# Towards Generic Database Management System Fuzzing

Yupeng Yang*, Yongheng Chen*, Rui Zhong^, Jizhou Chen*, and Wenke Lee*

\* Georgia Tech   ^ paloalto NETWORKS

# Background and Motivation

- Database Management Systems (DBMSs) are widely used for data storage, retrieval, and management.

- Both relational (SQL) DBMSs and non-relational (NoSQL) DBMSs have wide adoption in real world for the diverse requirements of various applications.



......

The security and robustness of these prevalent and critical systems are vital!

# Background and Motivation

- Fuzzing can be used to test software systems by injecting random inputs to them.

- Fuzzers targeting **SQL DBMSs** have proven useful and effective over the years.
  - SQLSmith, Squirrel, SQLancer…

- **NoSQL DBMSs** lack an effective fuzzing solution
  - Existing SQL DBMS fuzzers have challenges migrating to NoSQL DBMSs.
  - Generic fuzzers (e.g., AFL) struggle to generate valid inputs to DBMSs.

# Challenges and Limitations

- We discover three major challenges when designing a fuzzer that extends to NoSQL DBMSs.

- C1: It is hard to generalize.

- C2: Semantics can change based on the context.

- C3: Loose data dependencies.

# C1: It is hard to generalize

- Semantic correctness is vital for exploring deep DBMS logic.
- NoSQL DBMSs have **diverse** interfaces, and their semantics vary drastically.

| DBMS | Input Format | Examples |
|---|---|---|
| redis | *key-value commands* | `HSET key field value [field value ...]` |
| AgensGraph | *ASCII-art (Cypher)* | `MATCH (p:person {name: 'Tom'})-[r:knows*1..2]->(f:person)`<br>`RETURN f.name, r[1].fromdate;` |
| mongoDB | *JSON documents* | `db.products.insertOne( { item: "card", qty: 15 } );` |

Georgia Tech

# C2: Semantics can change based on the context

```
 1  // Partial grammar rules:
 2  createtbl_stmt:
 3    'CREATE TABLE'
 4    tbl_name '('
 5        ...
 6    ')';
```
*Grammar*

```
 7
 8  // The test case:
 9  > CREATE TABLE t1(
10      c1 date
11  );
```
*Parsing*

```
12
13
14  // Bind "TABLE define"
15  // to `tbl_name`.
16
17  // When we traverse to the
18  // node `t1`, we know a
19  // TABLE t1 is defined.
```
*Static Constraint*

- Existing works bind "static semantics" to the syntax structures.
- This works well for modeling common SQL semantics.

Georgia Tech.

# C2: Semantics can change based on the context

However, for NoSQL, semantics often change based on the context.

------------------------------------------------------------

- One syntax structure can have different semantics in different syntactic contexts.

```
MATCH (n:L) WHERE (n)-[]->() RETURN n.x;
            │                │
            ▼                ▼
      define                use
            identifier
                                        A cypher query
```

```
1   // Partial grammar rules:
2   match_clause:
3       MATCH pattern_part where_part;
4
5   pattern_part:
6       node_pattern;
7
8   where_part:
9       WHERE node_pattern;
10
11  node_pattern:
12      '(' identifier ')' | ...;
13
```

- Data types can depend on other values in the context.

