

## Distributed Crawler Application

### *General architecture specification, phase A*

The DC application is multithreaded server class software. Inter process interactions uses ZMQ TCP sockets connections server and client type. In process interactions uses ZMQ inproc sockets connections servers and clients. Asynchronous inproc resources access as well as requests processing provided by MOM-based transport provided by ZMQ sockets technology. No system or POSIX access controlling objects like mutexes or semaphores used. The only one kind of processing item is message. The several messages containers types like queues, dictionaries, lists and so on can be used to accumulate and organize processing sequences. The key-value DB engine is SQLite.

There are two listened by DTM application TCP ports:

1. Administration.
2. Client.

**Administration** – handles requests from tool client applications. Requests of that type fetches statistical information from all threading objects as well as can be used to tune some configurations settings of DTM application directly at runtime period.

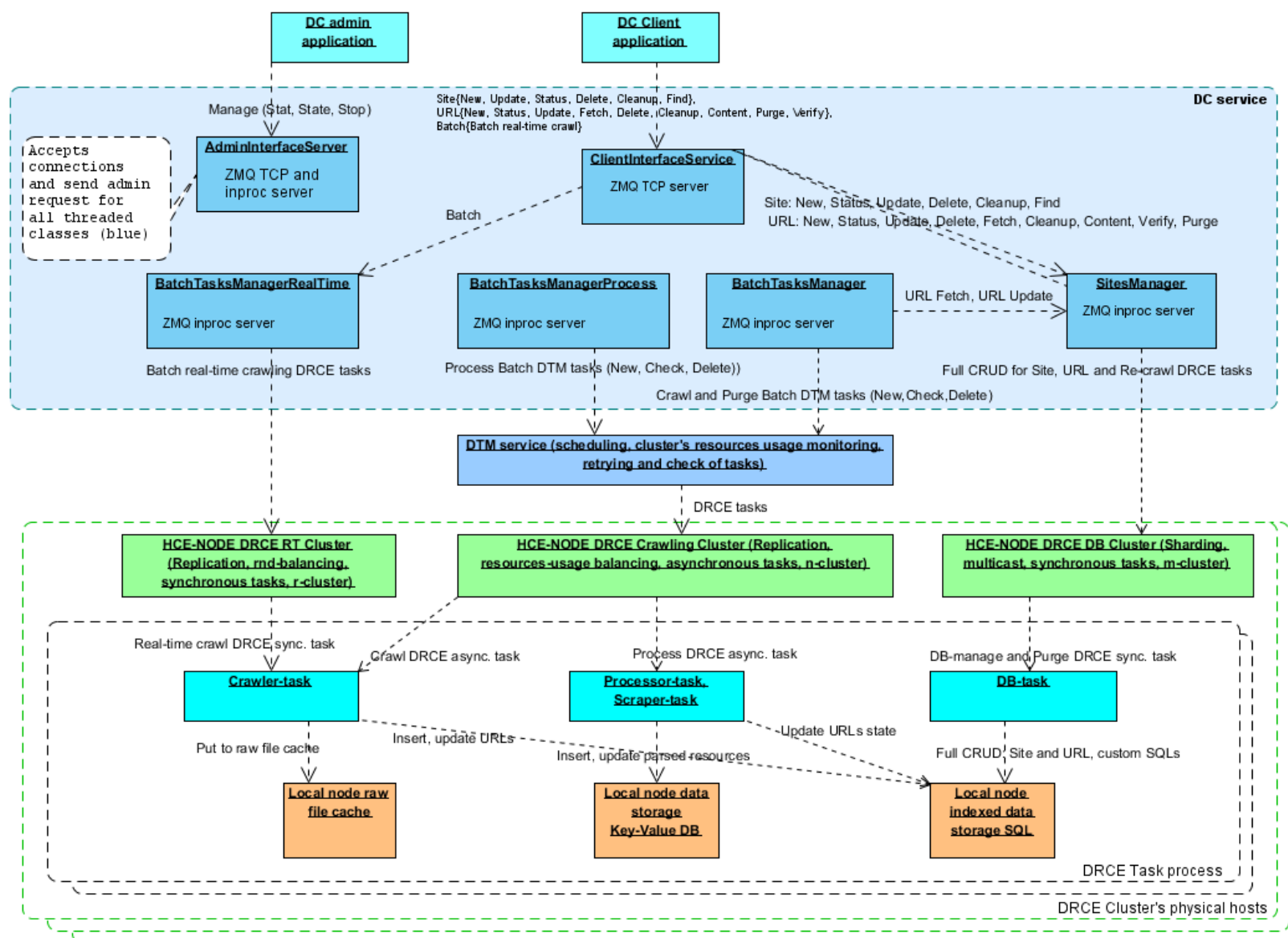
**Client** – handles requests from client side API or tool application. Requests of that type are main functional queries for tasks and related with HCE DRCE interaction protocol.

All requests are represented by json format of message. The structure of requests and main tasks logic principles are separated from concrete execution environment. This gives potential possibility to use some several execution environments engines, but not only HCE DRCE cluster.

### Object model

All active networking usage objects (blue primitives on architecture UML diagram) are based on threaded class object. Also, all threaded objects uses list of transport level items of two types: client and server. Both are based on ZMQ sockets, server type can be TCP or inproc. TCP used for external connections, inproc – for all internal inter-threading events interactions

## Application architecture model

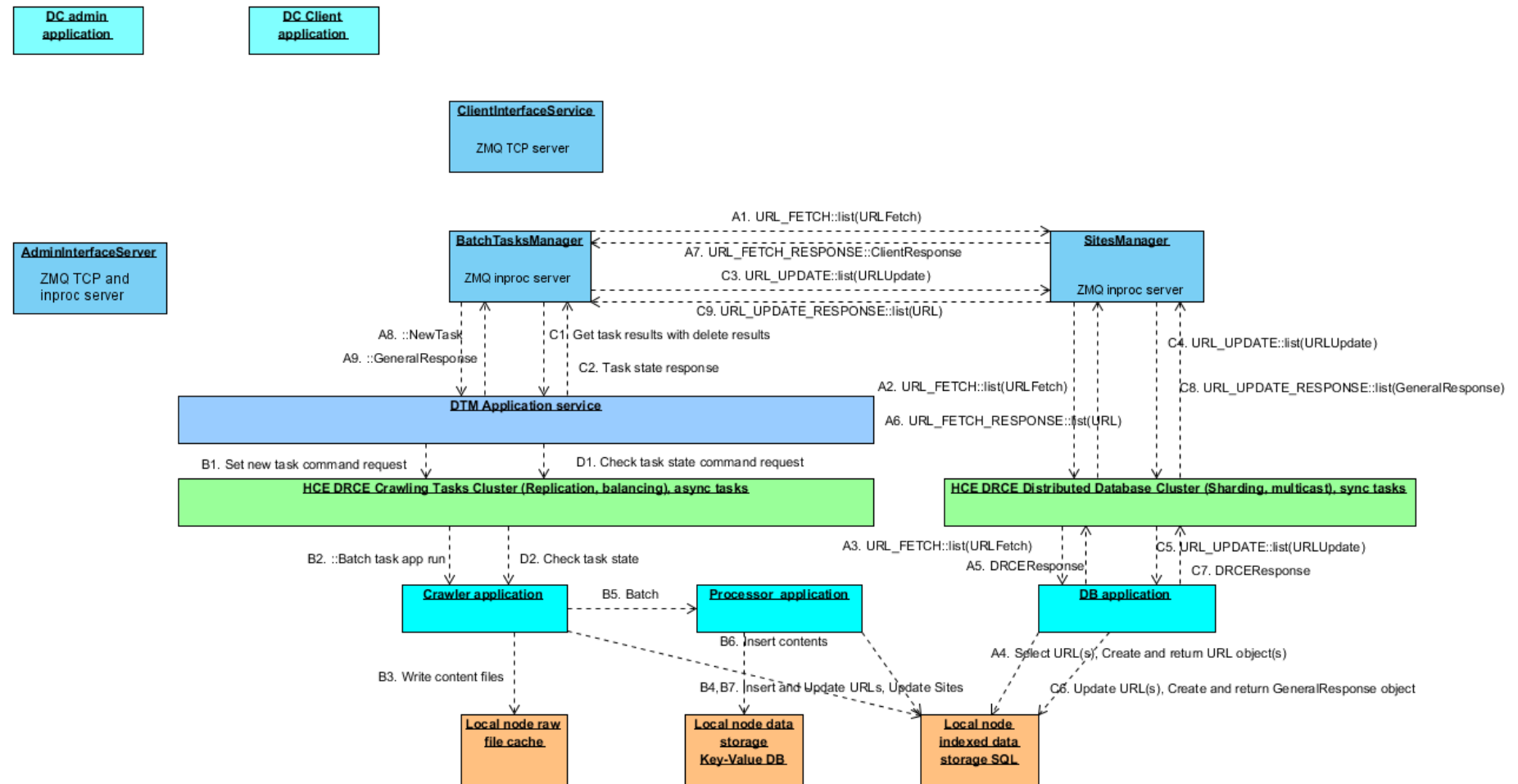


## Operations

Operations are supported inside the service as isolated actions and available via client interface. There are different objects like Site, URL and Batch can be involved in to the algorithm of operation execution as well as it can to have more deep and long results than immediately returned in response.

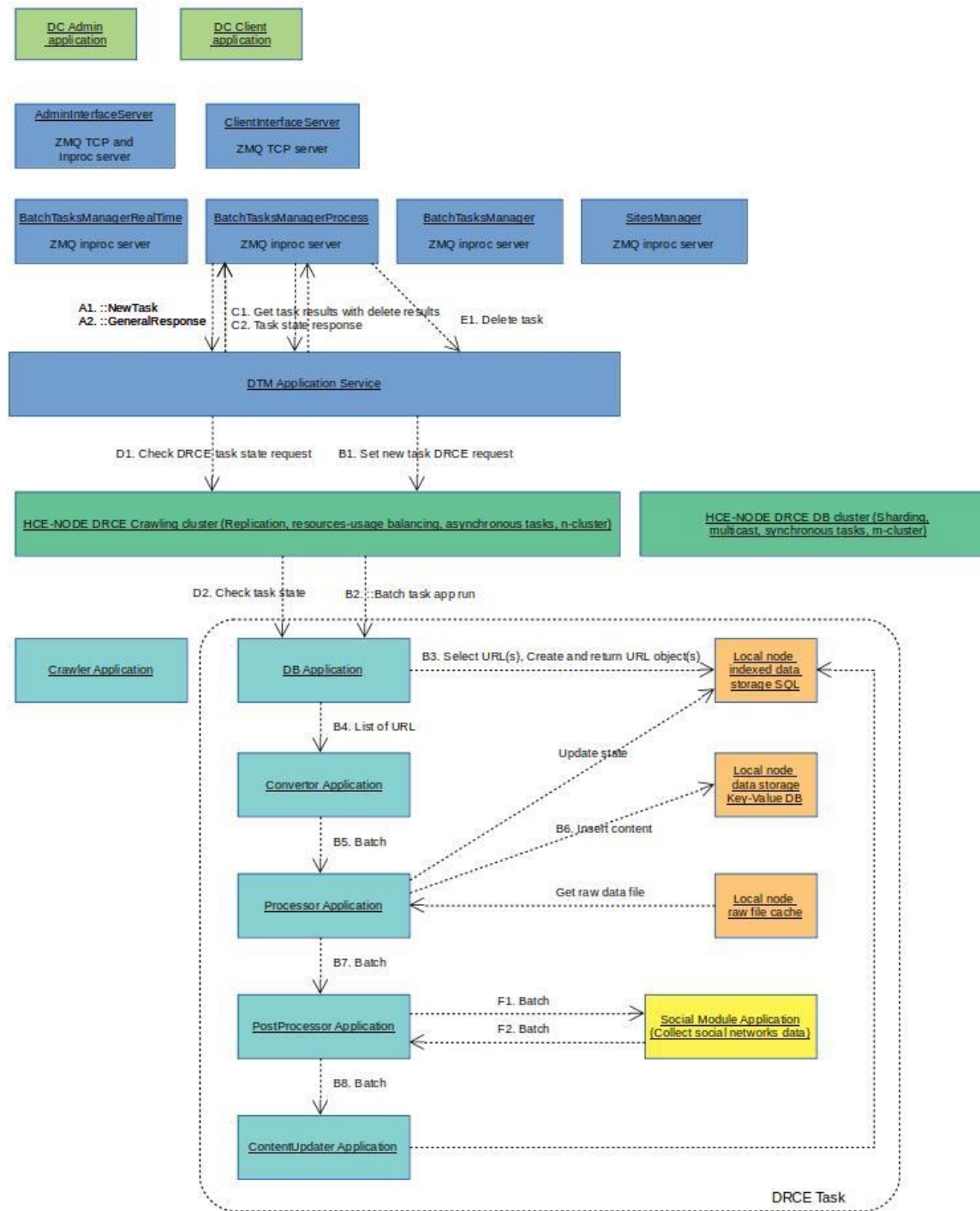
## The regular crawling

The regular crawling is periodic process that creates the Batch object by requesting the lists of URLs object from the SiteManager. The list of URLs is unique and only not present in another active batch tasks are accepted. The Batch objects sending in to the DTM service as tasks and managed by the tasks queue. After finish of task execution URLs from the Batch are updated in state CRAWLED(4) at all data hosts. Before task execution all URLs from the Batch are send to insert as new with state SELECTED\_TO\_CRAWL(2) to all data hosts – the URLs propagation method. The DTM service set the Batch task on n-type cluster and its router uses resource usage balancing mode to choose proper data host to execute. This way URLs crawled only on one data host from several, but collected on all with state CRAWLED(4) to prevent farther redundant crawling and processing.



