# codeq
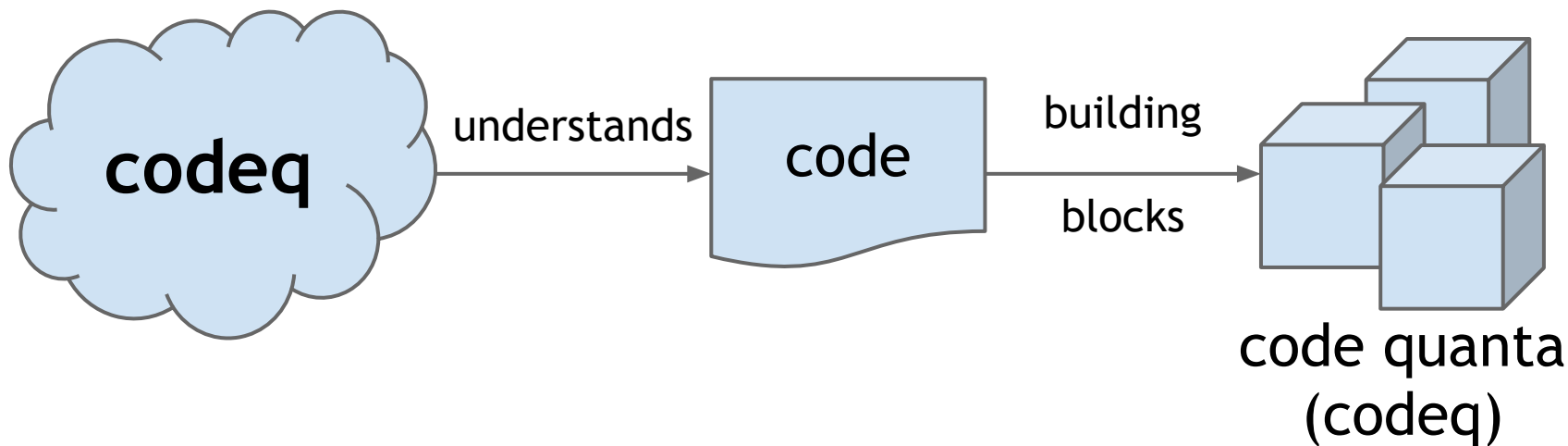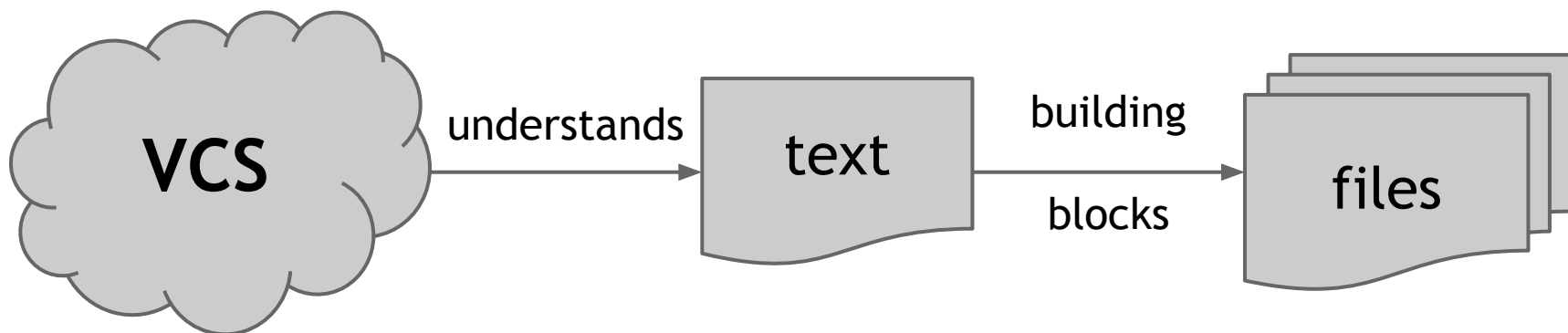## making a