
**Language resource management —
Corpus query lingua franca (CQLF) —
Part 1:
Metamodel**

*Gestion des ressources linguistiques — Corpus query lingua franca
(CQLF) —*

iTeh STANDARD PREVIEW
Partie 1: Métamodèle
(standards.iteh.ai)

[ISO 24623-1:2018](https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018)

<https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018>



iTeh STANDARD PREVIEW (standards.iteh.ai)

ISO 24623-1:2018

<https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018>



COPYRIGHT PROTECTED DOCUMENT

© ISO 2018

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Fax: +41 22 749 09 47
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

	Page
Foreword.....	iv
Introduction.....	v
1 Scope.....	1
2 Normative references.....	1
3 Terms and definitions.....	1
4 Aims.....	4
5 Assumptions.....	4
6 CQLF Metamodel.....	4
7 Conformance.....	7
Annex A (informative) Example CQLF conformance statements.....	8
Bibliography.....	12

iTeh STANDARD PREVIEW (standards.iteh.ai)

[ISO 24623-1:2018](https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018)

<https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018>

Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular the different approval criteria needed for the different types of ISO documents should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation on the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see the following URL: www.iso.org/iso/foreword.html. (standards.iteh.ai)

This document was prepared by Technical Committee ISO/TC 37, *Language and terminology*, Subcommittee SC 4, *Language resource management*.

A list of all parts in the ISO 24623 series can be found on the ISO website. Additional parts on single-stream and multi-stream ontology architectures are planned to be developed in the future.

Introduction

A range of standards relating to language resource management, with the Linguistic Annotation Framework (ISO 24612) at the centre, have been developed. These standards are mostly designed to regulate the representation aspect of language data – they look at the data from the point of view of preparation and curation. This document complements this perspective by that of the end-user, that is to say, from the point of view of processing and querying.

The corpus linguistic community has, by now, developed several corpus query languages (QLs), and there is a particularly large number of them if “dialects” and forks are included. There are two main reasons for this abundance. Firstly, there are socio-economic and organizational factors, with separate query systems having been created by isolated projects with un-coordinated funding, many of them eventually developing their own set of followers. Secondly, query systems are typically sensitive to the format of the data and are often designed with a specific purpose in mind. For example, systems for querying parallel audio and transcription streams with multiple speakers have different characteristics from systems designed to query purely textual data with a single layer of morphosyntactic description. Dependency and hierarchical annotations demand yet another set of solutions. All of this results in the richness of alternatives or near-alternatives on the one hand, and in the lack of interoperability among the variants on the other. As a consequence, a “wrong” choice at the beginning of a project can bury months of research by exposing inadequacies in the initial decision after the project has become mature enough to move to new extended functionality and towards addressing more complex information needs.

This document codifies, in a modular way, the best existing practices followed in the design of corpus query languages. Its theoretical aim is to provide a basis for the investigation of the relationships between language resource architecture and corpus query language properties. The practical aim of the Corpus Query *Lingua Franca* (henceforth COLF) is to provide linguists and language technology practitioners with a clear and coherent basis for making informed choices concerning data architectures and the query languages appropriate to them.

[ISO 24623-1:2018](https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018)

<https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018>

iTeh STANDARD PREVIEW
(standards.iteh.ai)

ISO 24623-1:2018

<https://standards.iteh.ai/catalog/standards/sist/ce21b803-f548-42c0-9d55-d3f2af64597a/iso-24623-1-2018>

Language resource management — Corpus query lingua franca (CQLF) —

Part 1: Metamodel

1 Scope

This document describes the abstract metamodel designed to accommodate any corpus query language (QL) and providing a basis for coarse-grained classification. The metamodel consists of several components referred to as CQLF classes, levels, and modules, and is illustrated with examples from the Single-stream class (where a single data stream is used to organize the relevant data structures). Within this class, this document discusses three CQLF levels (Linear, Complex and Concurrent), as well as their subdivisions into modules, dictated by functional and modelling criteria.

This document does not provide a way to specify further details beyond the above-mentioned divisions, and neither does it contain within its scope QLs designed to query more than one concurrent data stream, as in multimodal corpora or in parallel corpora (such QLs can still be classified according to the criteria suggested here for less expressive QLs).

