

Programming with Millions of Examples

**Alon Mishne, Hila Peleg, Sharon Shoham,
Eran Yahav, Hongseok Yang**

Components are everywhere

facebook.

HIBERNATE

spring
source



ANDROID

eclipse
PLUGIN CENTRAL

django

GRAILS

iOS 6

Struts

jQuery
write less, do more.

Components are Accessed Via APIs

```
class File {  
    File(String);  
    void open();  
    void close();  
    byte read();  
    void write(byte);  
}
```

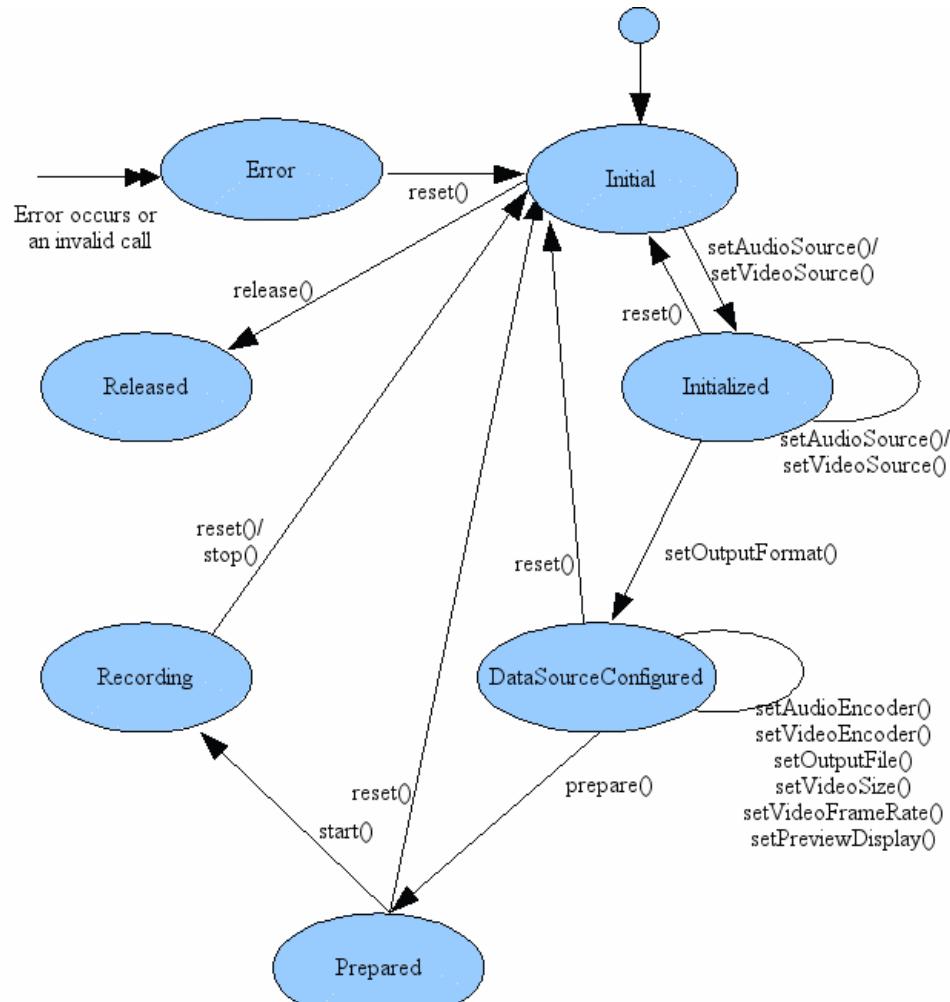
```
void writeBuff(File f, byte[] buff) {  
    f.open();  
    for (byte b : buff) f.write(b);  
    f.close();  
}
```

Component APIs are Complicated



There is only one thing more painful than learning from experience and that is not learning from experience.
– Archibald MacLeish

Component APIs are complicated



MediaRecorder state diagram

From the official Android documentation for MediaRecorder

(<http://developer.android.com/reference/android/media/MediaRecorder.html>)

Examples are prevalent

GitHub

Explore Features Enterprise Blog Sign up Sign in

Search android.media.MediaRecorder Search

We've found 10,873 code results

10,873 Sort: Best match

Repositories 10,873 Code Issues 2 Users

yumekochan/VoiceTalkTest – MediaRecorderTest.java Java

```
import android.media.MediaRecorder;
```

matteofini/MediaNote – MyRecorder.java Java

```
import android.media.MediaRecorder;
```

eldoth/gesprek – UsuarioServiceImpl.java Java

```
import android.media.MediaRecorder;
```

Sort: Best match

Languages

Java	7,858
HTML	1,600
XML	75
PHP	44
JavaScript	7
Scala	6
Haxe	2
Python	2
INI	2
Erlang	2

Advanced Search Cheat Sheet

Stats: > 3.5M users, over 8M repositories, growing daily

Examples are prevalent

The screenshot shows a Stack Overflow question page. The title of the question is "Retrieve column names from java.sql.ResultSet". Below the title, there is a snippet of Java code demonstrating how to retrieve column names using `ResultSetMetaData`. The code is as follows:

```
ResultSet rs = stmt.executeQuery("SELECT a, b, c FROM TABLE2");
ResultSetMetaData rsmd = rs.getMetaData();
String name = rsmd.getColumnName(1);
```

Below the code, a note states: "and you can get the column name from there."

Answers:

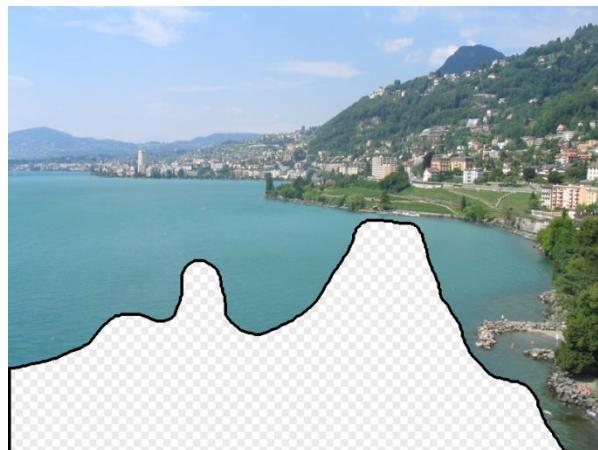
- 10 Upvotes, question by user 10: "With `java.sql.ResultSet` is there a way to get a column's name as a `String` by using the column's index? I had a look through the API doc but I can't find anything." (This answer is highlighted with a blue background).
- 19 Upvotes, question by user 19: "See `ResultSetMetaData`" (e.g. `ResultSet rs = stmt.executeQuery("SELECT a, b, c FROM TABLE2"); ResultSetMetaData rsmd = rs.getMetaData(); String name = rsmd.getColumnName(1);`)

Stats: > 2.3M registered users, > 5.6M questions, > 10M answers

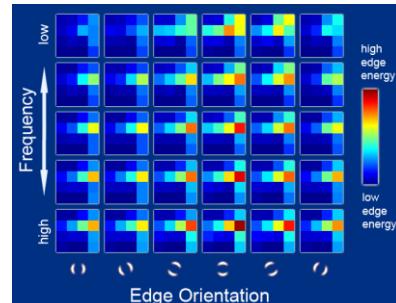
Challenge

How do we leverage the millions of examples to synthesize client code using the library?

Scene Completion from Millions of Photographs



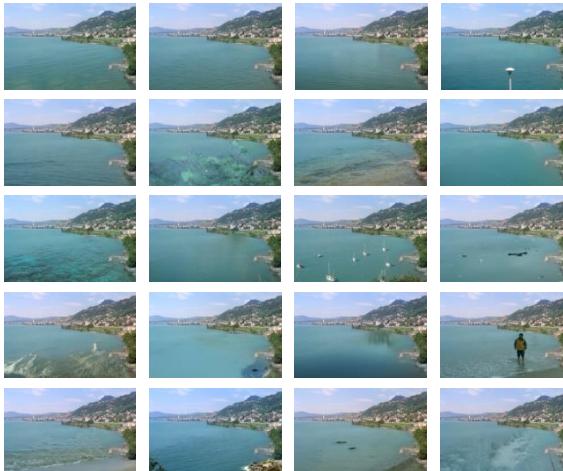
Query



Semantic Descriptor



Semantic Index



20 completions



Context matching
+ blending

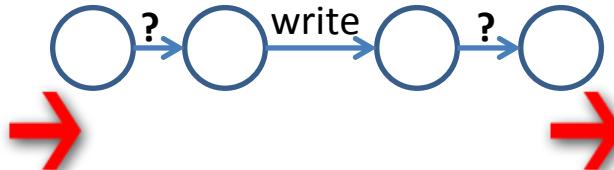


200 matches

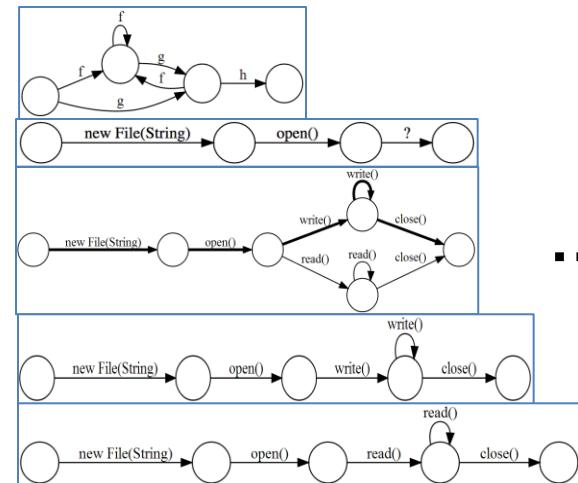
Code Completion from Millions of Examples

```
File f = ?  
f.write("73");  
f.?
```