VCS code-sensitive

an evaluation
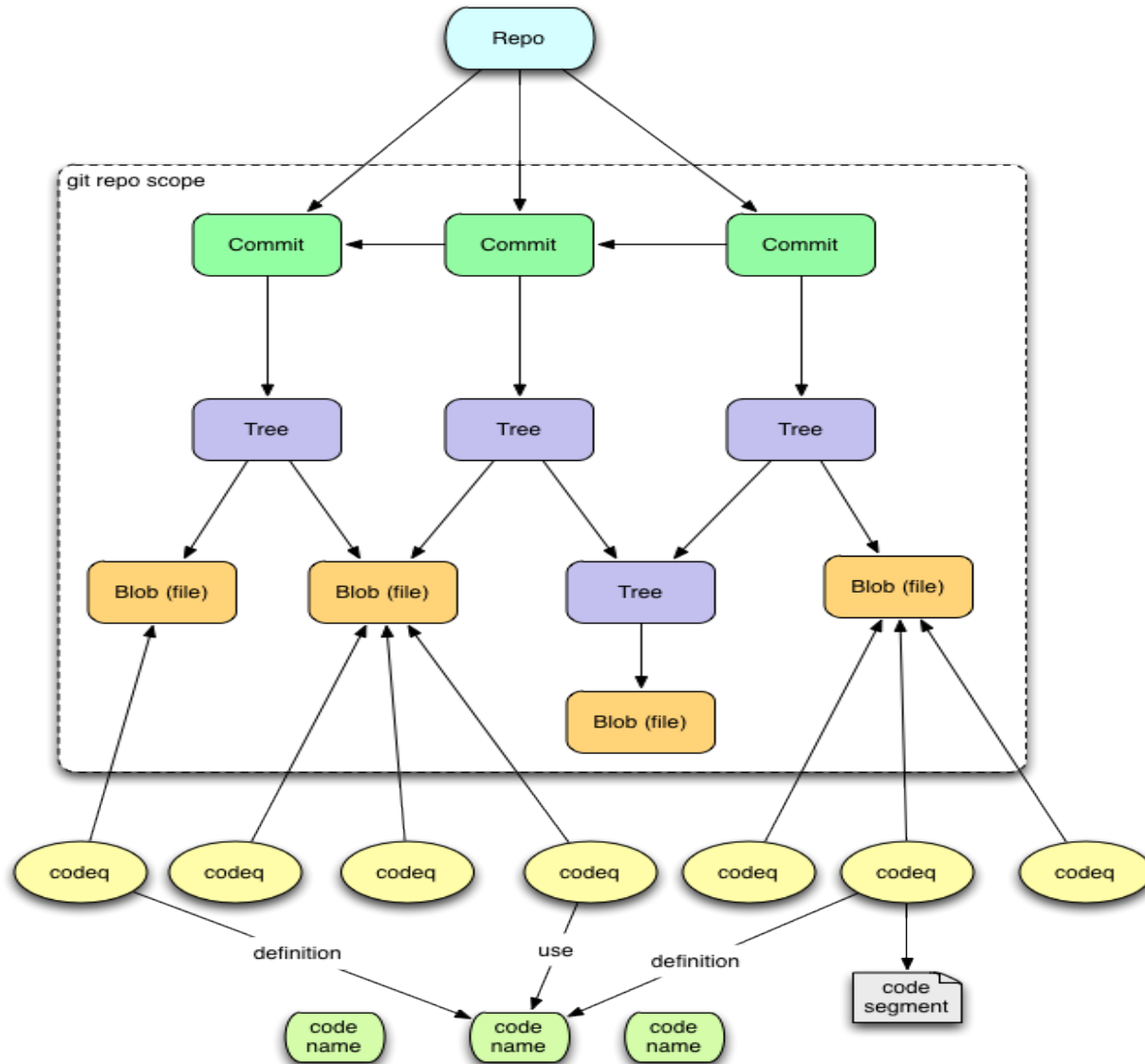
George {Kastrinis, Kollias}

