

BLAHP: A local batch system abstraction layer for global use

Primary Authors:

REBATTO, Davide¹ (david.rebato@mi.infn.it) - PRELZ, Francesco¹ - FIORENTINO, Giuseppe¹ - MEZZADRI, Massimo¹ - MARTELLI, Enzo¹ - MOLINARI, Elisabetta¹

Co-Authors:

DVORAK, F.² - KOURIL, D.² - KRENEK, A.² - MATYSKA, L.² - MULAC, M.² - POSPISIL, J.² - RUDA, M.² - SALVET, Z.² - SITERA, J.² - SKRABAL, J.² - VOCUS, M.² - AVELLINO, G.³ - BECO, S.³ - CAVALLINI, A.³ - MARASCHINI, A.³ - PACINI, F.³ - PARRINI, A.³ - SCARCELLA, C.³ - SOTTILARO, M.³ - TERRACINA, A.³ - MONFORTE, S.⁴ - PAPPALARDO, M.⁴ - ANDREZZI, S.⁵ - CHECCHI, M.⁵ - CIASCHINI, V.⁵ - FERRARI, T.⁵ - GIACOMINI, F.⁵ - LOPS, R.⁵ - RONCHIERI, E.⁵ - VENTURI, V.⁵ - GUARISE, A.⁷ - PATANIA, G.⁷ - PIRO, R.⁷ - WERBROUCK, A.⁷ - ANDREETTO, P.⁶ - BORGIA, A.⁶ - DORIGO, A.⁶ - GIANELLE, A.⁶ - MARZOLLA, M.⁶ - MORDACCHINI, M.⁶ - SGARAVATTO, M.⁶ - ZANGRANDO, L.⁶

¹INFN Milano – ²CESNET – ³Datamat – ⁴INFN Catania – ⁵INFN CNAF – ⁶INFN Padova – ⁷INFN Torino

Overview

In current, widely deployed management schemes, intensive computing farms are locally managed by **batch systems** (e.g. *Platform LSF*, *PBS/Torque*, *BQS*, etc.). When approached from the outside, at the global – or *grid* – level, these local resource managers (LRMS) are seen as services providing at least a basic set of job operations, namely **submission**, **status retrieval**, **cancellation** and **security credential renewal**. The Batch-system Local ASCII Helper Protocol (**BLAHP**) was designed to offer a simple abstraction layer over the different *LRMS*, providing uniform access to the underlying computing resources. In order to preserve the simplicity and portability of the scheme and the robustness of the implementation, the functionality in the abstraction had to be carefully limited. The daemon, originally developed for the **EGEE gLite Condor¹-based Computing Element**, is going to be used by Condor also outside the gLite framework. It is also a component of **CREAM²**, the Web Services oriented Computing Element for gLite.

Batch system abstraction

Submission must be considered successful only if the job is actually accepted and enqueued by the LRMS. Therefore a double-check through the LRMS log file is performed.

The LRMS server can be easily overloaded if flooded by status requests. To avoid direct queries to the server, a job state machine, fed by the log file, has been implemented.

No assumption are made on job cancellation. BLAHPD relies on LRMS capability for this operation, and trusts the LRMS cancellation command's result.

LRMSes offer no way of dispatching files to the workdir of a running job, so a small daemon is sent along with the job and started by the job wrapper. It accepts GSI secured connections, used to send the fresh proxy file.

