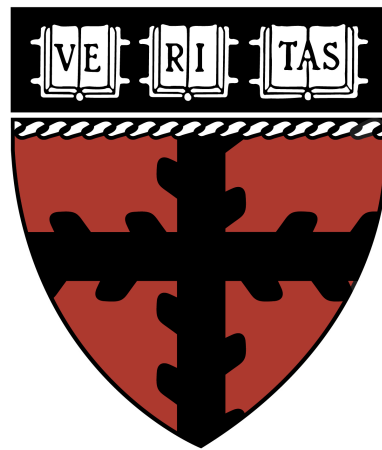


Toward a Verified Relational Database Management System

Gregory Malecha, Greg Morrisett,
Avraham Shinnar, **Ryan Wisnesky**

POPL 2010



Overview

- We built a **verified RDBMS**.
 - DB2, Oracle, BerkeleyDB, MySQL, SQLite,...
- Code available at:
`http://ynot.cs.harvard.edu`

Motivation

- Data management is ubiquitous and difficult; correctness is important.
- RDBMS components have clean specifications with deep underlying theory.
- We can verify compilers and operating systems.

Our RDBMS Functionality

- Data storage
- Data update
- Query optimization
- Single-threaded B+ Tree execution

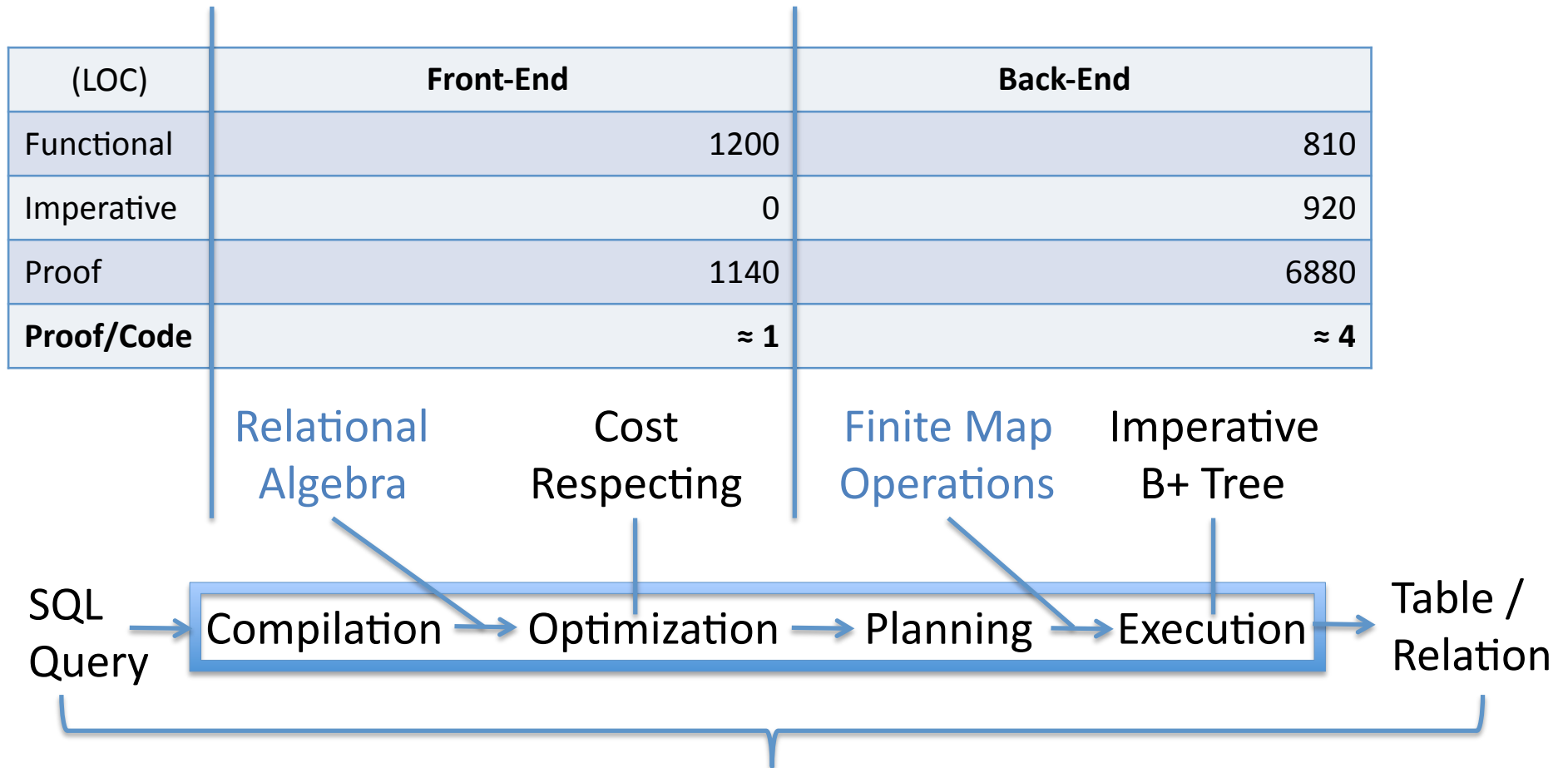


Mechanically
verified

RDBMS State of the Art

- Data storage
- Data update
- Query optimization
- Concurrent execution
 - **A**tomicity
 - **C**onsistency
 - **I**solation
 - **D**urability

Our RDBMS Pipeline



The table the SQL query denotes and the table the RDBMS returns are equal (partial correctness).

Verification Methodology

- For Purely Functional Code:
the Coq proof assistant
- For Imperative Code:
the Ynot extensions to Coq (*Nanevski et al, ICFP '08*)
 - Hoare Type Theory (*Nanevski et al, ICFP '06*)
 - Separation logic (*Reynolds, LICS 02*) (*O'Hearn et al, POPL '04*)
 - Proof Automation (*Chlipala et al, ICFP '09*)

Lessons Learned

Challenge

- Reasoning about queries and relations
- Modular, first-class imperative finite maps
- Reasoning about B+ Trees in separation logic

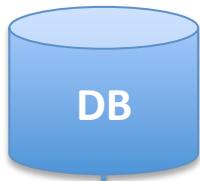