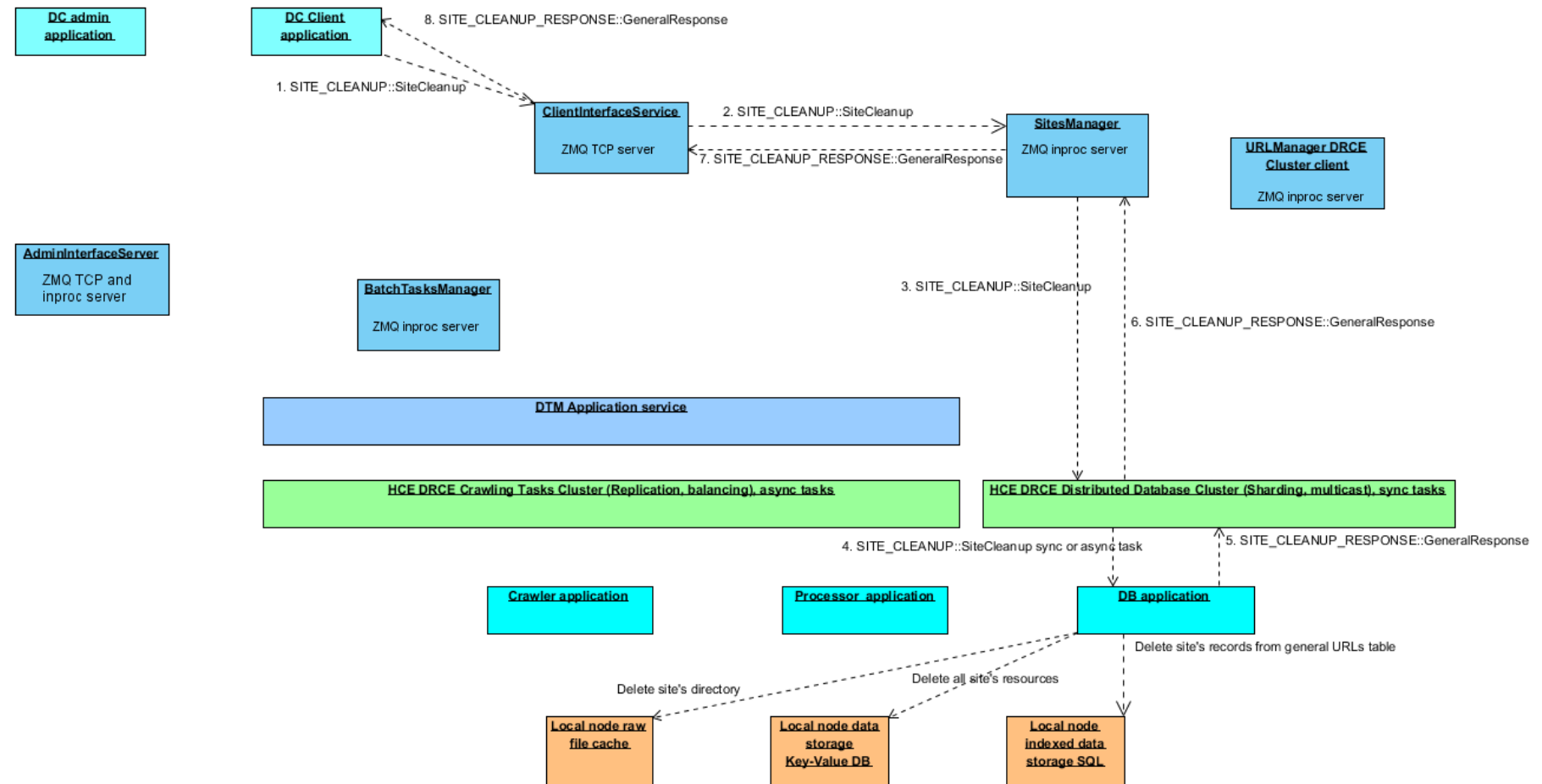
### The processing

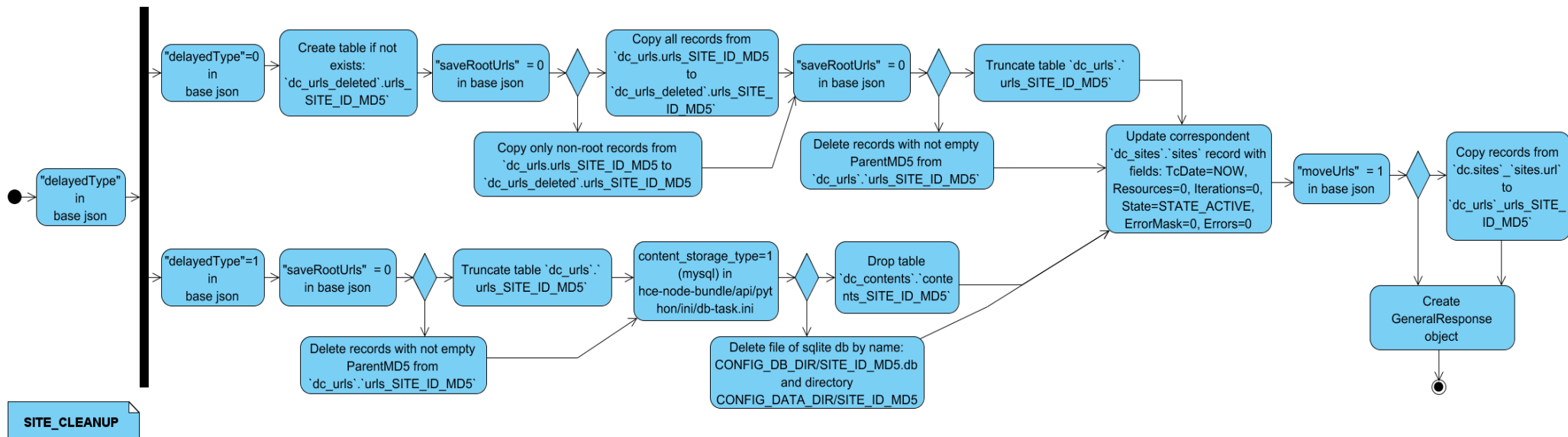
The processing is periodic process that creates the URLFetch object and sets it as DTM task to execute. The execution has three steps: 1) execution of the URLFetch operation, 2) convert the URLs list to the Batch object and 3) execution of the processing operation. The result of the processing depends on the processing method or algorithm, typically it is a scraping. The processed resulted content are stored depends on the method, for the scraping it is a local storage. The URLs objects inside a DC service that are processed changed status to PROCESSED(7).



## Site cleanup

The Site object cleanup operation performs completely deletes all related URLs objects and contents and free space. Can be performed with delete of root URLs and without. If uses delete method that inherited from default.

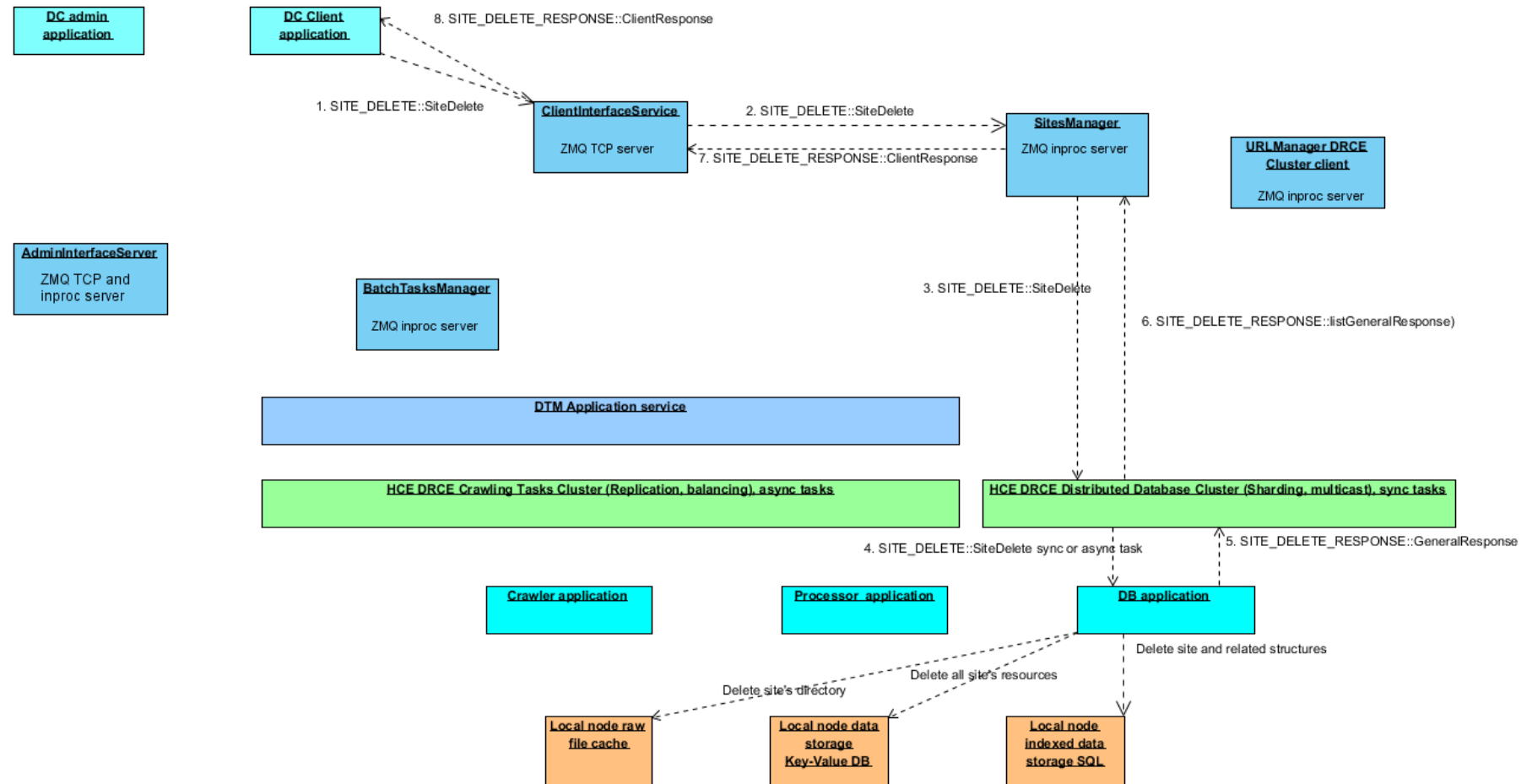


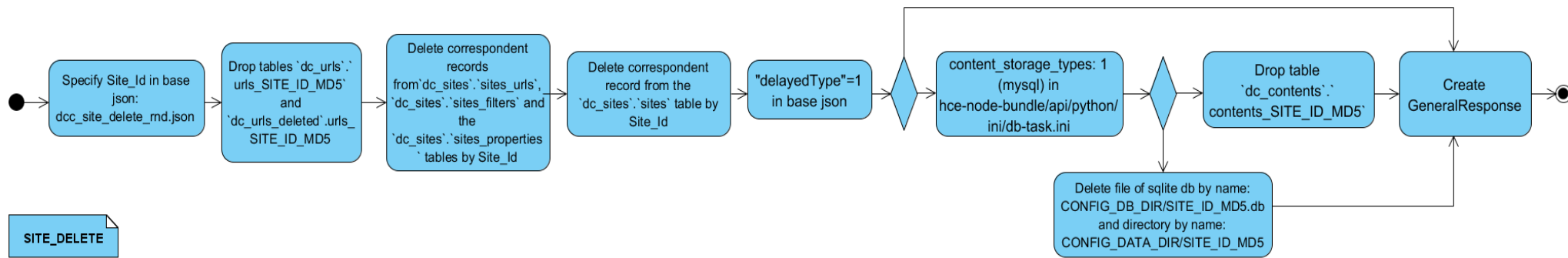




## Site delete

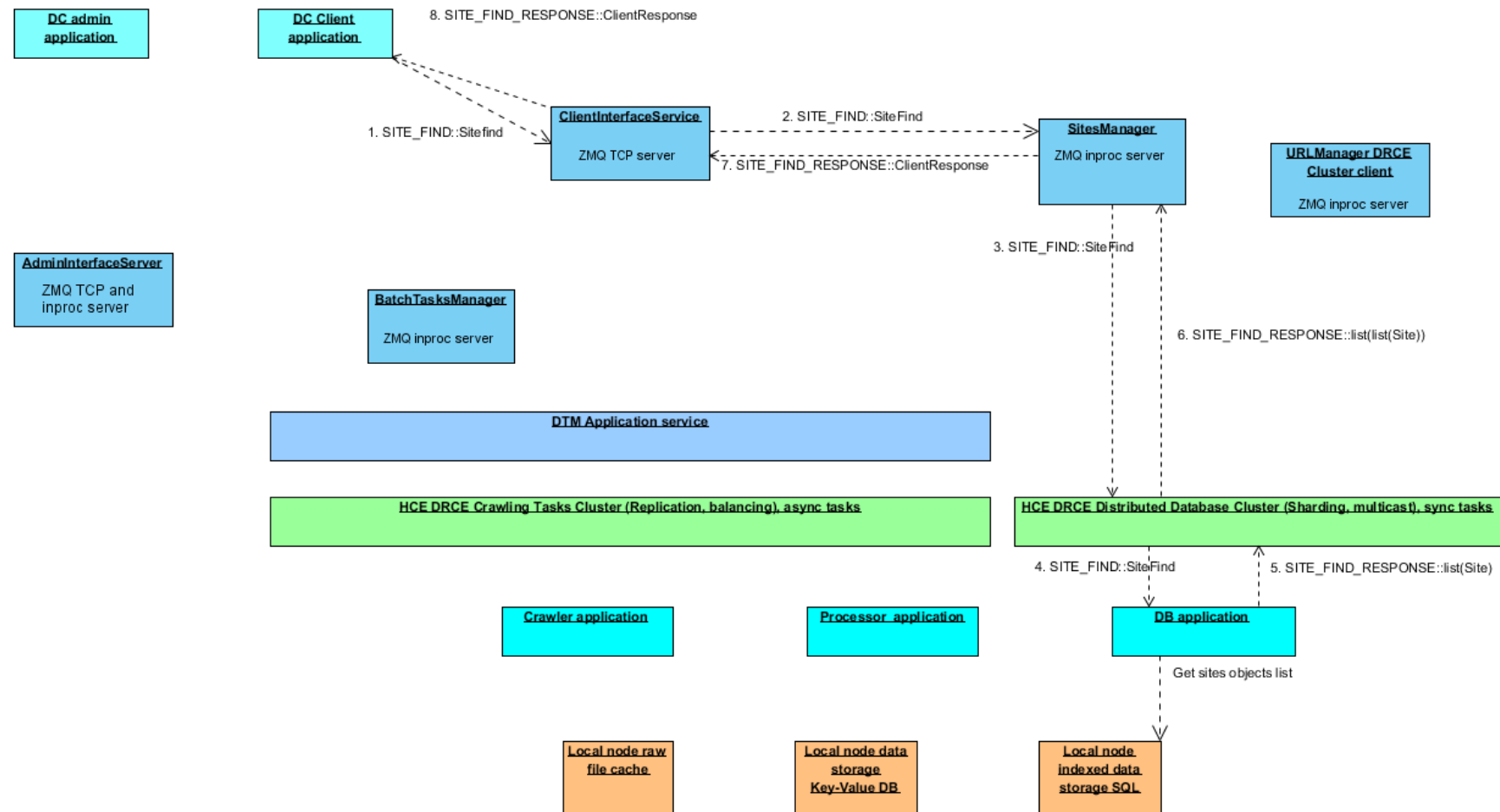
Site delete removes the Site representation from the service storages and completely cleans all related structures, delete all accumulated data and lists of items like URLs or documents.

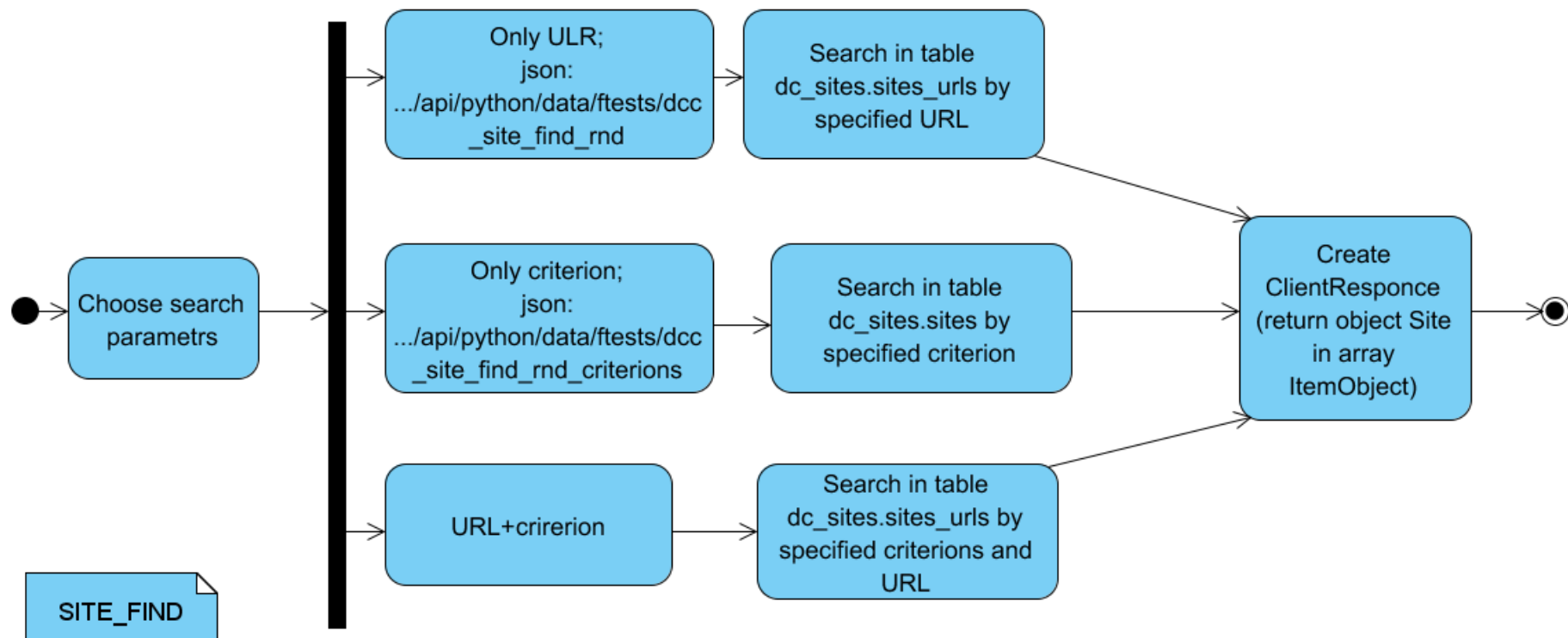




## Site find

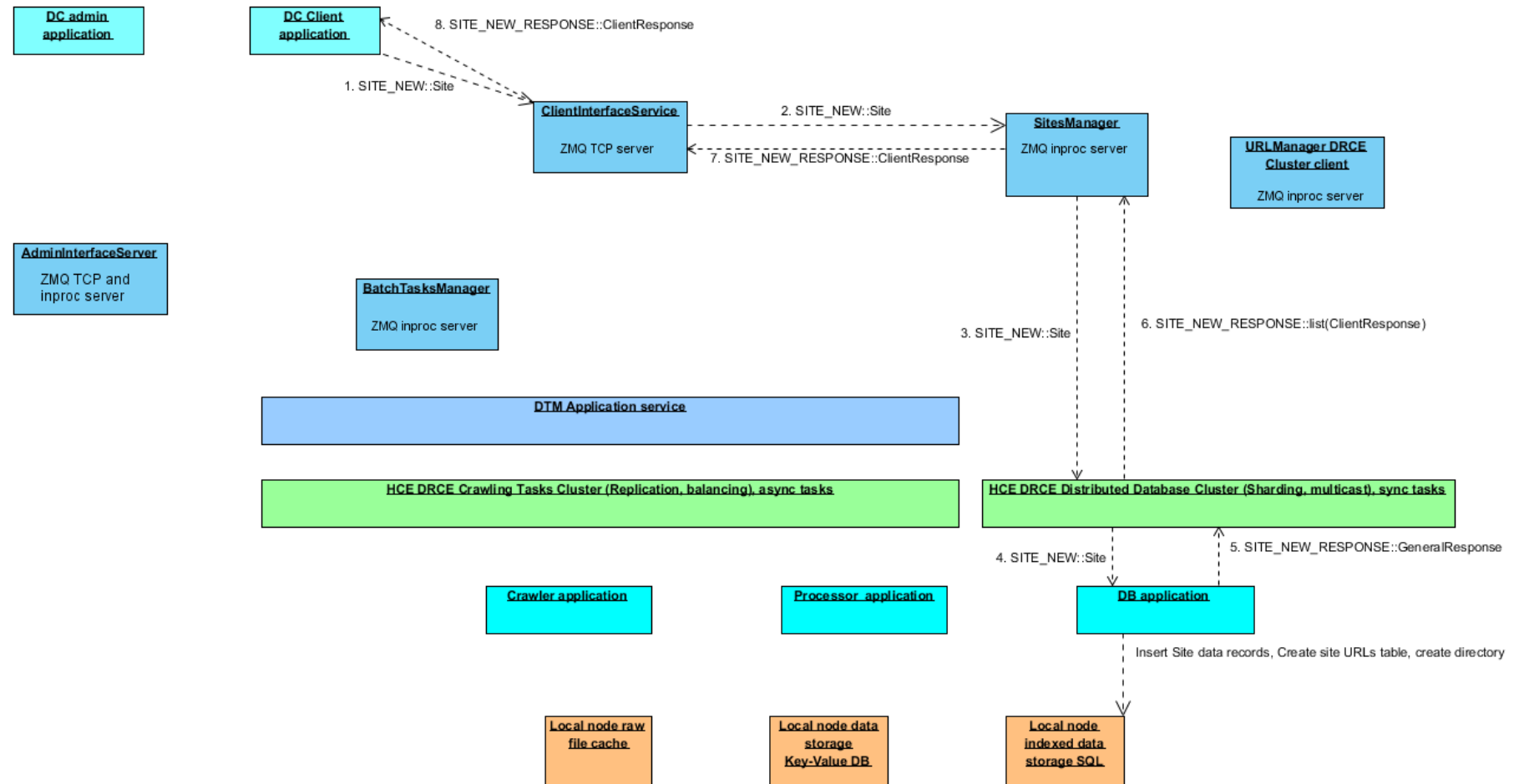
Site find performs search action with defined criterions or URL pattern. If URL pattern used – root URLs of all sites in the service compared. In results list of the Site objects returned or empty list if no one site matched.