ABSTRACT RESOURCE MANAGER SERVICE

JOB SUBMISSION

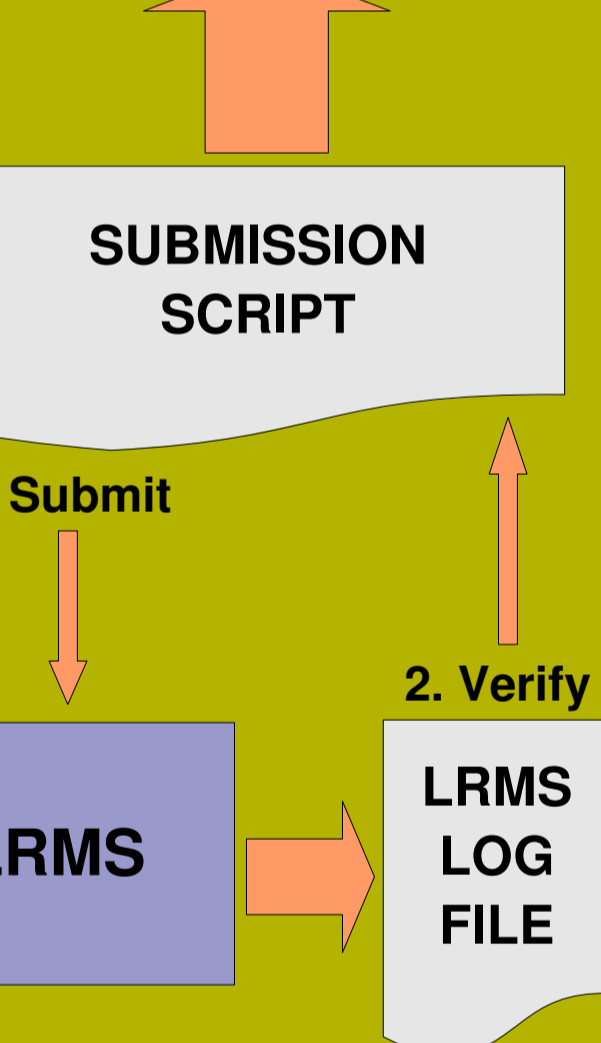
JOB STATUS

JOB CANCELLATION

CREDENTIAL RENEWAL

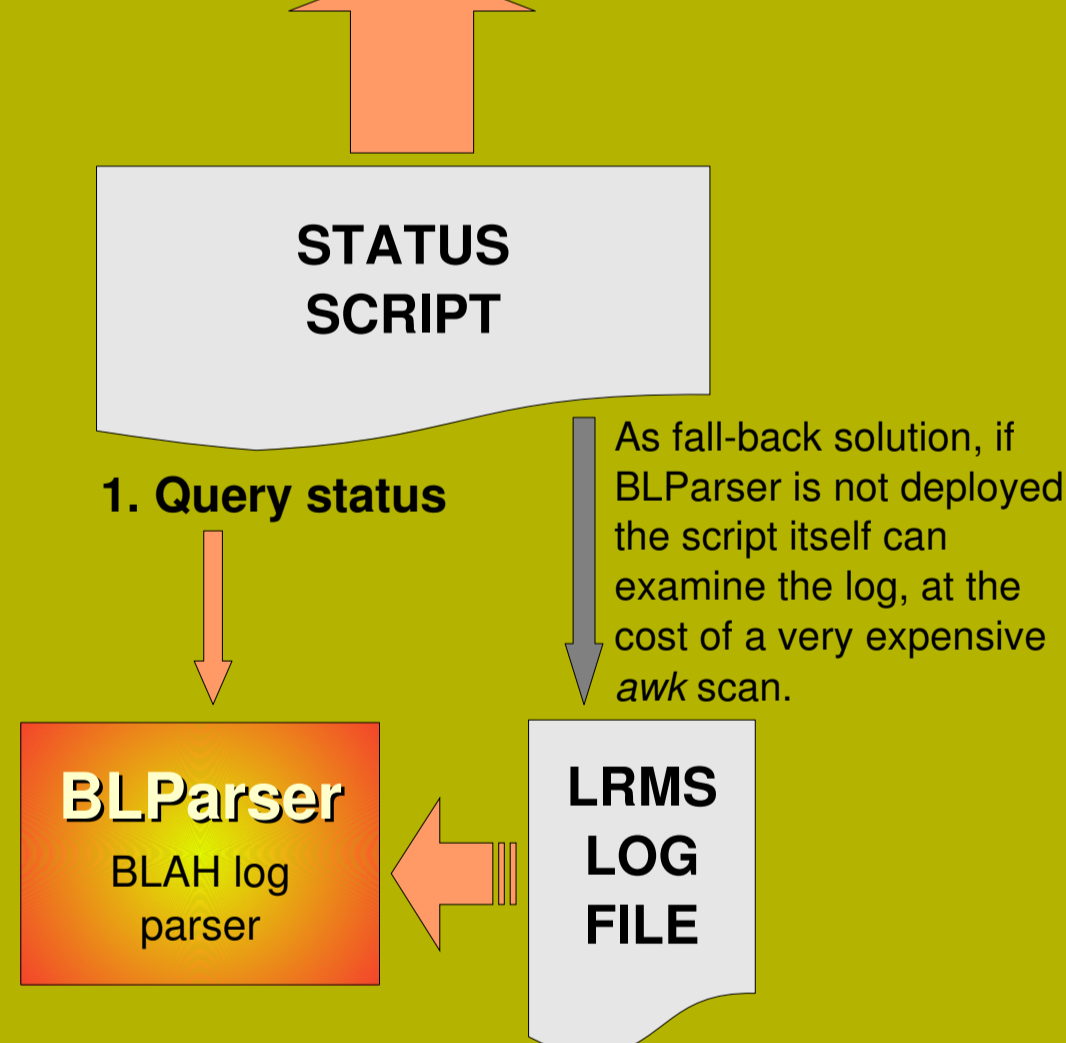
3. Return the job ID

SUBMISSION SCRIPT



2. Return the job status

STATUS SCRIPT



BLAHPD

BLParser

BLAH log parser

1. Query job status and location

Job is running

Yes

BPRClient

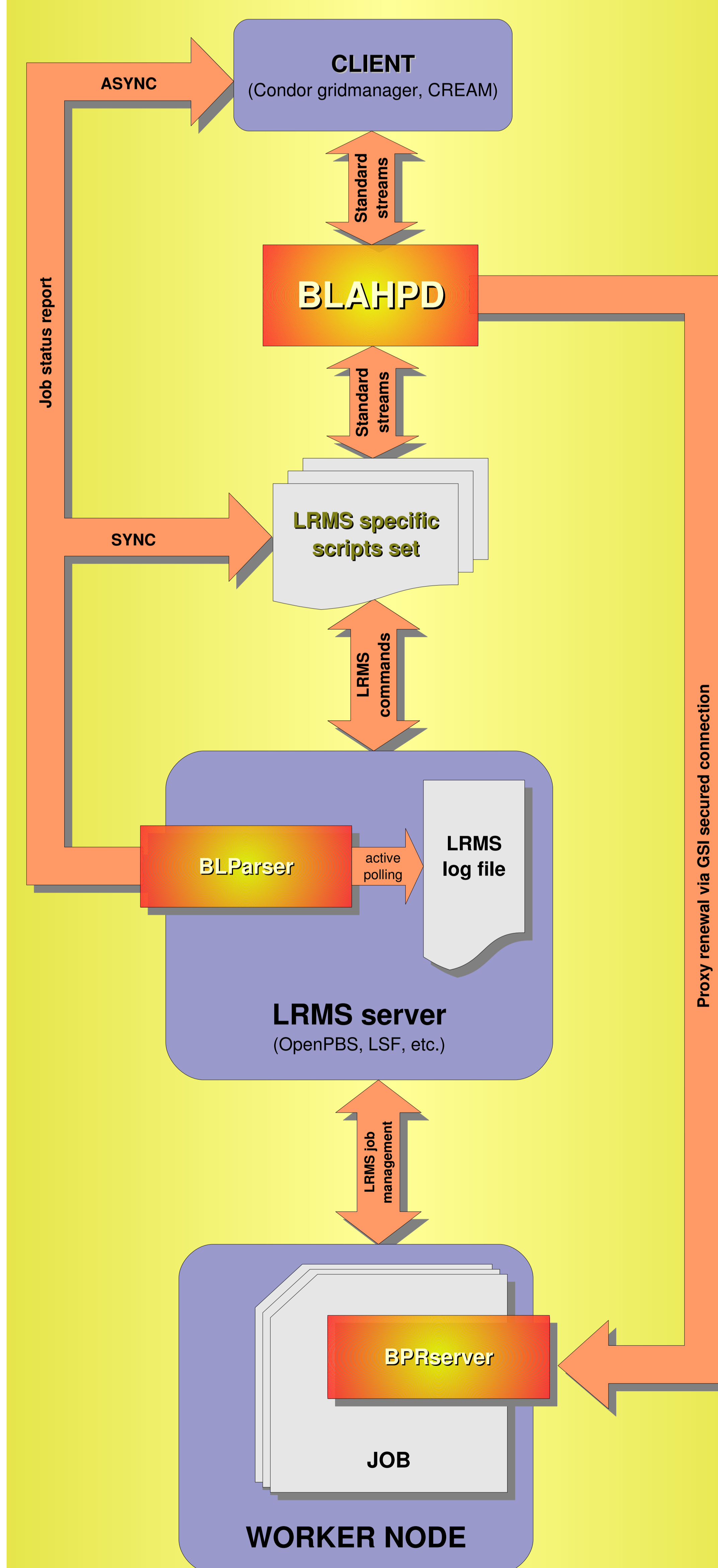
2. Send proxy over GSI channel

BPRServer

No

No operation needed, proxy renewed in place

BLAHP Architecture



BLAHP commands³

COMMANDS

Provide the list of supported commands

VERSION

Return BLAHPD version

ASYNC_MODE_ON

Enable async notification of unread results

ASYNC_MODE_OFF

Disable async notification of unread results

RESULTS

Read and empty the result buffer

BLAH_JOB_SUBMIT

Submit a job

BLAH_JOB_STATUS

Retrieve jobs' status

BLAH_JOB_CANCEL

Cancel a job