Query



Semantic Descriptor



Semantic Index

```
File f = new File();  
f.open();  
//while(y.data < 42) {  
//    f.write("73");  
//    y = y.next;  
//}  
f.close();
```

```
File x = new File();  
x.open();  
while(y.data < 42) {  
    x.write(y);  
    y = y.next;  
}  
x.close();
```

1 of 20 completions

Context matching
+ blending

file f = new File(); f.open(); //while(y.data < 42) { // f.write("73"); // y = y.next; //} f.close();	file f = new File(); f.open(); while(y.data < 42) { f.write("73"); y = y.next; } f.close();	file f = new File(); f.open(); while(y.data < 42) { x.write(y); y = y.next; } x.close();	file f = new File(); f.open(); while(y.data < 42) { y.write(x); y = y.next; } x.close();	file f = new File(); f.open(); while(y.data < 42) { x.write(y); y = y.next; } x.close();
file f = new File(); f.open(); //while(y.data < 42) { // f.write("73"); // y = y.next; //} f.close();	file f = new File(); f.open(); while(y.data < 42) { f.write("73"); y = y.next; } f.close();	file f = new File(); f.open(); while(y.data < 42) { x.write(y); y = y.next; } x.close();	file f = new File(); f.open(); while(y.data < 42) { y.write(x); y = y.next; } x.close();	file f = new File(); f.open(); while(y.data < 42) { x.write(y); y = y.next; } x.close();

200 matches

Your wish is my command,
do you know **what is your wish?**

Discovery

Android location services

```
LocationManager x = ?  
Location y = ?  
y.getLastKnownLocation(?)
```

```
LocationManager manager =  
    (LocationManager) context.getSystemService(LOCATION_SERVICE);  
String provider = manager.getBestProvider(criteria, true);  
boolean gpsEnabled =  
    locationManager.isProviderEnabled(GPS_PROVIDER);  
boolean networkEnabled =  
    locationManager.isProviderEnabled(NETWORK_PROVIDER);  
...  
Location latestLocation = null;  
if (networkEnabled)  
    latestLocation = manager.getLastKnownLocation(GPS_PROVIDER);  
else if (gpsEnabled)  
    latestLocation = manager.getLastKnownLocation(NETWORK_PROVIDER);  
...  
return latestLocation;
```



Query

1 of 20 matches

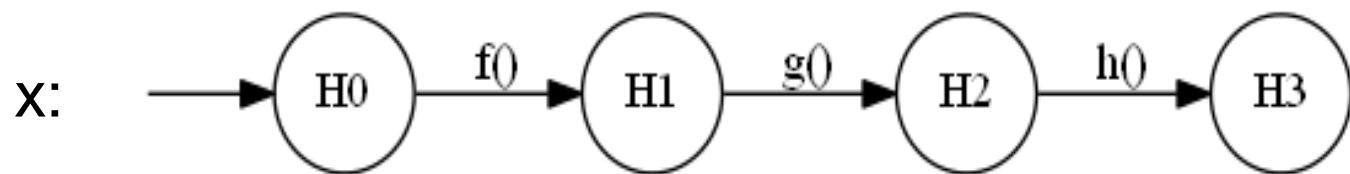
Naïve specification does not lead to “common” solution

Semantic Index: what do we want?

- Find semantically similar programs
 - Compare program behaviors
 - *Represent temporal ordering of method invocations across multiple objects*
- Represent **partial information** due to partial programs
- Recover information using **consolidation**
- Feasible index construction and comparisons
- Precise enough to distinguish good from bad

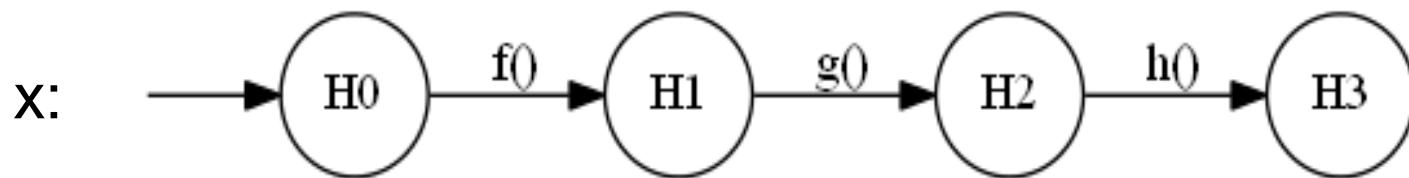
Concrete history

```
public void method(Something x) {  
    x.f();  
    x.g();  
    x.h();  
}
```



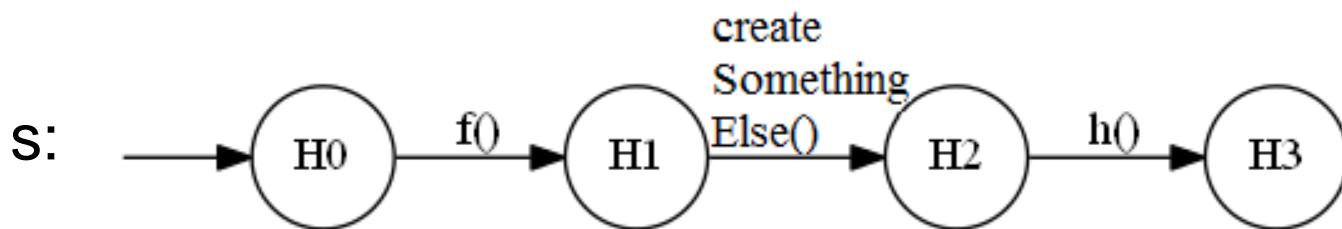
Aliasing

```
public void method(Something x) {  
    x.f();  
    BaseOfSomething y = x;  
    y.g();  
    y.h();  
}
```



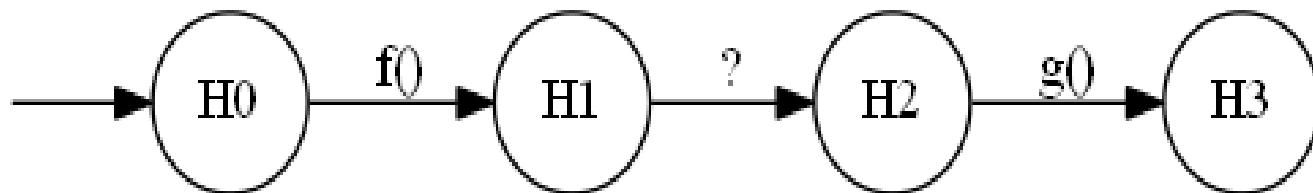
Creation context

```
public void method(Something x) {  
    x.f();  
    SomethingElse s = x.createSomethingElse();  
    s.h();  
}
```



Dealing with the unknown

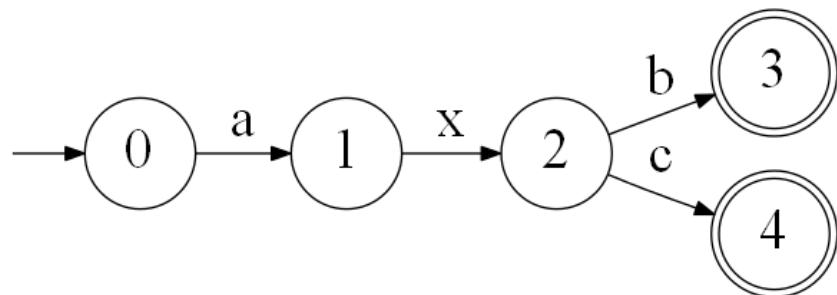
```
public void method1() {  
    Something x = new Something();  
    x.f();  
mystery(x);  
    x.g();  
}
```



Representing Histories: DSA

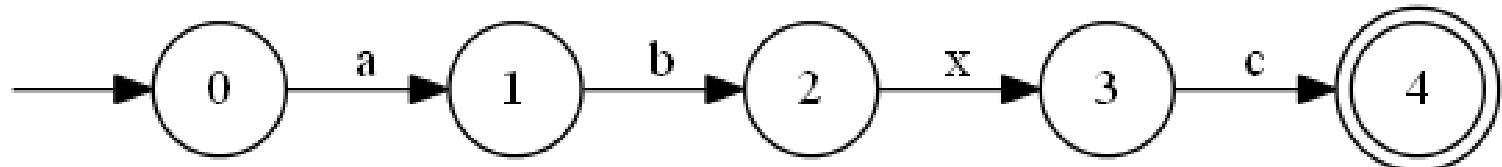
A Deterministic
Symbolic Automaton is
a tuple $(\Sigma; Q; \delta; \iota; F; \text{Vars})$

- Σ is a finite alphabet
- Q is a finite set of states
- δ is the transition relation,
 $Q \times (\Sigma \cup \text{Vars}) \rightarrow Q$
- $\iota \in Q$ is the initial state
- $F \subseteq Q$ is the set of final
states
- **Vars is the finite set of
variables**