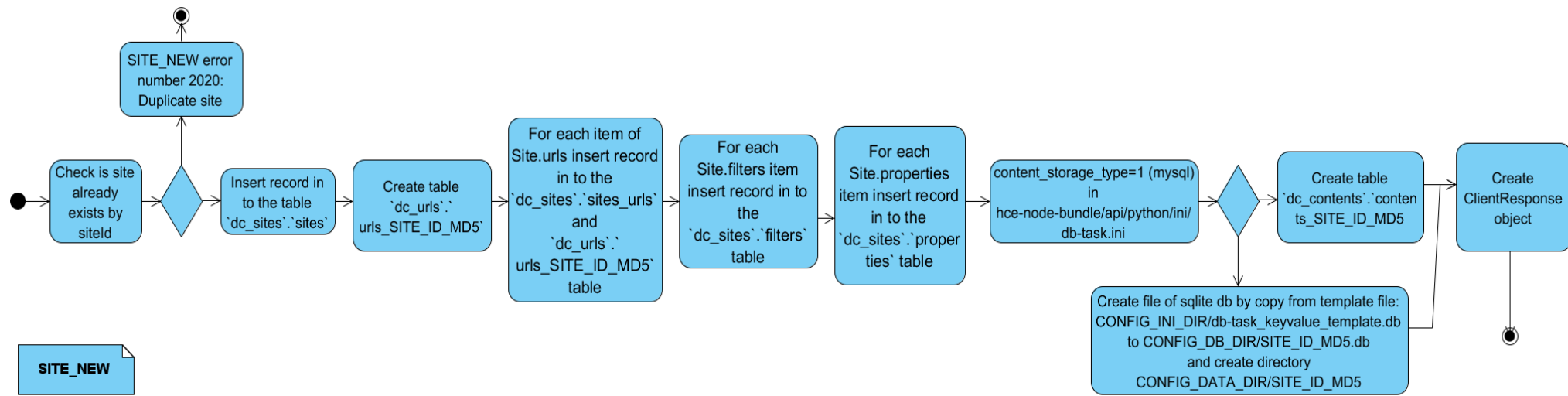




## Site new

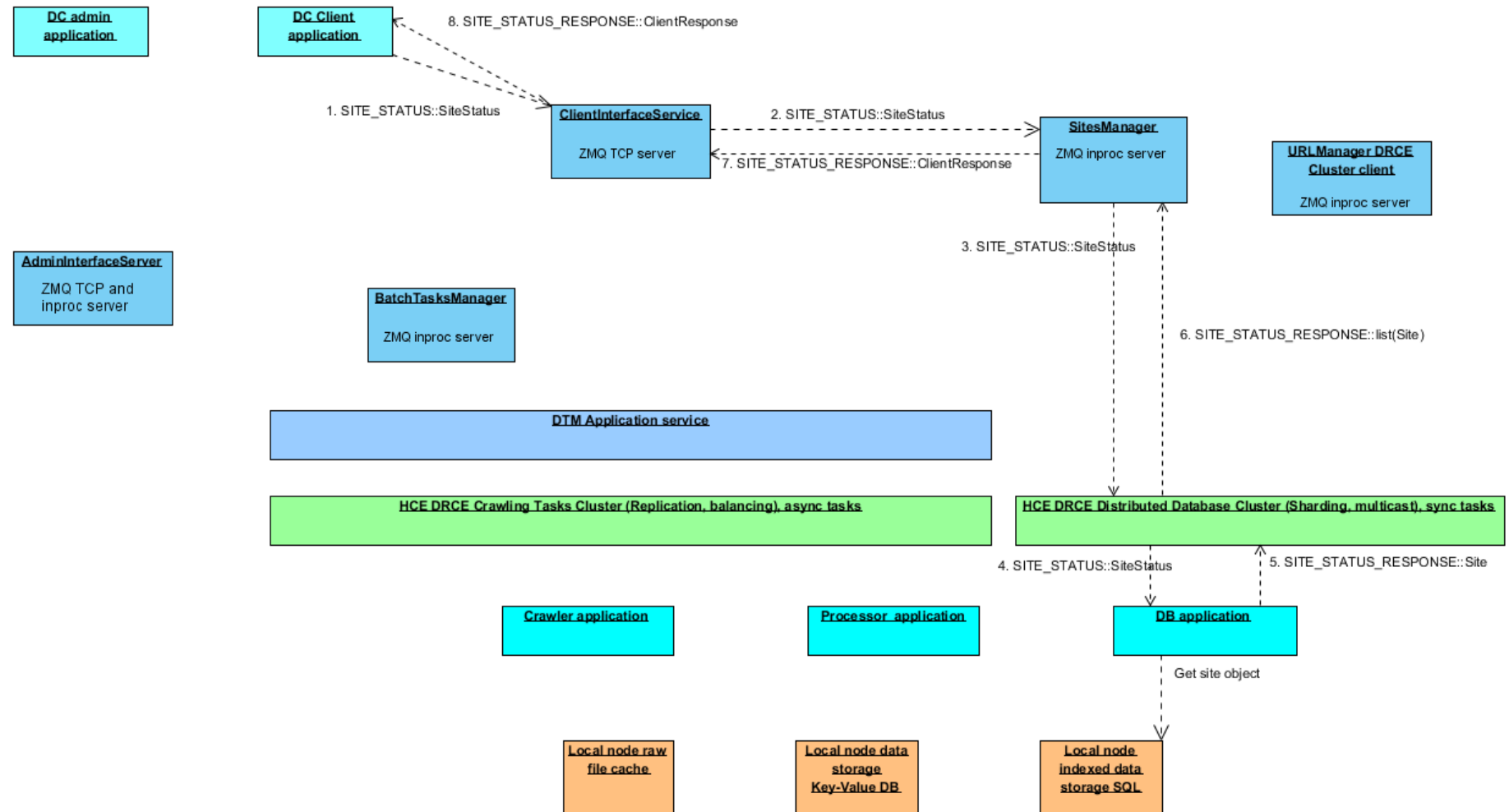
Creates new Site object representation inside the service and initializes all related structures and lists of related objects. Site can start to take a part in the regular crawling or another periodic processes depends on initial values of its fields.

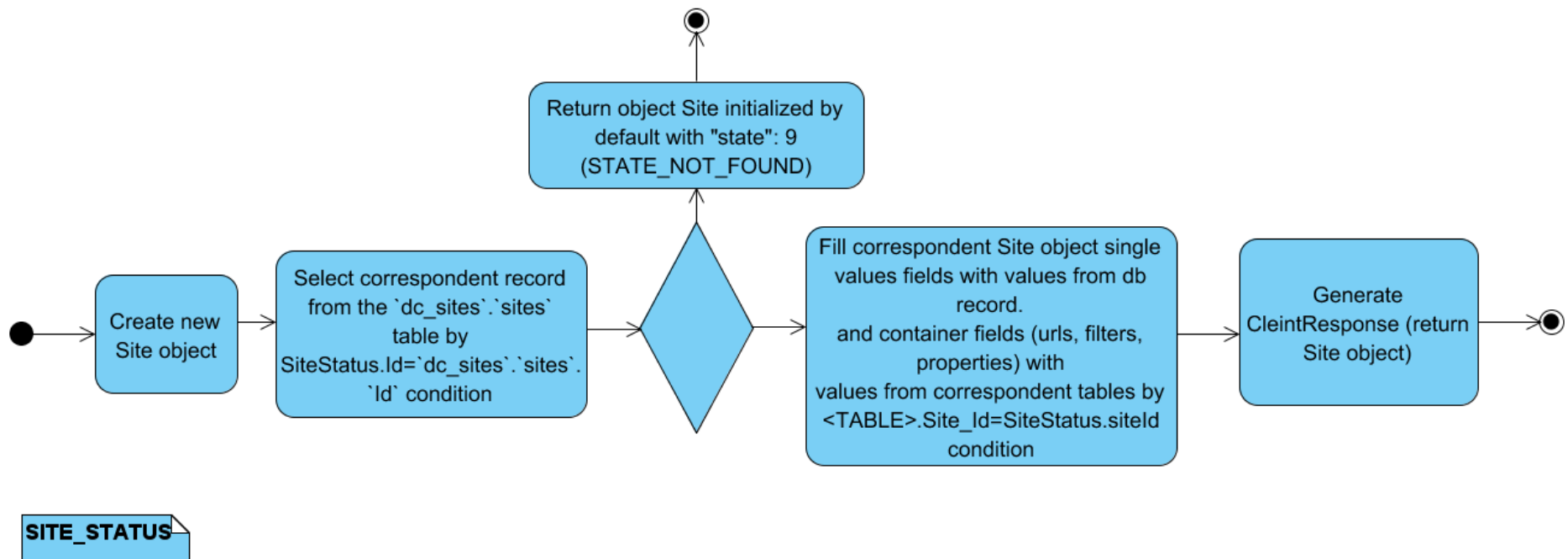




## Site status

This operation performs simple check of presence of the Site object in the system by the site Id and returns it. If site not present – error code and error message set. Operation can be used to get the Site object fast way without involving the search process.

