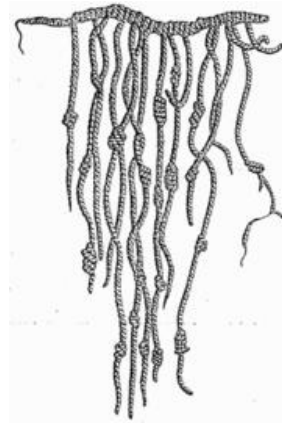




DATA WAREHOUSE MANAGEMENT

WHITEPAPER QUIPU version 1.1



OPEN SOURCE DATAWAREHOUSING

Contents

Introduction	4
QUIPU.....	4
History	4
Data warehouse architecture	6
Starting points.....	6
Modeling.....	7
Core components of QUIPU	7
Enterprise DW Architecture.....	9
Data Vault	11
Basics of the Data Vault	11
Hubs.....	11
Hub Satellites	12
Links.....	13
Link Satellites.....	14
Reference tables.....	14
Historical Data Archive.....	15
Model Generation	17
Data Vault Generation.....	17
Data Vault generation step 1: Import model	18
Data Vault generation step 2: Generate Data Vault proposal ('Analyze staging')	18
Data Vault generation step 3: Evaluate and modify generated Data Vault proposal ('Prepare staging')	20
Data Vault generation step 4: Regenerate Data Vault proposal	21
Data Vault generation step 5: Modify / add metadata.....	22
Script generation.....	22
QUIPU Technology	25
Client Server Architectuur	25
Back-end: Groovy / Grails.....	26
Front-end: C# / .Net.....	26
Database support.....	27
QUIPU deployment options.....	27
Future developments	29
QUIPU OS functionality	29
Commercial modules.....	29

More information..... 33

- About QOSQO33
- Licence model33
- Source code availability33

Introduction

QUIPU

This Document is not a manual

This document delivers some background information on QUIPU, the Open Source data warehouse generation software developed by QOSQO. It aims at creating understanding of the position of QUIPU in the BI landscape and development cycle. It is *not* a manual, and does not replace release notes. The document may not always be 100% in sync with the latest release of the software, so please refer to the release notes for definitive information.

Version 1.1 release provides only core functionality

The Version 1.0 release of QUIPU offered only a subset of all the functions and features we think are essential, important or just pretty cool to have in such a product. But still we feel this release has proven to be a time saver for many data warehouse architects and has assisted them in developing reliable data warehouses based on solid modeling principles.

In the course of the first year after its original release in July 2010 several upgrades have been made to bring the current version to 1.1. These developments have added some important features to QUIPU and also affected some of the views expressed in this document. Hence an update to the whitepaper.

Open Source

After the initial release of QUIPU in the Open Source, QUIPU has found its way to its intended users and organizations small and large many of which have reported to benefit from it. We continue to invite the feedback from the user community, as it will be the basis for further development from our side. We may have already recognized the need for certain features, but we would love to be surprised with unexpected use cases that require features currently not present.

Community

We also hope and expect that gifted developers will step forward and pick up some of the work, willing to share with the world. We will actively support any developer that wishes to contribute with documentation and guidance. More information to this can be found in the user documentation or on sourceforge (<http://sourceforge.net/projects/quipu/>)

Commercial modules

In the past year it became apparent that more advanced functionality is required to satisfy the needs of our larger customers. We have decided to develop some of these functions as paid modules to complement the core functionality in the Open Source application. This paper describes changes to the basic architecture of QUIPU to allow add-on modules and explains what type of future developments will remain Open Source and what functionality we will develop as commercial modules.

Work in progress

Clearly the product is not finished. This document too should be considered as 'work in progress'. Expect frequent updates of the software and this document, as we implement new features and fix bugs found in active use.

History

Data Vault modeling

Consultants of our sister company Nippur as well as from QOSQO have been actively participating in data warehouse and business intelligence projects for many years. They have contributed or managed projects with large multinational companies based in the Netherlands as well as abroad. They have used the Data Vault modeling technique for many years now as best practice for mod-

From best practice to software development

eling data in a way that fits business needs best.

The Data Vault model is particularly strong in ensuring lineage of data warehouse data back to their source systems, thus enabling a complete audit trail. At the same time the Data Vault model is very easy to adapt to changing requirements and extremely scalable. The strong architectural foundation of Data Vault and the simplicity of its design open the way to automation. Our consultants have created quite a few customer specific generator scripts to create and populate a data warehouse.

Early automation efforts

In some cases we created tailor-made fully automated data warehouse management solutions. Complete source and target models, ETL code specifications, scheduling and documentation are managed and maintained using a single repository. Our lead architect of QUIPU was involved in the design and realization of one such effort of a large Netherlands based Bank / Insurance company, which started around 2001. Currently, QOSQO is responsible for the support of this software.

Development of QUIPU

This specific solution, combined with our best practices at other customers sparked the idea to develop a data warehouse management software solution back in 2008.

QUIPU: an ancient Inca recording device

The Quipu or khipu (sometimes called talking knots) was a recording device used in the Inca Empire and its predecessor societies in the Andean region. A quipu usually consisted of colored spun and plied thread or strings from llama or alpaca hair. It could also be made of cotton cords. The cords contained numeric and other values encoded by knots in a base ten positional system. Quipus might have just a few or up to 2,000 cords.



Quipucamayocs (Quechua khipu kamayuq, "khipu-authority"), the accountants of Tawantinsuyu, created and deciphered the quipu knots. Quipucamayocs could carry out basic arithmetic operations such as addition, subtraction, multiplication and division. They kept track of mita, a form of taxation. The Quipucamayocs also tracked the type of labor being performed, maintained a record of economic output, and ran a census that counted everyone from infants to "old blind men over 80." The system was also used to keep track of the calendar.

[source and images: wikipedia,<http://en.wikipedia.org/wiki/quipu>]

Data warehouse architecture

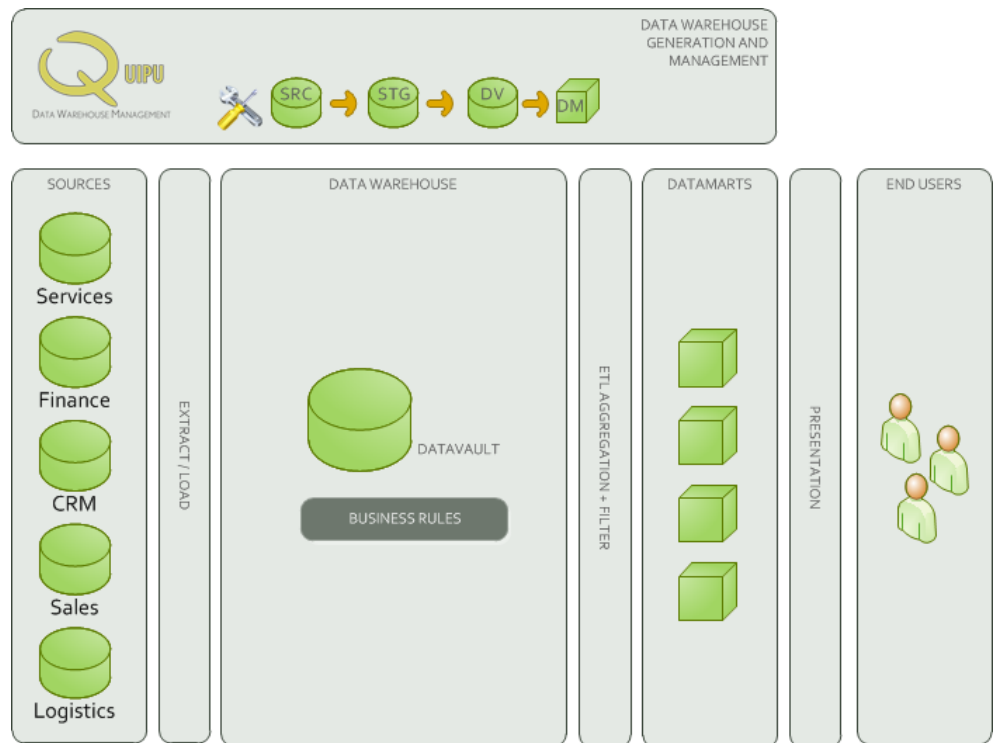
Starting points

Layered data warehouse architecture

'Classic' data warehouse (DW) architectures consist of several layers. The first one is the source layer, representing the data to be used for reporting and analysis (business intelligence) purposes. The second layer is the actual data warehouse where selected (or in rare cases: all) source data is offloaded from source systems and collected in the 'memory' of the organization. Subsets of data are then created to satisfy needs of specific functional areas or domains in the third layer. This layer is called the data mart layer. The fourth layer presents the relevant information to end users, the presentation layer.

