

Violating dataflow and getting away with it

Mistral Contrastin

Andrew Rice

Mistral.Contrastin@cl.cam.ac.uk

dodisturb.me

@madgen_

Datalog syntax in this talk

```
advisor("Euler", "Lagrange").  
advisor("Lagrange", "Fourier").  
advisor("Lagrange", "Poisson").
```

Fact

Predicate

Atom (atomic formula)

```
academicAncestor(X, Y) :- advisor(X, Y). ← Clause  
academicAncestor(X, Y) :- academicAncestor(X, Z), advisor(Z, Y).
```

Head

Body

Variable Constant

```
?- academicAncestor(X, "Fourier").
```

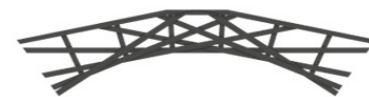
Query

Frustrating Datalog: programmer's misery

```
lt100(X) :- X < 100.  
?- in(X, 1, 50), lt100(X).
```

Frustrating Datalog: programmer's misery

```
pure(E) :- ! vars(E, _).  
?- exp(E), pure(E).
```



Frustrating Datalog: programmer's misery

?- exp(E_1), exp(E_2), independent(E_1, E_2).

independent(E_1, E_2) :- pure(E_1).

independent(E_1, E_2) :- pure(E_2).

independent(E_1, E_2) :-
vars(E_1, V_1), vars(E_2, V_2), ! intersect($V_1, V_2, _$).

Frustrating Datalog: compiler writer's misery

```
?- exp(E1), exp(E2),  
   (pure(E1);  
    pure(E2);  
    vars(E1, V1), vars(E2, V2), ! intersect(V1, V2, _)).
```

Overview

- ▶ Dataflow restrictions (well-moding & range restriction)
- ▶ Information flow (modes & adornment)
- ▶ Dataflow repair
- ▶ Synthesising control (reorder atoms)

Dataflow restrictions

Well-moding

❌ 'i' is not bound to a value.

```
❌+ 1  predicate p(int i) { 1 <= 1 and i <= 10 }
⚠️  2  predicate q(int i) { i in [1..10] }
```

$p(I) :- I \leq 1, I \leq 10.$

$q(I) :- \text{in}(I, 1, 10).$

- ▶ Not every predicate is created equal
- ▶ \leq (as implemented in QL) is *extralogical*; *in* is *logical*
- ▶ Foreign predicates, parameters indexed in a database, semantically one way functions, etc.

Information flow

Mode = Binding requirement

- ▶ **+** means “needs to be bound”
- ▶ **?** means “we don’t care”
- ▶ Static property independent of where the predicate occurs
- ▶ $\text{in}^{?++}(\text{I}, \text{Lower}, \text{Upper})$
 $\text{A} \leq^{++} \text{B}$
 $\text{sha256}^{+?}(\text{Content}, \text{Hash})$
 $\text{plus}^{++?, +?+, ?++}(\text{A}, \text{B}, \text{C})$

Mode = Binding set (?)

- ▶ In QL, **binding sets** *seem* to be the equivalent notion
- ▶ `plus++?,+?+,?++(A, B, C)`
- ▶ `bindingset[a, b]`
`bindingset[b, c]`
`bindingset[a, c]`
`predicate plus(int a, int b, int c) {`
 `a + b = c`
`}`

Adornment = active binding

- ▶ **b** for bound
- ▶ **f** for free
- ▶ Dynamic property depending on the context
- ▶ ? – $p_{fb}(X, 1)$, $q_{bfb}(X, Y, \emptyset)$.

Well-modable: Intuition

▶ ?- $\text{lt}100(X)$.

$\text{lt}100(X) :- X < 100$.

▶ ?- $\text{in}(X, 1, 50)$, $\text{lt}100(X)$.

$\text{lt}100(X) :- X < 100$.

Well-modable: Formalisation

Need a b for every $+$

► ?- $\text{lt}100_f(X)$.

$\text{lt}100_f(X) :- X <_{fb^{++}} 100$.

► ?- $\text{in}_{fbb}(X, 1, 50)$, $\text{lt}100_b(X)$.

$\text{lt}100_b(X) :- X <_{bb^{++}} 100$.

Well-modable vs well-moded

- ▶ Semantically clause heads and atoms are related

?- $\text{in}_{fb}(\text{X}, 1, 50)$, $\text{lt100}_b(\text{X})$.
 $\text{lt100}_b(\text{X})$:- $\text{X} <_{bb^{++}} 100$.

- ▶ But the evaluator treats each clause independently

?- $\text{in}_{fb}(\text{X}, 1, 50)$, $\text{lt100}_b(\text{X})$.
.....
 $\text{lt100}_f(\text{X})$:- $\text{X} <_{fb^{++}} 100$.