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO 24611, *Language resource management — Morpho-syntactic annotation framework (MAF)*

ISO 24612, *Language resource management — Linguistic annotation framework (LAF)*

ISO 24615-1, *Language resource management — Syntactic annotation framework (SynAF) — Part 1: Syntactic model*

3 Terms and definitions

For the purposes of this document, the following terms and definitions apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- IEC Electropedia: available at <https://www.electropedia.org/>
- ISO Online browsing platform: available at <https://www.iso.org/obp>

3.1 annotation

information added to *primary data* (3.9), independent of its representation

[SOURCE: ISO 24612:2012, 2.3, modified — "linguistic" at the beginning of the definition was deleted.]

3.1.1

concurrent annotations

multiple, potentially conflicting *annotations* (3.1) describing, entirely or partly, the same *character span* (3.2) or an overlapping sequence of character spans

Note 1 to entry: Concurrent annotations may be expected to conflict in several ways: content-wise (with different tags for the same character span), structure-wise (assuming different structural arrangements within the targeted character spans), and also in terms of segment edges (which is typically due to structurally conflicting claims concerning the encompassing character spans). Concurrent annotations typically come from different sources (e.g. tools or human annotators) or result from different settings (e.g. different parsing models or segmentation rules) within a single tool. When encoded in XML, concurrent annotations are typically expressed by means of stand-off techniques.

3.1.2

dependency annotation

annotation (3.1) that encodes the dependency relations between *character spans* (3.2)

Note 1 to entry: An example of a dependency relation (see ISO 24615-1:2014, 3.5) is one between a verb and its subject or direct object, between an attributive adjective and its head noun, or between a preposition and the head of its dependent noun phrase. Dependency relations may be defined at the word-level alone, or may involve higher-level syntactic constructs, in which case it is possible to speak of mixed hierarchical-dependency annotations.

3.1.3

hierarchical annotation

annotation (3.1) that encodes the relationship of dominance (often also precedence) necessary to define syntactic trees over *character spans* (3.2)

Note 1 to entry: Annotating hierarchical relationships requires only the relation of dominance to be indicated. Precedence is typically implicit in the ordering of character spans.

3.1.4

segmentation annotation

annotation (3.1) that delimits linguistic elements that appear in the *primary data* (3.9)

Note 1 to entry: These elements include (1) continuous segments (appearing contiguously in the primary data), (2) super- and sub-segments, where groups of segments will comprise the parts of a larger segment (e.g. contiguous word segments typically comprise a sentence segment), (3) discontinuous segments (linking continuous segments) and (4) landmarks (e.g. time stamps) that note a point in the primary data. In current practice, segmental information may or may not appear in the document containing the primary data itself.

[SOURCE: ISO 24612:2012, 2.5]

3.1.5

simple annotation

annotation (3.1) that constitutes a single information package whose interpretation is not dependent on other annotations

Note 1 to entry: This definition is intended to distinguish the simplest (“tabular”) kind of annotation from more complex relational structures (providing hierarchical, dependency, or alignment information); simple annotations are the only kind of annotations present at the linear level of complexity.

3.1.6

stand-off annotation

annotation (3.1) that can be layered over *primary data* (3.9) but is separated from the data stream that it targets

Note 1 to entry: Stand-off annotations refer to specific locations in the primary data, by addressing the character offsets, elements or coordinates to which the annotation applies. They can be serialized as separate documents, but do not have to be. Multiple stand-off annotation documents for a given type of annotation can refer to the same primary document (e.g. two different part of speech annotations for a given text). It is also possible to construct hierarchies of stand-off annotation layers, where layer n can reference layers $0..n-1$.

[SOURCE: ISO 24612:2012, 2.7, modified — The definition and note were modified.]

3.2

character span

sequence of characters, identified by start and end offsets, to which an annotation may be applied

Note 1 to entry: This definition is a relaxed version of the definition in ISO 24615-1:2014, 3.16, the difference lying in the use of “may be applied” over “is applied”. Compare also the definition of “region” in ISO 24612:2012, 2.10.