QUIPU generates models and load code

QUIPU aims at fulfilling the generic core components in the data warehouse and data mart layers of the architecture. It is particularly strong in generating, maintaining and populating (database) structures to store historical data, both transactional data and reference data. Apart from the data structures QUIPU also generates the (SQL) code to load the data correctly in the data warehouse. QUIPU is designed to operate in a multitude of environments.



NOTE: QUIPU generates the basis Data Vault structures and ETL to load the data. QUIPU does *NOT* provide a solution to the business rules that are often required to integrate data from different sources or derive performance indicators from the 'raw' transaction data. As such QUIPU does not replace ETL tooling or rules engines but rather complements them.

Modeling

Data Vault - efficient, expandable historical record

QUIPU is based on the Data Vault model. See <http://danlinstedt.com/about/> for more information and background.

The key strengths of this model are:

- it allows capturing an historical record of any data ('transaction' as well as 'master' data)
- it is optimized for (near) real-time updates (allows massive parallel processing, late arriving transactions etc.)
- it isolates changes in the model and is thus flexible (can be adapted to changing business needs with relative ease)

Reasons for a Data Mart layer

QUIPU typically creates more tables than the original 'point-in-time' 3NF source models. As such it creates a model quite opposite of a star-schema and not particularly strong in query-performance on the dominant RDBMS systems (Oracle, MS SQLServer, DB2, MySQL, etc.).

This is the technical reason why on most platforms a Data Mart layer is required.

The data is stored in 'raw' (source) format in the data warehouse, so all of the original data is available. This raw format of the data is good for operational reporting, but not optimal for other reporting and analysis purposes, as there are no derived data (e.g. KPI's) available and integration of data from different sources can be hard.

This is the business reason for creating datamarts.

Summarizing: the Data Mart layer contains a subset of the data and is optimized for query performance: data is integrated, aggregated and pre-calculated, often stored in a star-schema.

Core components of QUIPU

No business rules in QUIPU

As QUIPU's added value lies with the generic modules, it implies all *solution specific* modules (mainly the data transformations that may be required) will have to be realized outside of the QUIPU product.

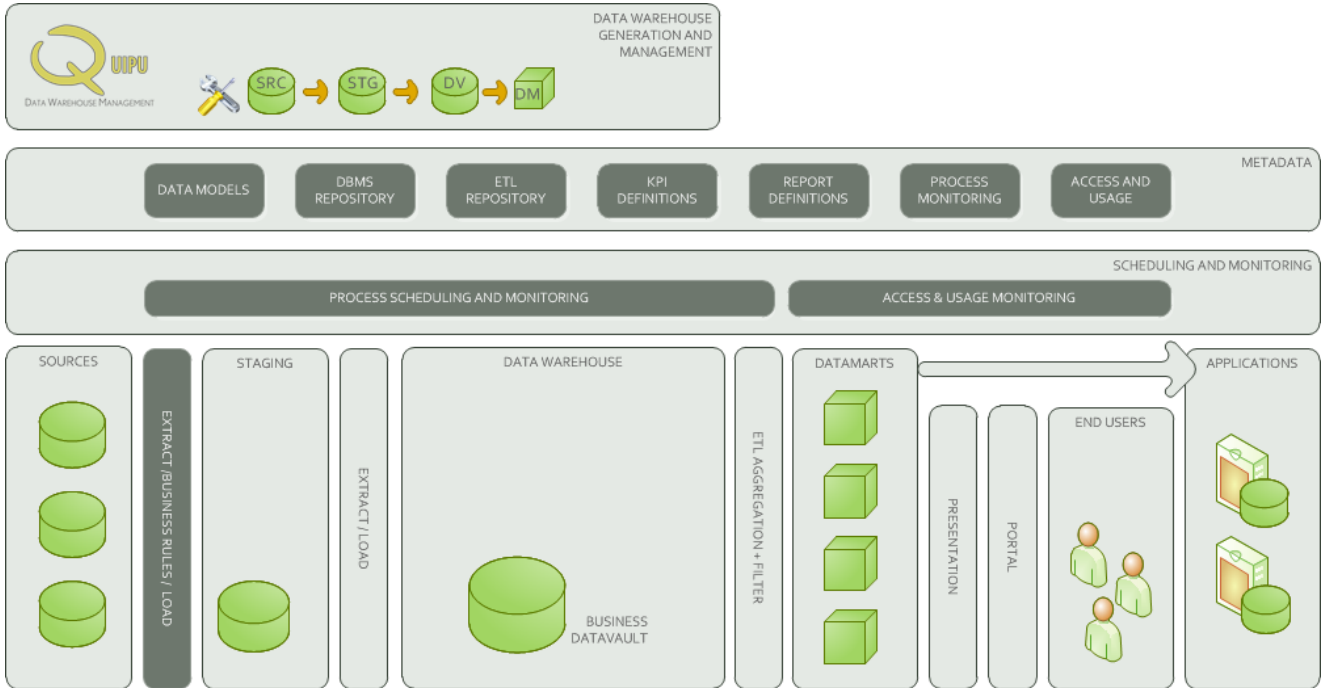
In most situations:

- the raw source data is not fit to fulfil the needs of the business. Sometimes complex business rules are required to integrate data from different sources (e.g. match external customer or supplier data to internal data structures).
- hierarchies are added to the source data structures, and complex KPI's are calculated for different levels of the hierarchy.

All of these derivations need to be implemented outside QUIPU, either using ETL tools or BI front-end tools.

Business Data Warehouse

In its most basic form QUIPU controls a single, integrated data warehouse based on a single, harmonized business model. The business model provided to QUIPU can be completely time-variant -capturing changing reference data- and thus describe the business entities and their relations as they are defined at any point in time.



Business Data Warehouse

3NF business data model

The business data model should be in 3NF (third normal form) and can be created using one of many database modeling tools. The business model cannot be imported in QUIPU directly (in the current version). The easiest way to get the business data model into QUIPU is by physically implementing the -empty- database in one of the many RDBMS'es supported by QUIPU and using QUIPU's reverse engineering function.

Once the model is available in QUIPU, the following functions are available:

Generate Staging area

- Generate Staging area. The source data will have to be delivered in a format that is more or less dictated by the business model into a staging area. The tables are generated by QUIPU, as well as template SQL statements that can fill the staging area from the source tables - assuming these tables exist physically in the same RDBMS system-. It is the task of an external ETL tool to load the data in the staging area if the source does not exist in this exact form or lives in a different database platform.

Generate Data Vault model

- Generate Data Vault model. A simple straightforward translation will be generated through analyzing the source model's tables, keys and relations. This model can be optimized through direct interactive manipulation (or in future versions via export and import functions).

Generate Load functions

- Generate Load Functions. The data delivered into the staging area are loaded into the Data Vault tables. Extensive validation and logging assures that the data warehouse always retains a valid state. All data elements in a load are marked with a unique load id (or alternatively a source id and timestamp), allowing various housekeeping functions on the data loaded (e.g. archiving data or removing faulty loads). QUIPU also generates code to load the staging area, but this should be regarded as demo code as in practice the staging area is often loaded

Generate views

directly from the source systems. Generation is based on templates that can be adapted to exploit the features of various database management systems.

- Generate various views on the Data Vault. As a first step towards supporting BI, QUIPU generates views that reproduce the 3NF source model (thus hiding all the complexities of the Data Vault) for query purposes. As the Data Vault stores historical data that may not be present in the original 3NF model some views deliver the data *as it currently is available in the source* (current values). A sample query is also provided, that demonstrates how a specific time and date can be specified so the data warehouse can provide the exact data that was available in the source at that specific moment (a 'point-in-time view').

NOTE Data Marts are required to deliver data to the business. Functionality to assist the construction of Data Marts is not yet available in the Version 1.1 release of QUIPU, but currently in Beta, and scheduled to be released as part of Version 1.2.