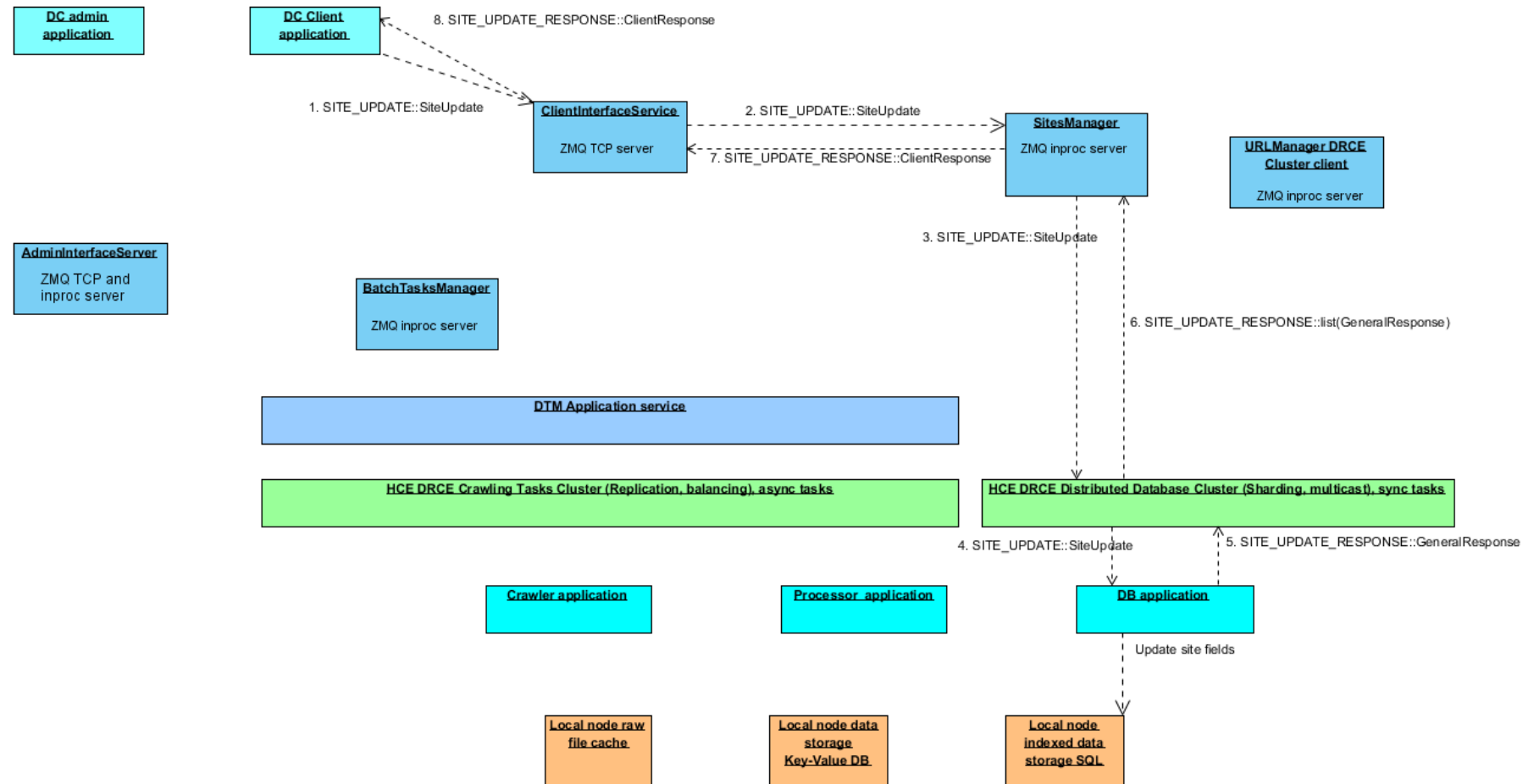


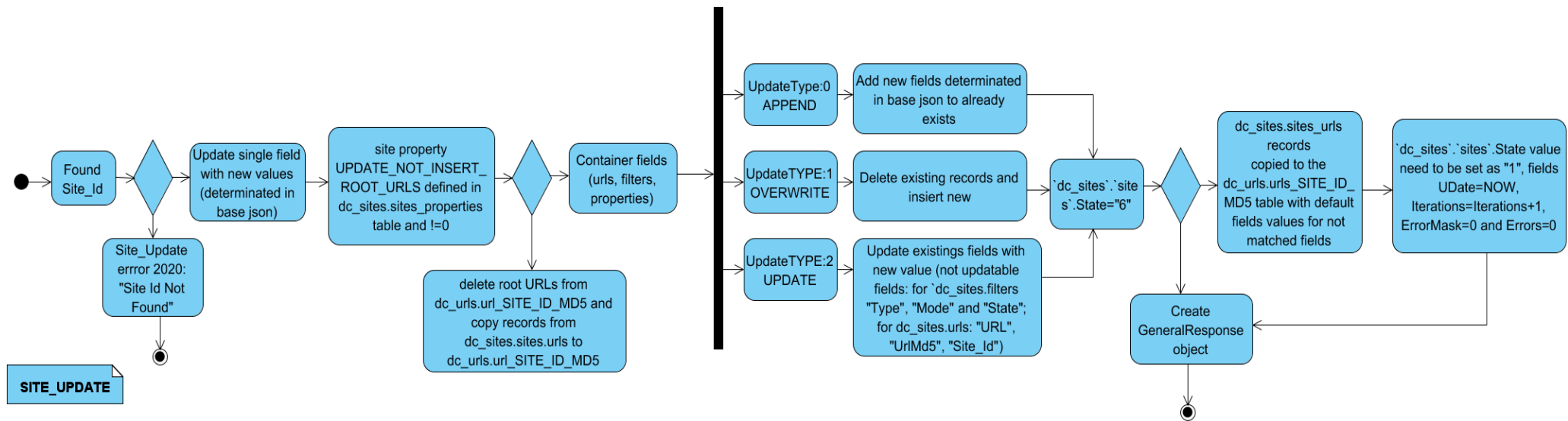




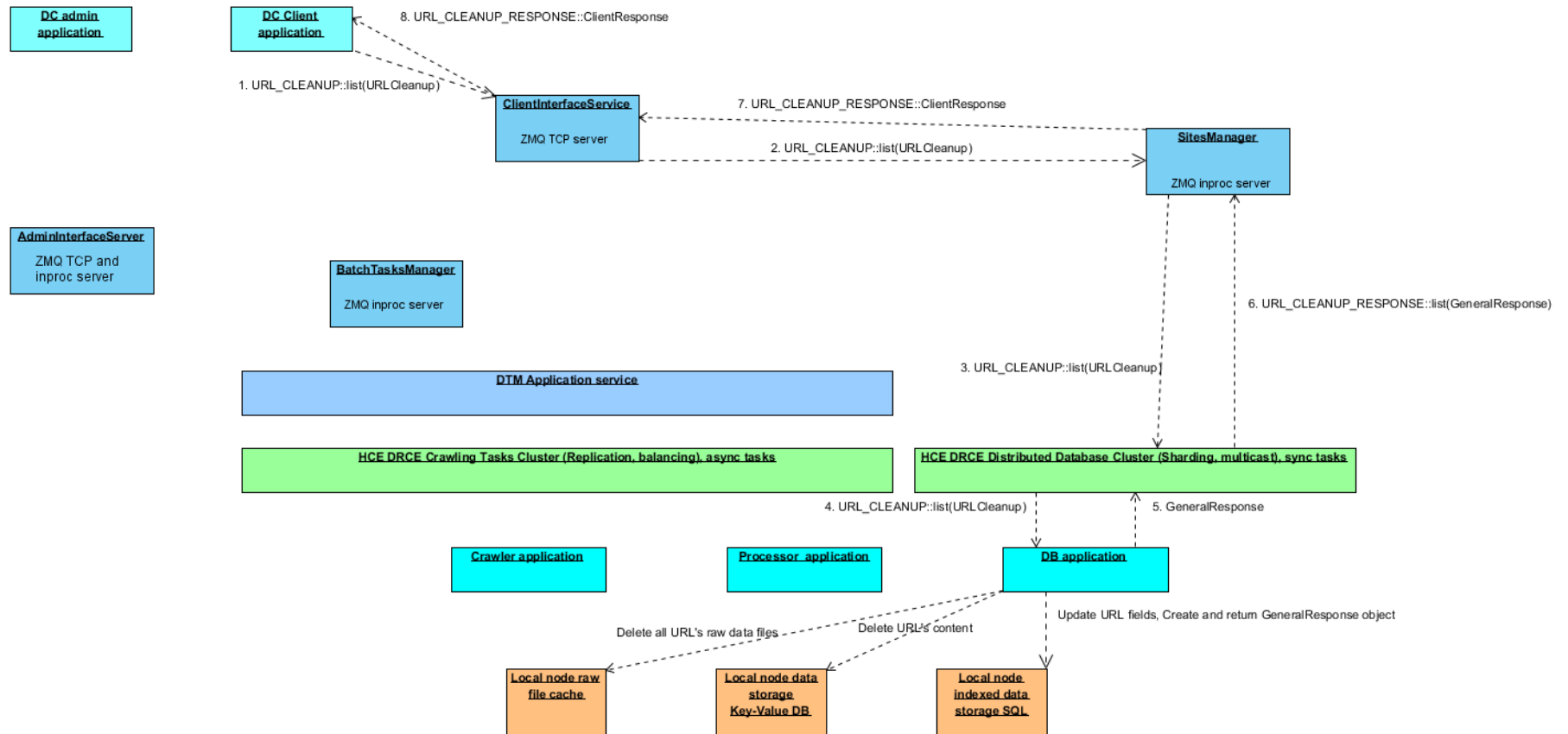
## Site update

This operation updates the Site object representation inside the service and/or related data structures. For multi-item structures like lists of objects like root URLs, filters or templates three different actions possible – append, overwrite and update. Note that operation is complex and can be done partially in case of some database or another kind errors.