Solution

- Tuples as heterogeneous lists
- Schema-indexed queries
- Axiomatic interfaces
- Hoare Type Theory
 - Key: higher order iterator
- Generic, higher order traversal
- Multiple invariants

Challenge #1

- Reasoning about queries and relations

DB Model



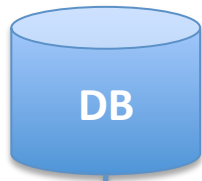
Section Assignments

Student	Year	Section #
John Doe	Freshman	3
Jane Doe	Senior	2

Room Assignments

Course	Room #
CS 51	115
CS 252	319
CS 152	323

DB Model



Section Assignments

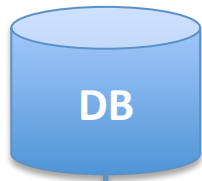
Student	Year	Section #
John Doe	Freshman	3
Jane Doe	Senior	2

Room Assignments

Course	Room #
CS 51	115
CS 252	319
CS 152	323

(* nameless, ordered *)
Definition Schema : Type :=
list Type.

DB Model



Section Assignments

Student	Year	Section #
John Doe	Freshman	3
Jane Doe	Senior	2

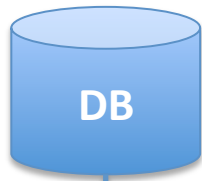
Room Assignments

Course	Room #
CS 51	115
CS 252	319
CS 152	323

(* nameless, ordered *)
Definition Schema : Type :=
list Type.

Fixpoint Tuple (T: Schema) :
Type :=
match T with
| Nil \Rightarrow unit
| Cons a b \Rightarrow a * Tuple b
end.

DB Model



Section Assignments

Student	Year	Section #
John Doe	Freshman	3
Jane Doe	Senior	2

Room Assignments

Course	Room #
CS 51	115
CS 252	319
CS 152	323

(* nameless, ordered *)
Definition Schema : Type :=
list Type.

Fixpoint Tuple (T: Schema) :
Type :=
match T with
| Nil \Rightarrow unit
| Cons a b \Rightarrow a * Tuple b
end.

(* unordered, no duplicates *)
Definition Table (T: Schema) :
Type := FSet (Tuple T).