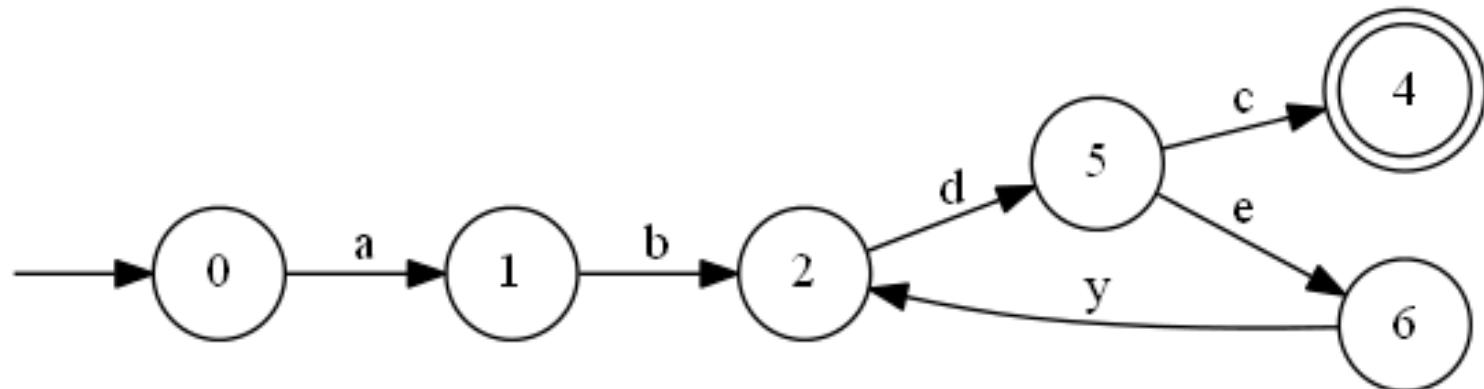


Assignment

An assignment σ maps a variable x in context $sw1$, $sw2$ (sw_1, x, sw_2) to a non-empty symbolic language



$$\sigma(\epsilon, x, \epsilon) = d(eyd)^*$$

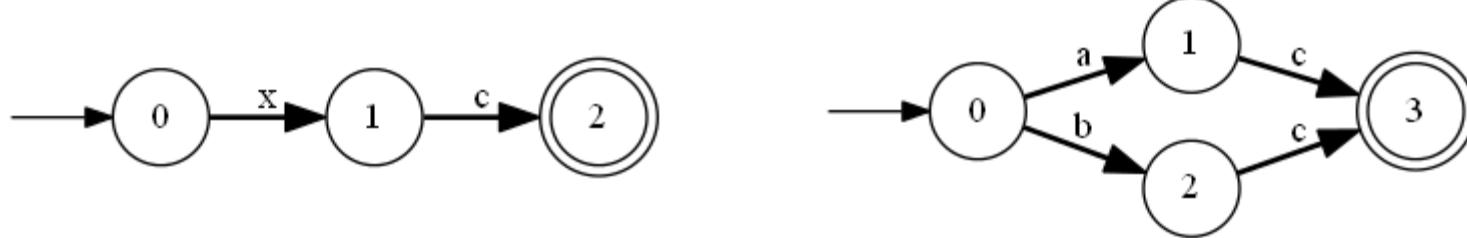


Desired Operations

- Analysis of individual example
 - Extend DSA with new call (known / unkown)
 - Join (later)
- Operations on multiple examples
 - Consolidation
 - Query matching

Order Between DSAs

- A natural way to define order between automata is language inclusion
- Won't work (directly) for symbolic automata:



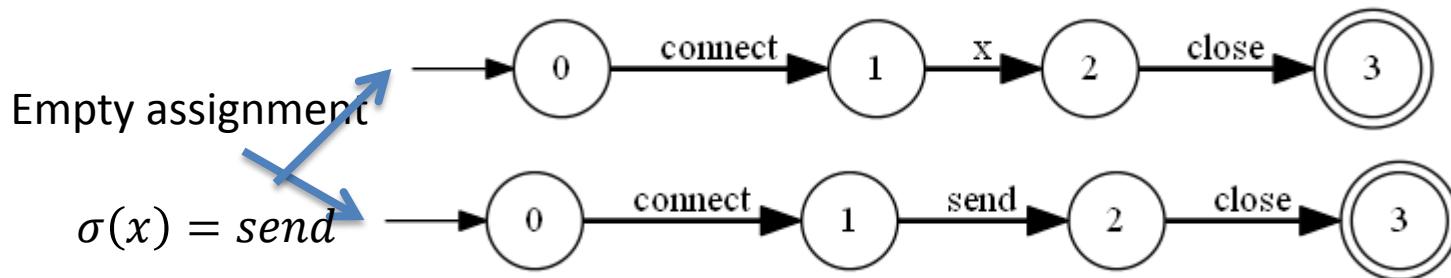
$$\{xc\}$$



$$\{ac, bc\}$$

Partialness

- A word is *less partial* (or *more complete*) than another if it represents a more concrete scenario

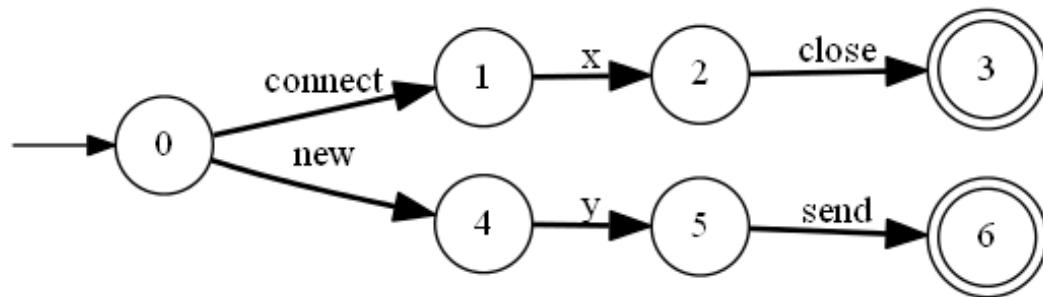


- Formally: w_1 is more partial than w_2 if for each assignment σ_2 to w_2 there is an assignment σ_1 to w_1 s.t. $\sigma_1(w_1) = \sigma_2(w_2)$

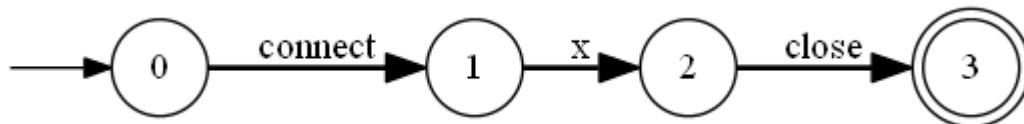
Partial Order

- We define the order over DSAs to capture both axes:
 - Precision: the natural concept of language inclusion
 - Partialness: of the individual words
- Intuitively: a DSA is smaller if it is more precise and more partial

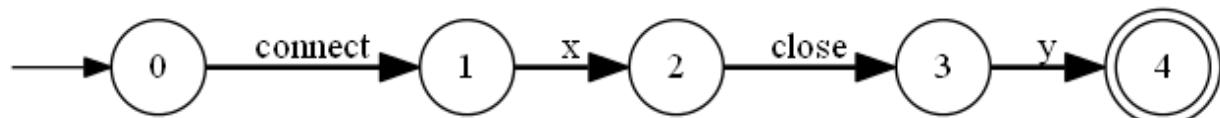
Precision vs. Partialness



The less *precise* you are, the more behaviors you describe



And the more *partial* you are, the more behaviors you describe...



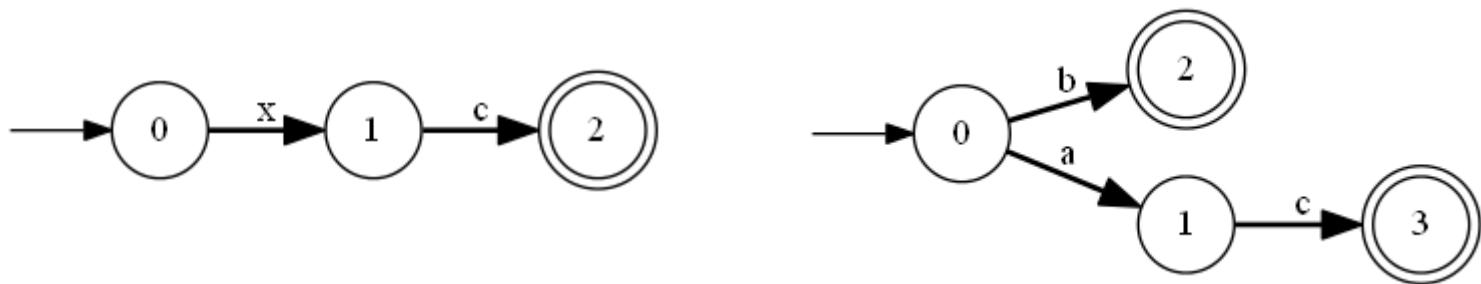
\leq

DSA Order

$A_1 \leq A_2$ if

$$\forall \sigma_2 \exists \sigma_1 L(\sigma_1(A_1)) \subseteq L(\sigma_2(A_2))$$

where σ_1 and σ_2 are concrete assignments

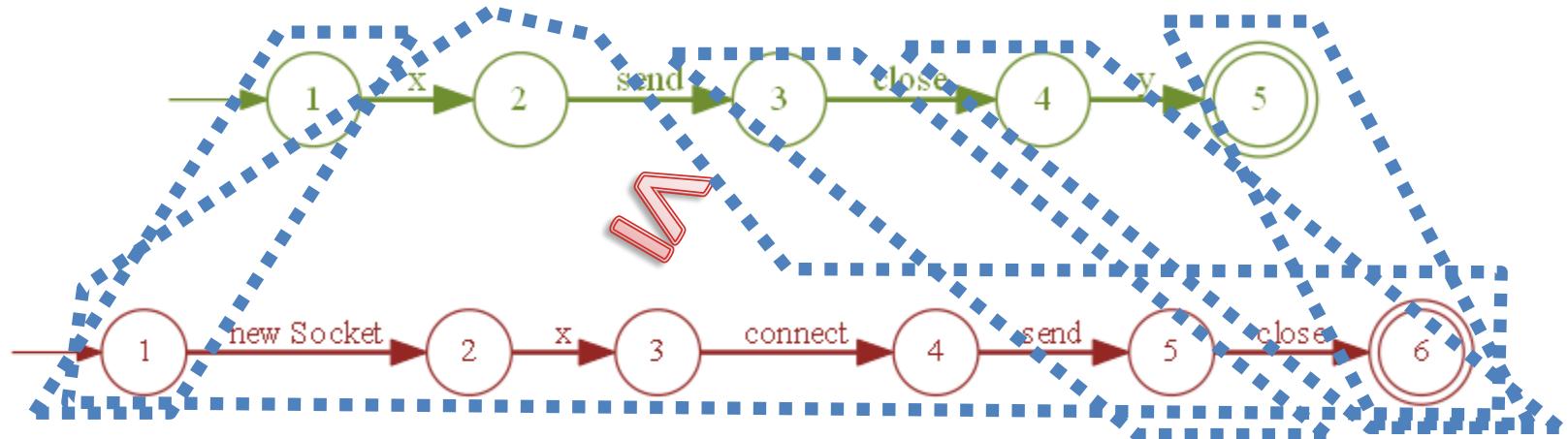


$$\sigma(x) = a, \text{ so } \{ac\} \subseteq \{ac, b\}$$

Calculating Inclusion via Simulation

- Want to determine whether $A_1 \leq A_2$
- Adapt DFA simulation to DSAs: **symbolic simulation**
 - Find pairs of **one state from A_1** (smaller) and a **set of states from A_2** (larger) that witness structural inclusion
 - Collect possible candidates using outgoing transitions
 - Shows us the matching word for each word in A_1
- DFA simulation already captures the notion of precision
- DSA simulation adds the notion of partialness
 - Symbols can “**swallow**” parts of the other DSA