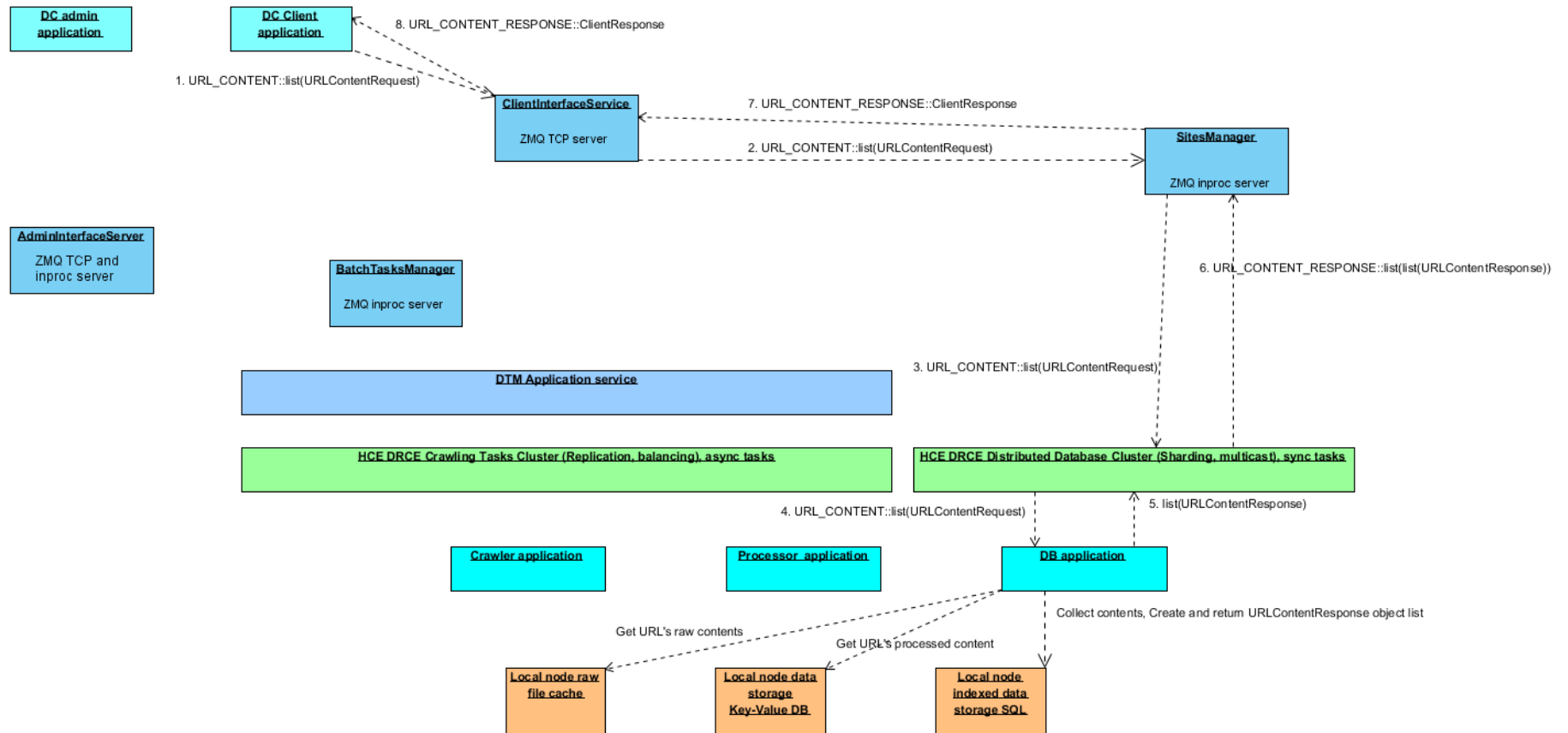




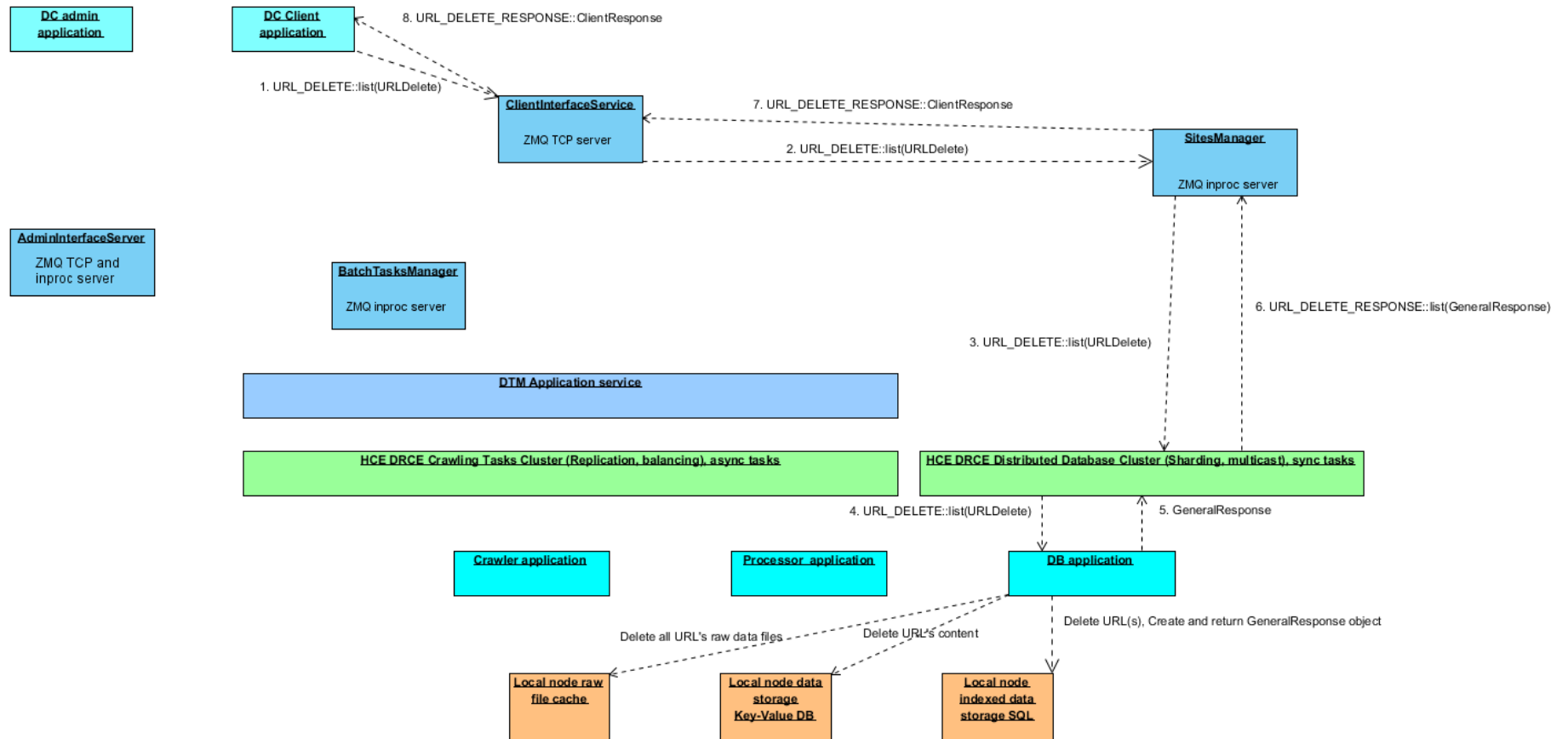
## URL cleanup



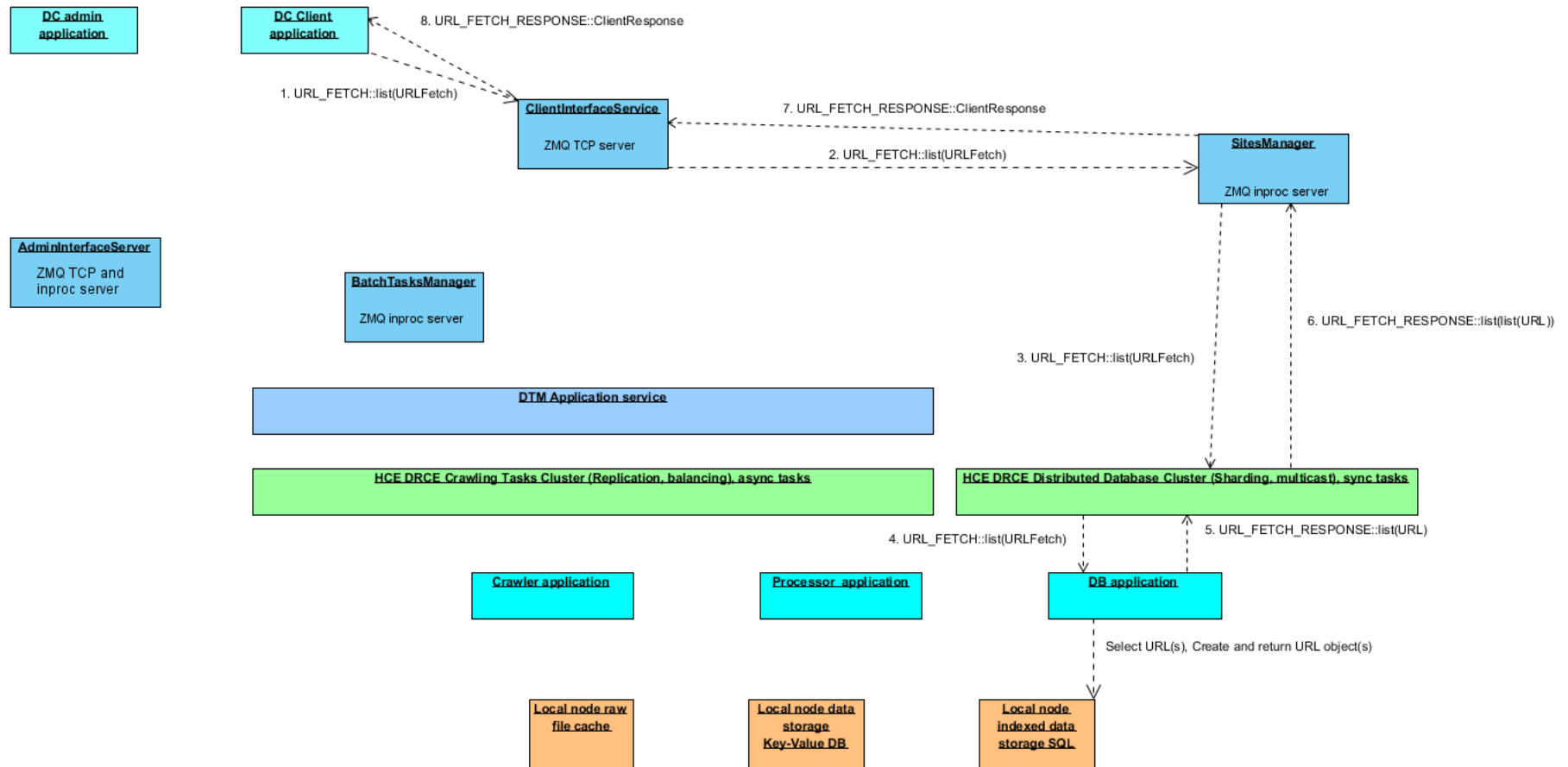
## URL content



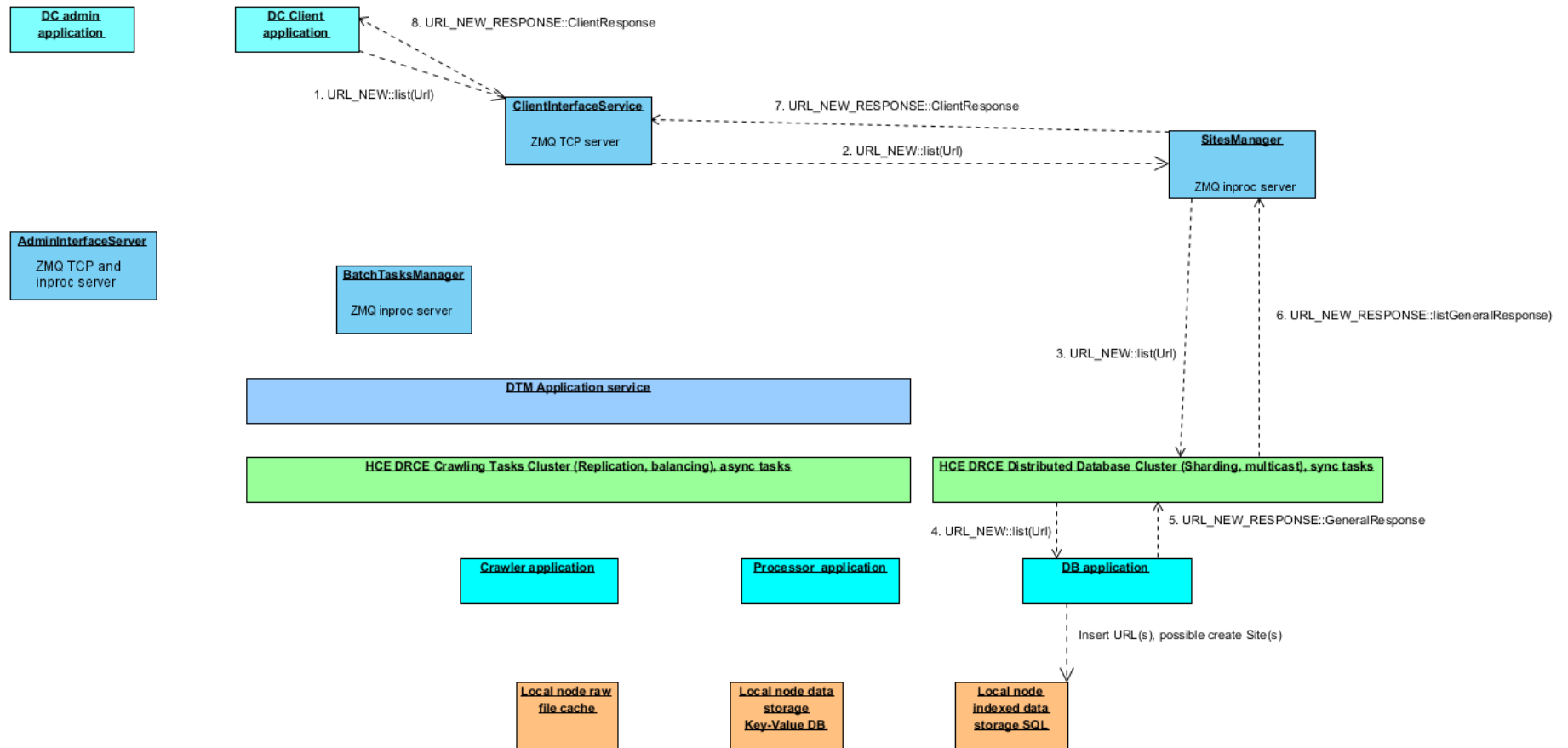
## URL delete



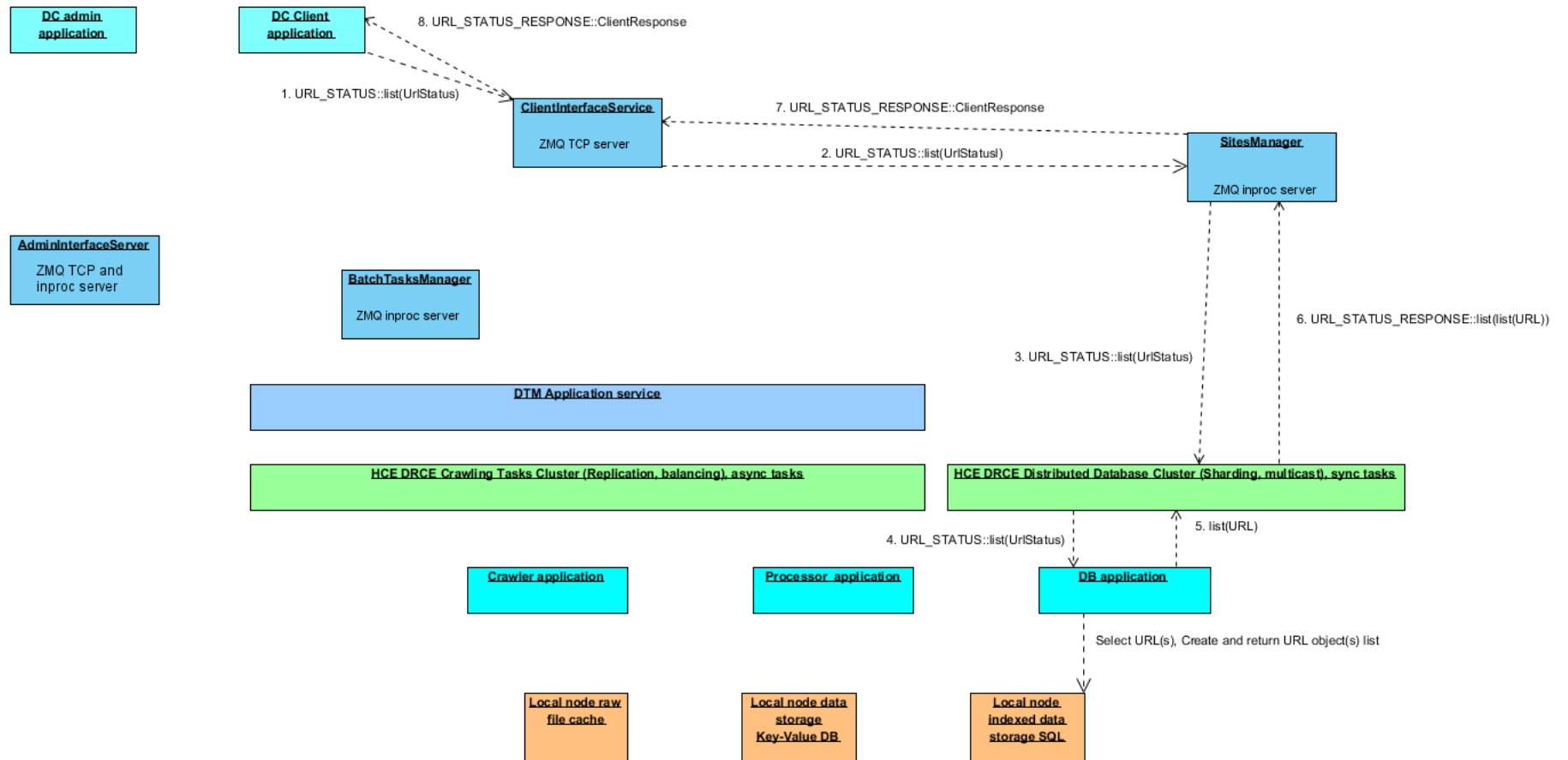
## URL fetch



## URL new

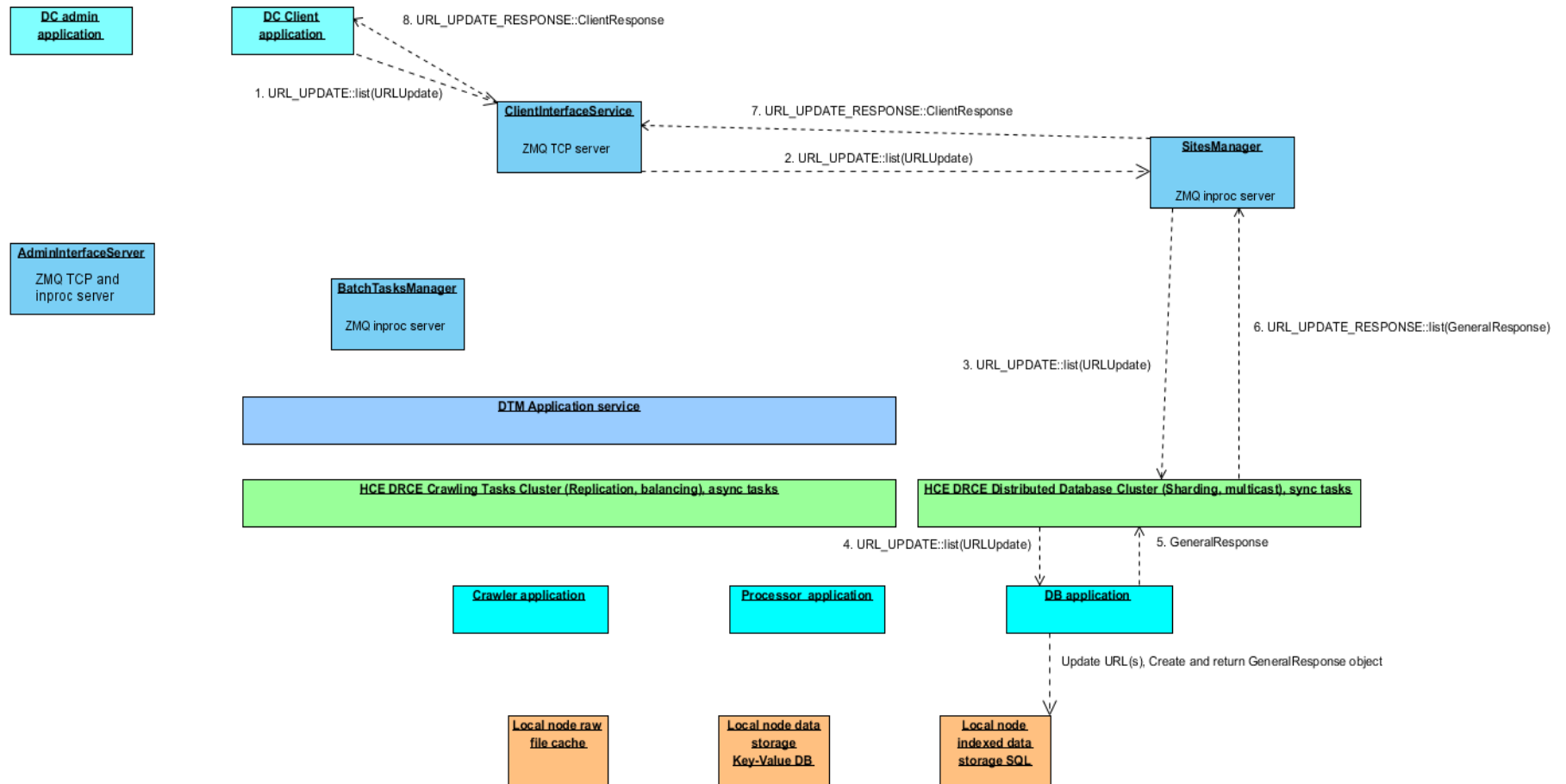


## URL status

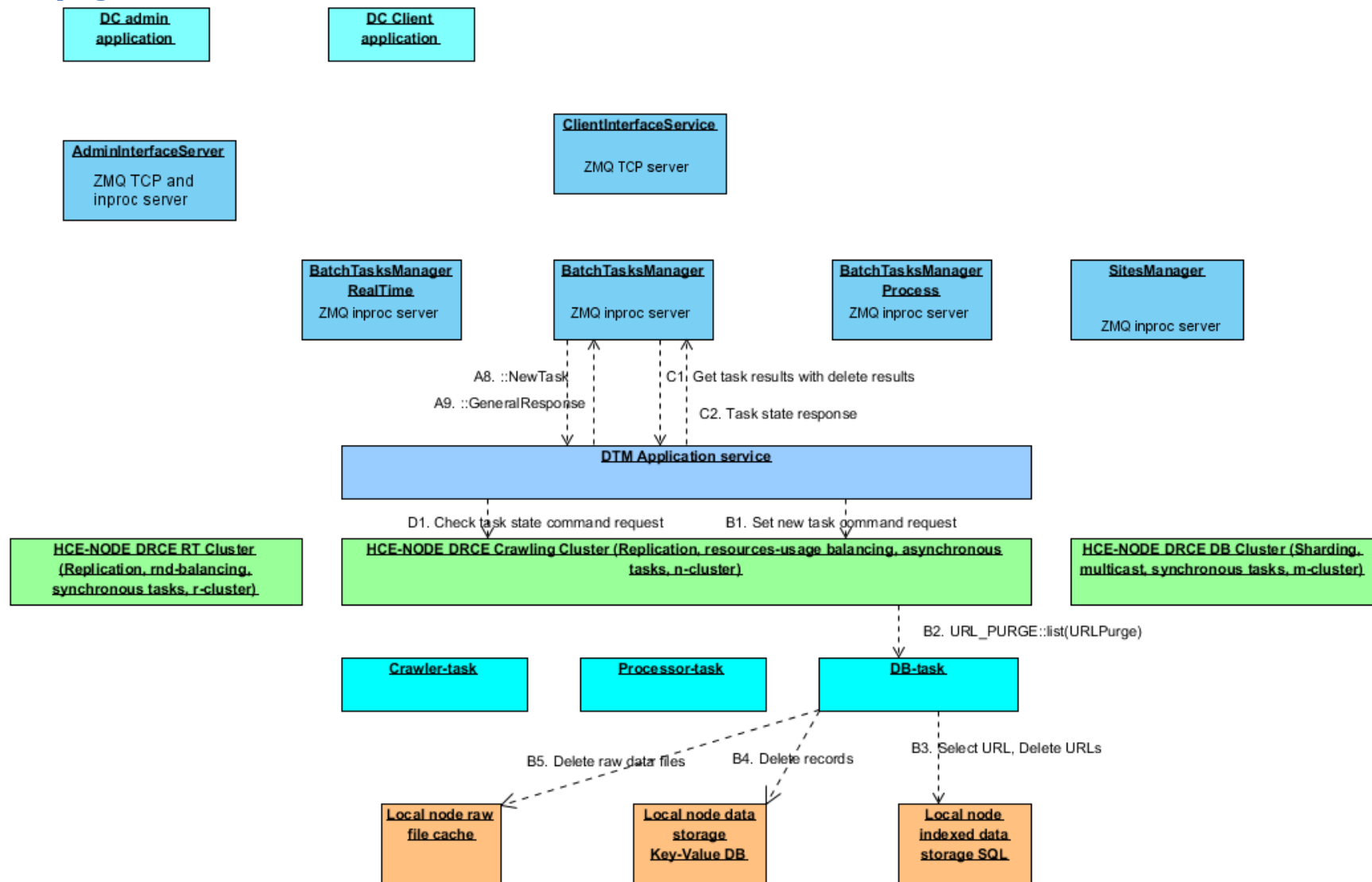




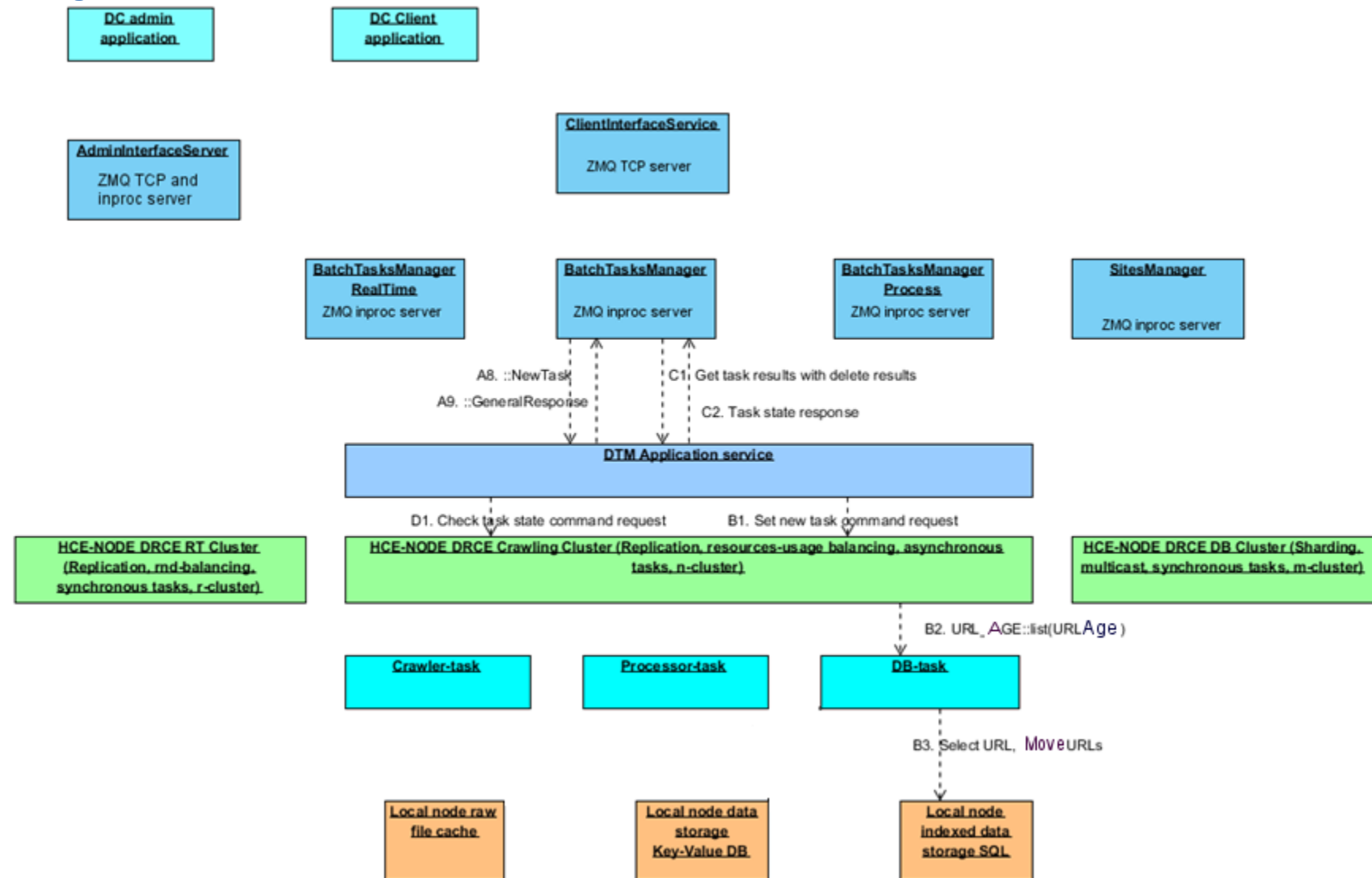