3.3

character span containment

relation obtaining between *character spans* (3.2) of *primary data* (3.9) in which character span A contains character span B if the initial offset of span A is equal to or higher than that of span B, and the final offset of span A is smaller than or equal to that of span B

Note 1 to entry: The relation of character span containment is used for stating a relationship between two or more character spans or simple annotations, without the need to utilize tree-based concepts and mechanisms. Instead of tree traversal, operators such as *contains*, *in* or *within* are typically used for character span containment queries.

3.4

corpus query language

formal language designed to retrieve specific information from (large) language data collections, and thereby incorporate certain abstractions over commonly shared data models that make it possible for the user (or user agents) to address parts of those data models

3.5

CQLF class

top-level division in the CQLF data model

iTeh STANDARD PREVIEW
(standards.iteh.ai)

Note 1 to entry: The CQLF Metamodel distinguishes two classes: Single-stream (where the annotation structure is built upon a single data stream, typically a character stream) and Multi-stream (corresponding to e.g. multi-modal corpora or parallel corpora).

3.6

CQLF implementation

query language that has been analysed with respect to the criteria described by the CQLF Metamodel, and thus has been “located” in the proposed feature matrix as “conformant with CQLF”

3.7

CQLF level

part of the matrix of QL properties, defined in terms of the general features of the assumed corpus data models, and consequently the set of properties of a corpus query language that is used to address these features

Note 1 to entry: The CQLF Metamodel distinguishes three levels of complexity within the Single-stream class: Linear, Complex and Concurrent.

3.8

CQLF module

subcomponent of a CQLF level, defined with reference to a specified data-model characteristic

Note 1 to entry: CQLF Metamodel currently distinguishes three modules within CQLF Level 1, Linear (plain-text, segmentation, and simple annotation), and three modules within CQLF Level 2, Complex (hierarchical, dependency, and containment).

3.9

primary data

electronic representation of language data

3.10 token

non-empty contiguous sequence of graphemes or phonemes in a document

[SOURCE: ISO 24611:2012, 3.21, modified — The note was deleted.]

4 Aims

The CQLF Metamodel is intended to establish a frame and a basis for establishing the potential extent and the limits of interoperability between different corpus query systems. It aims to provide a single matrix of a few well-defined properties in which any corpus QL can be located for the purpose of coarse-grained comparison with the others. Further parts of the standard elaborate on these properties and flesh out the relationships among them. A long-term goal of CQLF is, additionally, to help reduce the gap between end users with a linguistic or literary background and powerful search environments.

While CQLF as a whole might be expected to mediate between individual corpus QLs as an interlingua, such initiatives raise a host of problems, ranging from low-level technical descriptions (e.g. the inability to preserve information when translating between regular expressions on the one hand and wildcards on the other) to issues of epistemology (“Does the result of the query reformulated in QL 2 address exactly the information need expressed in QL 1?”). The immediate goal of CQLF is therefore more modest: to serve as the target space within which QLs can be located with respect to their basic properties. It can thus also serve as a measure of compatibility and interoperability, but without the added claim to provide QL-to-QL mappings. A robust CQLF-based bi-directional mapping system in action (with the epistemological burden appropriately controlled) would be an interesting challenge in the long run, but it makes a lot of sense to start smaller, and in the spirit of other standards, lay the ground for a pivot-based system with monodirectional mapping from various corpus QLs into a representation defined by a superset of their individual properties.

The metamodel presented here circumscribes the outer limits of QL compatibility. CQLF may be used as a set of guidelines to be applied in the development of a new QL or in the enrichment of an existing QL with new functionality. It is likely that information about the extent and points of conformance of a given QL with CQLF will also be useful to corpus linguists for the purpose of identifying the QL suitable for the task that they are faced with.

5 Assumptions

The metamodel described here builds on models defined in other standards. In particular, its infrastructure includes:

- the general data model for corpus and annotation description defined by ISO 24612 (LAF), together with the more detailed models defined by ISO 24611 (MAF) and ISO 24615-1 (SynAF);
- a common repository of data content and data containers (see ISO 12620).