Enterprise DW Architecture

Enterprise Requirements

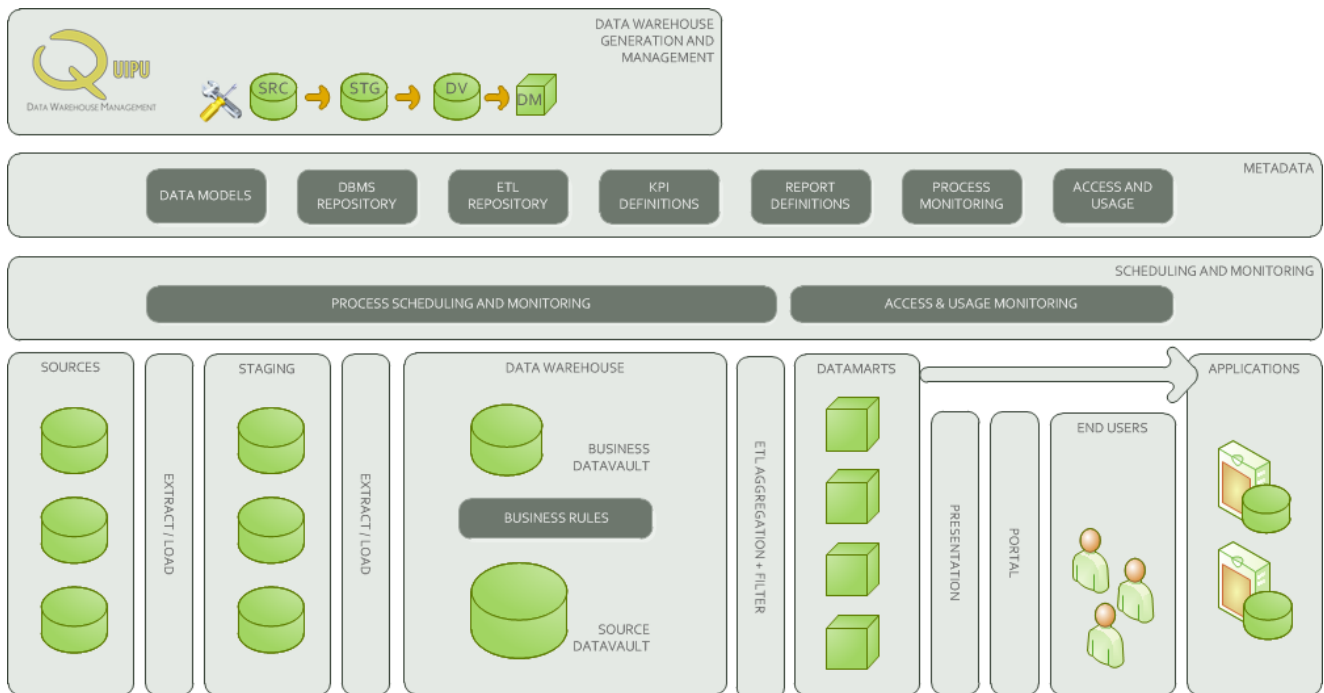
QUIPU is designed to support the Enterprise Data Warehouse concept. It can certainly be used in less complex environments but is designed with the requirements of larger enterprises in mind. Thus, it is capable of capturing and integrating information simultaneously from multiple sources and scaling to large data volumes.

BI stack

QUIPU can integrate with Open Source or commercial products, to offer a complete solution. The solution then includes a modeling tool, ETL and a BI front-end stack.

Near-real-time

QUIPU is aware of the shifting role of the DW in today's businesses. More and more the DW is considered a key component in the information architecture, supporting not only the strategic information requirements of senior management, but also the tactical and even operational requirements at the shop floor. The latter require frequent data updates (near real-time).



Extended DW Architecture

Implement a multi-layered version

With QUIPU it is possible to implement a multi-layered version of the Data Vault.

Thus QUIPU supports the creation of data marts on the Business Data Vault, but also directly on the Source Data Vault.

'Source Data Warehouse'

The 'raw' source data is captured in a 'source Data Warehouse'. These are generated as Data Vaults and are directly derived from the (logical) data models of the sources. They are ideal for operational, single-source reporting and analysis.

Where business keys can be recognized they can be the first step towards building the 'business Data Warehouse' as they can bridge data from multiple source systems.

In QUIPU hubs are separated physically and logically. This means shared hubs need to be created virtually, by defining views over multiple physical hubs. Future versions of QUIPU will support shared hubs more transparently.

'Business Data Warehouse'

The business needs a unified model to report from: the 'business Data Warehouse'. It is built for integrated, business wide reporting and analysis.

The Business Data Warehouse complements the data in the Source Data Warehouse with integrated data (through cleansing and mapping algorithms) and derived data (calculated using business rules). It is based on a consistent data model, specifically designed to support the business.

In most cases the integration requires mapping of business keys and transformation of the source data. Although it would be possible to distribute the transformed data directly into data marts, this is often not desirable:

- if the transformed data is considered business data in its own right: even when the transformation rules change, the data needs to remain available as it was historically derived (transformed).
- if transformations need to be implemented only once in order to safe-guard consistency. It can be difficult to enforce that the same business rule is correctly applied when building multiple data marts. It is often much easier to achieve consistency when the result is stored in a second Data Warehouse layer: the Business Data Warehouse.

Data Vault

Basics of the Data Vault

The Data Vault concept has been developed by Dan Linstedt. A good starting point for learning more on the Data Vault is the dedicated wiki page http://en.wikipedia.org/wiki/Data_Vault_Modeling.

The Data Vault structure consists of three basic table types: Hubs, Links and Satellites. To satisfy specific requirements QUIPU additionally supports Reference tables.

Hubs

List of business keys

Hubs contain a simple list of *business keys*. They *must* be meaningful for the business and *uniquely identify* an entity (or object) recognized by the business (like customer, product, supplier, invoice, etc.).

Sometimes -ideally- these business keys can be found in source systems as primary keys in databases, or unique indices.

The business keys - and thus the hubs- are the anchor points of the entire model!

Business keys must be unique

It is important to stress the requirement for *global and historical uniqueness* of the business key.

This means that the business key should not be used in different parts of the company to mean different things (often the case in larger organizations, where e.g. *different* customers with the *same key* may be created in different countries).

It also means business keys *should never be reused* (often seen in older systems, where a limited set of keys is re-used after a certain period of time). In both cases the Data Vault will treat the data it stores from these different business entities as belonging to the same entity. This will confuse users of the data warehouse, as it probably confuses the business in the operational processes.

So it is extremely important to find and analyze the business keys before constructing the Data Vault. Analysis should include inspection of source data (data profiling) and verification with business domain experts.

Inserts only

The hub does *not* contain any other business data. As the business keys themselves do not change, the hub is not timevariant. Records will *never ever be updated or deleted*. The table will only see record insertions (which guarantees a very efficient loading process).

Hub structure

The hubs contain:

- the business key
- the source system that first reported the business key to the data warehouse
- the timestamp when it was first recorded

- optionally a meaningless id (may perform better when business keys are compounded and/or very large)
- some technical fields (like the meaningless ID generated by the data warehouse if chosen and an audit id identifying the process that created the record).

Hub Satellites

At least one satellite per hub per source

When the hub only contains the business keys, the relevant data *belonging to* a business key is stored in the satellites. Each hub has at least one satellite that records at least the business validity of the business key.

Source system clones can feed into a single satellite

As a principle a separate satellite for each source system should be created, recording the data that this specific source holds on the business key, over time. Each of these satellites contains the 'voided' attribute, that indicates the business key is no longer valid in a particular source system.

Generally a single implementation of a transaction system should be considered a source system. So an organization that has implemented an ERP system for each of its division should consider each as a separate source.

However, under certain conditions it may be possible -and then often also desirable- to consider multiple transaction systems as a single source. This is the case when:

- all systems share the same logical and technical (data) design (multiple clones of the same system)
- master data management is implemented across all systems, assuring business keys are uniquely defined and shared across all systems.

Under these conditions all systems can feed into the same source Data Vault, possible greatly reducing the number of tables and load processes and at the same time simplifying querying the data.

Still one needs to realize that this simplification may come at a price. For example parallel loading from multiple source systems becomes more complex or even impossible. There are very likely to be other consequences.

At least one record per business key

In this minimal case the satellite would only contain technical fields, and minimally a single record per business key that is known by this source (creation/first seen).

When the business key is deleted in the source another record is inserted indicating the business key is no longer valid. This demonstrates that in the satellite the history is stacked: new data is recorded by inserting new records with a Start Timestamp.

Hub entries can be generated from transaction tables

Ideally hub entries (Business Keys) are generated when a source table is loaded containing the master data on this Business Key, *before the transactions are processed*. Thus, the information on the BK is available when transactions arrive.

This can not always be guaranteed, mostly because of timing issues (e.g. master data may be loaded only once a day, whereas transactions are loaded as they arrive during the day).

All relevant data elements of a business key

In the Data Vault this is not an issue, as the Hub will be updated when the transaction arrives. The corresponding record in the hub-satellite will be absent until the master data is loaded.

Normally there is more interesting data available on business key. The relevant source data elements are all stored in one or more satellites. By default for each source table containing data related to a business key a separate satellite will be created. This keeps the model close to the source model and facilitates load procedures.

Inserts Only

As changes in the source data are recorded in *new* entries in the satellites also this table can be refreshed with only insert statements.

End-dating

Although logically not required, for performance reasons on specific database management systems it may be necessary to end-date a record when a new record is created in the satellite. An additional End Timestamp can be generated. This is the only reason for a record to be updated in the satellite.

Satellite Structure

The satellites contain:

- the unique key of the hub it refers to (either the business key or meaningless id)
- the 'start' timestamp: from what time is the data valid (usually the data was recorded in the data warehouse)
- optionally an 'end' timestamp: until what time is the data valid.
- the source system that delivered the data to the data warehouse
- some technical fields (like the audit id identifying the process that created the record, the voided flag that marks the data as deleted in the source).