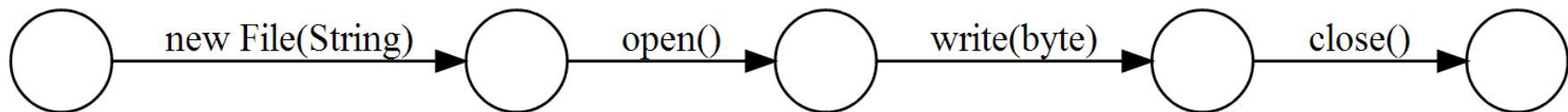
Simulation Example



- ✓ $(1,\{1\}), \quad (4,\{6\}),$
 $(2,\{1,2,3,4,5,6\}), \quad (5,\{6\})$
 $(3,\{5\}),$

Index Construction: Client-Side Static Analysis

```
1  File f = new File(filename);  
2  f.open();  
3  f.write(z);  
4  f.close();
```



- Interprocedural
- Alias-aware with points-to information

Example

How should I use a
java.nio.channels.SocketChannel?

Analyzing a Single Code Sample

```
example1() {
```

```
    SocketChannel sc = SocketChannel.open();
```

```
    sc.configureBlocking(false);
```

```
    sc.connect();
```

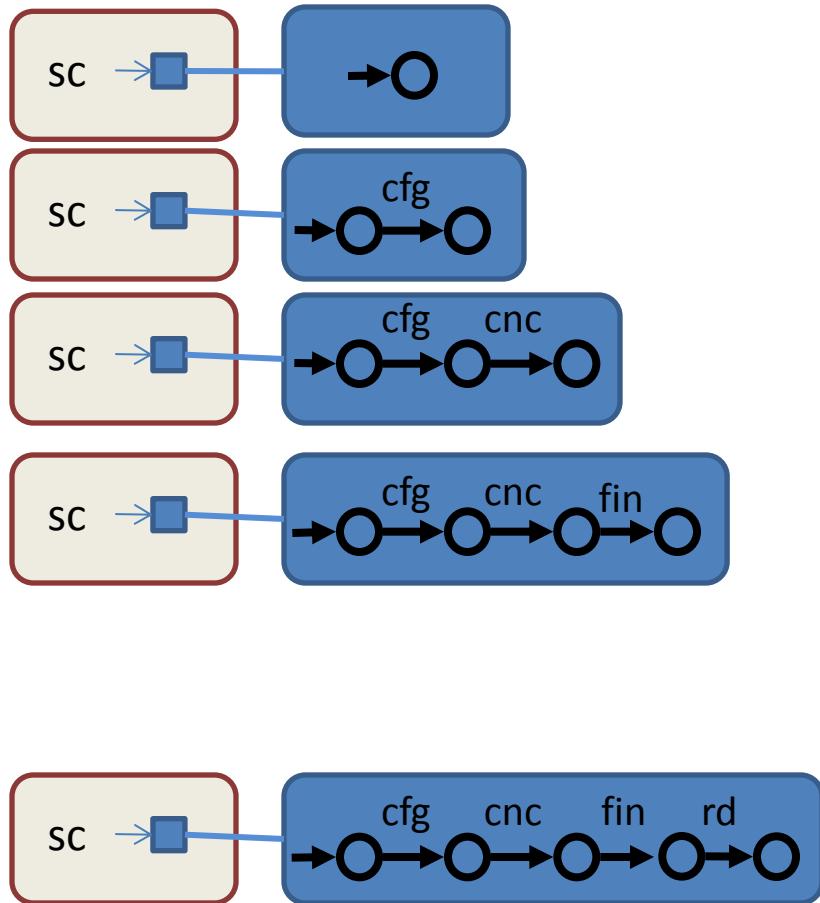
```
    sc.finishConnect();
```

```
    ByteBuffer dst = ...;
```

```
    sc.read(dst);
```

```
}
```

Heap History



```

SocketChannel createChannel (...) {
{
    SocketChannel sc =
        SocketChannel.open();
    sc.configureBlocking(false);
    return sc;
}

```

```

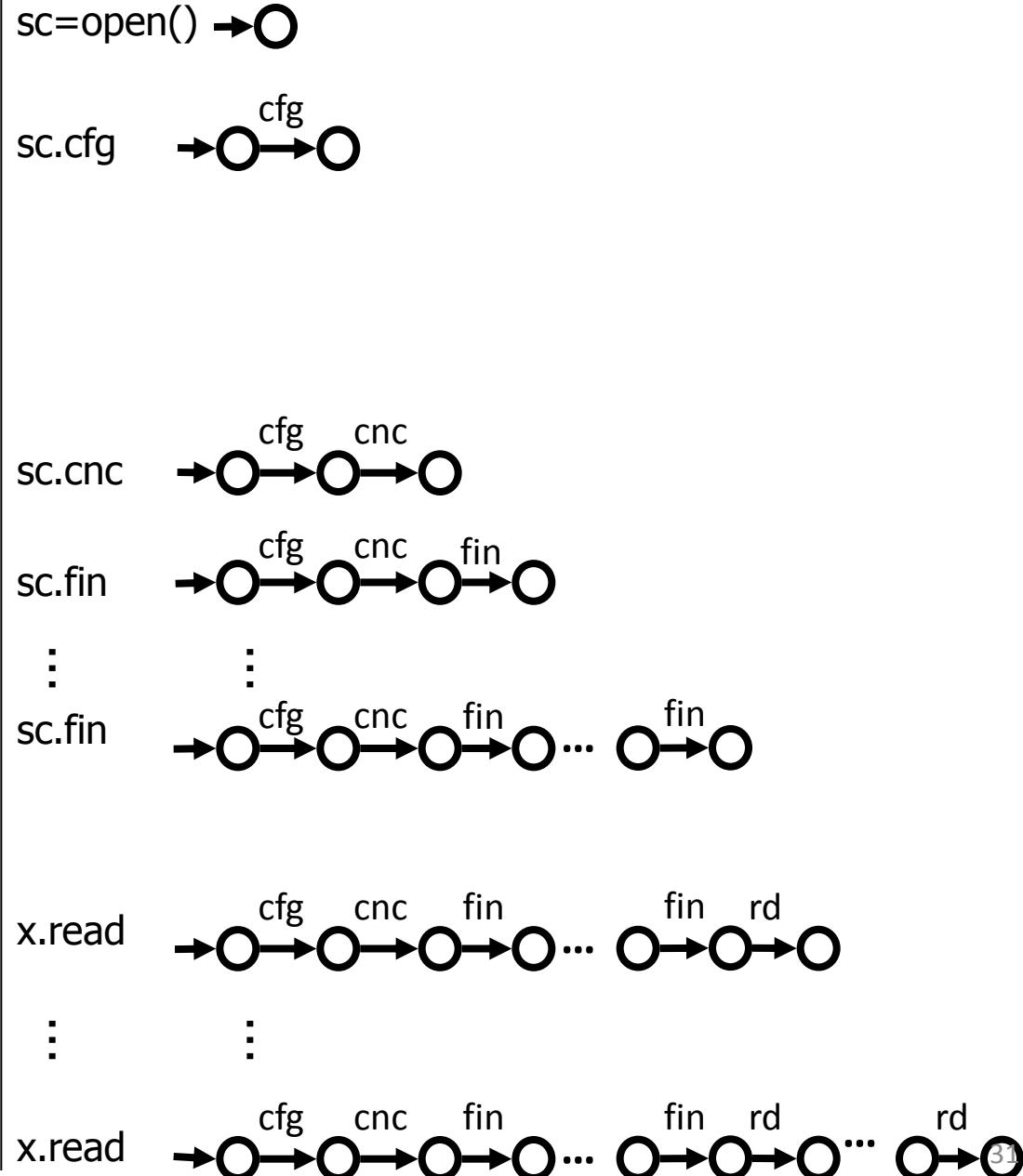
void example() {
    Collection<SocketChannel> chnls =
        createChannels();
    for (SocketChannel sc : chnls){
        sc.connect(new ...);
        while (!sc.finishConnect()) { ... }
        if (?) { receive(sc); }
        else { send(sc); }
    }
    closeAll(channels);
}