BLAH_JOB_HOLD

Suspend a job (if supported by LRMS)

BLAH_JOB_RESUME

Resume a suspended job

BLAH_JOB_REFRESH_PROXY

Refresh job credentials

BLAH_SET_GLEXEC_DN

Use user's credential to issue LRMS commands

BLAH_SET_GLEXEC_OFF

Stop using user's credential

BLAHP implementation details

- The **blahpd daemon** runs on the submission machine of the batch system. It translates BLAHP commands into abstract LRMS actions. Most job management commands are served in dedicated threads: upon completion the results are queued in a (thread-safe) buffer and can be retrieved by the client later on.
- The LRMS specific **scripts set** is where the abstraction effort is concentrated. The scripts for different LRMS must have the same semantic and syntax, offering a uniform interface toward the core daemon. This architecture allows to make BLAHPD support new batch systems without going through the daemon code.
- The **BLParser** (BLAH Log Parser) implements a state machine, updated by a constant watching of the LRMS log file. It accepts synchronous queries via a socket connection, and can asynchronously notify job status changes to subscribed clients (e.g. CREAM).
- The **BPRServer** (BLAH Proxy Renewal Server) is a small daemon, sent along with the job to the worker node. The server is started by the job wrapper when the job start running, and when an incoming connection is detected, it answers with the jobID. The connecting client, if the jobID is correct, open a GSI secured channel and sends the fresh proxy, which replaces the old one.

¹ Condor Project: <http://www.cs.wisc.edu/condor/>

² Cream web page: <http://grid.pd.infn.it/cream/field.php>

³ For a detailed description of BLAH protocol and syntax, see http://egee-jra1-wm.mi.infn.it/egee-jra1-wm/ce_blahp.shtml