Links

Links record relations between business keys

Business entities interact in many ways, and these interactions are reflected in (logical) relations between business keys in source systems. In the Data Vault these relations are recorded in Link tables.

Examples of these relations are all kind of structures and classifications in the master data (e.g. chart of accounts, product structure, organizational structure, customer classifications, etc.) and of course the facts registered in the source (e.g. a sale of a product to a customer by a sales representative on a specific date).

Similar to the hub the link table only contains the business keys or meaningless id's, when these are preferred. All data linked to the relation are stored in a satellite table (link satellite).

Link structure

The link contains:

- the unique key of the hubs it refers to (either the business key or meaningless id): at least 2.
- the timestamp: when was the data recorded in the data warehouse
- the source system that first reported the relation to the data warehouse
- some technical fields (like the audit id identifying the process that created the record).

Link Satellites

Similar to Hub satellites

Link satellites share all the characteristics of a hub satellite.

The difference is that a link satellite contains data describing a relationship. For links between master data elements (e.g. the organizational structure) these satellites contain little or no data and provide only information on the validity of the relationship.

But when facts are registered in the relation (e.g. a sale) the link satellite will contain all the details of the fact (e.g. number of products sold, the unit price, total value, etc.)

Reference tables

Reference tables can lead to an excessive number of links

A reference table is generally speaking a very simple, rather static table that is often referenced from many other tables. It is often used in source systems to define a restricted value set (domain) for a certain attribute. E.g. the list of currencies, countries, colors available for a product, etc.

By default (and following strict Data Vault modeling principles) a hub and satellite are created for each of these tables, as well as link tables to all hubs that reference these tables.

Querying these tables can then become rather complex, as in all cases the satellites of the reference tables need to be joined in (via the link tables) when the current value is requested for a given business key.

Also the model can become more cluttered due to the abundance of link tables that are generated.

Links are not generated for reference tables

To alleviate this problem it is possible to mark a source table as a reference table. The result is that a hub and satellite will be created for the reference table as normal, but *no link tables* between the reference hub and the referring hubs are generated.

Joins to reference tables must be added in the queries

The reference *to* the reference table must of course still be stored somewhere. This is now a function of the satellite of the referring entity.

Of course the data mart architect will have to understand what reference tables exist, and create the appropriate joins in his queries.

History is kept for reference tables

In the QUIPU implementation of reference tables, with a 'regular' hub and hub-satellite all history of (changes to) the reference table are kept. A further simplification would be to simply overwrite changes to the reference table, as they are very likely to occur only sporadically. Of course dropping history violates a basic concept of the data warehouse (always store all history) and potentially poses an auditing risk.

Reference tables without history are not (directly) supported by QUIPU.

Caution: use with care. The Reference table feature should be used with proper caution, as it breaks one of the key aspects of the Data Vault. In a strict Data Vault all relations are modeled as links. This makes it extremely easy to assess the impact of change. The existence of references from satellites to the hubs of reference tables breaks with this concept. This means changes to reference tables can impact a larger part of the Data Vault model.

Historical Data Archive

'Informal' business keys

In our practice working with Data Vault models for our clients we have learned that is not always practical or even possible to identify the 'true business key' for all business objects. Many systems work internally on technical keys. Business keys are often available, but seldom guarded as vigorously as the technical keys. This leads to a variety of data quality issues with these 'informal' business keys.

There may be any or all of the following problems:

- missing keys (no entry)
- duplicate keys for different objects (e.g. due to re-use of keys)
- multiple keys for the same object (e.g. due to spelling faults)
- missing keys in the master data system (new keys are entered in slave systems)
- expired keys in slave systems

Bad business keys break business processes

Problems with business keys are a problem in operations, as it introduces ambiguity in the business processes. The resulting confusion is the cause of many inefficiencies in the business.

Bad business keys lead to inconsistent BI

But apart from that they seriously reduce the capability of BI to produce reliable integrated reports spanning multiple systems and processes. Reports on the same business facts from different systems will produce different results!

Master Data Management is a pre-requisite

Master Data Management practices need to be introduced and enforced to tackle these problems. In most cases that will mean operational procedures and systems need to be adapted and possibly a dedicated MDM system needs to be implemented.

The business cannot wait for such efforts to bear fruit before they can get access to their vital management information.

Historical Data Archive (HDA) built on technical keys

What can be done when business keys are not available? The next best thing is to store the historical data linked to the technical keys of the sources. This at least assures that a true record of the historical changes of the source can be kept.

This can easily be achieved using QUIPU, by simply following the default process. When QUIPU analyzes a source model it will always look for business keys as primary keys or unique indices of the source tables. If it finds such a key it will suggest to use this key as the business key. Of course in practice it often is a technical key, not a business key.

Risk: Disconnected source Data Vaults

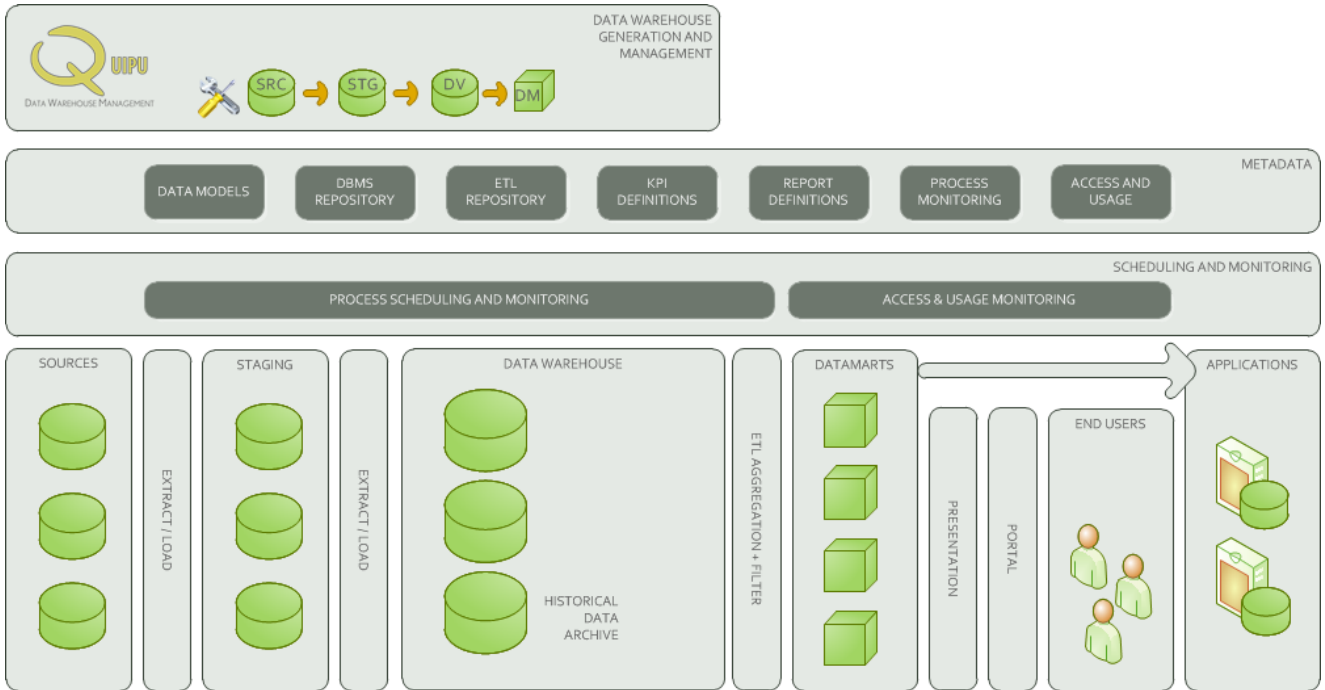
In the most extreme case (where all source systems have only disconnected technical keys) this approach results in a set of disconnected source data vaults, one for each source.

The value of a hub in this environment can be disputed, as the primary role of the hub is to serve as a connection point between systems/processes. When a hub contains technical keys that are only known in a single system this value evaporates.

HDA: satellites only

It is therefore logical to dispel the the hubs (and links) and only implement satellites.

The resulting environment is of course no longer a Data Vault, but is sometimes known as a Historical Data Archive (HDA). You could also label it as a Persistent Staging Area.



Historical Data Archive

Building the Data Warehouse on the HDA

The advantage of implementing the HDA akin to Data Vault satellites is that it is still possible to later find/define the true business keys and then start building a 'Data Vault'-like data warehouse by adding hubs and links. In some cases it may be possible to link these new hubs directly to the existing satellites.