## URL update



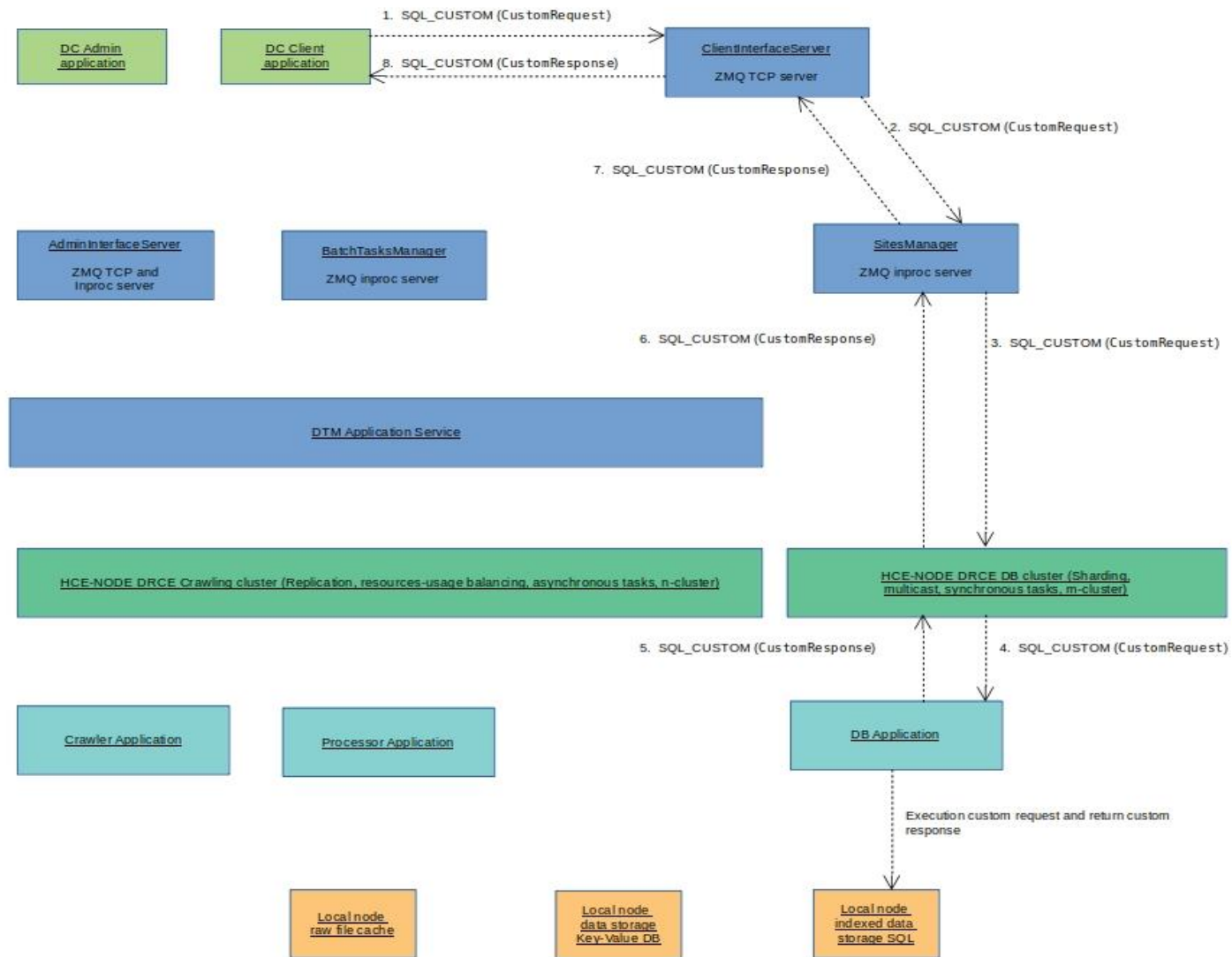
## URL purge



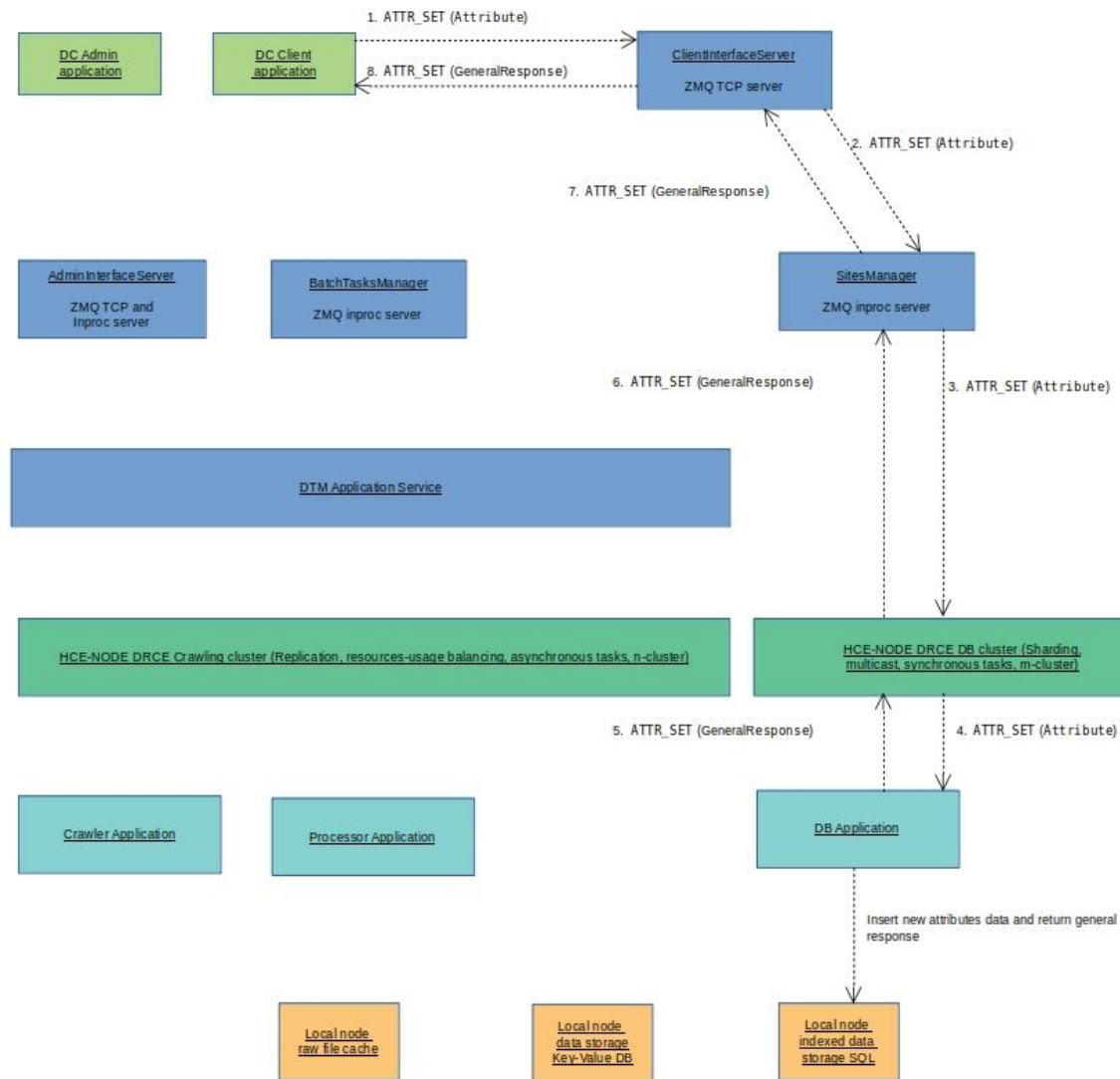
## URL age



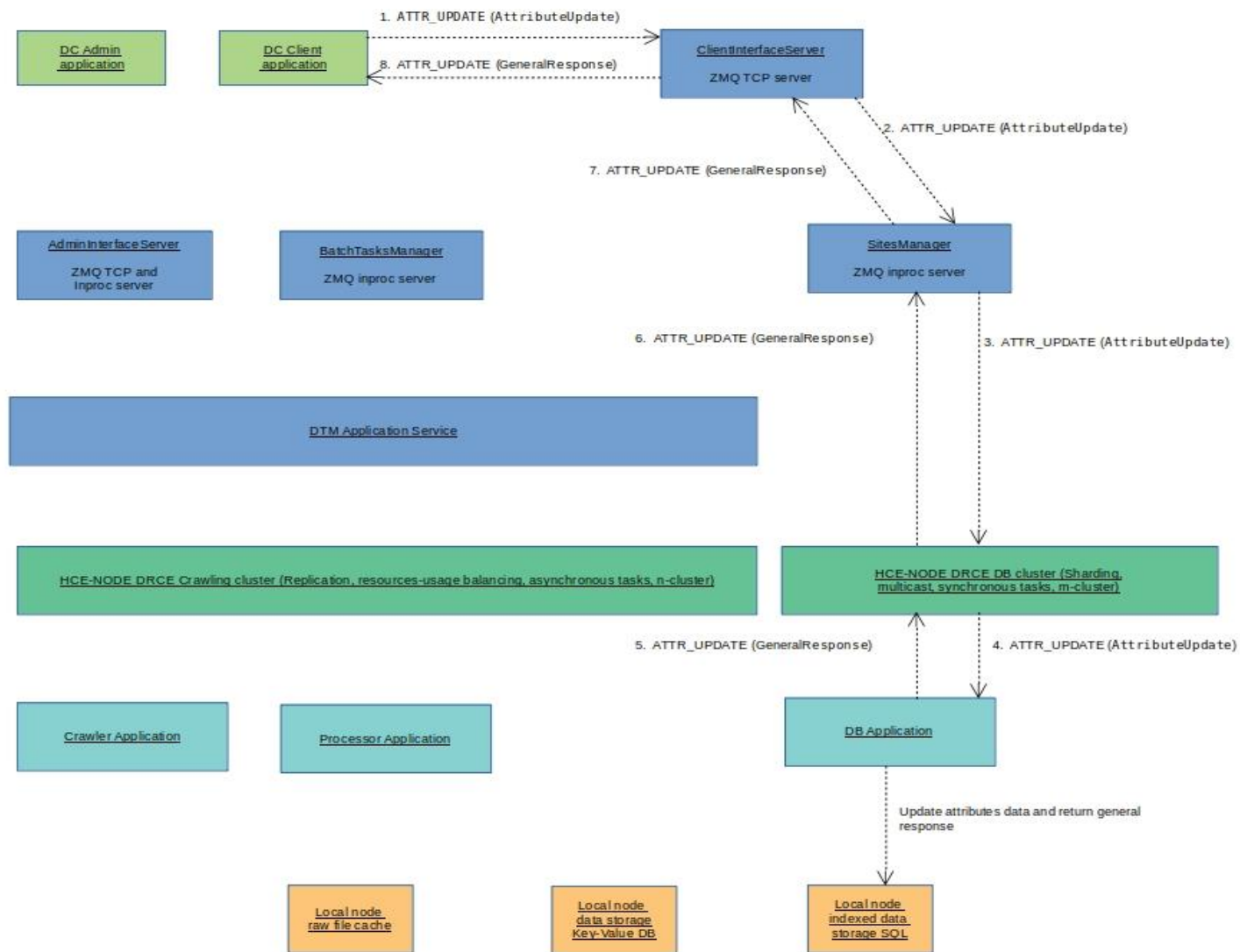
## SQL Custom



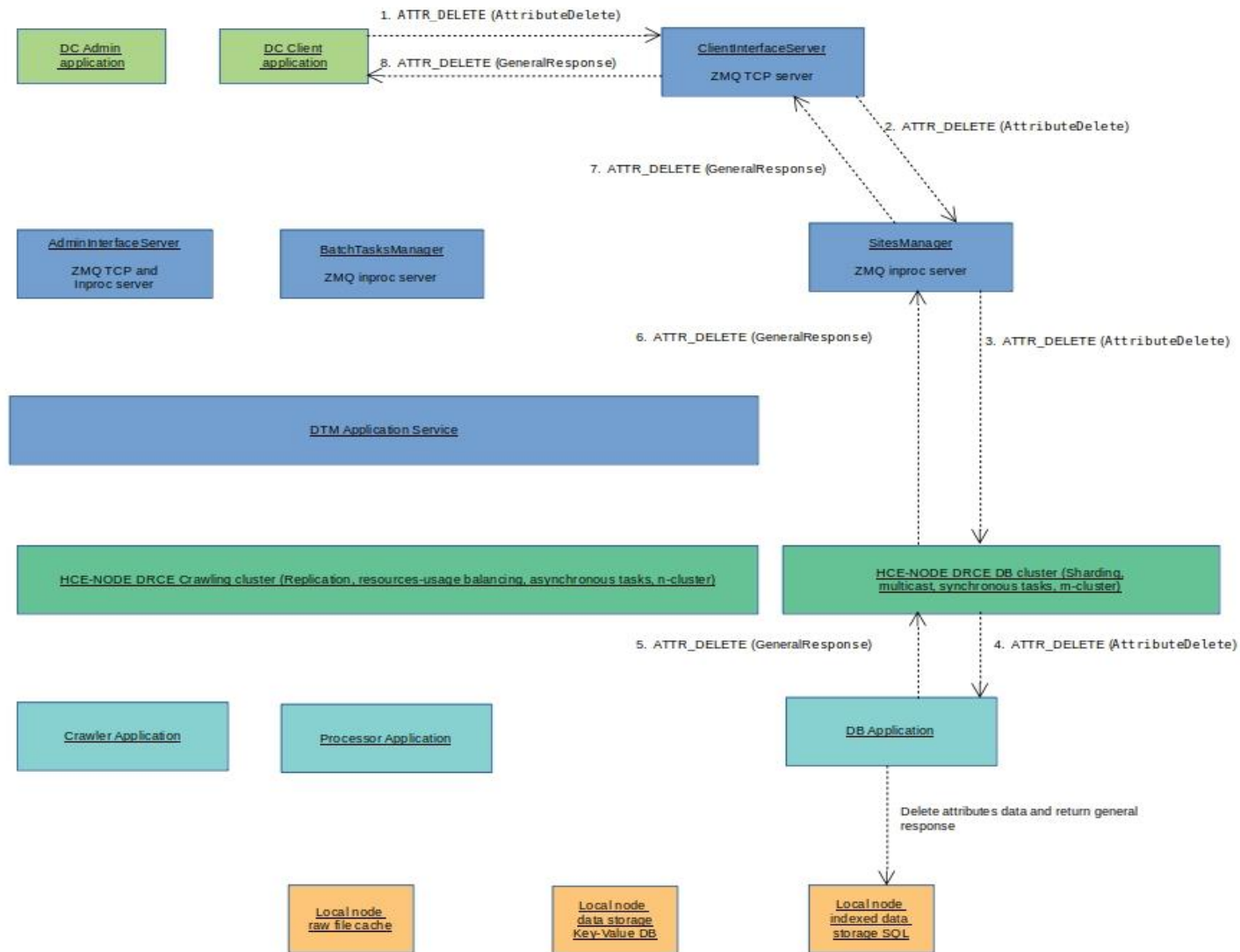
## Attribute set



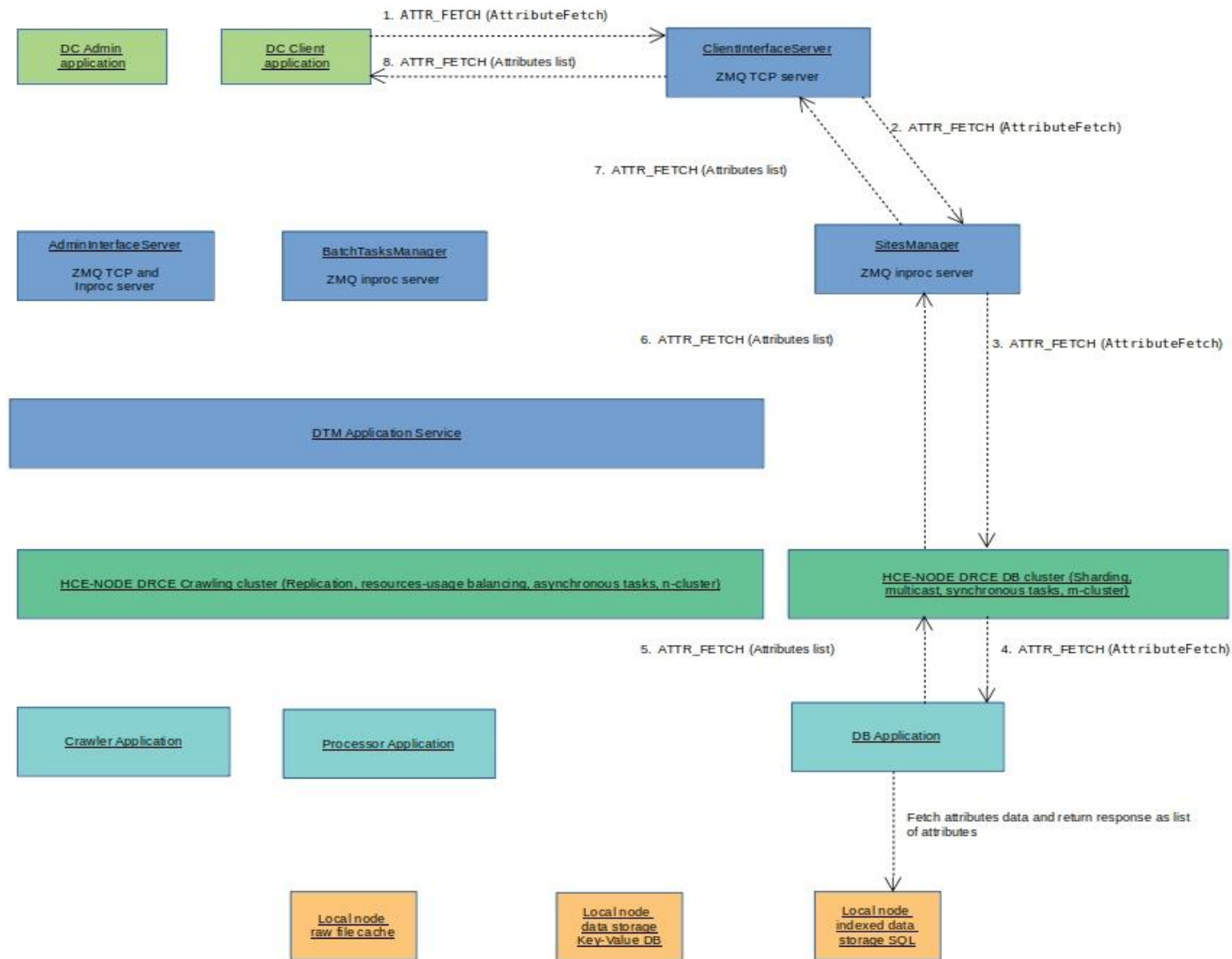
## Attribute update



## Attribute delete



## Attribute fetch





### **Thread classes**

AdminInterfaceServer, ClientInterfaceService, BatchTasksManager and SitesManager.