Query Abstract Syntax

Inductive RAExp	:	Type :=
var : $\forall (v: \text{name}), \text{RAExp}$		
union :		
RAExp \rightarrow RAExp \rightarrow RAExp		
select : $\forall (t: \text{Schema}),$		
RAExp \rightarrow (Tuple $t \rightarrow \text{bool}$) \rightarrow RAExp		
...		
product :		
RAExp \rightarrow RAExp \rightarrow		
RAExp .		

Query Abstract Syntax

Inductive RAExp (G: Context) : Schema → Type :=
| var : ∀ (v: name), RAExp G (G v)
| union : ∀ (t: Schema),
 RAExp G t → RAExp G t → RAExp G t
| select : ∀ (t: Schema),
 RAExp G t → (Tuple t → bool) → RAExp G t
| ...
| product : ∀ (t t': Schema),
 RAExp G t → RAExp G t' →
 RAExp G (t ++ t').

Optimization

Definition denote $T(q: \text{RAExp } T)$
: $\text{FSet } (\text{Tuple } T) := \dots$

For legibility, ignore context and environment.

Optimization

Definition denote T (q : RAExp T)
: FSet (Tuple T) := ...

Definition rewrite T : Type :=
RAExp $T \rightarrow$ RAExp T .

Optimization

Definition $\text{denote } T \text{ (} q : \text{RAExp } T \text{)}$
 $: \text{FSet (Tuple } T \text{)} := \dots$

Definition $\text{rewrite } T : \text{Type} :=$
 $\text{RAExp } T \rightarrow \text{RAExp } T.$

Definition $\text{semantics_preserving } T \text{ (} r : \text{rewrite } T \text{)} :$
Prop $:= \forall (q : \text{RAExp } T), \text{denote } q = \text{denote } (r \ q).$

Optimization

Definition $\text{denote } T \ (q: \text{RAExp } T)$
: $\text{FSet } (\text{Tuple } T) := \dots$

Definition $\text{rewrite } T : \text{Type} :=$
 $\text{RAExp } T \rightarrow \text{RAExp } T.$

Definition $\text{semantics_preserving } T \ (r: \text{rewrite } T) :$
Prop $:= \forall (q: \text{RAExp } T), \text{denote } q = \text{denote } (r \ q).$

Definition $\text{optimization } T : \text{Type} :=$
{ r : $\text{rewrite } T$
; $\text{pf}_1 : \text{semantics_preserving } r$
; $\text{pf}_2 : \text{cost_respecting } r$ }

Challenge #2

- First-class, modular imperative finite maps.


Mutable ADTs in Ynot





Class myADT : Type := {

handle : Type;  Implementation type

model : Type := ...;  Model / “Ghost state” type

rep : handle → model → heap → **Prop**;  Representation invariant

“ghost state”


myOp : $\forall (m : \text{model}) (h : \text{handle}),$
IO ($\text{rep } h \ m$)  Pre-condition
 (fun r : T \Rightarrow rep h m_{post})  Post-condition
 Return value