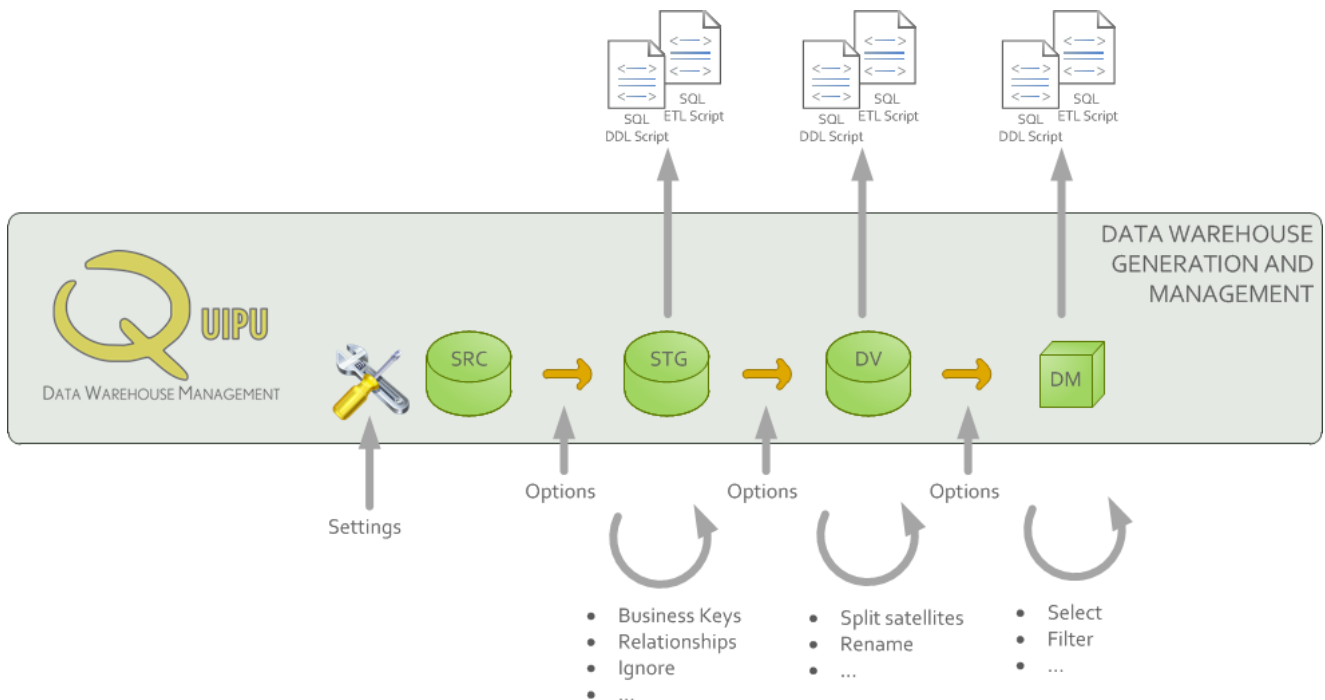
Model Generation

Data Vault Generation

Five steps to Data Vault generation

The generation of a Data Vault from a source system follows a five step approach:

- Import or reverse engineer a source model, generating the staging area as a (near) 1 to 1 copy
- Generate a Data Vault proposal from analysis of the staging model
- Evaluate the result and influence the generation process by:
 - (Re-)Define the business keys
 - (Re-)Define the relations between source tables
 - Ignore source columns that contain invalid or irrelevant data
- Regenerate the Data Vault model
- Final customization
 - Rename tables and/or columns if required to get an easier to understand model and add meta data to document the model
 - Split satellites, move attributes between satellites to group logically related attributes and/or separate frequently updated attributes from more static attributes



Data Vault generation step 1: Import model

Reverse engineer RDBMS

Reverse engineering is supported for a large number of relational database servers and file based 'databases' (Excel, XML, CSV, etc.). Requirements are:

- availability of a JDBC connector
- a connection between the QUIPU back-end and the database (or file) server
- sufficient rights to read the database catalog.

The reverse engineer module will import tables, columns, relationships and indices.

No direct support for database modeling tools

Direct import of the model from a database modeling tool (e.g. ERWin) is currently not supported. The easiest way to transfer a model from such a tool to QUIPU is via generating a DDL script in the tool and using it to create a dummy database in any of the many supported database systems. This database can then be reverse engineered by QUIPU.

Data Vault generation step 2: Generate Data Vault proposal ('Analyze staging')

Rule based generation

The Data Vault is generated from the source model based on a fixed decision tree implemented in the code.

The first rule is that the table must have a unique key. It consists of a single or multiple columns. And a each column can contain a technical or business key. If a unique key does not exist, the table is ignored.

Further rules:

A source table is a Hub unless proven otherwise: the general principle in the generation of the Data Vault is that all tables in the source become hubs, unless it can be determined (using a set of rules) that the table should be treated as a link or satellite table. This is generally the case when the principles ruling the Data Vault can be recognized in the source system.

A clear example is the many-to-many tables found in most relational systems that register transactions (e.g. a sale may be registered in a table with as primary key a reference to the product table, the customer table and a timestamp). Such a table should be recognized and translated to link and link satellite construction.

Analysis options

It happens that one (or more) of the business keys in a table is not a real business key. For instance: a table contains data with a business key and extra date column as its unique key. Strictly speaking, such tables should translate into link-satellites, with the extra (in this case: date) column representing a second Business Key placed in its own Hub. But often these extra columns have little business value and are NOT linking pins to other parts of the data model. The resulting Data Vault is then overly complex and hard to query.

QUIPU offers two distinct options to prevent this situation:

Multi-active satellite: Allow creating satellites with a different cardinality then its hub. The table is

mapped to a satellite on the hub, but the satellite contains the extra key field(s).

Allow peg legged links: Allow links that link only to one hub. Now the extra key field(s) are in the link table, as one would expect, but the corresponding hub is not implemented. All data on this link is kept in a link-satellite. See http://en.wikipedia.org/wiki/Data_Vault_Modeling for more information.

Log of the results of the decision tree

The result of evaluation of each node in the decision tree is recorded in a log file and is shown in the GUI. This allows the architect to evaluate the results and make the necessary changes to the source model to arrive to the desired result within the limits of the decision tree.

And of course that also enables a qualified software engineer to modify the QUIPU source code to implement the desired behavior.

Business keys derived from primary or alternate keys

The Data Vault decision tree uses a primary key or unique index constraint to find the candidates for business keys.

In many cases the primary key is a simple code that is actually meaningful to the business, providing a natural business key.

The source system may use meaningless keys, however, that have little or no value to the business. The real business key is then often defined as an alternate key and/or a unique index is defined on the real business key.

Source system blues

QUIPU will always create a Data Vault that is technically correct. It will however not always be the most efficient model, nor will it always be the model that you would like to have.

The reason for this is that the optimal Data Vault model can not be derived by analyzing the source system *structure* as it is registered in the catalog of the RDBMS.

There are plenty of potential problems, including but not limited to:

- Many source systems are implemented without foreign key relationships guarded by the RDBMS. The database consistency is then assured by rules in the software. Without foreign key relations, it is not possible for QUIPU to create link tables. NOTE: these links can be added manually in QUIPU to aid automated generation.
- Table and column names in the source may be generated, or for other reasons be impossible to interpret
- Tables are often 'overloaded'. Multiple logical entities are mapped into a single table.
 - Example 1: A column in the table indicates to what entity the row belongs. This is often found in source systems as implementation for reference tables, consisting of a code and description.
 - Example 2: Table are not fully normalized, so a 'parent' table (e.g. product group) is redundantly stored in a 'child' table (e.g. product)..
- Columns may be used for other purposes then their names imply
- Logical entities may be split over multiple physical tables (e.g. transaction per month)

All of these issues (and many more) will lead to the generation of a sub-optimal Data Vault. Issues

related to the *structure* can -within limits- be resolved within QUIPU by creating or modifying relations, renaming tables and columns, defining business keys. But issues related to the *contents* of tables (like overloading) cannot currently be resolved by QUIPU. If you want to resolve these, you will have to create some ETL to split or combine the physical tables and find a way to offer the logical 3NF model to QUIPU (e.g. by creating a dummy database with the correct structure).

Data Warehouse Architect is in control

Ultimately the data warehouse architect needs to check the result of the Data Vault generation process and manually intervene if necessary.

NOTE: The initial release for QUIPU will generate a separate Data Vault for each source model offered. So QUIPU will *not* integrate different source model into a single, linked Data Vault.

Data Vault generation step 3: Evaluate and modify generated Data Vault proposal ('Prepare staging')

Including / Excluding columns or tables

Specific tables and/or columns can be excluded from the generation process by setting them to 'ignore'. This is particularly important when generating a source data warehouse from a source that contains tables or columns that are not used (properly) by the business.

History

By default all columns are set to track history after generating the Data Vault. However for some columns you might want to disable history tracking and include only the last provided data. This may be the case for highly volatile, operational data (like process status-information), or on the opposite side of the spectrum: for very stable, descriptive information (like the descriptive name for organizations, equipment, etc.). In those cases you can be sure that people will only ever be interested in the actual values.