```

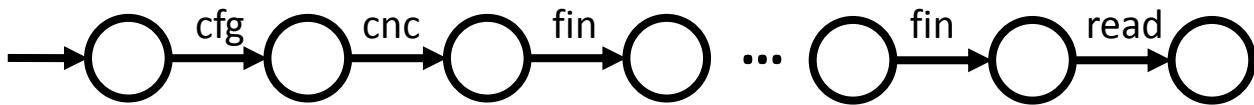
```

void receive(SocketChannel x) {
...
    while (numBytesRead >= 0) {
        numBytesRead = x.read(dst);
        fos.write(dst.array());
    }
...
}

```



History Abstraction

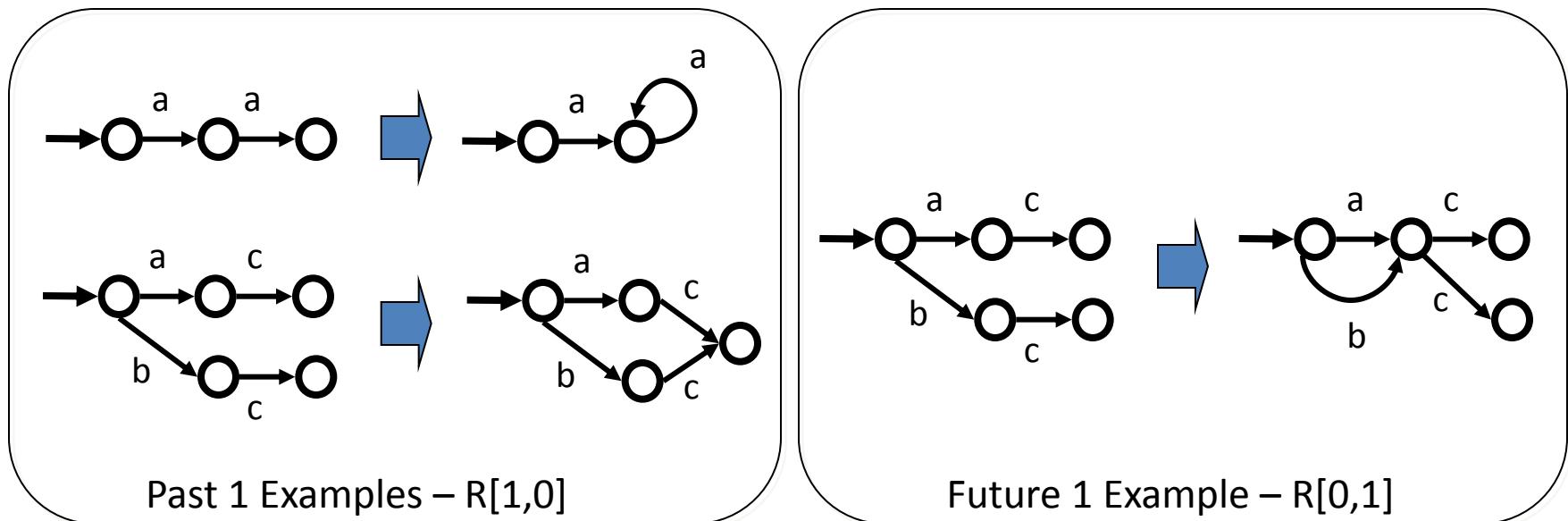


- Abstract history
 - DSA over-approximating unbounded histories
- **Quotient-based abstractions** for history
 - DSA states which are equivalent w.r.t. a given **equivalence relation R** are merged

History Abstraction

- Past-Future Abstraction

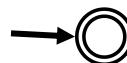
$(q_1, q_2) \in R[k_{in}, k_{out}]$ if q_1 and q_2 share both an incoming sequence of length k_{in} and an outgoing sequence of length k_{out}



Abstract Semantics

- Initial abstract history
 - empty sequence automaton
- When an API method is invoked
 - history extended: append event and construct quotient

sc = open



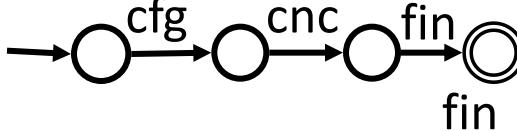
sc.config



sc.connect

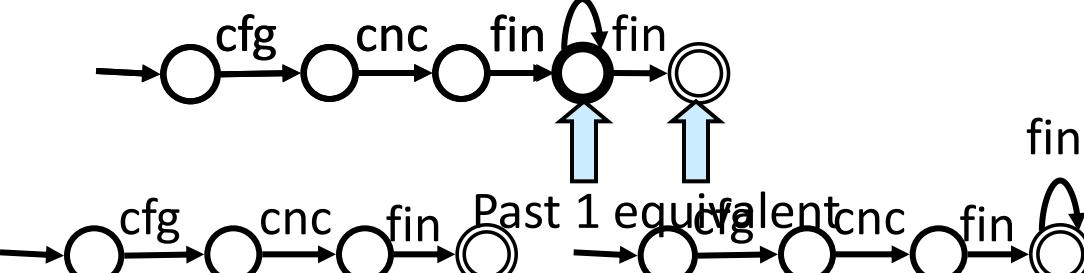


while (!sc.finCon) {



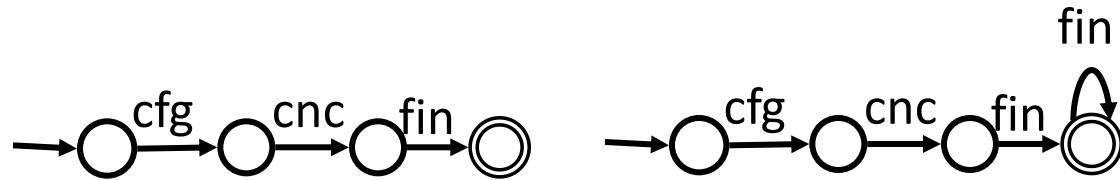
fin

} //endof while



Join

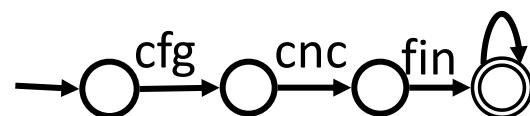
endof while



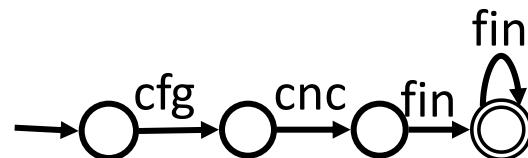
union



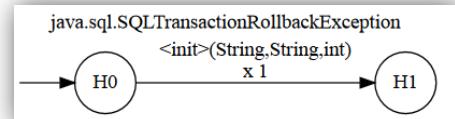
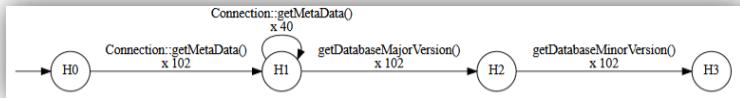
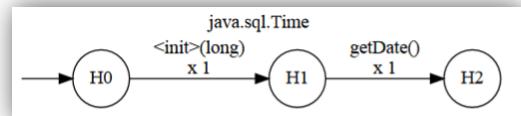
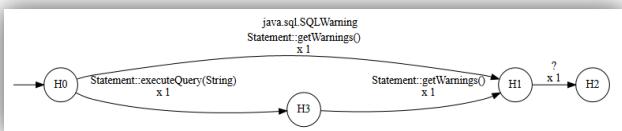
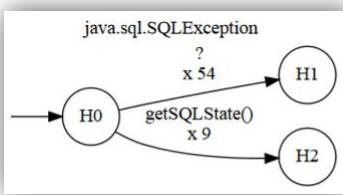
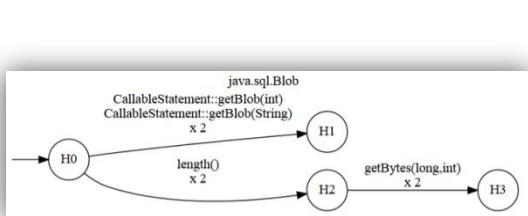
fin