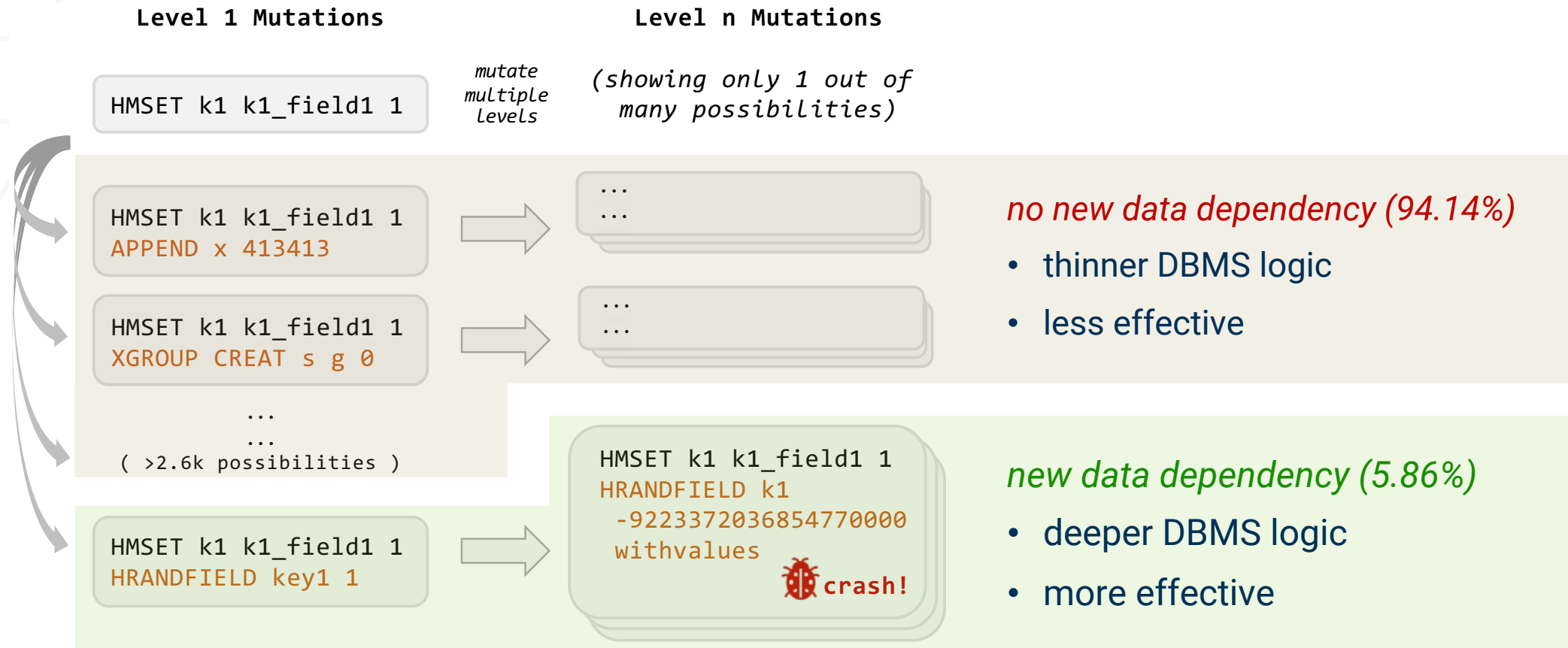Georgia Tech

# C2: Semantics can change based on the context

However, for NoSQL, semantics often change based on the context.

---

- One syntax structure can have different semantics in different syntactic contexts.

- Data types can depend on other values in the context.

```
> HSET k1 k1_field1 "Hello"                    redis commands
                          ↘ context
                      ASCII string

> HSET k2 k2_field1 "123"
                          ↙ context
                      Numeric string

> HINCRBY k1 k1_field1 1  ✗ Only a numeric string is valid.
(error) value not an integer
```

Georgia Tech

# C3: Loose Data Dependencies

Random mutations tend to generate loose data dependencies.

**Level 1 Mutations**

**Level n Mutations**

*mutate multiple levels*

*(showing only 1 out of many possibilities)*

```
HMSET k1 k1_field1 1
```

```
HMSET k1 k1_field1 1
APPEND x 413413
```

```
...
...
```

```
HMSET k1 k1_field1 1
XGROUP CREAT s g 0
```

```
...
...
```

```
...
...
( >2.6k possibilities )
```

```
HMSET k1 k1_field1 1
HRANDFIELD key1 1
```

```
HMSET k1 k1_field1 1
HRANDFIELD k1
  -9223372036854770000
  withvalues
                  🐞crash!
```

*no new data dependency (94.14%)*
- thinner DBMS logic
- less effective

*new data dependency (5.86%)*
- deeper DBMS logic
- more effective

*Random Mutation Running Examples (for redis)*

# Our Solution

We propose three approaches to tackle the three challenges.

- *Semantics Abstraction*
  - C1: Non-generic

- *Context-sensitive Constraint Resolution*
  - C2: Context-based Semantics

- *Dependency-guided Mutation*
  - C3: Loose Data Dependency

We implemented our approaches into a generic fuzzing framework, BuzzBee, that can fuzz **both** SQL and NoSQL DBMSs effectively.

# Semantics Abstraction -> C1

To generalize, we model common DBMS operations at a highly abstract level using three basic data operations: *Define*, *Use*, and *Invalidate*.

Next, we constrain the abstract semantics
- When to *Define, Use,* or *Invalidate* (scope constraints)
- What *type* to *Define, Use,* or *Invalidate* (type constraints)

*Constraints:*

The semantic rules to avoid a DBMS execution error.