# Marketing Idea

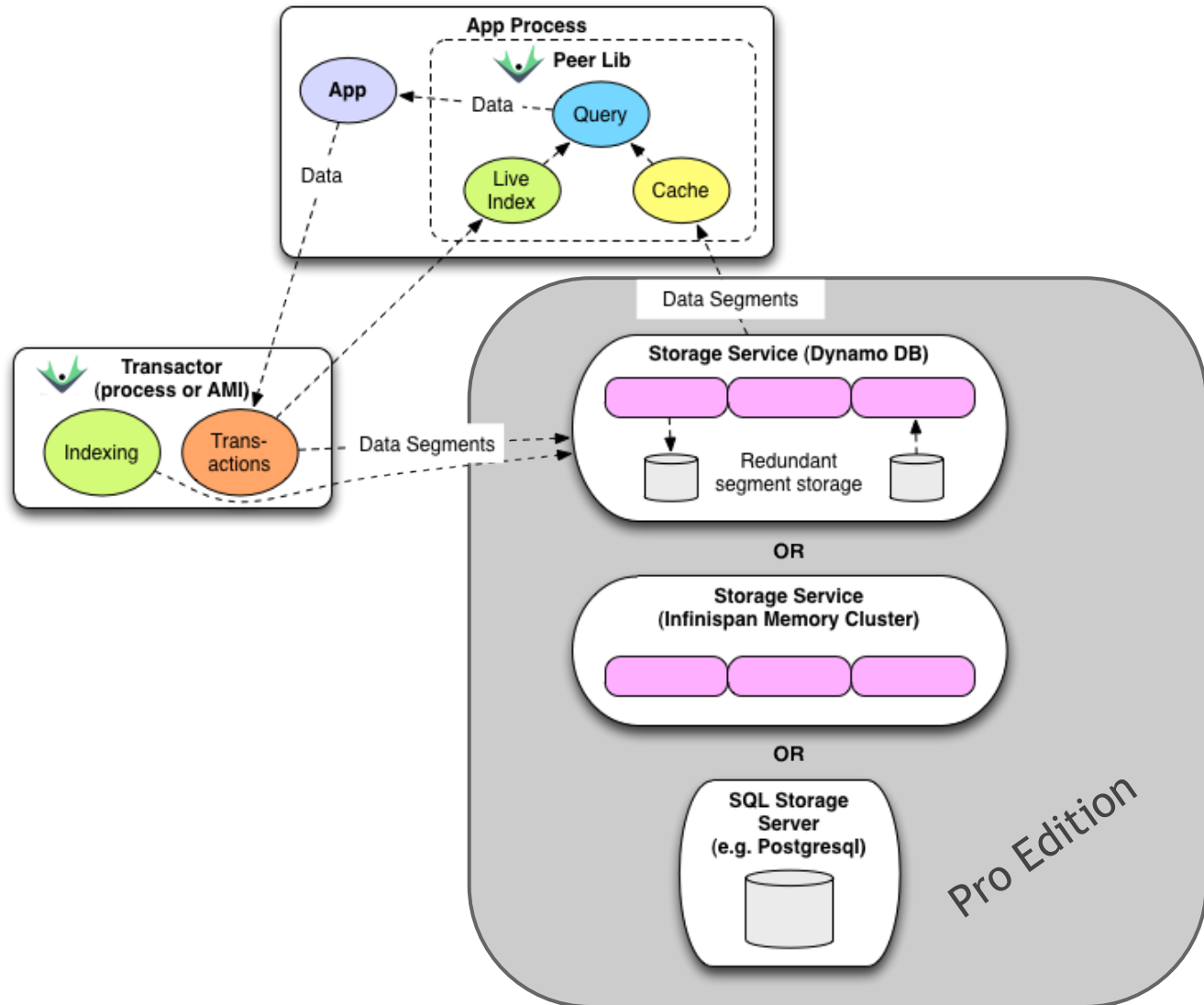- **Import** a code repo in a DB
  - DB = Datomic (uses a Datalog variant)
- **Analyze** the source files
- Issue **queries** on code history

VCS —understands→ text —building blocks→ files

codeq —understands→ code —building blocks→ code quanta (codeq)

