other operators
- add CMSSW release indexer
- 0.4.11
 - adjust abstract data-service and mongocache to use DAS compliant header if it is supplied by DAS compliant API, e.g. Tier0.
 - added cmsswconfigs data-service
 - work on xml_parser to make it recursive. Now it can handle nested children.
 - Fix problem with multiple look-up keys/API, by using api:lookup_keys dict. This had impact on storage of this information within das part of the record. Adjust code to handle it properly
 - added map for Tier0 monitoring data-service
 - fix problem with id references for web interface
 - fix problem with None passed into spec during parsing step
- 0.4.10
 - added new mapping for Phedex APIs
 - work on aggregator to allow merged records to have reference to their parent records in DAS cache, name them as cache_id
 - improve DAS admin interface:
 - * show and hide various tasks
 - * DAS tasks (query db, clean db, das queries)
 - * Add digest authentication to admin interface, based on `cherrypy.tools.digest_auth`
 - allow to use multiple aggregators at the same time, e.g. site=T1_* | count(site.id), sum(site.id), avg(site.id)
 - enable aggregators in DAS core
 - migrated from CVS to SVN/GIT
 - added AJAX interface for DAS query look-up in admin interface
 - bug fix in core to get status of similar queries
 - validate web pages against XHTML 1.0, using <http://validator.w3.org/check>
- V0.4.9
 - update admin interface (added query info)
 - integrate DAS lexer in to DAS parser
 - add new class DASLexer, which is based on [PLY]
 - remove >, <, >=, <= operators from a list of supported ones, since they don't make sense when we map input DAS query into underlying APIs. The API usually only support = and range operators. Those operators are supported by MongoDB back-end, but we need more information how to support them via DAS <-> API callback
 - work on DAS parser to improve error catching of unsupported keywords and operators
 - split apart query insertion into DAS cache from record insertion to ensure that every query is inserted. The separation is required since record insertion is a generator which may not run if result set is empty

- synchronized expire timestamp in DAS cache/merge/analytics db's
- V0.4.8
 - fix pagination
 - display DAS key for all records on the web to avoid overlap w/ records coming out from multiple data-providers (better visibility)
 - protect DASCACHEMgr with queue_limit configurable via das.cfg
 - found that multiprocessing is unreliable (crash on MacOSX w/ python version from macports); some processes become zombies. Therefore switch to ThreadPool for DAS cache POST requests
 - added ThreadPool
 - work on DBS2 maps
 - make monitoring_worker function instead of have it inside of DASCACHEMgr
 - re-factor DASCACHEMgr, now it only contains a queue
 - switch to use <major>.<minor>.<release> notations for DAS version
 - switch to use dot notation in versions, the setup.py/ez_tools.py substitute underscore with dash while making a tar ball
- V04_00_07
 - re-factor DAS configuration system
 - switch to pymongo 1.5.2
 - switch to MongoDB 1.4
 - added admin web interface; it shows db info, DAS config, individual databases and provide ability to look-up records in any collection
- V04_00_06
 - added support for proximity results
 - resolve issue with single das keyword provided in an input query
 - dynamically load of DAS plugins using __import__ instead of eval(klass)
 - first appearance of analytics code
 - fix issue with data object look-up
 - switch to new DAS QL parser
- V04_00_05
 - re-wrote DAS QL parser
 - move to stand-alone web server (remove WebTools dependency)
 - adjust web UI
- V04_00_04
 - choose to use flat-namespace for DAS QL keys in DAS queries
 - added aggregator functions, such as sum/count, etc. as coroutines
 - added “grep” filter for DAS QL
 - extended dict class with _set/_get methods

- re-wrote C-extension for dict_helper
 - added wild_card parameter into maps to handle data-service with specific wild_card characters, e.g. *, %, etc.
 - added ability to handle data-service HTTPErrors. The error records are recorded into both DAS cache and DAS merge collection. They will be propagated to DAS web server where admin view can be created to view them
- V04_00_02, V04_00_03
 - bug fix releases
- V04_00_01
 - minor tweaks to make CMS rpms
 - modifications for init scripts to be able to run in stand-alone mode
- V04_00_00 - incorporate all necessary changes for plug-and-play - modifications for stand-alone mode

Release V03 series

Major change in this release was a separation of DAS cache into independent cache and merge DB collection. The das.cache collection stores *raw* API results, while das.merge keeps only merged records.