QUIPU allows switching history tracking off for particular columns.

NOTE: switching off history breaks the Data Vault. You will no longer have a full audit trail for this data in your data warehouse. So use it wisely!

Adjusting the business key

The QUIPU analysis routine may not always be able to correctly find the Business Key(s) from the structural information it gets from reverse engineering the database. You can add and remove columns to the business key to correct these situations.

Columns that are ignored cannot be part of the business key.

Every table needs to have a business key. If not the data vault generation process will automatically set the table to "ignore".

We don't recommend using technical (or generated) columns as business keys. While this will result in a technically correct Data Vault model, you can encounter problems when creating datamarts.

Reference tables

A reference table is a table that will be included when generating a Data Vault, but all foreign key relationships connected to this table will be ignored. This will result in a reduced amount of link tables in the generated data vault. Tables can be marked as reference table in the QUIPU GUI.

Add / remove relationships

QUIPU may not always get the correct relationships from the reverse engineering routine. These relationships are essential to the Data Vault generation routine, as QUIPU will map them to hub-satellites, links and/or link-satellites. It is therefor essential that QUIPU allows relations to be edited.

Relationship groups

Each relationship should have a sensible, descriptive name, as this name is used when the relationship is transformed to a link as the link name.

As a default QUIPU will generate a link table between 2 hubs for each relationship it finds. Very often it is more sensible to combine several relationships into a single link table, linking multiple hubs together.

This can be achieved by defining a 'Relationship group'.

Data Vault generation step 4: Regenerate Data Vault proposal

Implementation options

When the logical model has been defined to satisfaction, it is time to start considering the - physical- implementation options. These options have no impact on the data content of the data warehouse, but may have an impact on the performance in a particular database system or affect the possibilities for scheduling, monitoring and data maintenance tasks.

Model level implementation options

Record source and timestamp: Choosing this option will add a record source and timestamp column to each table in the staging area (default).

Audit ID: Choosing this option will add an audit ID column to each table in the staging area. The audit ID can point to a record in an audit table with record source and timestamp information, but could also contain additional information about the loading process (e.g. on extraction date, # of records in the export, operator, etc.).

NOTE: QUIPU will load these fields into the data warehouse for later reference, but will NOT fill these fields (in the generated ETL code to load the staging area). Loading the staging area is outside the scope of QUIPU.

Model/Table level implementation options

Generate surrogate keys: QUIPU can provide each data vault hub and link with a generated primary key (hub_id or link_id). Subsequently this key will be used in the relationships that connect these tables to each other and to their satellites. When this option is disabled the business key of the staging table will be used as primary key in the Datavault. Using generated keys can greatly reduce the size of the data warehouse and improve query performance, especially when business keys are multi-field and large (strings). It comes at a cost, however: loading the data warehouse becomes (potentially a lot) slower, as all business keys need to be translated to the correct surrogate key. Queries may become harder to build and read, as the surrogate keys have no meaning to the business.

Use end-dating: By default all satellites in the generated data vault will get an extra column called "load_date" that will contain the date of inserting. This column is required and sufficient for keeping the history. It is therefor mandatory and part of the primary key.

Retrieving the correct historical record from such a table may be difficult for a relational database: a complex query may be required that could also perform badly. QUIPU offers the option to add another column called "load_date_end" that can be used to line up all records in chronologic order. This field will be filled when a new version of the record is entered in the database, thus 'closing' the previous record. It can be used to simply single out the correct record with a SQL construct (load_date_start <= date < load_date_end).

Data Vault generation step 5: Modify / add metadata

Modify the generated Data Vault

When the Data Vault is generated, there may be cause to modify some aspects; to simplify the model or improve performance.

Split a satellite

The fields in a satellite may not change at the same rate. Also, some fields may always change as a set, not individually.

Example: An employee table may contain info on name, sex, date and place of birth and address details. The personal identification attributes will (in general) never change. The address details may change, but if they do the individual attributes (street, city) often change together.

This is a typical case where splitting the satellite could result in a better model: attributes that 'belong together' and 'change together' are placed in a separate satellite. This makes the model easier to read and interpret when at the same time requiring less disk space and offering better performance.

If you want to split a satellite, right-click on a satellite and select "split satellite" and give the new satellite a descriptive unique name. This will create a new satellite with only system columns and the business key. Now you can drag and drop data columns (from a satellite linked to the same hub or link) into the new satellite. Quipu will take care of (re)creating mappings.

Merge satellites

Sometimes it makes sense to combine multiple satellites in a single one. Thus, one can consolidate information from multiple source tables in a single table in the data warehouse. The model becomes simpler at the (potential) cost of storage requirements.

NOTE: Combining satellites from different source systems is not recommended. Different source systems have very often different characteristics in terms of completeness, correctness, timeliness etc. It is wise to store the data apart, so these characteristics can be easily assigned to a complete satellite table in the data warehouse (as opposed to columns in a combined satellite).

Rename objects

You can rename all tables and columns in QUIPU, so you can apply any naming convention. This will never alter any mappings.

Script generation

Template based generation

QUIPU uses templates to generate scripts to generate databases and load code. These templates can be customized to support specific features of the various RDBMS's. So for each SQL dialect a template could be provided. The initial release comes with a template for standard ANSI SQL 92 intermediate code. Additional templates will be released when they become available through development at QOSQO or in the community.

DDL

QUIPU generates Data Definition Language scripts to create:

- The staging area of the data warehouse
- The Data Vault area of the data warehouse
- 3NF-views on the Data Vault (mimicking the original source model)

ETL

QUIPU generates the load code for the staging area and the Data Vault.

Delta vs Full Load

The staging area load code assumes that the source database is situated in the same RDBMS, and the specified connection to this RDBMS has sufficient rights in the source database. This is of course often not the case. In these situations the code may serve as a sample demonstrating the logic to fill the staging area.

The Data Vault load code implements a delta detection mechanism. This means the code will function properly on both full data sets and delta data sets (only containing the changes in the source), *and you are not required to indicate whether you offer full or partial data sets.*

The problem with delta data sets is that QUIPU will not be able to detect *deletions* in the source (unless deleted records are marked by an attribute and offered in the delta set). This is the *one and only reason* for QUIPU to offer an option indicating that the data set offered is a full data set: *any record not in the full data set is considered deleted from the source.*

Staging Load table

QUIPU will end-date the current satellite record(s) and create a new record marked as deleted. Also, any hub-entry not found in the source will be marked as deleted.

In practice, delta and full loads are often required intermittently for a specific table. It is common practice to process the smaller delta data sets on a daily basis, accepting that deletes are not recorded. Once every week or month a full data set is processed to record the deletes that occurred in this period.

To facilitate this practice QUIPU implements a special table in the staging area with a record for each source table. The data warehouse manager can indicate for each source table whether a full load or a partial load is required.

QUIPU generates code that will check this table and automatically perform the full or partial load as specified in this table.

Splitting the load process

It is possible to split the load process in multiple jobs. For instance in many cases it makes sense to load the master data only once a day (with the full load option), but load the transaction data multiple times a day (to get more actual data and prevent very large load jobs at the end of the day).

This works well, with the one consequence that for NEW master data only the business key will be available until the all master data is refreshed.

Re-starting load jobs

All of this is not directly supported by QUIPU, but can be achieved by either splitting up the generated ETL scripts or building a scheduling solution that directly addresses the QUIPU repository.

When the same data set is offered multiple times (either full or delta sets), no change will be found, and thus no changes will be recorded in the Data Vault.

This means interrupted data loads can be restarted without any adverse affects.

All tables in a single export can be loaded in any particular order (and even in parallel), although it makes sense to load master data before loading transaction data.

Please note that offering jobs out-of-order (older extracts after newer extracts) will cause incorrect results in the Data Vault!

NOTE: that the ETL scripts will only implement mapping from source tables to target tables. No transformations can be specified.

End-dating

There is an option in QUIPU to add end-dates to the satellites. The end-date allows simpler SQL (and often: faster) queries on the Data Vault, using the SQL construct "WHERE load_date_start <= point-in-time < load_date_end"

End-dating the satellites means that the each record in the satellite needs to be updated with the start-date of a later record. If no later record exists, the end-date gets a date very far in the future (31-12-9999).

End-dates can be updated when new entries arrive, but this slows down the data load processes, which can be a problem in certain scenario's (e.g. for near-realtime updates for operational BI). That is why QUIPU implements end-dating as a separate process that should be executed *after* loading the data.

Views

Some views will be generated by QUIPU that will present an easier to interpret (and query) representation of the Data Vault.

The views will come close to rebuilding the 3NF data model that was used for the Data Vault generation (unless manual interventions have changed the model).

As the Data Vault stores historical versions of all data, and the 3NF model in general does not, the views will need a point-in-time to indicate which historical view is requested.