IO Monad

A Finite Map ADT

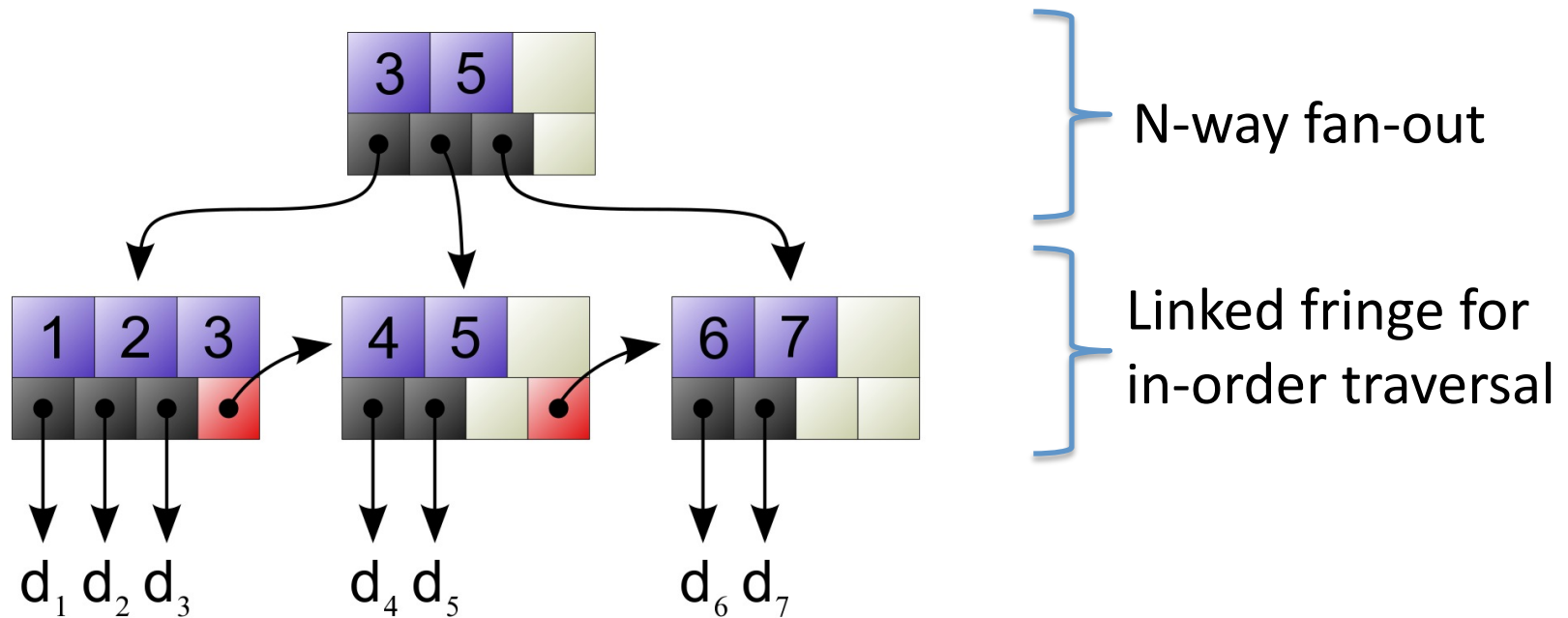
```
Class Fmap (K V: Type) : Type := {  
  
  handle : Type;  
  
  model   : Type := list (K*V);  
  
  rep      : handle → model → heap → Prop;  
  
  add      : ∀ (m: model) (k: K) (v: V) (h: handle),  
    IO (      rep h m )  
    (fun _: unit ⇒ rep h ((k,v)::m));  
  
  lookup   : ... ;  
  iterate  : ... ;  
  ...
```

Challenge #3

- Reasoning about B+ Trees in separation logic.