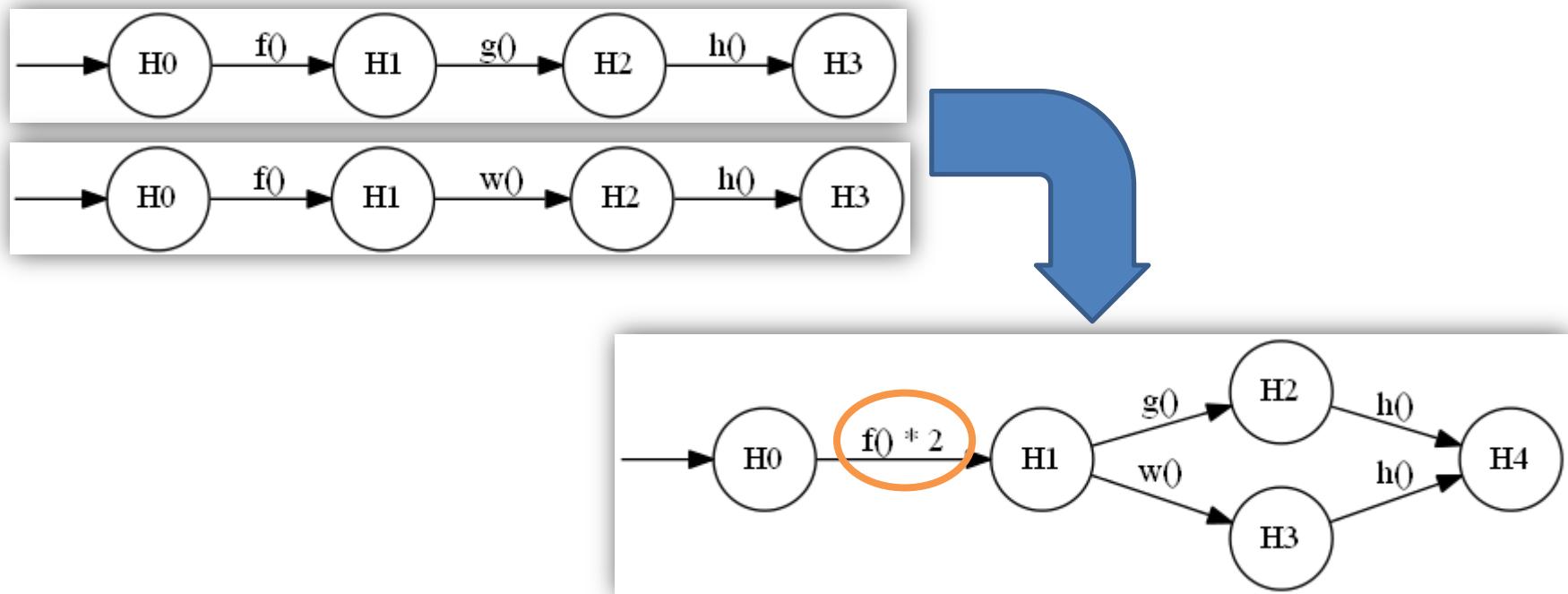
quo



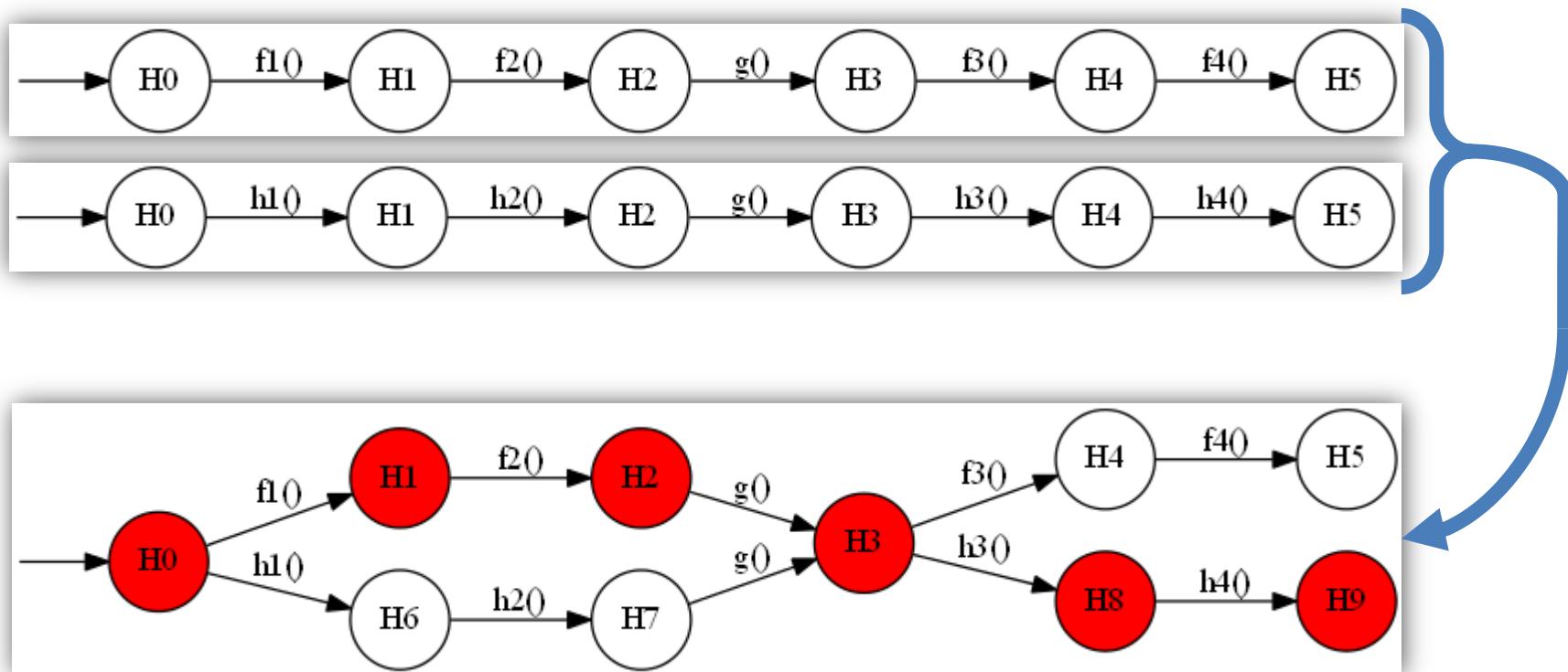
Consolidating the Index



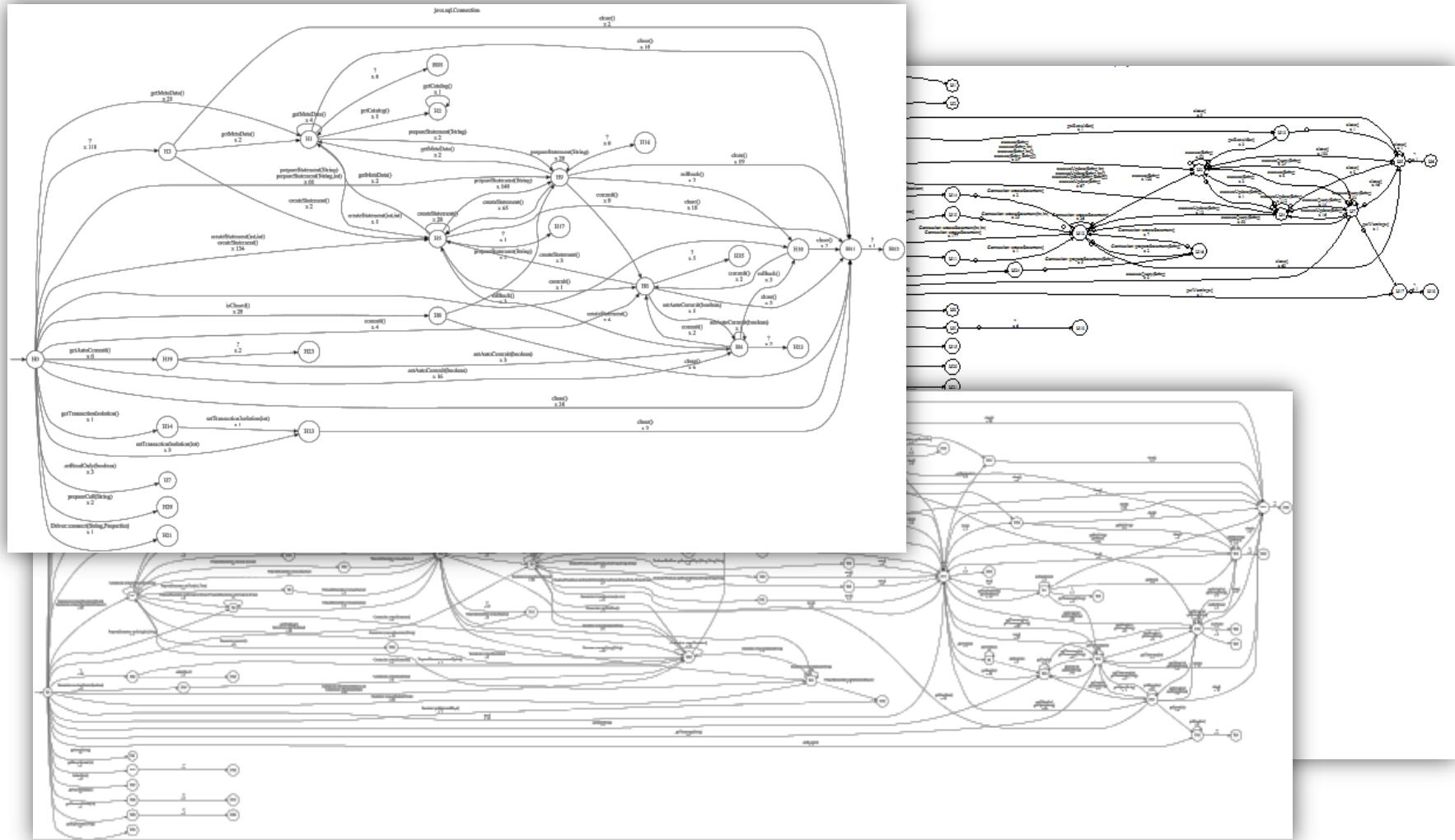
Merge Same Type Together



Merge All Samples of Same Type

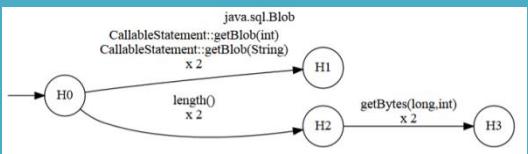


Merging all Samples of Same Type

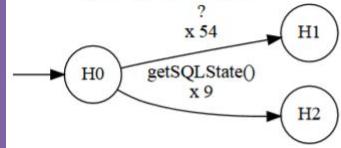


Merge by Use Case

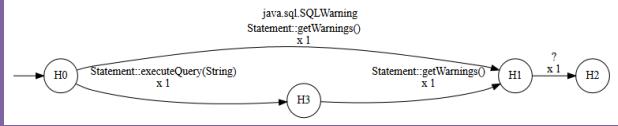
Use case 3



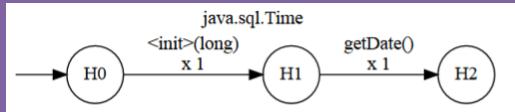
java.sql.SQLException



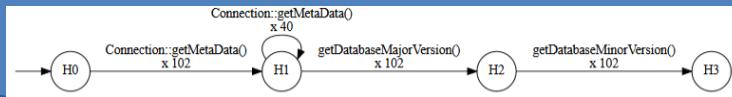
Use case 1



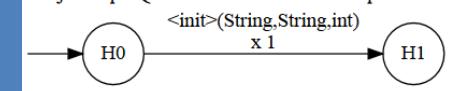
java.sql.Time



Use case 2

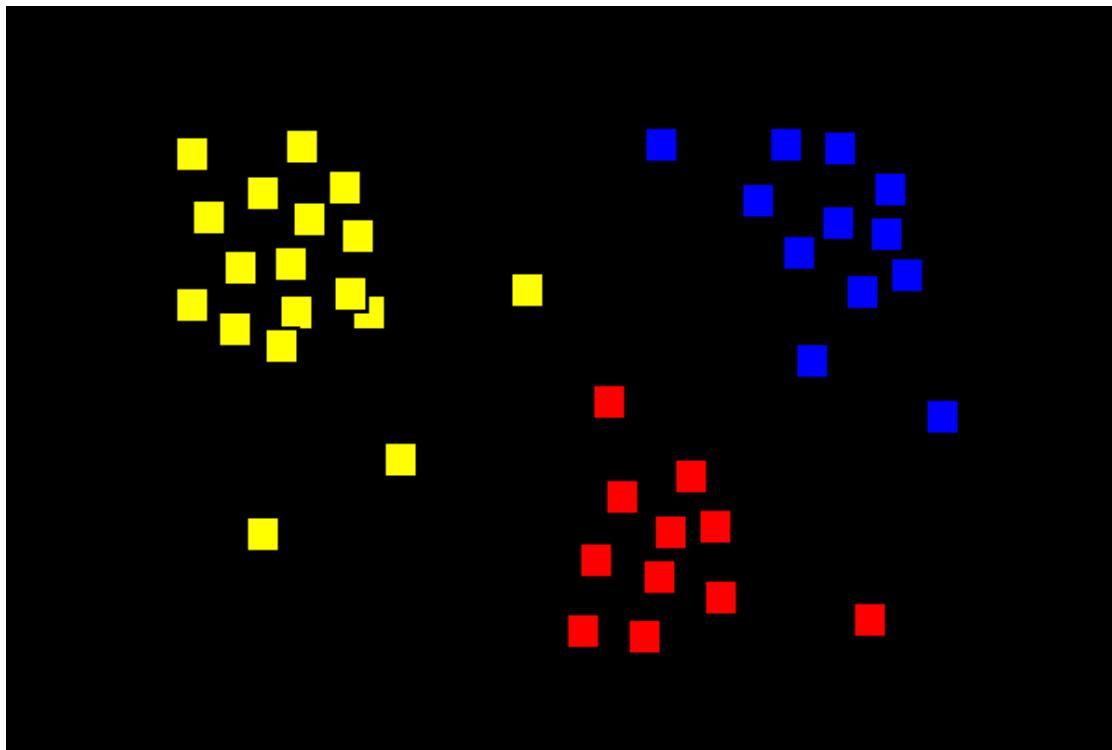


java.sql.SQLTransactionRollbackException