We design an *Annotation System* to let users annotate the abstract semantics and constraints on the input grammar.



```
1  hset:
2      HSET key① (field② value)+;
```

*Input Grammar*

```
1  ①: "default": {                    scope constraint
2         operation: Define,          type constraint
3         args: { type: "HSET key" }
4  }
```

*Annotation*

# Semantics Abstraction - Internals

- We design an IR to carry the syntactic and semantic information specified by the user.
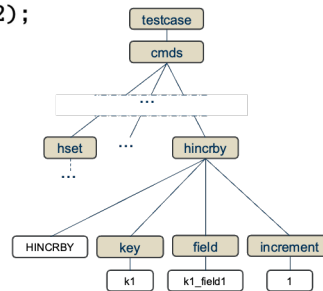- We maintain scope trees and symbol tables to track the data.

```
1  HSET k1 k1_field1 "Hello"
2  HSET k2 k2_field1 "123"
3  HSET k3 k3_field1 "456"
4  // DEL k1
5  HINCRBY k1 k1_field1 1
```
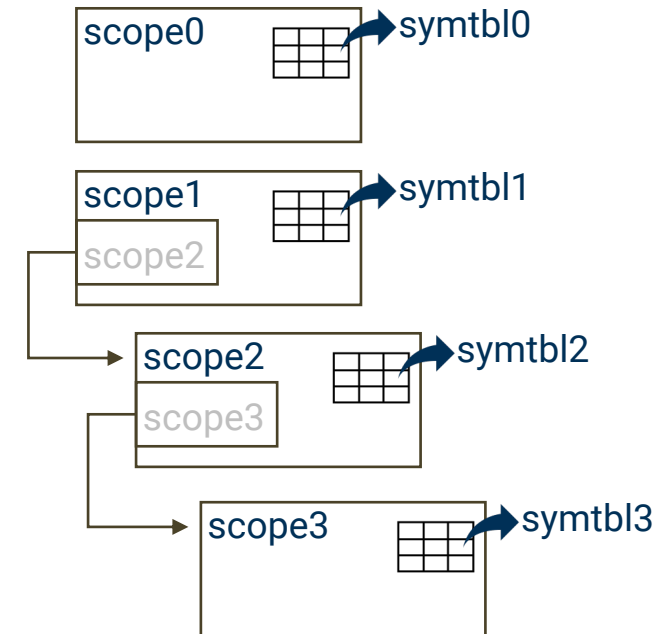
*Test Case*

*Grammar*          *Annotation*

```
13  node3(id=3, type=hset, text=NULL, annotation={},
14      children={&node4, &node5, &node7, &node9}, parent=&node2);
15  ...
16  node33(id=33, type=hincrby, text=NULL, annotation={},
17      children={&node34, &node35, &node37, &node39},
18      parent=&node32);
19  node34(id=34, type=terminal, text="HINCRBY", annotation={},
20      children={}, parent=&node33);
21  node35(id=35, type=key, text=NULL,
22      annotation={"default":{operation:Use,
23                              args:{type:"HSET key"}}},
24      children={&node36}, parent=&node33);
```

*An IR Program*

*Scope Trees and Symbol Tables*

# Context-sensitive Constraint Resolution -> C2

To achieve context sensitivity, we design two features for the *Annotation System* so that users can specify constraints based on the context.

- *Context Query Language (CQL)* for simplicity – targeting common semantics
- *Custom Resolvers* for expressiveness – targeting complex semantics

Georgia Tech

# Context Query Language (CQL)

- CQL is a lightweight language to fetch information from the context.
- To fetch certain information, we need to know:
  - **where** to fetch (which part of the context do we care about?)
  - **what** to fetch (what property of that part are we interested in?)

*where*

```
1  cql:
2      navigator* property ;
3
4  navigator:
5      '.parent'
6      | '.child' arg | '.lsib' arg | '.rsib' arg
7      | ... ;
8
9  arg:
10     '(' num ')' ;
11
```

*what*

```
12 property:
13     '@text' | '@id' | '@sym_type' | ... ;
```

*Grammar of CQL*

```
.lsib(1)@text

.parent.rsib(1)@id

.parent.rsib(1).child(0)@id
```

*CQL Examples*

```
args: {
  type: "HSET numeric field
         of {.lsib(1)@text}"
}
```