# Codeq Overview
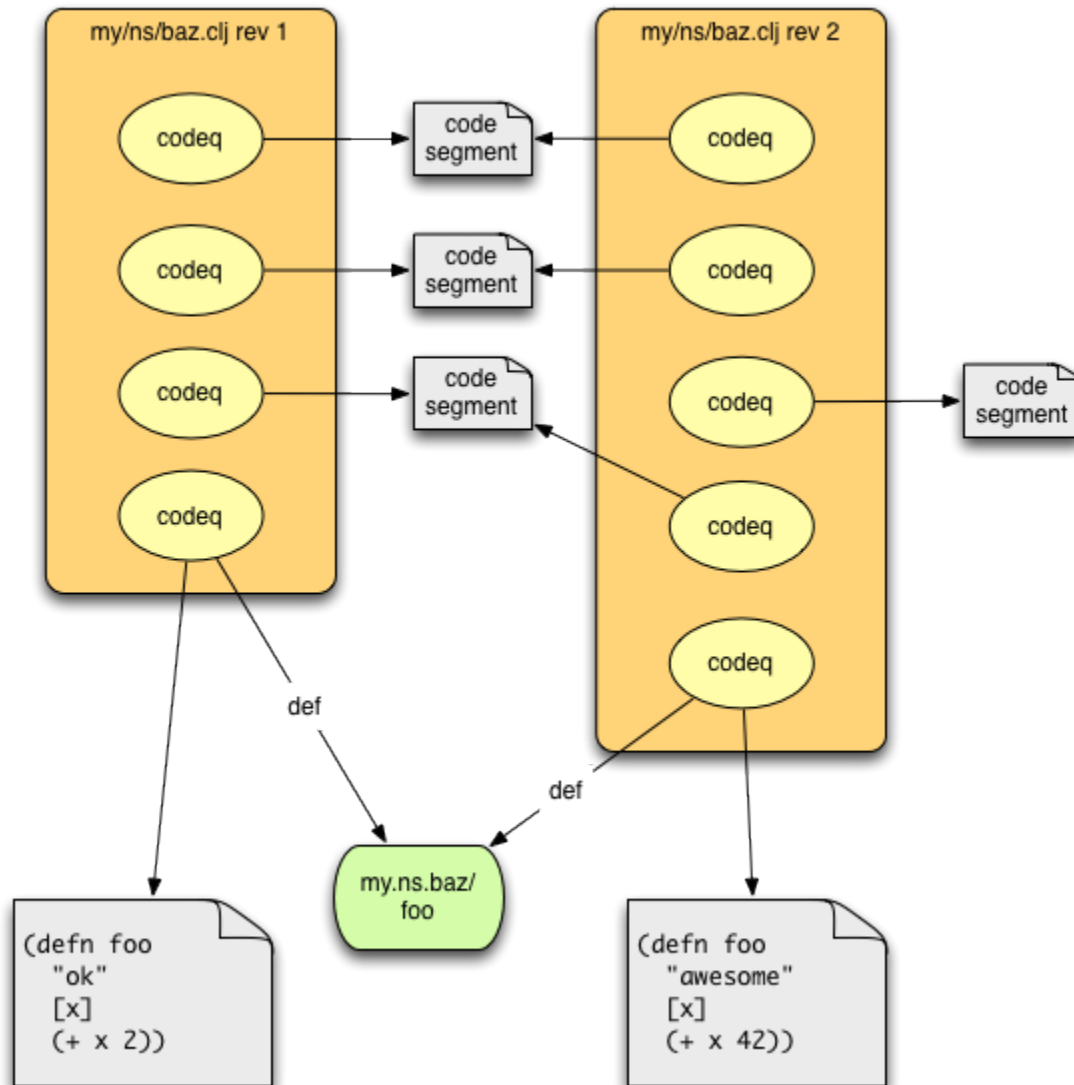
# Storing in Datomic

# The long story

# Importing git repos

- Parser for git repos structure

  - Finds the general file hierarchy

  - Imports all commits sequentially

  - 1:1 mapping of **commits** to **Datomic transactions**

    - Each Datomic transaction **annotated** with codeq commit's info

# Analyzing repo code

- **Language-specific** analyzers
  - Break each (source) file to **codeqs**
  - Determine level of **granularity**
  - Label each codeq with **semantic** information
  - 1:1 mapping of **analyses/modified files** to **Datomic transactions**
  - Written in Clojure

# Analyzing source file

# Codeq

- **Code Q**uantum (segment)

- Encodes the **name**, **code, location** & **semantic** information
  - The semantics depend on analyzer

- Must be uniquely identifiable
  - e.g. Global namespacing

# Datomic

- The underlying db

- Rules in a **Datalog** variant

- **Universal** Schema
  - **immutable** facts (*Datoms*) stored over time in the form:

*<entityID, attribute, value, transactionID, add/retract>*

# Datomic Versioning

> *conn.transact(list(":db/**add**", 1, ":intern/firstName", "George"))*

&lt;**1**, :intern/firstname, "George", **2**, add&gt;

&lt;**2**, :db/txInstant, 2013-06-14 14:05:04.0000, add&gt;

> *conn.transact(list(":db/**add**", 1, ":intern/firstName", "Kostas"))*

&lt;**1**, :intern/firstname, "George", **3**, retract&gt;

&lt;**1**, :intern/firstname, "Kostas", **3**, add&gt;

&lt;**3**, db/txInstant, 2013-06-14 14:06:04.0000, add&gt;

> *conn.transact(list(":db/**retract**", 1, ":intern/firstName", "Kostas"))*

&lt;**1**, :intern/firstname, "Kostas", **4**, retract&gt;

&lt;**4**, :db/txInstant, 2013-06-14 14:07:04.0000, add&gt;

We can get a DB value as of any specific point in time!
(using a specific transaction's timestamp)
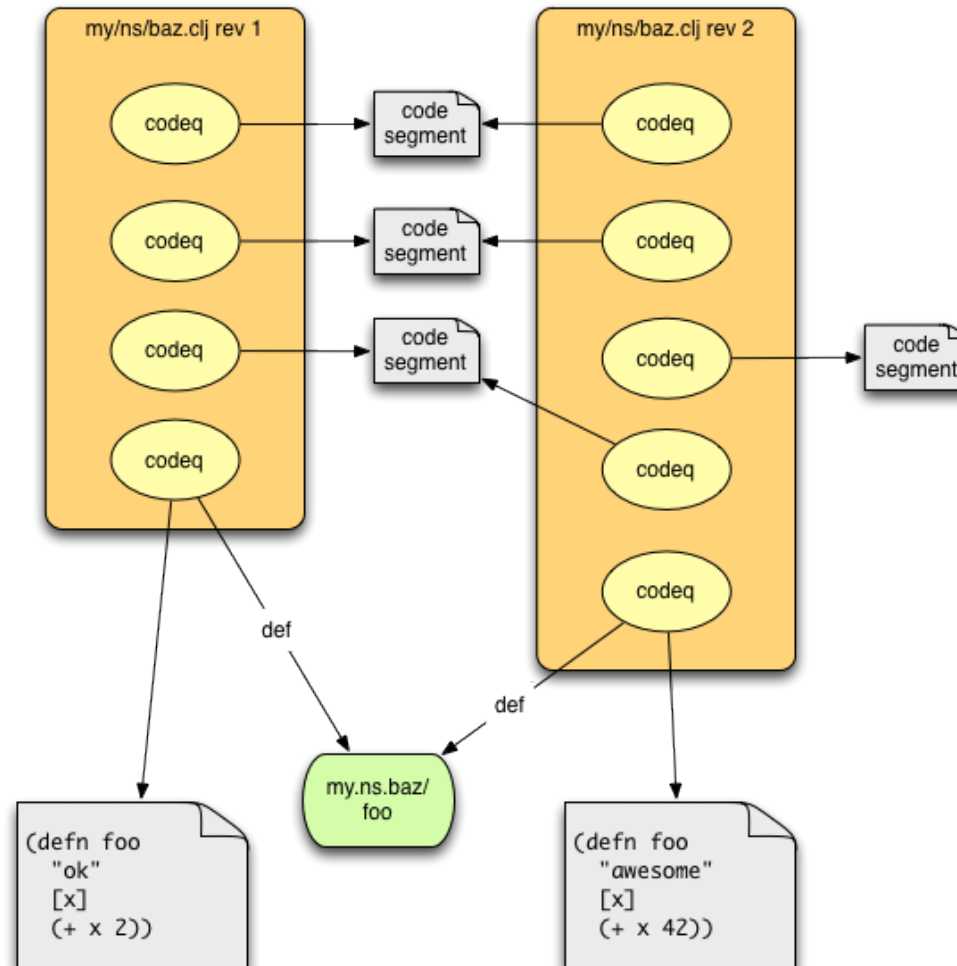
# Codeq Schema



codeq schema

# Codeq Versioning

- Reflects git's structure
  - Keeps **parent** commits, **date**, **authors**,…
- Needs a single total order
  - In effect, import **only the active** branch
  - Datomic supports linear history
- Must import **in order**
- Ignore any missing **dependencies**
  - User must import them manually