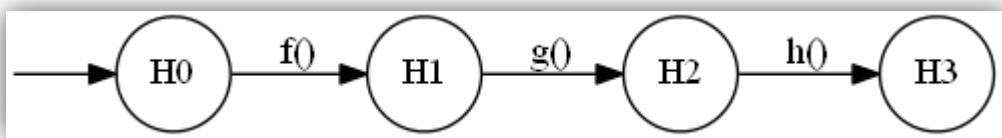
But how?

Clustering

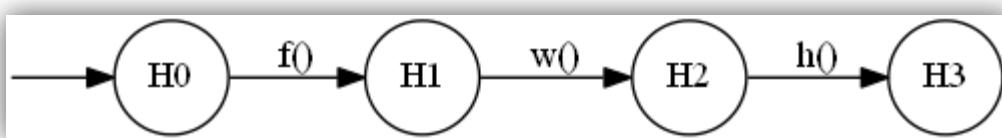


Example distance function

Different use-cases typically use different methods



f()	g()	w()	h()
1	1	0	1



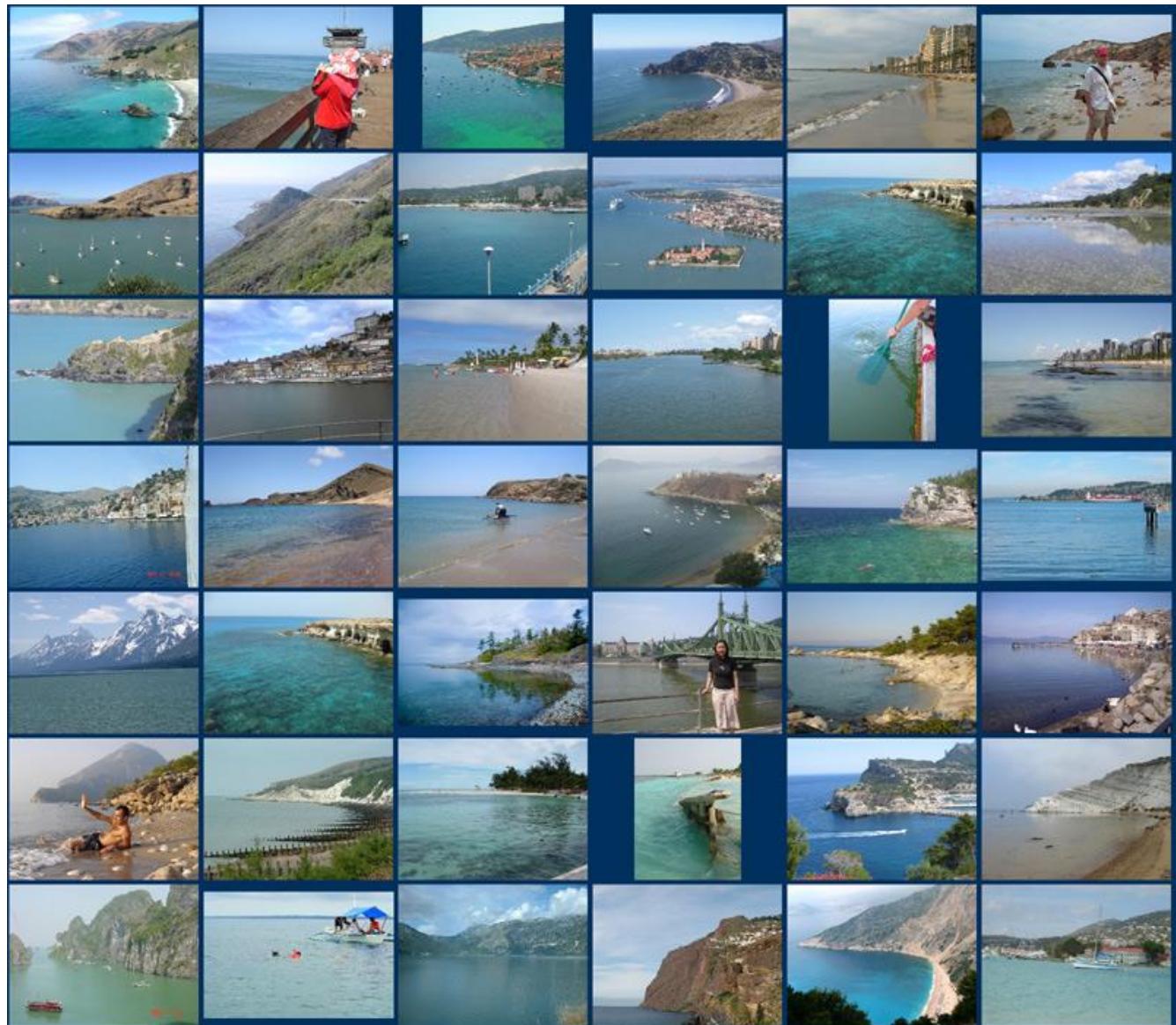
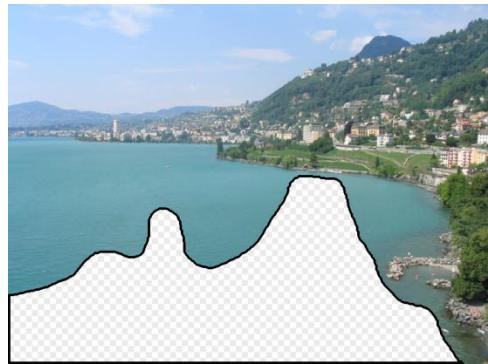
f()	g()	w()	h()
1	0	1	1

$$\text{Distance} = \sqrt{(1-1)^2 + (1-0)^2 + (0-1)^2 + (1-1)^2} = \sim 1.41$$

Clustering Results

	# samples	% samples	size
▶ java.sql.Statement	396	18	40
	239	60	5
	77	19	7
	32	8	28
	27	7	5
	9	2	8
	6	2	6
	6	2	7
▶ java.sql.Timestamp	62	3	22
	29	47	3
	13	21	9
	7	11	7
	7	11	10
	6	10	12

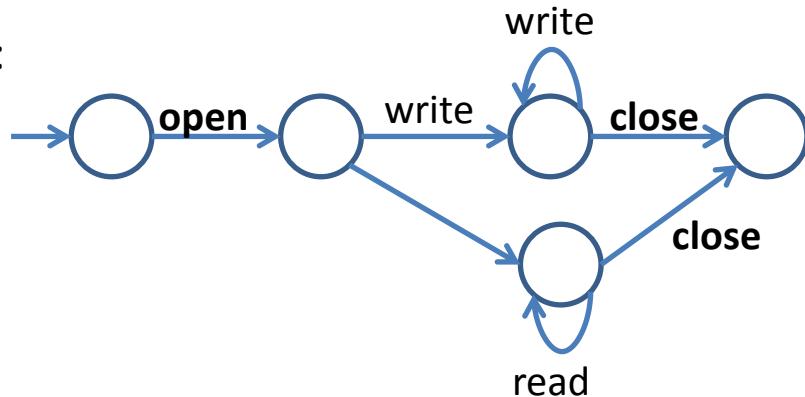
Query Matching



Query Matching



Index (partial):