B+ Trees

- Generalized Binary Search Trees



B+ Tree Invariant

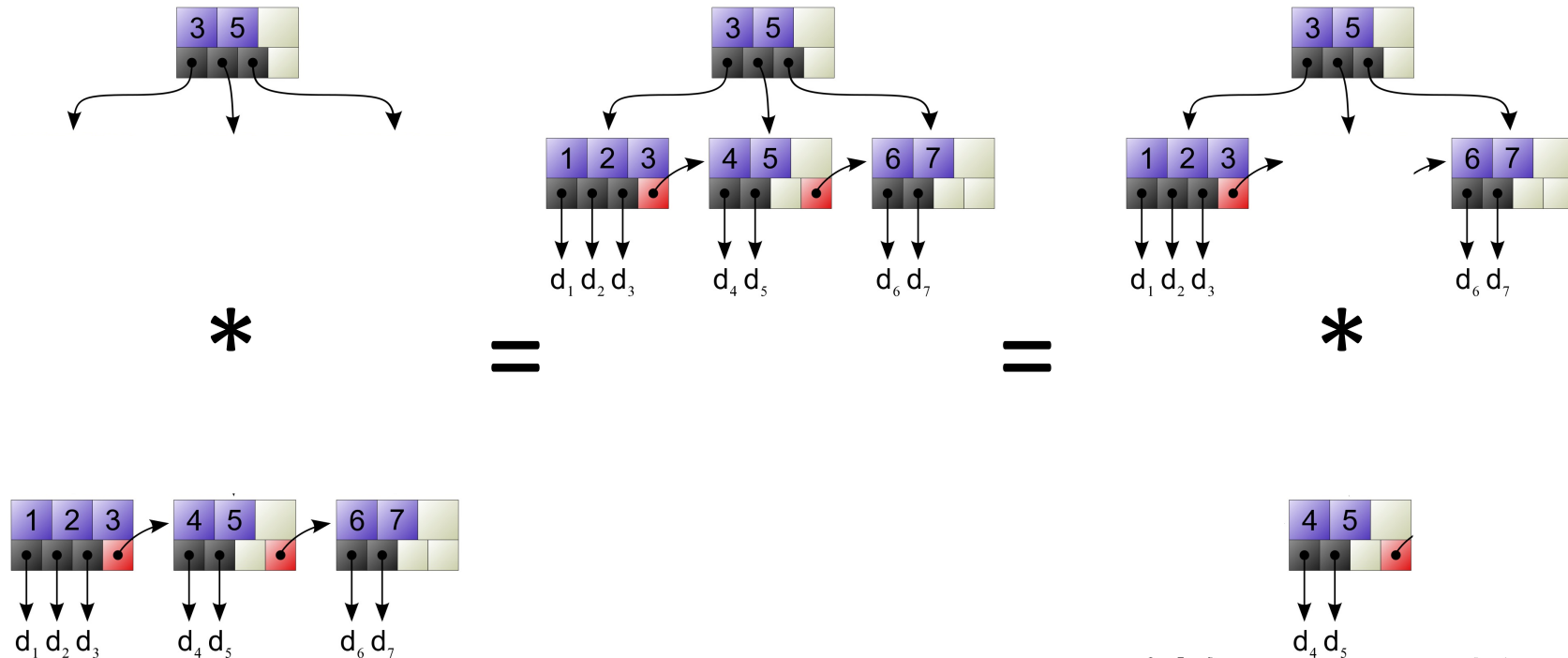
$\text{repTree } 0 \ r \ \text{optr} \ (p', ls) \iff$
 $[r = p'] * \exists \text{ary}. r \mapsto \text{mkNode } 0 \ \text{ary} \ \text{optr} *$
 $\text{repLeaf } \text{ary} \ |ls| \ ls$

$\text{repTree } (h + 1) \ r \ \text{optr} \ (p', (ls, \text{nxt})) \iff$
 $[r = p'] * \exists \text{ary}. r \mapsto \text{mkNode } (h + 1) \ \text{ary} \ (\text{ptrFor } \text{nxt}) *$
 $\text{repBranch } \text{ary} \ (\text{firstPtr } \text{nxt}) \ |ls| \ ls *$
 $\text{repTree } h \ (\text{ptrFor } \text{nxt}) \ \text{optr} \ \text{nxt}$

$\text{repLeaf } \text{ary} \ n \ [v_1, \dots, v_n] \iff$
 $\text{ary}[0] \mapsto \text{Some } v_1 * \dots * \text{ary}[n - 1] \mapsto \text{Some } v_n *$
 $\text{ary}[n] \mapsto \text{None} * \dots * \text{ary}[\text{SIZE} - 1] \mapsto \text{None}$

$\text{repBranch } \text{ary} \ n \ \text{optr} \ [(k_1, t_1), \dots, (k_n, t_n)] \iff$
 $\text{ary}[0] \mapsto \text{Some } (k_1, \text{ptrFor } t_1) *$
 $\text{repTree } h \ (\text{ptrFor } t_1) \ (\text{firstPtr } t_2) \ t_1 * \dots *$
 $\text{ary}[n - 2] \mapsto \text{Some } (k_{n-1}, \text{ptrFor } t_{n-1}) *$
 $\text{repTree } h \ (\text{ptrFor } t_{n-1}) \ (\text{firstPtr } t_n) \ t_{n-1} *$
 $\text{ary}[n - 1] \mapsto \text{Some } (k_n, \text{ptrFor } t_n) *$
 $\text{repTree } h \ (\text{ptrFor } t_n) \ \text{optr} \ t_n *$
 $\text{ary}[n] \mapsto \text{None} * \dots * \text{ary}[\text{SIZE} - 1] \mapsto \text{None}$