# Example Queries

- Find all the files referenced by a commit
- Find all the codeqs referenced by a commit
- Find all the commits including a file
- Find all the commits including a codeq
- Find all commits including a specific string
- Find all codeqs including a specific string
- …

# Real Example Query

Find all the different definitions of the function
*datomic.codeq.core/test* when **firstly introduced** (each def)

# Real Example Query (LogiQL)

Find all the different definitions of the function
*datomic.codeq.core/test* when **firstly introduced** (each def)

```
_[src] = mindate <-
agg<<mindate = min(date)>>(
    name[n] = "datomic.codeq.core/test",
    clj:def[n] = cq,
    codeq:codeSegm[cq] = cs, code:text[cs] = src,
    codeq:file[cq] = f, file:commits(f, c),
    commit:date[c] = date
).
```

# Real Example Query (Datomic)

Find all the different definitions of the function
*datomic.codeq.core/test* when **firstly introduced** (each def)

```
(d/q '[:find ?src (min ?date)
       :in $ % ?name
       :where
       [?n :code/name ?name]
       [?cq :clj/def ?n]
       [?cq :codeq/code ?cs]
       [?cs :code/text ?src]
       [?cq :codeq/file ?f]
       (file-commits ?f ?c)
       (?c :commit/authoredAt ?date)]
     db  rules "datomic.codeq.core/test")
```

```
[["(defn test\n  ..."
  #inst "2012-10-06T18:07:38.00"]
 ["(defn test\n  ..."
  #inst "2012-10-03T02:29:19.00"]
 ["(defn test\n ..."
  #inst "2012-10-03T19:09:39.00"]
 ["(defn test\n ..."
  #inst "2012-09-28T21:55:25.00"]
 ["(defn test\n ..."
  #inst "2012-10-02T21:38:38.00"]]
```

# Evaluation

# Language Support

- Must write an **analyzer** for each language
  - Only the **Clojure** one is provided
    - Handles only **top-level** Clojure structures (functions, namespace definitions, etc)
    - Doesn't parse top-level Call-Sites (totally **ignored**)
    - Very **simple** (~40 LOC)
  - There is a half-baked **3rd party** Java Analyzer

# Language Support - Thoughts

level of code understanding

- text
*(simply another CVS)*

- syntax
*(what codeq currently does)*

- semantics
*(might need expensive analysis of whole program)*

# Activity

- Main support only from their team

- Sporadic commits
  - Last useful commit, months old
  - Not a single issue (out of 8) addressed

- Weak user community

# Installation

- Download **Datomic** free edition

- Start Datomic's **Transactor**

- Build **codeq** with Clojure specific **Leiningen** build tool

- Run codeq over a **git repo**
  - sends the code segments to the Transactor

- Start querying!

# Usage Experience - Pros

- It is **simple**

- It just works

- Quite **fast**
  - Uses directly git capabilities

- **Incremental** support

- Makes both **git repos** & the actual **code history** queryable

# Usage Experience - Cons

- It is **simple**
- More of a **proof-of-concept** than an actual tool
- Naive Clojure Analyzer
- Tightly tied to Datomic
  - On an implementation level
  - On a concept level, any db would do
- Codeq redundancy even if unchanged
- Codeq name-changes end its history

# **M**ining **S**oftware **R**epositories

**commits & metadata** -> co-changing files -> frequent patterns between files

**commits** -> **AST (XML)** -> evolutionary change data

**file structure** & config files (filtering) -> detection of clones (similar code segments)

**source files** -> code duplication, code quality

**git API** -> queryable database (like codeq)

**crash reports** & **bug reports** & **commits** -> crash correlations, bug localization

**bytecode** -> class / packages graph & dependencies

source code & AST -> usage of dynamic features

commits & AST & **POS system** -> documentation from identifiers

**bug tracking system** & **IR** (information retrieval) techniques -> bug localization

**binaries** -> infer usage of libraries and other binaries