Using unknown elimination we can get for example (assignment σ):

$$\sigma(var) = \begin{cases} \text{open}, & var = x \\ \text{write}^* \cdot \text{close}, & var = y \end{cases}$$

```
Public void print(String s) {  
    File f = new File("x");  
    String path = f.getAbsolutePath();  
    f.write(s);  
    f.close();  
    record(path);  
}
```

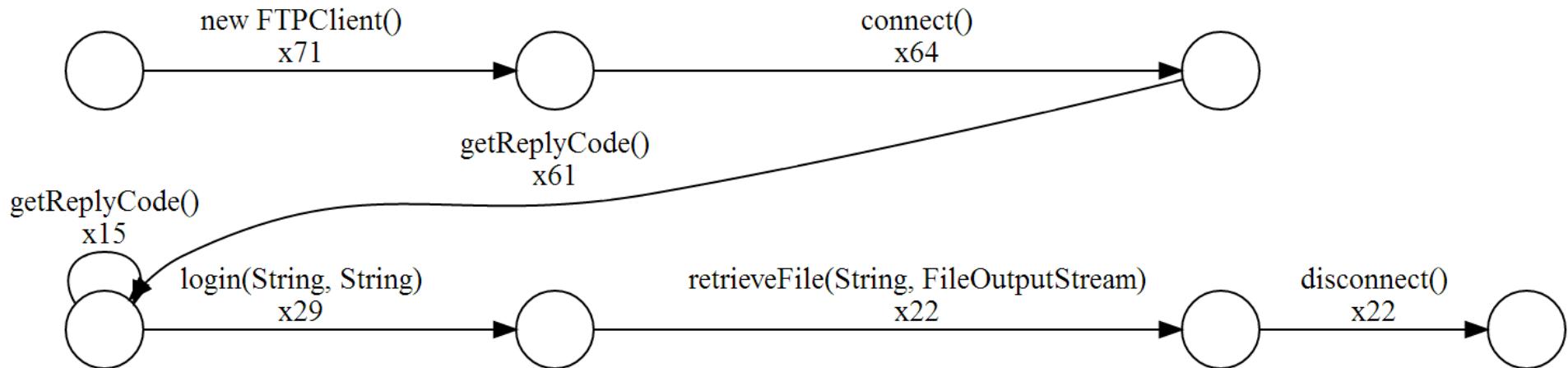
```
File f = ?  
f.write(_);  
f.?
```

```
Public void print(String s) {  
    File f = new File("x");  
    if (s != null)  
        f.write(s);  
    else  
        f.write("error");  
    f.close();  
}
```

:

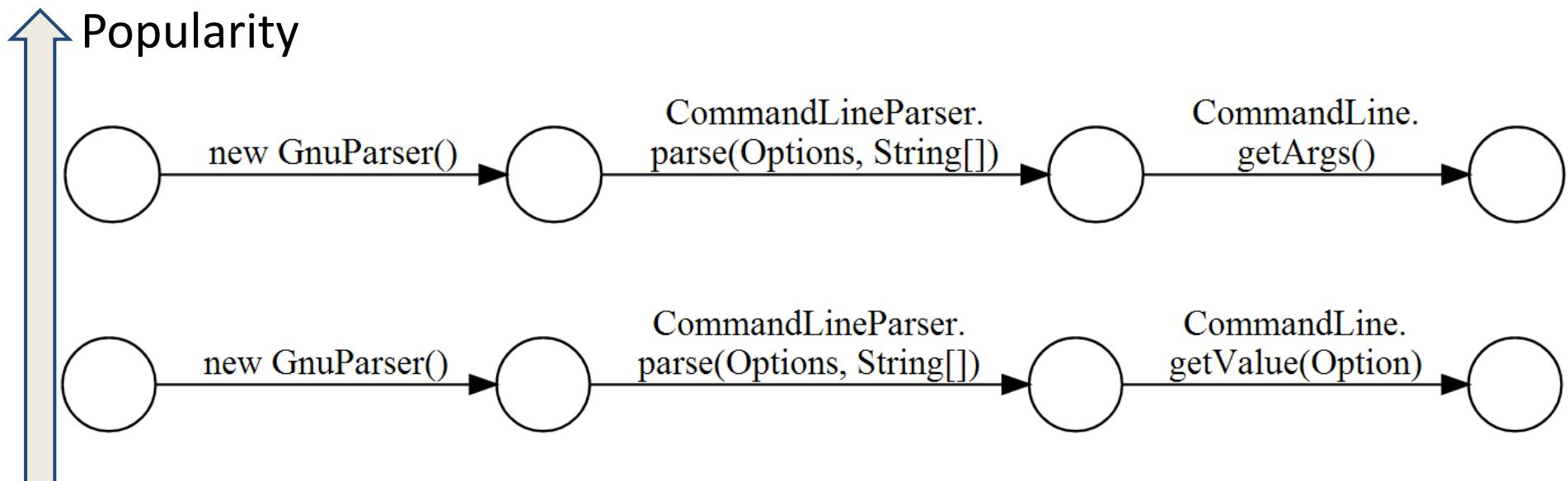
Example: Complete Sequence

```
FTPClient c = ?;  
c.login("username", "password");  
c.?;  
c.retrieveFile("file path", out);
```



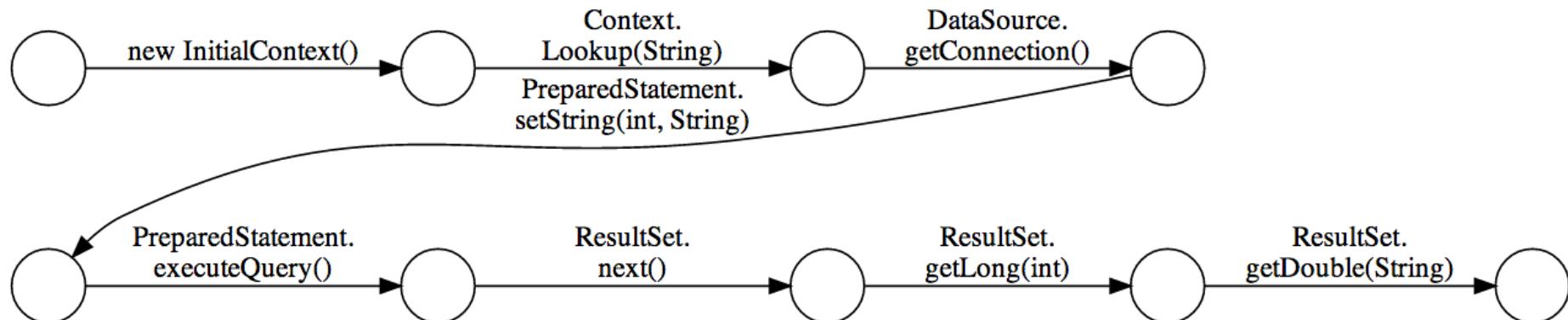
Example: Query Over Multiple Types

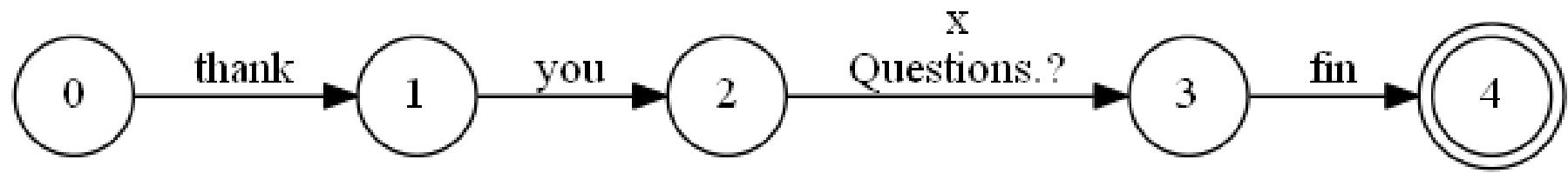
```
GnuParser p = new GnuParser();  
p.?;  
String s = p.?;
```



Example Result: Long Query

```
InitialContext ic =  
new InitialContext();  
ic.?;  
PreparedStatement ps = ic.?;  
ResultSet rs = ps.executeQuery();  
rs.?;  
double d = rs.getDouble("column");
```

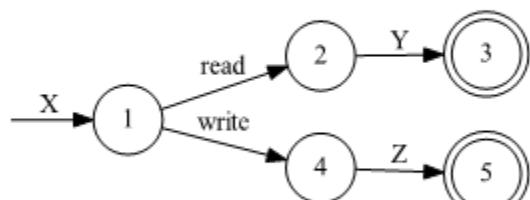




References

- [SAS13] Symbolic Automata for Static Specification Mining
- [OOPSLA12] Typestate-based semantic code search over partial programs
- [ISSTA07] Static Specification Mining using Automata Based Abstractions
- [ISSTA06] Effective Typestate Verification in the presence of Aliasing
- [Bodik-PLDI05] Jungloid mining: helping to navigate the API jungle
- [Kuncak-CAV11] Interactive Synthesis of Code Snippets
- [Perelman-PLDI12] Type-directed completion of partial expressions

context



(a)

$(*, X, \text{read}) \mapsto \text{open}, \text{canRead}$
 $(*, X, \text{write}) \mapsto \text{open}$
 $Y \mapsto \text{close}$
 $Z \mapsto \text{close}$

(b)

- (a) Symbolic automaton representing the query for the behavior around read and write methods and (b) the assignment with contexts to its symbolic transitions which answers the query.

Simulation Example