- V03_00_04
 - minor changes to documentation structure
- V03_00_03
 - added DAS doc server
 - added sphinx support as primary DAS documentation system
- V03_00_02
 - work on DAS cli tools
- V03_00_01
 - bug fixes
- V03_00_00
 - separate DAS cache into das.cache and das.merge collections

Release V02 series

This release series is based on MongoDB. After a long evaluation of different technologies, we made a choice in favor of MongoDB.

- added support for map/reduce
- switch to pipes syntax in QL for aggregation function support
- switch DAS QL to free keyword based syntax

Release V01 series

Evaluation series. During this release cycle we played with the following technologies:

- Memcached
- CouchDB
- custom file-based cache

At that time DAS QL was based on DBS-QL syntax. During this release series we added DAS cache/web servers; made CLI interface.

1.15 References

Indices and tables

- `genindex`
- `modindex`
- `search`

- [CMS] <http://cms.web.cern.ch/cms/>
- [LHC] <http://en.wikipedia.org/wiki/Lhc>
- [DAS] <https://github.com/dmwm/das/>
- [Python] <http://www.python.org/>, version 2.6
- [CPF] <http://www.cherrypy.org/>, version 3.1.2
- [YUI] <http://developer.yahoo.com/yui/>, version 2
- [Gen] <http://www.dabeaz.com/generators/>
- [YAML] <http://en.wikipedia.org/wiki/Yaml>, <http://pyyaml.org/wiki/PyYAMLDocumentation>, version PyYAML-3.08
- [PYLINT] <http://pypi.python.org/pypi/pylint>
- [MongoDB] <http://www.mongodb.org>, version 1.4.1
- [MongoDBOverview] http://www.paperplanes.de/2010/2/25/notes_on_mongodb.html
- [Pymongo] <http://api.mongodb.org/python/>, version 1.6
- [Cheetah] <http://www.cheetahtemplate.org/>, version 2.4.0
- [Cherrypy] <http://www.cherrypy.org/>, version 3.1.2
- [Sphinx] <http://sphinx.pocoo.org/>, version 0.6.4
- [IPython] <http://ipython.scipy.org/>, version 0.10
- [Memcached] <http://memcached.org>
- [Couchdb] <http://couchdb.apache.org>
- [JSON] <http://en.wikipedia.org/wiki/JSON>
- [CJSON] <http://pypi.python.org/pypi/python-cjson>
- [PLY] <http://www.dabeaz.com/ply/>
- [Gevent] <http://www.gevent.org/>
- [YAJL] <http://lloyd.github.com/yajl/>
- [CURL] <http://curl.haxx.se/download.html>
- [PyCurl] <http://pycurl.sourceforge.net/>

[PyYAML] <http://pyyaml.org/>

[VENV] <https://raw.githubusercontent.com/pypa/virtualenv/master/virtualenv.py>

[GIT] <http://git-scm.com/>

d

DAS.keywordsearch, [43](#)
DAS.keywordsearch.entity_matchers, [47](#)
DAS.keywordsearch.entity_matchers.kwd_chunks,
 [48](#)
DAS.keywordsearch.entity_matchers.string_dist_levenstein,
 [48](#)
DAS.keywordsearch.metadata, [44](#)
DAS.keywordsearch.metadata.das_output_fields_adapter,
 [45](#)
DAS.keywordsearch.metadata.das_ql, [45](#)
DAS.keywordsearch.metadata.schema_adapter2,
 [44](#)
DAS.keywordsearch.presentation, [48](#)
DAS.keywordsearch.presentation.result_presentation,
 [48](#)
DAS.keywordsearch.rankers, [48](#)
DAS.keywordsearch.tokenizer, [45](#)
DAS.tools.create_das_config, [49](#)
DAS.tools.das_bench, [49](#)
DAS.tools.das_client, [50](#)
DAS.utils.cern_sso_auth, [51](#)
DAS.utils.das_config, [52](#)
DAS.utils.das_option, [52](#)
DAS.utils.jsonwrapper, [53](#)
DAS.utils.logger, [54](#)
DAS.utils.query_utils, [55](#)
DAS.utils.regex, [56](#)
DAS.web.das_codes, [42](#)