The simple data model for Single-stream architectures (as opposed to, e.g., multi-modal corpora or parallel corpora) recognizes a minimum of two layers at which text can be queried: the layer of characters and the layer of labelled abstractions over characters (interpreted by the annotator or tools). The sequences in the former layer can be defined by means of character-based regular expressions, whereas sequences in the latter can be defined using ISO/IEC 14977 (EBNF).

An additional assumption made in the present specification is that tokenization is merely a special kind of segmentation annotation, upon which hierarchies of other annotation layers can be built.

6 CQLF Metamodel

CQLF is designed to be a modular construction with several components. Each component is characterised with respect to some aspect of the data models describing corpus objects that are the target or context of queries. A schematic view of the components of a CQLF Metamodel is presented

in [Figure 1](#). The top-level components are referred to as CQLF Classes and correspond to the major division into data models built upon a single data stream vs. those which use more than one data stream (be it binary or text-based) in parallel. This document illustrates an instantiation of the metamodel for the Single-stream class (to be introduced below), which consists of three CQLF Levels that correspond to the major kinds of data organization in linguistic corpora.

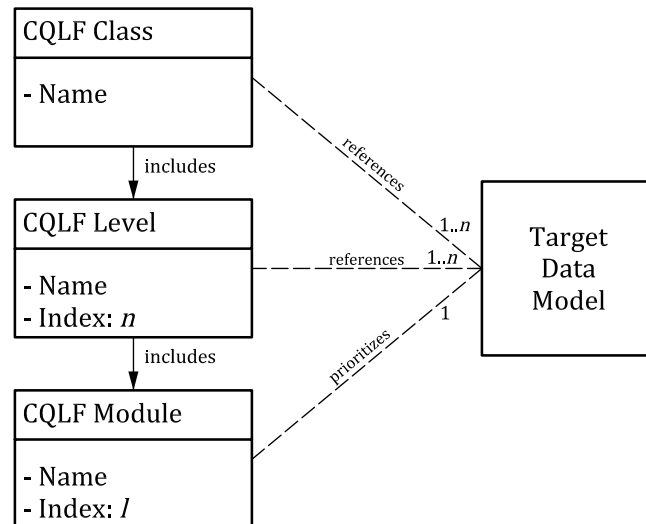


Figure 1 — Components of the CQLF Metamodel

The *index* variables of CQLF Level and Module are convenience details that allow for easier reference. The identification of levels and modules, while grounded in formal properties of data models and in functional characteristics of query languages, is also partially utilitarian, where it neglects some distinctions that could otherwise be recognised, in order to provide a simpler mechanism for determining conformance with the overall model and for stating the most important similarities and differences between corpus QLS.

The relationships sketched in [Figure 1](#), together with the basic divisions made on the basis of the existing corpus QLS, yield a basic taxonomy presented in the diagram in [Figure 2](#).

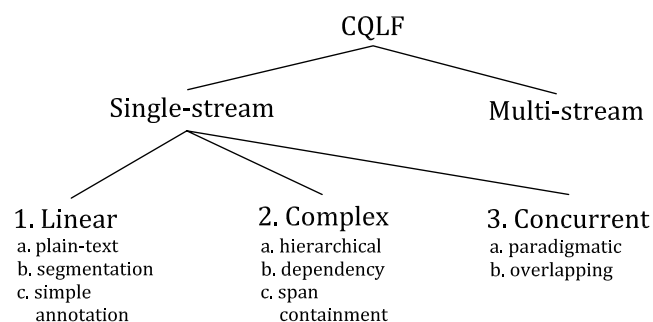


Figure 2 — Partial CQLF taxonomy of query language properties

Within the Single-stream class, each consecutive level introduces a more complex search dimension. The level system is based on the distinction between different major types of data organization and consequently different types of annotations (modelled by various parts of the LAF family of standards), as well as queries that correspond to them.

- **Level 1 (Linear)** addresses, in any combination, plain-text search (1a) as well as search in segmented data (segmentation annotations; 1b), and in simple annotations (1c) attached to particular segments:
 - at this level, annotations (if present) form a single layer of objects that exhaustively or partially describe the primary data stream;