### **Site**

The structural unit that represents crawling item from user side. Has properties applied for each URL fetched from page of this site.

### **Resource**

The structural unit that represents results of HTTP request downloaded from site by URL. Mainly it is http resource like html page, image or any kind of MIME type content.

### **URL**

The structural unit that represents the resource URL for the site. Used as base of main resource identifier for md5 checksum calculation.

### **HTTP Request**

The object that used to do HTTP querying and resources downloading algorithms and operations on data nodes.

### **HTTP Response**

The object that used to represent and to parse downloaded resources on data nodes.

### **Client interface service**

This class accepts the client requests events from TCP connection and sends them to the proper target threading class. This class works mostly as a proxy between the target functional threaded class and clients. It connects to three target threaded classes as shown on application architecture schema.

### **Batch Tasks Manager**

It is a main class that implements algorithms of continues crawling processing iteration, interacts with DTM application service, sets crawling tasks to execute and checks state of that tasks periodically. To start the crawling task it collects some URLs from SitesManager and creates the crawling task Batch. The crawling task batch send as main data set and input data of the main crawling application that works on target data node in the HCE DRCE Cluster.

## The Crawler application

The application that acts as web-crawler, uses crawling batch as a main job order, executes HTTP requests and stores resources in the local file system cache, updates state of resources in the local SQL DB. Located and acts on data node.

This is main crawling business logic and algorithms application. Steps of run processing:

1. Gets the Batch serialized object from stdin.
2. For each BatchItem from Batch.urls list (md5 string URL identifier) – gets URL record from MySQL db table:
  - If BatchItem.siteId=="" the table is `dc\_urls`.`urls\_0`.
  - If BatchItem.siteId!= "" then table is "`dc\_urls`.`urls\_" + BatchItem.siteId.

If record not found then insert new URL record in to the correspondent table with values from the batchItem.urlObj field contains the URL object filled from the db located on the source host by URL\_FETCH operation.

Farther, the table that contains correspondent record will be named just "urls", but depends on site type (general or not) it must be treated as `dc\_urls`.`urls\_0` or as `dc\_urls`.`urls\_%SITE\_ID\_MD5%` where is the "%SITE\_ID\_MD5%" it is a macro name that need to be replaced with value detected in some way.

2a. Reset `urls`.`ErrorMask` to accumulate only newly error bits for current url.

3. From `dc\_sites`.`sites` read correspondent record. If BatchItem.siteId=="" siteId="0".
4. If SITE.State is not "Active" then skip this BatchItem and process next.
5. If SITE.Errors> SITE.MaxErrors or SITE.Resources>SITE.MaxResources then **accumulate(add to existing bits)** corresponding bit of SITE.ErrorMask, increment SITE.Errors, set TcDate=NOW(), set SITE.State="Suspended" and update all this fields of this site in correspondent site's table detected above. Then process next BatchItem.
6. Get site property "HTTP\_HEADERS" from `dc\_sites`.`sites\_properties` for this site. If empty or no records – then use default headers from correspondent crawler-task\_headers.txt file from **config ini file**.
7. Get site property "HTTP\_COOKIE" from `dc\_sites`.`sites\_properties` for this site. If empty or no records – act the same way as above, file name crawler-task\_cookie.txt.
8. Prepare HTTP request using URL.URL, headers and cookie.
9. Update urls.Status="Crawling", urls.Crawled++, urls.Batch\_Id for correspondent url's record in MySQL db.
10. Wait URL.RequestDelay, msec.
11. Using prepared HTTP headers make HTTP request and get robots.txt content. If content found and not empty – parse them and store the set of rules in temporary container (list or map) for this batch URLs set. When the batch processing finished – this set of rules need to be removed. Check URL in robots.txt rules and if not succeeds – then set urls.State="Error", corresponding bit in ErrorMask fields for sites(**add to existing bits**) and urls(**set only current error bits**) tables, urls.Status="Crawled", and another related Site's and URL's fields like urls.UDate and process next BatchItem.
12. Using prepared request – make HTTP request and wait on response urls.HTTPTimeout msec.

13. If response got – detect: HTTP response code, Charset, MIME Content-Type, and so on downloaded resource's properties and HTTP response fields. If timeout reached – set the corresponding ErrorMask bit for urls.ErrorMask field, and other fields the same way as in case of robots.txt rule not matched case.
14. Update fields of this URL record urls.Status="Crawled", urls.State="Error" (if some error), urls.ContentType, urls.Charset, urls.ErrorMask, urls.TotalTime, urls.CrawlingTime, urls.UDate=NOW(),urls.HTTPCode to proper value.
15. If some other error like zero response, wrong response, HTTP error, content size greater thanMaxResourceSize and so on – set corresponding ErrorMask bit, update urls record.
16. If any kind error trapped that lead to zero content size – remove item from the Batch.items and process next item.
17. Check files directory and create if not found if BatchItem.siteId!="":
  - CONFIG\_DATA\_DIR/SITE\_ID\_MD5/PathMaker(RESOURCE\_ID\_MD5).getDir()/
 and if BatchItem.siteId=="":
  - CONFIG\_DATA\_DIR/0/PathMaker(RESOURCE\_ID\_MD5).getDir()/
18. Write HTML content file (if MIME Content-Type requires – converted to utf-8):
  - RESOURCE\_ID\_MD5\_ASCIIIME.bin
19. Write HTTP response header file:
  - RESOURCE\_ID\_MD5\_ASCIIIME.headers.txt
20. Write HTTP request file:
  - RESOURCE\_ID\_MD5\_ASCIIIME.request.txt
21. Collect URLs from HTML "A", "IMG", "IFRAME", "FRAME" and so on, complete list see here ["http://stackoverflow.com/questions/2725156/complete-list-of-html-tag-attributes-which-have-a-url-value/"](http://stackoverflow.com/questions/2725156/complete-list-of-html-tag-attributes-which-have-a-url-value/) and canonize them to have only fully qualified URLs uses "HTTP" protocol only (later another protocols support will be added).
22. For each collected URL check if already exists by URLMD5 field and if not – insert record to the urls table detected on step 2) (general or corresponds to this site), set fields with proper values of Site\_Id, URL and so on. The Type, RequestDelay, HTTPTimeout set from correspondent Site fields. Fields that not defined leave default value. The URL candidate to insert need to be verified as belonged to this site by usage `dc\_sites`.`site\_filters` table records by `Site\_Id`. The record has `Pattern` field defines the regexp pattern to check and the `Type` field defines check result treat as exclude this URL from farther processing (0) or include (1). The URL candidate to insert in to the database need to be checked by all filters for this site and in case of satisfy (no one pattern will exclude it) need to be inserted. If some check is not success – the URL skipped and not inserted in to the database for farther processing.
23. Serialize the resulted Batch object (possible, some items were removed) and print serialized string content to the stdout. Finish with zero exit status value in case of no critical errors and with 1 in another case. The error that can be treated as critical related with some DB structure, file I/O operations like permissions, wrong input object and so on...

### *Site.ErrorMask and URLs.ErrorMask specification*

Bit #	Description
<b>Crawling errors</b>	

0	Wrong URL
1	Timeout
2	HTTP error
3	Empty content
4	Wrong MIME type
5	Connection error
6	Code page convert error
7	Bad redirection
8	Size error
9	Authorization error
10	File operation error, write file, create directory and so on
11	Robots.txt rule not matched.
12	HTML_PARSE_ERROR
13	BAD_ENCODING
14	SITE_MAX_ERRORS
15	SITE_MAX_RESOURCES
16	RAW_CONTENT_NOT_STORED
17	MAX_ALLOW_HTTP_REDIRECTS
18	MAX_ALLOW_HTML_REDIRECTS
<b>Processing errors</b>	
32	MaxErrors limit reached
33	MaxResources limit reached
34	Unsupported content-type for processing
35	Error raw content decoding
36	Scraper processing error
22	
23	

24	
25	
26	
27	
28	
29	
30	

## Processor-task application

This application represents sets of data processing algorithms that can be used to process locally stored data. Located and acts on data node. This is main an after crawling data processor business logic and algorithms application. Steps of run processing:

1. Gets the Batch serialized object from stdin.
2. For each BatchItem from Batch.urls list (md5 string URL identifier) – gets URL record from table:
  - If BatchItem.siteId=="" the table is `dc\_urls`.`urls\_0`. Site Id="0".
  - If BatchItem.siteId!= "" then table is "`dc\_urls`.`urls\_" + BatchItem.siteId.
3. From `dc\_sites`.`sites` read correspondent record. If BatchItem.siteId=="" Site Id="0".
4. If SITE.State is not "Active" then skip this BatchItem and process next.
5. If SITE.Errors> SITE.MaxErrors or SITE.Resources>SITE.MaxResources then set corresponding bit of SITE.ErrorMask, increment SITE.Errors, set TcDate=NOW(), set SITE.State="Suspended" and update all this fields of this site. Then process next BatchItem.
6. Form `dc\_urls`.`urls\_SITE\_ID\_MD5` table read corresponded record by the `URLMd5` field value.
7. From `dc\_sites`.`sites\_properties` read all records belongs this site by the `Site\_Id` field. For all records with field `Name`="PROCESS\_CTYPES" check is value of field `dc\_urls`.`urls\_SITE\_ID\_MD5`.`ContentType` the same as one of `dc\_sites`.`sites\_properties`.`Value`. If not, skip this URL, set correspondent `ErrorMask` as "Unsupported content-type for processing", update URL record and go to process next item from the batch.
8. Update URL.Status="Processing", URL.Processed++ for correspondent record.
9. Call processing algorithms methods and create resulted object.
10. Serialize resulted object to the json format.
11. If BatchItem.siteId!= "" – connect to the site file db using path:
  - CONFIG\_DB\_DIR/SITE\_ID\_MD5.db
 If BatchItem.siteId== "" – connect to the site file db using path:
  - CONFIG\_DB\_DIR/0.db
12. Insert or update the record, key is URL\_MD5\_ID, value is json string, CDate is NOW().
13. Update fields of this URL record URL.Status="Processed", URL.State="Error" (if some error), URL.ErrorMask=URL.ErrorMask & "Error mask" (if some error), URL.TotalTime=URL.TotalTime+time, URL.ProcessingTime=time, URL.UDate=NOW() to proper value.
14. Update `dc\_sites`.`sites`.`Contents` + 1 for current siteId.
15. Serialize the resulted Batch object and print serialized string content to the stdout. Finish with zero exit status value in case of no critical errors and with 1 in another case.

## Db-task application

The application that acts as a main interface with DTM application to distributed DB units. The DB units manage the local databases of all objects and items defined on local data node level. Also, this application fetches locally stored resources. Located and acts on data node. The main SQL database schemas defined in the `dc_sites.sql` and `dc_sites_template.sql` dump files.

This is main data storages (SQL DB, key-value DB and files) operations business logic and algorithms application. Steps of run processing:

1. Gets the serialized object from stdin. This object it is tuple of (OPERATION\_CODE, OperationObject).
2. Qualify the object type by one of two operations groups:
  - **Site:** (SITE\_NEW, Site), (SITE\_UPDATE, SiteUpdate), (SITE\_STATUS, SiteStatus), (SITE\_DELETE, SiteDelete), (SITE\_CLEANUP, SiteCleanup)
  - **URL:** (URL\_NEW, list(URL)), (URL\_STATUS, list(URLStatus)), (URL\_UPDATE, list(URLOpdate)), (URL\_DELETE, list(URLDelete)), (URL\_FETCH, list(URLFetch)), (URL\_CLEANUP, list(URLCleanup)), (URL\_CONTENT, list(URLContent)), (URL\_PURGE, list(URLPurge)), (URL\_AGE, list(URLAge))
3. Step depends on operation code and object type. Object types and operation codes:
  - **SITE\_NEW::Site** – Insert new site operation. Steps:
    - i. Check is site already exists by Id, if yes – return error and move to step ix If no – move to step ii.
    - ii. Insert record in to the table `dc\_sites`.`sites`.
    - iii. Create table `dc\_urls`.`urls\_SITE\_ID\_MD5`.
    - iv. For each item of Site.urls insert record in to the `dc\_sites`.`sites\_urls` table.
    - v. For each item of Site.urls insert record in to the `dc\_urls`.`urls\_SITE\_ID\_MD5` table.
    - vi. For each Site.filters item insert record in to the `dc\_sites`.`filters` table.
    - vii. For each Site.properties item insert record in to the `dc\_sites`.`properties` table.
    - viii. If 'Content storage types: 0' in hce-node-bundle/api/python/ini/db-task.ini then move to ix, if 'Content storage types: 1' then skip ix and x, move to xi
    - ix. Create file of sqlite db by copy from template file:
      - a) CONFIG\_INI\_DIR/db-task\_keyvalue\_template.db to CONFIG\_DB\_DIR/SITE\_ID\_MD5.db
    - x. Create directory CONFIG\_DATA\_DIR/SITE\_ID\_MD5
    - xi. Create table `dc\_contents`.`contents\_SITE\_ID\_MD5`.
    - xii. Create GeneralResponse object and fill error code and states if some errors.
    - xiii. Print serialized object to stdout.
  - **SITE\_UPDATE::SiteUpdate** – updates fields of site table and correspondent records of related tables in the `dc\_sites` database.
    - i. Check is site already exists by Id, if yes – return error 2020 and move to step iv, If no – move to step ii
    - ii. If site property UPDATE\_NOT\_INSERT\_ROOT\_URLS exists and !=0 then delete root URLs from dc\_urls.urls\_SITE\_ID\_MD5 and copy records from dc\_sites.sites\_urls to dc\_urls.urls\_SITE\_ID\_MD5. Otherwise skip this step

- iii. If field of the SiteUpdate object is not None and it is single value – then update it with new value determinated in base json. If it is container (urls, filters, properties) than behavior depends on the SiteUpdate.updateType field.
    - a) If it is UPDATE\_TYPE\_APPEND (updateType=0) than correspondent Site\_Id records need to be added to tables `dc\_sites`. `sites\_filters`, and so on
    - b) If it is UPDATE\_TYPE\_OVERWRITE (updateType=1) than correspondent Site\_Id records need to be removed and insert new to tables `dc\_sites`. `sites\_filters` and so on
    - c) If it is UPDATE\_TYPE\_UPDATE (updateType=2) then update existing fields with a new value
  - iv. If SiteUpdate.state=STATE\_RESTART (dc\_sites`. `sites`.State = 6) than dc\_sites.sites\_urls records need to be copied to the dc\_urls.urls\_SITE\_ID\_MD5 table with default fields values for not matched fields; `dc\_sites`. `sites`.State value need to be set as STATE\_ACTIVE, also, fields UDate=NOW, Iterations=Iterations+1, ErrorMask=0 and Errors=0.
  - v. Create GeneralResponse object and fill error code and statuses if some errors.
  - vi. Print serialized object to stdout.
- **SITE\_STATUS::SiteStatus** – used to get registered Site object status. Returns the Site object filled with values from DB tables. Steps:
  - i. Create new Site object.
  - ii. Select correspondent record from the `dc\_sites`. `sites` table by SiteStatus.siteId=`dc\_sites`. `sites`. `Id` condition. If record not found, filled Site.State=STATE\_NOT\_FOUND value and go to last step.
  - iii. Fill correspondent Site object single values fields with values from db record.
  - iv. Fill correspondent Site object container fields (urls, filters, properties) with values from correspondent tables by <TABLE>.Site\_Id=SiteStatus.siteId condition.
  - v. Print serialized Site object to stdout.
- **SITE\_FIND::SiteFind** – used to find registered site and return the Site object filled with values from DB tables
  - i. Create new Site object
  - ii. Select correspondent record from the `dc\_sites`. `sites` table by SiteStatus.siteId=`dc\_sites`. `sites`. `Id` condition. If record not found, return empty itemObject
  - iii. Fill correspondent Site object single values fields with values from db record.
  - iv. Fill correspondent Site object container fields (urls, filters, properties) with values from correspondent tables by <TABLE>
  - v. Create GeneralResponse object and fill error code and statuses if some errors.
  - vi. Print serialized object to stdout.
- **SITE\_DELETE::SiteDelete** – used to delete registered site and remove all data from local storages. Steps:
  - i. Delete correspondent record from the `dc\_sites`. `sites` table by SiteStatus.siteId=`dc\_sites`. `sites`. `Id` condition.