A

ApiDef (in module DAS.keywordsearch.metadata.schema_adapter2), 44

are_wildcards_allowed() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter class method), 44

avg_std() (in module DAS.tools.das_bench), 49

C

check_auth() (in module DAS.tools.das_client), 50

check_glidein() (in module DAS.tools.das_client), 50

check_result_field_match()
(DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 44

cleanup_query() (in module DAS.keywordsearch.tokenizer), 46

compare_dicts() (in module DAS.utils.query_utils), 55

compare_specs() (in module DAS.utils.query_utils), 55

compare_str() (in module DAS.utils.query_utils), 55

ConfigOptionParser (class in DAS.tools.create_das_config), 49

convert2pattern() (in module DAS.utils.query_utils), 55

convert_time() (in module DAS.tools.das_client), 50

D

DAS.keywordsearch (module), 43

DAS.keywordsearch.entity_matchers (module), 47

DAS.keywordsearch.entity_matchers.kwd_chunks (module), 48

DAS.keywordsearch.entity_matchers.string_dist_levenstein (module), 48

DAS.keywordsearch.metadata (module), 44

DAS.keywordsearch.metadata.das_output_fields_adapter (module), 45

DAS.keywordsearch.metadata.das_ql (module), 45

DAS.keywordsearch.metadata.schema_adapter2 (module), 44

DAS.keywordsearch.presentation (module), 48

DAS.keywordsearch.presentation.result_presentation (module), 48

DAS.keywordsearch.rankers (module), 48

DAS.keywordsearch.tokenizer (module), 45

DAS.tools.create_das_config (module), 49

DAS.tools.das_bench (module), 49

DAS.tools.das_client (module), 50

DAS.utils.cern_sso_auth (module), 51

DAS.utils.das_config (module), 52

DAS.utils.das_option (module), 52

DAS.utils.jsonwrapper (module), 53

DAS.utils.logger (module), 54

DAS.utils.query_utils (module), 55

DAS.utils.regex (module), 56

DAS.web.das_codes (module), 42

das_configfile() (in module DAS.utils.das_config), 52

das_readconfig() (in module DAS.utils.das_config), 52

das_readconfig_helper() (in module DAS.utils.das_config), 52

DASOption (class in DAS.utils.das_option), 52

DASOptionParser (class in DAS.tools.das_client), 50

dasql_to_nl() (in module DAS.keywordsearch.presentation.result_presentation), 48

DasSchemaAdapter (class in DAS.keywordsearch.metadata.schema_adapter2), 44

debug() (DAS.utils.logger.PrintManager method), 54

decode() (DAS.utils.jsonwrapper.JSONDecoder method), 53

decode_code() (in module DAS.web.das_codes), 42

decode_mongo_query() (in module DAS.utils.query_utils), 56

dump() (in module DAS.utils.jsonwrapper), 54

dumps() (in module DAS.utils.jsonwrapper), 54

E

emit() (DAS.utils.logger.NullHandler method), 54

encode() (DAS.utils.jsonwrapper.JSONEncoder method), 53

encode_mongo_query() (in module DAS.utils.query_utils), 56

entities_for_inputs() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 44

error() (DAS.utils.logger.PrintManager method), 54
extract_value() (in module DAS.tools.das_client), 50

F

fescape() (in module DAS.keywordsearch.presentation.result_representation), 48
flatten() (in module DAS.keywordsearch.metadata.das_output_fields_adapter), 45
flatten() (in module DAS.keywordsearch.metadata.das_q1), 45

fullpath() (in module DAS.tools.das_client), 51
funcname() (in module DAS.utils.logger), 54

G

gen_passwd() (in module DAS.tools.das_bench), 49
get_api_param_definitions()
(DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 44

get_connection() (DAS.tools.das_client.HTTPSClientAuthHandler method), 50

get_connection() (DAS.utils.cern_sso_auth.HTTPSClientAuthHandler method), 51

get_data() (in module DAS.tools.das_client), 51
get_data() (in module DAS.utils.cern_sso_auth), 51

get_from_configparser() (DAS.utils.das_option.DASOption method), 52

get_from_wmcore() (DAS.utils.das_option.DASOption method), 52

get_keyword_without_operator() (in module DAS.keywordsearch.tokenizer), 46

get_method() (DAS.tools.das_bench.UrlRequest method), 49

get_operator_and_param() (in module DAS.keywordsearch.tokenizer), 46

get_operator_synonyms() (in module DAS.keywordsearch.metadata.das_q1), 45

get_opt() (DAS.tools.create_das_config.ConfigOptionParser method), 49

get_opt() (DAS.tools.das_bench.NClientsOptionParser method), 49

get_opt() (DAS.tools.das_client.DASOptionParser method), 50

get_outputs_field_list() (in module DAS.keywordsearch.metadata.das_output_fields_adapter), 45