Standard views included are:

- the current status (using today as the point-in-time).
- the current status including all deleted data (using today as the point-in-time)
- the full historical record

QUIPU Technology

Client Server Architectuur

Internal structure

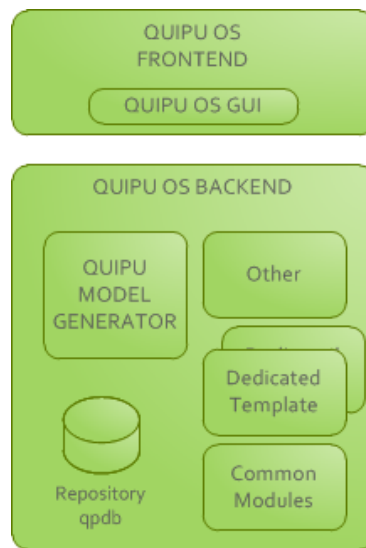
QUIPU is internally structured with a strict separation between components and (front- and back-end) layers based on open interfacing standards in order to:

- isolate the generic, data warehouse and data vault specific logic from the specific optimizations and implementation details of the technology used
- allow concurrent development of new functionality
- allow new technologies to be integrated and optimized (through add-in components like plugins, external libraries etc.)

Thin client concept

QUIPU has been designed with a thin client in mind. This means all core functions are implemented in the backend, and the client only implements the GUI component of the application.

Template based generation



The generation of DDL and ETL code is template based. This allows QUIPU to generate optimized code for a great variety of target (database) systems. The standard templates produce ANSI SQL code that will run unchanged on many RDBMS systems.

In the future optimized templates can be added to QUIPU. This is something we will turn to the community of QUIPU users to co-develop and co-support, as the variety out there is overwhelming and it is impossible for us to set up test-environments for each.

In essence it is an open system allowing users to change or add templates.

Back-end: Groovy / Grails

Java Virtual Machine

A key technology decision for QUIPU is to have the backend implemented on the Java Virtual Machine.

The Java Virtual Machine that executes the Java (and Groovy) code is the premier platform for portable applications (that will run on almost any hardware platform and operating system) and is also well known for its capabilities for web applications.

Running the QUIPU backend on the Java Virtual Machine assures that this key component of QUIPU can run in all conceivable technical platforms, from stand-alone pc's to enterprise server environments.

The latter is specifically relevant as for more advanced implementations production systems need direct access to the QUIPU backend. The QUIPU backend must then be able to run in the same protected environment as these production systems.

Groovy programming language

The drive to develop Groovy comes from the perception that Java is not so very productive due to the enormous (coding) overhead it requires for very mundane tasks. Groovy is much more permissive and requires much less code to achieve the same functionality.

Grails Framework

GRAILS stands for Groovy on Rails, and borrows a lot of the concepts and mechanisms of the well known Ruby on Rails environment: a framework that bundles a number of key Java libraries and puts convention over configuration.

Groovy -as programming language- allows the programmer to build nearly all perceivable functionality. But of course that would require a lot of time and effort. For many generic tasks Java libraries have been built and many of those are in the public domain. Grails has selected and bundled a number of these (e.g. Hibernate to interact with databases, or Spring to generate some basic GUI functionality) into the Grails package. A common problem associated with the use of such libraries is the need to configure them correctly. The more powerful the package, the more options can be configured in configuration files and the more difficult it becomes to make them work.

Grails solves this problem by pre-configuring the libraries to assume the programmer adheres to a set of well described *conventions*. These conventions consist of rules on how to structure the code (in files) and where to put them (a set of directories with a specific purpose each) as well as naming conventions. When the programmer adheres to these conventions virtually all need to configure the libraries is removed. Of course, the configuration files are still available to the programmer if certain -default- behavior needs to be tweaked.

Front-end: C# / .Net

Thin C# client

The front-end for QUIPU has been developed in C#, and requires the Microsoft .NET runtime modules to be installed on a pc with Windows installed. The front-end has deliberately been kept as thin as possible, so most functionality of QUIPU is implemented in the back-end. Front-end and back-end communicate through standard http calls. This makes it entirely feasible to develop alternative front-ends (e.g. using AJAX, Silverlight, etc.) or complimentary front-ends (e.g. script-

ing certain processes to execute without interaction).

Database support

Broad RDBMS support for Quipu Repository

QUIPU stores all the information it manages in a (relational, 3NF) repository. This repository can exist in a very large number of database management systems (courtesy of Hibernate). See Hibernate documentation for a list of all supported databases and versions.

This broad support of database systems allows the QUIPU repository to be installed in the target database environment, which is required for the (future) monitoring and scheduling functionality. QUIPU configures itself upon installation: only an empty database needs to exist with a valid connection and proper authorization. QUIPU will then automatically create an empty repository.

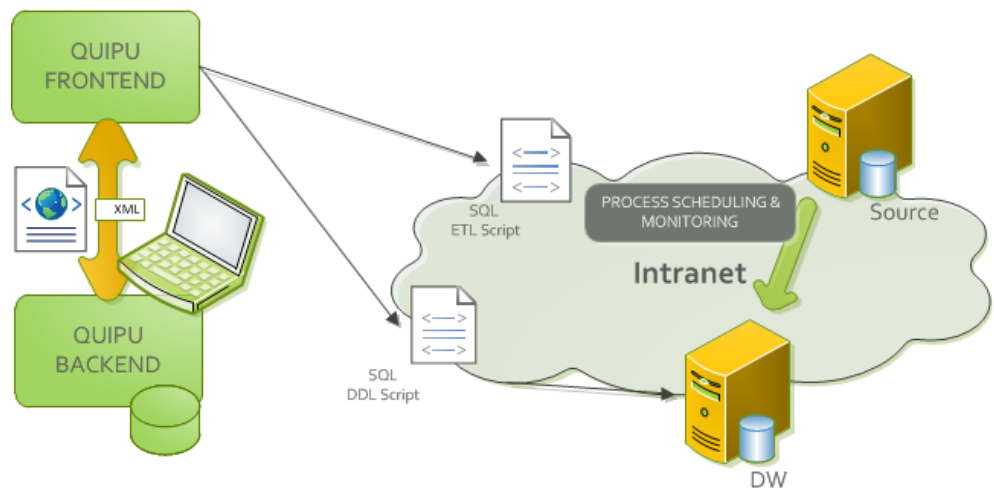
'Out of the box' QUIPU contains a file based RDBMS (HSQLDB) so no separate RDBMS is required to install and run QUIPU.

QUIPU deployment options

Standalone / Prototyping setup

In this minimal scenario QUIPU front-end and back-end are installed on the same machine (PC). QUIPU is used to generate the DDL scripts and ETL scripts, that can be stored in .TXT files using the Front-End.

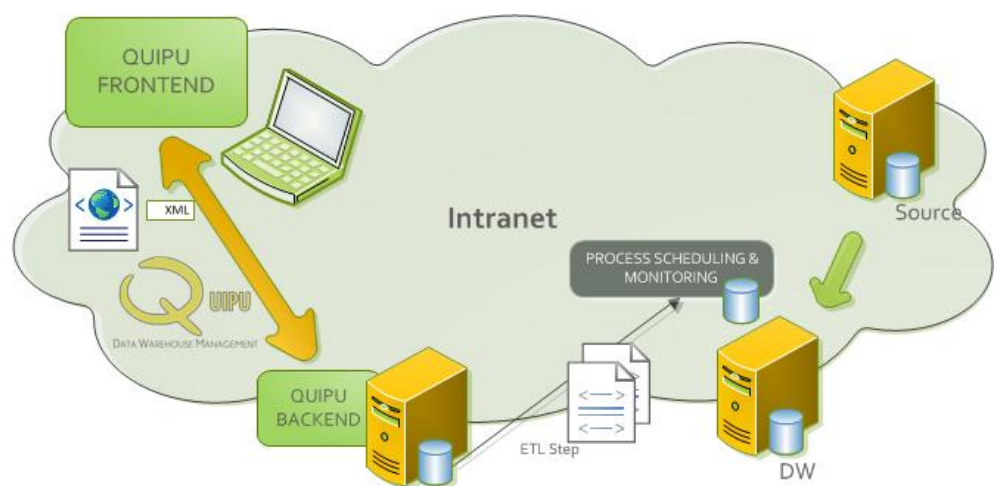
For each physical datamodel (schema) a script file will be generated that wil create or fill the entire model.



It is up to the architect to deploy these scripts in the production environment, to create the DW structures and then fill these structures with data.

Integration in server environment

A more advanced option is to install the QUIPU back-end on a machine that can be -reliably and securely- accessed by the process scheduling and monitoring tool (often a component of the ETL software).



In this setup it is possible for the ETL software to access individual ETL step generated by QUIPU and execute them separately. This allows for more elaborate monitoring and exception handling.

The setup can be made fully generic, so few or no modifications are required in the process scheduling and monitoring setup are required when new version of the ETL are generated by QUIPU.

