

# DataKnots: A Framework for Building Domain Specific Query Languages

Clark C. Evans; Kyrylo Simonov, PhD

DataKnots<sup>1</sup> is an extensible data processing framework. It lets collaborative research teams build their own domain specific query languages (DSQLs) that reflect their conceptual models and vocabularies. DataKnots is based upon an algebraic framework, Query Combinators<sup>2</sup>, which formalizes how query components can be defined and combined in a consistent manner. While DataKnots includes standard query operations (group, filter, sum, etc.), they are not given special treatment over domain specific operations. One can add new query components by encapsulating existing queries, lifting Julia language subroutines to queries, and authoring novel transformations using a pipeline construction interface. DataKnots is a platform for creating an ecosystem of mutually interoperable DSQLs, so that collaborative research groups can easily customize their query system to fit the kinds of data sources and analysis methods they use.

To show that ergonomic, high-performance DSQLs could be rapidly constructed, we built a DSQL inspired by the Clinical Quality Language (CQL) and used it to implement CMS124v7 “Cervical Cancer Screening”. Specifically, we constructed a processing pipeline that converts JSON encoded Fast Healthcare Interoperability Resources (FHIR) to its clinical quality measure (CQM) score. The pipeline loads JSON to an in-memory FHIR representation, converts to an intermediate Quality Data Model (QDM) form, then implements CMS124v7 logic to calculate the CQM score. This layered approach separates concerns, giving us a place to put encoding logic specific to FHIR that doesn’t belong in an CQM, as well as a place to isolate electronic health record vendor differences. This CMS124v7 measure, inclusive of necessary FHIR and CQM dependencies, took 27 working days to implement<sup>3</sup>. Support for additional measures can now be added with incremental effort.

We benchmarked using synthetic patient data. Computation of a single patient averaged 76ms with a single core on a i7-4770 desktop computer, while the reference implementation averaged 607ms. For aggregate computations over batches of patients, our bottleneck is memory usage, with a high-water mark of 11mb per patient. JSON parsing is expensive (17ms/patient). These benchmarks motivate future work on a custom JSON parser, suitable to our vectorized representation, that extracts FHIR lazily based upon the exact fields needed for a given computation.

CQL is a highly-targeted DSQL, created specifically for implementing clinical decision logic. Using DataKnots, we were able to construct a DSQL that matches CQL with comparable functionality and ergonomics. Moreover, we were able to represent not only a CQM calculation, but the entire processing pipeline, including conversion from JSON to FHIR and conversion of FHIR to QDM. For this DSQL, DataKnots was extended with relevant data types, such as a datetime `interval`, and corresponding operators, such as `and_previous` and `during`. These extensions are treated no differently than built-in operations such as `filter`. Unlike CQL, whose syntax and semantics are specifically designed around clinical decision logic requirements, our syntax and semantics are universal across any application of DataKnots, reducing training costs and enabling sharing of query components.

```
define "PapTest Within 5 Years":
  ("Pap Test with Results" PapTestOver30YearsOld
   with ["Patient Characteristic Birthdate"] BirthDate
   such that Global."CalendarAgeInYearsAt" (
     BirthDate.birthDatetime,
     start of PapTestOver30YearsOld.relevantPeriod)>= 30
   and PapTestOver30YearsOld.relevantPeriod 5 years or
     less before end of "Measurement Period")
```

CMS124v7 fragment using CQL with QDM  
<https://ecqi.healthit.gov/sites/default/files/ecqm/measures/CMS124v7.html>

```
@define PapTestWithin5Years =
  let birthDate => PatientCharacteristicBirthdate.
    BirthDateTime,
    previous5years => interval(MeasurePeriod.end) .
      and_previous(5years)
  PapTestWithResults.
  filter(years_between(relevantPeriod.start, birthDate) >= 30
    && relevantPeriod.during(previous5years))
end
```

an equivalent using DataKnots – cms124.jl  
<https://github.com/rbt-lang/DataKnots4FHIR.jl/blob/master/doc/src/cms124v7.jl>

At the time of submission, our effort is not used in any production system. Notably, this approach is not at all limited to clinical decision support, indeed, we’ve prototyped a DSQL tailored to Observational Health Data Sciences and Informatics (OHDSI) cohort construction<sup>4</sup>. The DataKnots project is MIT/Apache licensed, exquisitely documented, and has extensive regression tests. We are actively searching for a pilot project.

1. Evans CC, Simonov K, DataKnots Query System for Julia (<https://github.com/rbt-lang/DataKnots.jl/>)
2. Evans CC, Simonov K, Query Combinators (<https://arxiv.org/abs/1702.08409>)
3. Evans CC, DataKnots4FHIR : Query Adapters for FHIR (<https://github.com/rbt-lang/DataKnots4FHIR.jl/>)
4. Evans CC, Simonov K, DSQLs for Medical Research (<https://www.biorxiv.org/content/10.1101/737619v2>)