get_result_field_title() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 44

get_titles_by_field() (in module DAS.keywordsearch.metadata.das_output_fields_adapter), 45

get_value() (in module DAS.tools.das_client), 51

H

handle() (DAS.utils.logger.NullHandler method), 54

https_open() (DAS.tools.das_client.HTTPSClientAuthHandler method), 50

https_open() (DAS.utils.cern_sso_auth.HTTPSClientAuthHandler method), 51

HTTPSClientAuthHandler (class in DAS.tools.das_client), 50

HTTPSClientAuthHandler (class in DAS.utils.cern_sso_auth), 51

info() (DAS.utils.logger.PrintManager method), 54

init() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 45

is_reserved_field() (in module DAS.keywordsearch.metadata.das_output_fields_adapter), 45

JSONEncoder (class in DAS.utils.jsonwrapper.JSONEncoder method), 53

J

JSONEncoder (class in DAS.utils.jsonwrapper), 53
JSONEncoder (class in DAS.utils.jsonwrapper), 53

K

keys_attrs() (in module DAS.tools.das_client), 51

L

levenshtein() (in module DAS.keywordsearch.entity_matchers.string_dist_levenstein), 48

levenshtein_normalized() (in module DAS.keywordsearch.entity_matchers.string_dist_levenstein), 48

list_result_fields() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter method), 45

load() (in module DAS.utils.jsonwrapper), 54

loads() (in module DAS.utils.jsonwrapper), 54

M

main() (in module DAS.tools.create_das_config), 49

main() (in module DAS.tools.das_bench), 49

main() (in module DAS.tools.das_client), 51

make_plot() (in module DAS.tools.das_bench), 49

N

natcasecmp() (in module DAS.tools.das_bench), 49

natcmp() (in module DAS.tools.das_bench), 50

natsort() (in module DAS.tools.das_bench), 50

natsort_key() (in module DAS.tools.das_bench), 50

natsorted() (in module DAS.tools.das_bench), 50

NClientsOptionParser (class in DAS.tools.das_bench), 49

NullHandler (class in DAS.utils.logger), 54

P

prim_value() (in module DAS.tools.das_client), 51

print_debug() (in module
DAS.keywordsearch.metadata.das_output_fields_adapter), 45

print_from_cache() (in module DAS.tools.das_client), 51

print_msg() (in module DAS.utils.logger), 54

print_summary() (in module DAS.tools.das_client), 51

PrintManager (class in DAS.utils.logger), 54

R

random_index() (in module DAS.tools.das_bench), 50

raw_decode() (DAS.utils.jsonwrapper.JSONDecoder
method), 53

read_configparser() (in module DAS.utils.das_config), 52

read_wmcore() (in module DAS.utils.das_config), 52

result_contained_errors() (in module
DAS.keywordsearch.metadata.das_output_fields_adapter), 45

result_to_dasql() (in module
DAS.keywordsearch.presentation.result_presentation), 48

runjob() (in module DAS.tools.das_bench), 50

S

set_cherrypy_logger() (in module DAS.utils.logger), 54

shorten_value() (in module
DAS.keywordsearch.presentation.result_presentation), 48

size_format() (in module DAS.tools.das_client), 51

spammer() (in module DAS.tools.das_bench), 50

T

test_operator_containment() (in module
DAS.keywordsearch.tokenizer), 47

timestamp() (in module DAS.utils.cern_sso_auth), 52

tokenize() (in module DAS.keywordsearch.tokenizer), 47

try_int() (in module DAS.tools.das_bench), 50

U

unique_filter() (in module DAS.tools.das_client), 51

UrlRequest (class in DAS.tools.das_bench), 49

urlrequest() (in module DAS.tools.das_bench), 50

V

validate_input_params() (DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter
method), 45

validate_input_params_lookupbased()
(DAS.keywordsearch.metadata.schema_adapter2.DasSchemaAdapter
method), 45

W

warning() (DAS.utils.logger.PrintManager method), 54

web_code() (in module DAS.web.das_codes), 42

wmcore_config() (in module DAS.utils.das_config), 52

word_chars() (in module DAS.utils.regex), 56

write_configparser() (in module DAS.utils.das_config),
adapter), 52

write_to_configparser() (DAS.utils.das_option.DASOption
method), 53

write_to_wmcore() (DAS.utils.das_option.DASOption
method), 53

X

x509() (in module DAS.tools.das_client), 51