*CQL in the Annotation*

Georgia Tech

# Context Query Language (CQL)

```
> HSET k1 k1_field1 "Hello"
> HSET k2 k2_field1 "123"
> ...
> HINCRBY k1 k1_field1 1
```
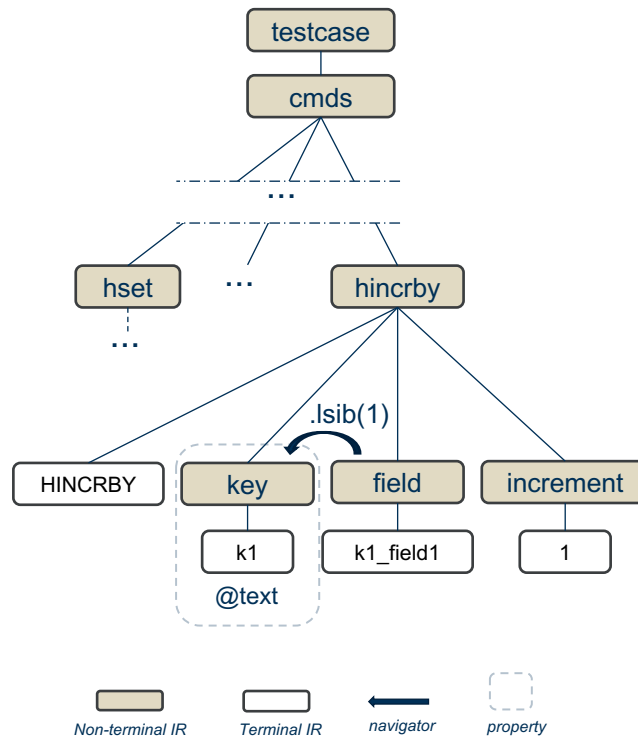
*Redis Test Case*

```
4  hincrby:
5      HINCRBY key③ field④ increment;
```

*Redis Grammar*

```
22  ④: "default": {
23      operation: Use,
24      args: {
25          type: "HSET numeric field
26                 of {.lsib(1)@text}"
27      }
28  }
```

*CQL in the Annotation*



*CQL Querying Process*

```
> HSET k1 k1_field1 "Hello"
> HSET k2 k2_field1 "123"
> ...
> HINCRBY k1 k1_field1 1
```

```
{
    operation: Use,
    args: {
        type: "HSET numeric field
               of k1"
    }
}
```

*Resolved Constraint*

# Custom Resolvers

Custom Resolvers are plugins to the *Annotation System*.

- Can be written in high-level languages like C++.

- Have access to all the context information visible to BuzzBee.

- Can express arbitrarily complex semantics to complement CQL.

Georgia Tech
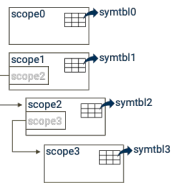
# Custom Resolvers



```
1  hset:
2      HSET key① (field② value)+;
```
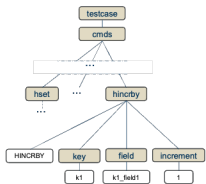
Redis Grammar

```
7  ②: "default": {
8        operation: Define,
9        args: {
10           type:
11               hset_field_type_resolver
12       }
13  }
```

Custom Resolver in Annotation

Symbols    IR Program

```
hset_field_type_resolver
> ... (custom code)
```

Custom Resolver

> HSET k1 k1_field1 "Hello"          type: HSET field of k1
> HSET k2 k2_field1 "123"
                                       type: HSET **numeric** field of k2
> ...
> HINCRBY k1 k1_field1 1
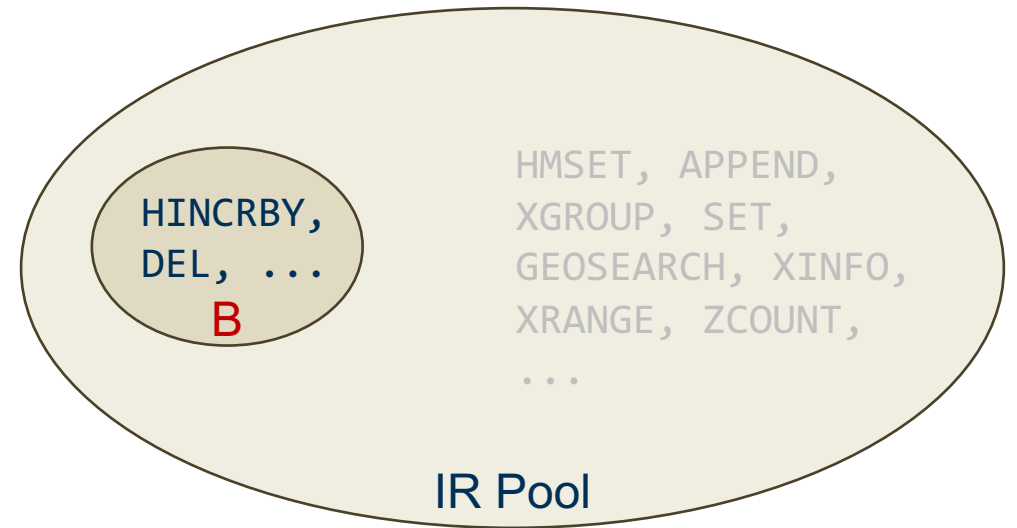
Redis Test Case

