

# Code Navigation with DXR

Jonas Finnemann Jensen

Mozilla

August, 2012

(Watch this presentation at [air.mozilla.org](http://air.mozilla.org))

# Outline

- Introduction
- Index Optimized String Matching <sup>1</sup>
- Search Semantics
- Demo
- Questions

---

<sup>1</sup>Boring technical details for nerds, you've been warned.

# Introduction to DXR

DXR is a web based source code index, featuring:

- Syntactic code search (substrings, regular expressions)
- Semantic code search (find method, subclasses, etc).
- Cross referencing (Jump to...)
- File and directory listings

DXR is a replacement for MXR.

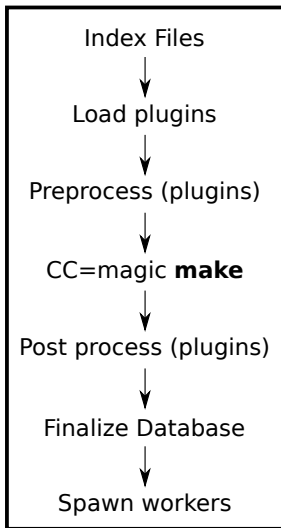
# DXR Architecture Overview



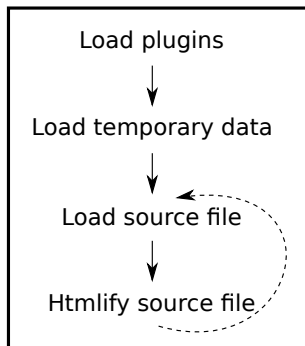
(We build HTML and database indexes offline).

# DXR Build Process

`./dxr-build.py --file ...`



`./dxr-worker.py ...`



```
$ hg clone mozilla-central  
$ ./dxr-build --file dxr.config  
$ tar -czf output.tar.gz ...  
$ scp output.tar.gz ...
```

# Index Optimized String Matching

## Problem

...WHERE text LIKE '%...%' and  
...WHERE text REGEXP '...'  
requires a **full table scan** in SQLite.

Enter trigrams...

## Definition (Trigram)

Given a text  $T$  a trigram of  $T$  is any substring of 3 characters in  $T$ .

e.g.  $\text{trigrams}('abcd') = \{'abc', 'bcd'\}$ .

## Solution

We create an index as a mapping  
 $\text{index} : \text{trigram} \rightarrow \text{doclist}$ , from  
trigrams to documents (ids).

To search for substring  $S$  we scan  
documents in the intersection  $C$  of  
doclists for trigrams of  $S$ .

$$C = \bigcap_{t \in \text{trigrams}(S)} \text{index}(t)$$

Notice that we still need to scan the  
text of every document in  $C$ .

ie. an inverted index using trigrams.

# The Inverted Index

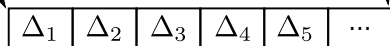
## Trigram integer encoding:

Trigram: 'aBc'  $\xrightarrow{\text{to lower}}$  'abc'  $\xrightarrow{\text{to int}}$  'a' + 'b' · 2<sup>8</sup> + 'c' · 2<sup>16</sup>

## Index b-tree (managed by sqlite)

```
CREATE TABLE %s_index (
  → trigram INTEGER PRIMARY KEY, /* Alias for rowid */
  doclist BLOB /* Binary blob */
);
```

## Delta List Encoding

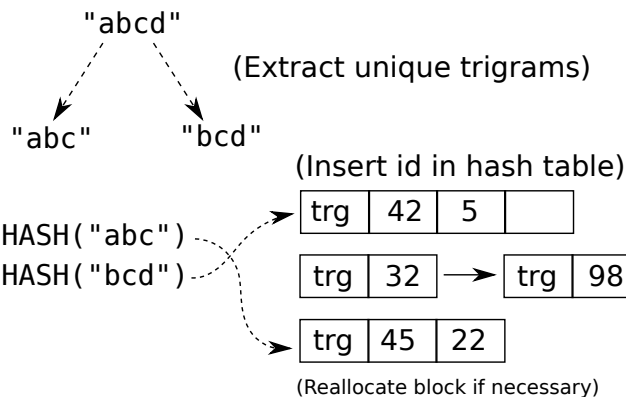


(Each  $\Delta$  uses a VarInt encoding)

$$id_n = \sum_{i=0}^n \Delta_i$$

## Building the Inverted Index

### INSERT



### COMMIT

Flush hash table to b-tree (ie. %s\_index)



# Extracting Trigrams From Regular Expressions

**Goal:** Extract required trigrams.

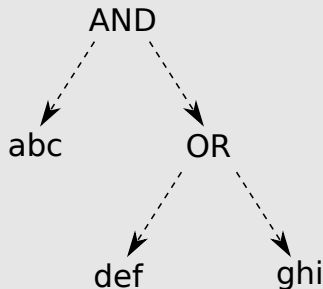
## Methods:

- (1) Recursive analysis of regular expression<sup>2</sup>  
(How Google Code Search Worked)
- (2) Analysis of underlying automaton<sup>3</sup>  
(As in recent WIP patch for PostgreSQL)

TriLite uses (1) as implemented by re2<sup>4</sup>.  
(An RE engine by guy who did Code Search)