This setup has been successfully deployed in combination with Open Source Pentaho components (Data Integration, formally known as Kettle).

Future developments

QUIPU OS functionality

QUIPU remains generation tool

The core function of QUIPU is the transformation of models (staging, Data Vault, data marts), the automatic generation of these models, their mappings and the ETL code to move data through the chain.

In the future the Open Source QUIPU tool will maintain this function, and the modeling capabilities will continue to expand and improve to fulfill the needs of businesses today.

Fill the gaps

There are certain key aspects of data warehouseing that are not fully supported as of yet. We will strive to add these functions in future releases. At the top of our list are:

- support for shared hubs, allowing model integration
- generation of data marts, including:
 - definition of shared dimensions and hierarchies
 - definition of KPI's

QUIPU is primarily designed to support the Data Vault model, but is not restricted to the 'pure' Data Vault.

QUIPU already supports some derivations of the Data Vault. In the end it is the responsibility of the data warehouse architect to choose an appropriate model for the data warehouse and up to QUIPU to optimally support this model.

Commercial modules

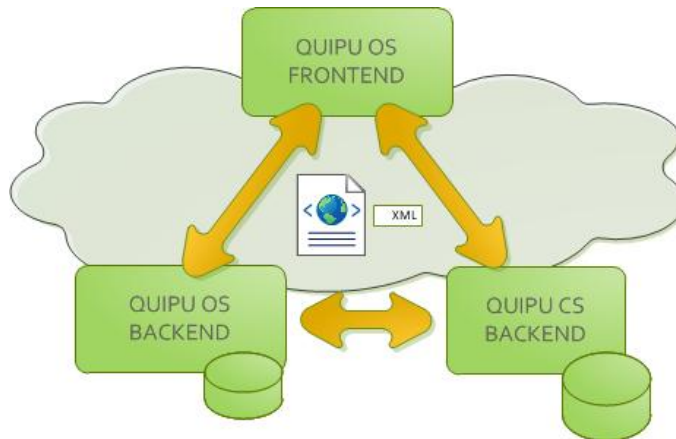
Deployment architecture changes

The technical deployment architecture of QUIPU was redesigned, in order to support both Open Source and Closed Source (commercial) modules in a single environment.

In the current release, QUIPU consists of a GUI module, built with C# on .NET, and a back-end module built with Groovy on Grails. The latter is distributed as a .war package -that can be deployed on web server- that communicates with the GUI via standard http calls.

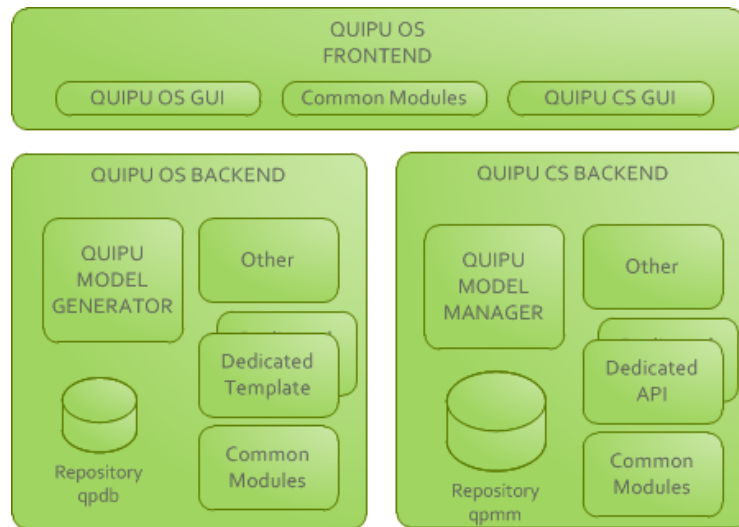


In future release (as of version 2.0) another server component is distributed containing all closed source modules. This component will implement a second web server, provides access to the closed source functionality.



Thus all QUIPU users will have the exact same software stack installed on their machines. Access to specific closed source modules will be based on a simple licensing mechanism.

The GUI module will contain all GUI functionality for both server modules, and communicate via https calls to both web servers. When a specific module is not licensed, the corresponding GUI components will be disabled. Both server components will also communicate directly with each other using web services.



Quipu Model Manager

The QUIPU Model Manager module will complement the functionality of the Open Source QUIPU primarily directed towards larger installations, where multiple project teams may work on expanding the data warehouse and modify existing data marts or design new ones.

Changes in the data warehouse structure (e.g. splitting a satellite) may require data to be moved in bulk from one table to another, with possibly recalculation of valid periods (start-end dates).

QMM is a substantial investment (development has started) that involves the creation of a new repository able to hold the additional (meta-)data elements and functions.

The first version of QMM will focus on change management on the various QUIPU models, and offer functions to:

- Check model versions in/out of a central repository
- Report on changes made between model versions
- Generate code to *modify existing* models in the data warehouse

As some of these functions can become very complex, the first release will probably focus on the more straightforward situations. We will then add more advanced routines gradually with each further release.

In later versions we will add functionality to move data from outdated to new structures in the data warehouse and support project teams.

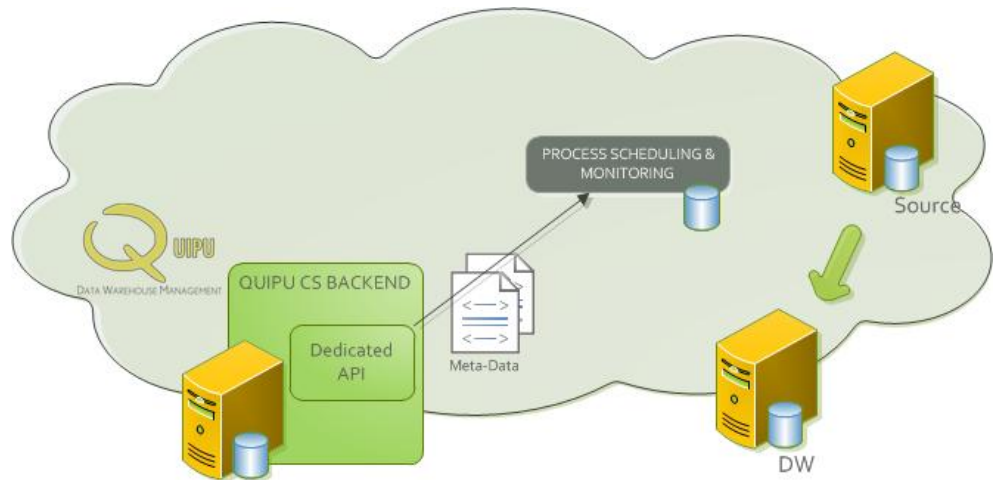
Integration with ETL/DW management tools

Larger organizations often own tools that assist them in designing and managing the operational data logistic processes. Tools like Informatica PowerCenter, SAS Data Integration Server, Oracle Data Integrator, etc.

These tools are built around a repository, that contains all the ETL procedures and schedules. Data flows can be modeled using dedicated design tools.

The drawback of building a Data Vault with these tools is that they are very 'model neutral' and thus have no fundamental understanding of the Data Vault. They therefore lack the intelligent support QUIPU offers when transforming source models to Data Vault models and then Data Vault

models to Data Marts.



QUIPU will not compete with the strong points of these tools (in terms of operational management of a data warehouse environment) but will instead seek to complement these tools by feeding the generated data models and ETL code into their repositories.

Clearly such connections will (greatly!) benefit only the customers that have made substantial investments in such tools. As the customer base for these developments is small we believe it is fair to develop these modules as commercial add-ons.

Support for BI tools

Many BI front-end tools have a propriety data access layer. Some tools implement an optimized physical data storage (e.g. Cognos cube, QlikView data file), others implement a virtual layer that retrieves the required data from an underlying data base (Business Objects Universe).

The maintenance challenge for BI competence centers is caused by the fact that these data access layers are affected by changes in the data warehouse (or data marts).

So, when new data elements are added to the data warehouse the data access layer needs to be adapted before the data becomes available in the front-end tool.

Specific modules will be developed to generate and re-generate components of the data access layer when necessary due to changes in the data warehouse or data marts.

Obviously these modules will be very specific for each target front-end environment and functionality may differ between modules.

More information

About QOSQO

QOSQO is founded in 2008 as a sister company of Nippur (founded in 2002). QOSQO is the leading company behind QUIPU, and provides data warehouse services for companies using QUIPU or other Data Vault modeled data warehouses.

Licence model

QUIPU is distributed under the GPLv3 licence model. More information can be found on <http://www.opensource.org/licenses/gpl-3.0.html>.

Source code availability

The QUIPU source code of all released versions is available for developers either via sourceforge.net website or upon request via email to support@QOSQO.nl.

QOSQO
datawarehouse services

[HTTP://WWW.QOSQO.NL](http://www.qosqo.nl)