Georgia Tech

# Dependency-guided Mutation -> C3

We add **guidance** to the *replacement* and *insertion* mutations.

```
> HSET k1 k1_field1 "Hello"
> HSET k2 k2_field1 "123"
> ...                          Mutation point A
> HINCRBY k1 k1_field1 1
```

1. Get all symbols available at A.
   - k1, k1_field1, k2, k2_field1

2. Favor B from the IR Pool, which can use the available symbols.



HINCRBY, DEL, ... **B**

HMSET, APPEND, XGROUP, SET, GEOSEARCH, XINFO, XRANGE, ZCOUNT, ...

IR Pool

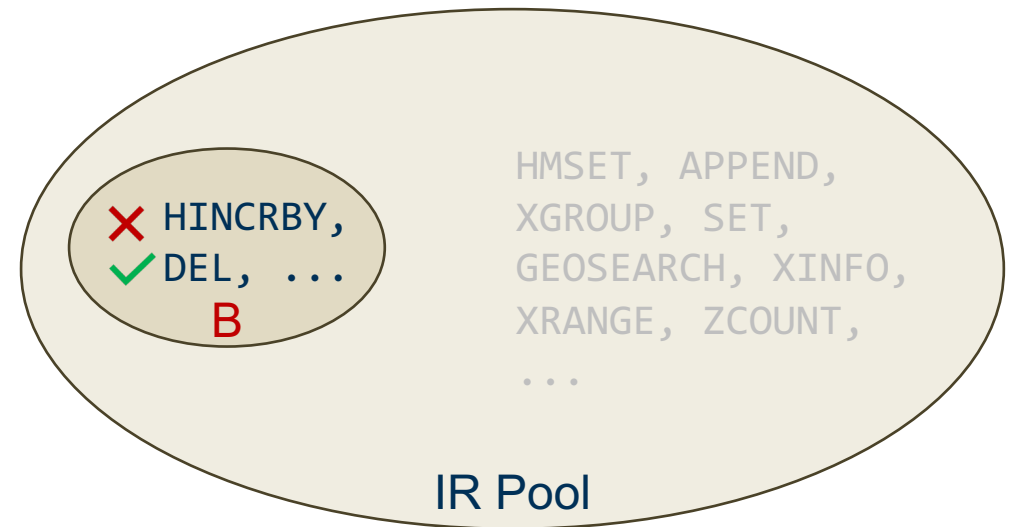IR Pool stores the mutation candidate IRs.

# Dependency-guided Mutation -> C3

We also introduce a finer-grained **prioritization** to cover more behaviors.

```
> HSET k1 k1_field1 "Hello"
> HSET k2 k2_field1 "123"
> ...                              Mutation point A
> HINCRBY k1 k1_field1 1
```

1. Get all symbols available at A.
   - k1, k1_field1, k2, k2_field1

2. Favor B from the IR Pool, which can use the available symbols.

3. Prioritizes IRs in B that do not exist in the test case.
   - DEL will be chosen over HINCRBY.



IR Pool

HMSET, APPEND, XGROUP, SET, GEOSEARCH, XINFO, XRANGE, ZCOUNT, ...

✗ HINCRBY,
✓ DEL, ...
B

# Implementation & Evaluation

- Implemented BuzzBee mainly in C++ and Python (9,130 LoC)

- Applied to **8** real-world DBMSs covering **4** major data models.
  - redis, KeyDB, RedisGraph, AgensGraph, MongoDB, ArangoDB, PostgreSQL, MySQL
- Discovered **40** bugs in the latest versions (with 4 CVEs).
- Outperformed generic fuzzers in NoSQL DBMSs
  - Up to 76.9% cov increase in NoSQL DBMSs
  - Discovered >30 bugs that generic fuzzers could not discover
- Achieved comparable results with SQL fuzzers
  - Achieved 92.7% cov of Squirrel
  - Found a similar # of bugs

# Thanks / Q&A

*Towards Generic Database Management System Fuzzing*

Yupeng Yang*, Yongheng Chen*, Rui Zhong^, Jizhou Chen*, and Wenke Lee*

\* Georgia Tech   ^ paloalto NETWORKS

Georgia Tech