Multiple B+Tree Views



During Iteration

During Deletion

Related Work

- **B+ Trees in Separation Logic**
 - Local Reasoning, Separation and Aliasing (*Bornat et al, SPACE 04*) [classical conj.]
 - Reasoning about B+ Trees with Operational Semantics (*Sexton et al, Elec. Notes. Theoretical CS 08*)
- **Verified DB Components**
 - Ph.D. Thesis (*Gonzalia, 2006*) [relations in Agda]
 - The Power of Pi (*Oury et al, ICFP '08*) [relational algebra in Agda]
 - A Generic Algebra for Data Collections Based on Constructive Logic (*Rajagopalan et al, LNCS '95*) [generic data in Nuprl]
- **Program Verification**
 - CompCert: Formal certification of a compiler back-end (*Leroy, POPL '06*)
 - A Verified Compiler for an Impure Functional Language (*Chlipala, POPL '10*)
 - seL4: Formal Verification of an OS Kernel (*Klein et al, SOSP '09*) [165k Proof]

Conclusion

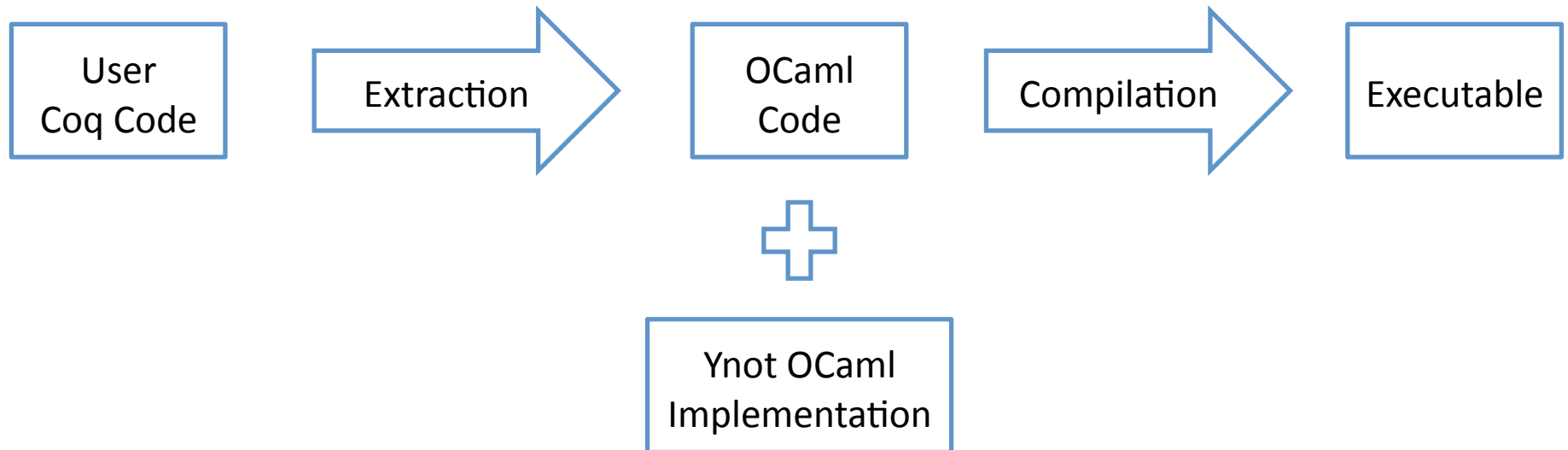
- Verified systems software is now viable.
- Verified RDBMSs are possible.
 - 3 Ph.D. students part-time 3-6 months
 - 30 minutes to verify (3GHz Pentium D, 1GB RAM)

... but still difficult, despite progress.

- Future Work
 - Concurrency and the ACID Properties
 - Dependent types vs traditional types

Questions?

Extraction



Cost Model

- Naïve implementation:
 - Set union: $O(n*m)$
 - Selection: $O(n)$
- No data statistics
- Conservative approximations
 - Selection preserves cardinality

First-class Sets

```
let user_schema : Schema :=  
  load_schema()  
in  
let user_data : Table user_schema :=  
  load_data(user_schema)  
in  
  process_data(user_schema, user_data)
```


LOC

- SQLite: 65k
- Compcert: 8k Code, 24k Proof, ratio ≈ 4
- seL4: 47.3k Code 165k Proof, ratio ≈ 3.5
- Chlipala: 6.6k Code, 1.8k Proof, ratio $\approx .3$