## Example

Regex: `/abc(def|ghi)/`



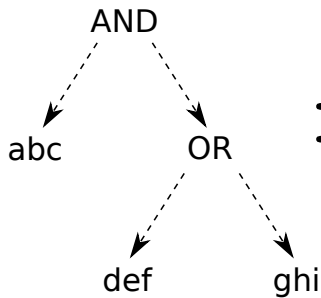
<sup>2</sup><http://swtch.com/rsc/regexp/regexp4.html>

<sup>3</sup><http://www.pgcon.org/2012/schedule/events/383.en.html>

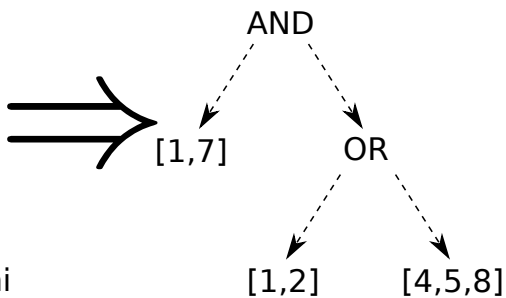
<sup>4</sup><http://code.google.com/p/re2/>

# Merging DocLists (1)

Expression Tree

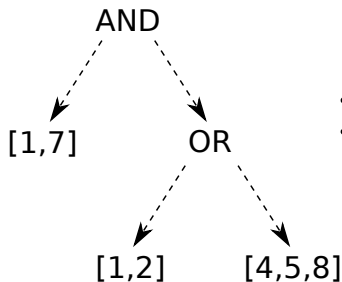


Fetch all doclists

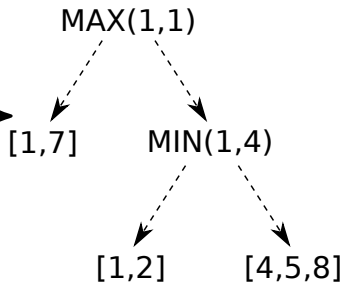


## Merging DocLists (2)

Fetch all doclists

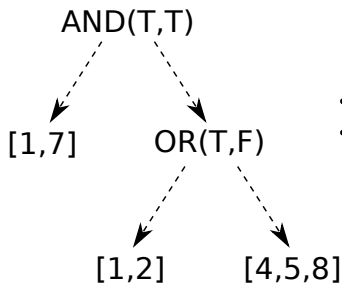


Find candidate id

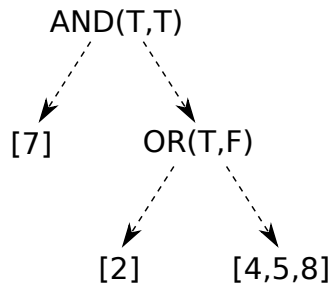


# Merging DocLists (3)

Test if 1 is a result

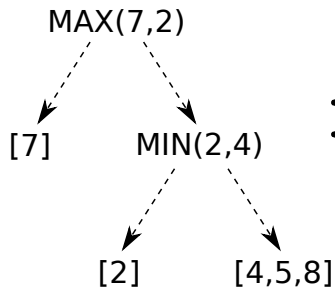


Advanced to > 1

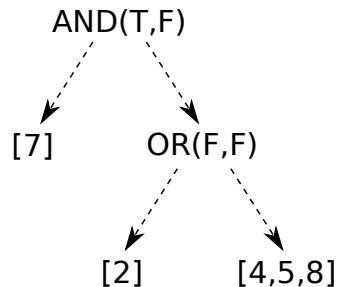


# Merging DocLists (4)

Find candidate id

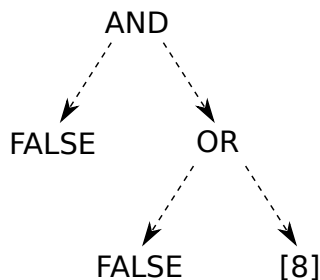


Test if 7 is a result

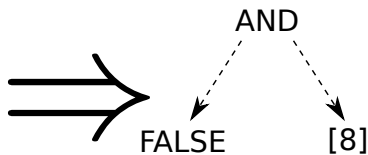


# Merging DocLists (5)

Advanced to > 7



Simplify on-the-fly (1)



# Merging DocLists (6)

## Simplify on-the-fly (2)

FALSE

This approach allows us to

- Skip more than one id at the time
- Terminate when there's no more solutions
- Merge doclists on-the-fly

# TriLite Module for Sqlite

## TriLite Features:

- Substring matching
- Regular expression matching
- Return start/end offsets of matches
- Scans the text after joins
- Protection from evil regular expressions.

TriLite is still in development, and **not ready** for production.



# Search Semantics

## Results $\subseteq$ Documents

(Only lines matched are returned)

Substrings are case sensitive.

Space denotes **AND**.

Dash **negates** a term (or argument).

Quotes searches for substrings with space.

**Example:** `"int main" Hello regexp:#(W|w)orld# -ext:h -ext:js`

## Special Parameters

regexp	regexp:/nsZip.*/
path	path:startupcache/
ext	ext:cpp
type	type:ptr
function	function:draw

For more see advanced query options.

Note: `regexp:` works like vim, pattern must start and end with the same character.

# Demonstration

# dxr.allizom.org

(allizom, mozilla spelled backwards)

# Questions

Ask away or drop by

# #static

at [irc.mozilla.org](https://irc.mozilla.org)

DXR: <https://github.com/mozilla/dxr>

TriLite: <https://github.com/jonasfj/trilite>