- ii. Delete correspondent records from `dc\_sites`.`sites\_urls`,  
`dc\_sites`.`sites\_filters` and the `dc\_sites`.`sites\_properties` tables by Site\_Id.
  - iii. Drop table `dc\_urls`.`urls\_SITE\_ID\_MD5`.
  - iv. Drop table `dc\_urls\_deleted`.`urls\_SITE\_ID\_MD5`.
  - v. If "delayedType" in base json is not determinate – then it default value is "1",  
move to vi; if "delayedType" value is "0" then skip steps vi-ix and move to x
  - vi. If 'Content storage type: 0' in hce-node-bundle/api/python/ini/db-task.ini then  
move to vi, if 'Content storage type: 1' then skip vi and vii; move to viii
  - vii. Delete file of sqlite db by name:
    - a) CONFIG\_DB\_DIR/SITE\_ID\_MD5.db
  - viii. Delete directory by name:
    - a) CONFIG\_DATA\_DIR/SITE\_ID\_MD5
  - ix. Drop table `dc\_contents`.`contents\_SITE\_ID\_MD5`.
  - x. Create GeneralResponse object and fill error code and statuses if some errors.
  - xi. Print serialized object to stdout.
- **SITE\_CLEANUP::SiteCleanup** – used to delete data of registered site from local storages.  
Steps:
  - i. If "delayedType" in base json is not determinate – then it default value is "1",  
move to ii; if "delayedType" value is "0" then skip ii, iii and move to iv
  - ii. Create if not exist table `dc\_urls\_deleted`.`urls\_SITE\_ID\_MD5`.
  - iii. If parameter "saveRootUrls" in base json = 0 then copy all records from  
`dc\_urls.urls\_SITE\_ID\_MD5` to `dc\_urls\_deleted`.`urls\_SITE\_ID\_MD5`; if  
saveRootUrls=1 then copy only non-root records from  
`dc\_urls.urls\_SITE\_ID\_MD5` to `dc\_urls\_deleted`.`urls\_SITE\_ID\_MD5`
  - iv. Parameter "saveRootUrls" in base json = 1 then delete records with not empty  
ParentMD5 from `dc\_urls.urls\_SITE\_ID\_MD5` (delete only non-root URLs) and  
skip step v; if it is 0 then move to step vi
  - v. Truncate table `dc\_urls`.`urls\_SITE\_ID\_MD5`.
  - vi. If "delayedType" in base json is not determinate – then it default value is "1",  
move to vii; if "delayedType" value is "0" then skip vii, viii, ix and move to x
  - vii. If 'Content storage type: 0' in hce-node-bundle/api/python/ini/db-task.ini then  
move to iv, if 'Content storage type: 1' then skip iv; move to v
  - viii. Delete directory CONFIG\_DATA\_DIR/SITE\_ID\_MD5 and sqlite db by name  
CONFIG\_INI\_DIR/db-task\_keyvalue\_template.db
  - ix. Drop table `dc\_contents`.`contents\_SITE\_ID\_MD5`.
  - x. Update correspondent `dc\_sites`.`sites` record with fields: TcDate=NOW,  
Resources=0, Iterations=0, State=STATE\_ACTIVE, ErrorMask=0, Errors=0.
  - xi. If parameter "moveUrls" in base json = 1 then copy records form  
`dc.sites`.`sites.url` to `dc\_urls`.`urls\_SITE\_ID\_MD5`; if it is 0 then skip this step  
and move to viii
  - xii. Create GeneralResponse object and fill error code and statuses if some errors.
  - xiii. Print serialized object to stdout.
- **URL\_NEW::list(URL)** – used to add list URLs to the system for sites or as not linked with  
sites, or with creation of new site if not found by Id or canonized domain name. Steps:
  - i. Create GeneralResponse object.
  - ii. For each URL object from list set fields that depends from the site if has None  
value (RequestDelay, HTTPTimeout, URLType):

- a) If `URL.siteId!=""` – select record from the ``dc_sites`.`sites`` table;
      - If record not found then behavior process step iii and get fields value from correspondent record;
      - If record found – get fields values from this record and process step iv.
    - b) If `URL.siteId==""` or `URL.siteId==None` then process step iii.
  - iii. Depend on the `URL.siteSelect` value:
    - a) If `SITE_SELECT_TYPE_EXPLICIT`:
      - If record not found or `URL.siteId==""` then insert record in to the ``dc_urls`.`urls_0`` table.
    - b) If `SITE_SELECT_TYPE_AUTO`:
      - Qualify `URL.URL` (use the: `scheme+netloc+"\"` operation to generate canonic url form), generate `SITE_ID_MD5` from qualified url and try to resolve site in ``dc_sites`.`sites`` table. If resolved, repeat step ii.
      - If site not resolved, create new site with fields values from `dc.EventObjects.Site` object and insert URL to the newly created table in ``dc_urls`.`urls_SITE_ID_MD5`` table. As single root URL use the `scheme+netloc+"\"` canonization and initialize new Site object with this url.
    - c) If `SITE_SELECT_TYPE_QUALIFY_URL`:
      - Qualify `URL.URL`, generate `SITE_ID_MD5` and try to resolve site in ``dc_sites`.`sites`` table. If not resolved, put URL to general DB table ``dc_urls`.`urls_0``.
  - iv. Check is URL exists for this site and if not – insert URL in to the correspondent urls table. If URL already exists set `GeneralResponse.states` item value as 2. If some db error – set it as 1.
  - v. Insert item in to the `GeneralResponse.states` with correspondent value.
  - vi. Print serialized results list object to stdout.
- **URL\_STATUS::list(URLStatus)** – used to get status of list of URLs. Steps:
  - i. Create results list.
  - ii. For each `URLStatus` object in the list: create new URL object, select record from table ``dc_sites`.`urls_URLStatus.siteId``, fill correspondent fields, put Site object in to the results list. The record selected by `URLStatus.url` field. If the `URLStatus.urlType==URL_TYPE_MD5` the ``dc_sites`.`urls_URLStatus.siteId`.`URLMd5`=URLStatus.url` condition used, else – the ``dc_sites`.`urls_URLStatus.siteId`.`URL`=URLStatus.url` condition used.
  - iii. Insert URL object in to the list filled with correspondent field values.
  - iv. Print serialized results list object to stdout.
- **URL\_UPDATE::list(URLOpdate)** – used to update list of URLs with new values. Only fields with not None values are updated. Steps:
  - i. Create `GeneralResponse` list.
  - ii. For each `URLOpdate` objects in list using `URLOpdate.url==<TABLE>.URLMD5` condition:
    - a) If `URL.siteId==""` – update record in the ``dc_urls`.`urls_0`` table;
    - b) If `URL.siteId==""` – update record in the ``dc_urls`.`urls_SITE_ID_MD5`` table.

- iii. Insert the `GeneralResponse.statuses` item with correspondent value.
- iv. Print serialized object to stdout.
- **URL\_DELETE::list(URLDelete)** – used to delete record correspondent to URL and data in file system and key-value DB. Steps:
  - i. Create `GeneralResponse` list.
  - ii. For each `URLDelete` objects in list:
    - a) If `URLDelete.siteId!=""`
      - delete record in the ``dc_urls`.`urls_SITE_ID_MD5`` by `URLDelete.siteId` with ``dc_urls`.`urls_SITE_ID_MD5`.`URLMd5`=URLDelete.url` condition;
      - Delete item in the key-value db file: `CONFIG_DB_DIR/"URLDelete.siteId".db` by `URLDelete.url` key;
      - Delete files in directory by mask: `CONFIG_DATA_DIR/URLDelete.siteId/PathMaker(URLDelete.url).getDir()/URLDelete.url*`
    - b) If `URLDelete.siteId==""`
      - delete record in the ``dc_urls`.`urls_0`` by `URLDelete.siteId` with ``dc_urls`.`urls_0`.`URLMd5`=URLDelete.url` condition;
      - delete item in the key-value db file: `CONFIG_DB_DIR/0.db` by `URLDelete.url` key;
      - delete files in directory by mask: `CONFIG_DATA_DIR/0/PathMaker(URLDelete.url).getDir()/URLDelete.url*`
  - iii. Decrement ``dc_sites`.`sites`.Resources` and ``dc_sites`.`sites`.Contents` in case of each record are deleted from ``dc_urls`.`urls_SITE_ID_MD5`` and from key-value db.
  - iv. Insert the `GeneralResponse.statuses` item with correspondent value.
  - v. Print serialized object to stdout.
- **URL\_CLEANUP::list(URLCleanup)** – used to update record correspondent to URL with value `State=0`, `Status=1` and delete raw data in file system and key-value DB. Steps:
  - i. Create `GeneralResponse` list.
  - ii. For each `URLCleanup` objects in list:
    - a) If `URLCleanup.siteId!=""`
      - update fields "state" and "status" of record record in the ``dc_urls`.`urls_SITE_ID_MD5`` by `URLCleanup.siteId` with ``dc_urls`.`urls_SITE_ID_MD5`.`URLMd5`=URLCleanup.url` condition;
      - delete item in the key-value db file: `CONFIG_DB_DIR/"URLDelete.siteId".db` by `URLDelete.url` key; delete files in directory by mask: `CONFIG_DATA_DIR/URLDelete.siteId/PathMaker(URLDelete.url).getDir()/URLDelete.url*`

- b) If `URLDelete.siteId==""` – update fields “state” and “status” of record in the ``dc_urls`.`urls_0`` table by `URLDelete.url=`dc_sites`.`sites`.`URLMD5`` condition; delete item in the key-value db file:
- `CONFIG_DB_DIR/0.db`
- by `URLDelete.url` key; delete files in directory by mask:
- `CONFIG_DATA_DIR/0/PathMaker(URLDelete.url).getDir()/URLDelete.url*`
- Set the `GeneralResponse.statuses` item with correspondent value.
- Print serialized object to stdout.
- URLFetch::list(URLFetch)** – used to fetch urls from correspondent tables. Steps:
- Create results list.
- If `URLFetch.sitesList` is empty or `None`, select all records from the ``dc_sites`.`sites`` table with condition `state="Active"` and fill `URLFetch.sitesList` with `Id` value.
- For each item in the `URLFetch.sitesList`:
- a) Select records from table ``dc_urls`.`urls_%URLFetch.siteId[i]`` with conditions using `URLFetch.criterions` items.
  - b) Create new `URL` object and fill fields’ values from selected record.
  - c) Insert `URL` object in to the results list
- Print serialized object to stdout.
- URLContentRequest::list(URLContentRequest)** – used to fetch resources corresponded to `URLContentRequest.siteId`. Steps:
- Create results list object.
- For each `URLContentRequest` object in list:
- a) Create new `URLContentResponse` object.
  - b) If `URLContentRequest.contentTypeMask` has `CONTENT_TYPE_PROCESSED` bit set – return the processed content from correspondent key-value db:
    - If `URLContentRequest.siteId!= ""` than use db file: `CONFIG_DB_DIR/"URLContentRequest.siteId".db`
    - If `URLContentRequest.siteId== ""` than use db file: `CONFIG_DB_DIR/0.db`
- and the `md5(URLContentRequest.url)` as key. Create new `Content` object, fill fields values and insert it in to the `URLContentResponse.processedContents` as item.
- c) If `URLContentRequest.contentTypeMask` has `CONTENT_TYPE_RAW_LAST`, `CONTENT_TYPE_RAW_FIRST`, `CONTENT_TYPE_RAW_ALL` bit(s) set – return correspondent raw content(s) from file by mask:
    - If `URLContentRequest.siteId!= ""` than use directory: `CONFIG_DATA_DIR/URLContentRequest.siteId/PathMaker(URLContentRequest.urlMd5).getDir()/`
    - If `URLContentRequest.siteId== ""` than use directory: `CONFIG_DATA_DIR/0/PathMaker(URLContentRequest.urlMd5).getDir()/`

Read content of correspondent raw data file, create new Content object, fill fields values and insert it in to the `URLContentResponse.rawContents` as item.

- iii. Insert `URLContentResponse` object in to the results list.
  - iv. Print serialized results list object to stdout.
4. Finish with zero exit status value in case of no critical errors and with 1 in another case.

## Sites manager

The main client interface for Site and URL objects operations

## Application configuration and start

### *Threaded classes instantiation sequence*

The threaded classes use inproc server and client connections. Inproc connections need to have server-side ready (bind and listen) to make connect operation from client side. This condition requires some instantiation sequence for threaded classes because they make connections at this time. The instantiation sequence:

1. AdminInterfaceServer.
2. URLManager.
3. SitesManager.
4. BatchTasksManager.
5. ClientInterfaceService.

## Command line arguments

### *DCDaemon dc-daemon.py*

- `--config, -c` – configuration file path name. Mandatory parameter, if application can't load config file it exits(1) with corresponding error message on console, no demonization, no application functionality started. If omitted application tries to open file with the same name as executable module and `".ini"` extension in sequence:
  - in the current directory,
  - in the `/etc/<application_binary_name>/` directory,
- `--help, -h` – displays brief usage information, including possible command line arguments names and values and exits(0).
- `--name, -n` – application instance name. Optional parameter, if omitted the same name as executable module file used. Overrides the same parameter in the configuration file "Application" section.

### *DCClient dc-client.py*

- `--config, -c` – the same as for the "dtmd" application.
- `--help, -h` – the same as for the "dtmd" application.
- `--command, -cmd` – task action, specifies the requested task operation, possible values {"SITE\_NEW", "SITE\_UPDATE", "SITE\_STATUS", "SITE\_DELETE", "SITE\_CLEANUP", "URL\_NEW", "URL\_STATUS", "URL\_UPDATE", "URL\_FETCH", "URL\_DELETE", "URL\_CLEANUP", "URL\_CONTENT"}. Optional if `--help` is specified, and mandatory – if not. A result of any action is file in json format. Structure corresponds with DRCE Functional Object response specification.
- `--file, -f` – data file for action, specified by `--task` option. Mandatory if `--task` is specified. The data file in json format, structure defined by DRCE functional object request protocol specification.
- `--fields, ff` – fields set json string that will be used after target object created from pre-created json file (`--file` parameter) or from regular construction as empty filled with None values all fields. More detailed see the `DC_public_client_API` document.

- `--merge` – specifies merge or not results received from different hosts in multi-host distributed cluster. By default if not specified by user the Boolean true or greater than zero value that is means – merge. Boolean false or zero value – means not merge and return the array of results from each node.

### *DCAdmin dc-admin.py*

- `--config, -c` – the same as for the “dtmd” application.
- `--help, -h` – the same as for the “dtmd” application.
- `--cmd` – command name, specifies the requested admin management operation. Possible values are {“STAT”, “SET”, “GET”, “STOP”}. “STOP” – lead to sequentially stop all threaded classes of application in opposite start sequence (specified in the “Server” section by the “instantiateSequence” option of configuration file) and application exit(0).
- `--fields` – comma separated fields list for requested operation. Mandatory, if “--cmd” is set and is not a “STOP” or “STAT” value. For “STAT” – specifies list of statistical indicators for class handler. If empty or not set – the empty list of fields used. For “GET” – specifies names of configuration options for all threaded classes and application that can be changed at runtime. For “SET” – specifies list of configuration options name:value pairs for all threaded classes and application that can be changed at runtime. All commands returned json file content. The structure depends on concrete request.
- `--classes` – comma separated list of threaded classes that will be used to execute request. Always optional for “STOP” command (always executed for all) and mandatory for all another commands. For “STOP” if omitted or empty – the list specified by the “instantiateSequence” option in “Server” section of configuration file.