- ▶ Well-modability is a **global** property, while well-modedness is a **local** one

**Dataflow repair:
Transforming well-modable to well-moded
and making clauses range-restricted**

Repairing basic ill-modedness

▶ $?- \text{in}_{fb}(\mathit{X}, 1, 50), \text{lt}100_b(\mathit{X}).$
.....
 $\text{lt}100_f(\mathit{X}) :- \mathit{X} <_{fb^{++}} 100.$



▶ $?- \text{in}_{fb}(\mathit{X}, 1, 50), \text{lt}100_b(\mathit{X}).$
.....
 $\text{lt}100_f(\mathit{X}) :-$
 $\text{in}_{fb}(\mathit{X}, 1, 50), \mathit{X} <_{bb^{++}} 100.$

Repairing basic range-restriction

- ▶ ?- exp(E_1), exp(E_2), independent(E_1, E_2).

independent(E_1, E_2) :- pure(E_1).

independent(E_1, E_2) :- pure(E_2).

...



- ▶ ?- exp(E_1), exp(E_2), independent(E_1, E_2).

independent(E_1, E_2) :- exp(E_2), pure(E_1).

independent(E_1, E_2) :- exp(E_1), pure(E_2).

...

Not so basic dataflow 1

▶ ? – $p\text{Clo}(X, 1)$.

$p\text{Clo}(X, Y) :- p^{?+}(X, Y)$.

$p\text{Clo}(X, Z) :- p\text{Clo}(X, Y), p^{?+}(Y, Z)$.

▶ Is this well-moded?

Not so basic dataflow 2

- ▶ ? – $p(X, 1)$, $r(Y)$, $p(X, Y)$.
 $p(X, Y) :- q(X)$.
- ▶ What flows into Y of p ?

What are we trying to find?

- ▶ For well-modedness, target is **body atom parameters**.
What flows into X of $X <^{++} 100$?

?- $\text{in}(X, 1, 50), \text{lt}100(X)$.

$\text{lt}100(X) :- X <^{++} 100$.

- ▶ For range restriction, target is **predicate parameters**.
What flows into second parameter of independent?

?- $\text{exp}(E_1), \text{exp}(E_2), \text{independent}(E_1, E_2)$.

$\text{independent}(E_1, E_2) :- \text{pure}(E_1)$.

What are we really trying to find?

- ▶ Enumerable domains for variables!
- ▶ QL makes this much easier with its **characteristic predicates**
- ▶ In pure Datalog, we must look hard for them

How to find these domains? (sketch)

1. Extract all local flow into a graph
2. Compute the shortest paths from data sources (constants, positive atoms) to data sinks (offending atom/predicate parameters)
3. Ensure data paths to offending sinks are **closed**
4. Extract domain into new predicate and prepend it to the problematic clause body

Implementation

- Feature extraction can be done in Datalog via MetaDL
- Flow path computation is about ~15 pure Datalog clauses
- Transformation cannot (yet) be done in Datalog!
- Similar analysis required for Magic Set Rewriting, potential infrastructure reuse

Mode inference and atom reordering

Adornments are contextual

Reminder: need a **b** for every **+**

- ▶ $\text{auth}_f(U) :- \text{hash}_{ff^{+?}}(P, H), \text{password}_{bf}(U, P), \text{valid}_{bb}(U, H).$ ❌
- ▶ $\text{auth}_f(U) :- \text{password}_{ff}(U, P), \text{hash}_{bf^{+?}}(P, H), \text{valid}_{bb}(U, H).$ ✅

Naïvely enumerate orderings?

- ▶ $\text{auth}_f(U) :- \text{hash}_{ff^{+?}}(P, H), \text{password}_{bf}(U, P), \text{valid}_{bb}(U, H).$ ❌
- ▶ $\text{auth}_f(U) :- \text{hash}_{ff^{+?}}(P, H), \text{valid}_{fb}(U, H), \text{password}_{bb}(U, P).$ ❌
- ▶ $\text{auth}_f(U) :- \text{password}_{ff}(U, P), \text{hash}_{bf^{+?}}(P, H), \text{valid}_{bb}(U, H).$ ✅
- ▶ $\text{auth}_f(U) :- \text{password}_{ff}(U, P), \text{valid}_{bf}(U, H), \text{hash}_{bb^{+?}}(P, H).$ ✅
- ▶ $\text{auth}_f(U) :- \text{valid}_{ff}(U, H), \text{password}_{bf}(U, P), \text{hash}_{bb^{+?}}(P, H).$ ✅
- ▶ $\text{auth}_f(U) :- \text{valid}_{ff}(U, H), \text{hash}_{fb^{+?}}(P, H), \text{password}_{bb}(U, P).$ ✅

It's a global reordering problem

- ▶ ?- auth_f(U)
auth_f(U) :- check_{ff}(U, P), password_{bb}(U, P).
check_{ff}(U, P) :- hash_{ff+?}(P, H), valid_{bb}(U, H). ❌
- ▶ ?- auth_f(U)
auth_f(U) :- password_{ff}(U, P), check_{bb}(U, P).
check_{bb}(U, P) :- hash_{bf+?}(P, H), valid_{bb}(U, H). ✅
- ▶ Combinatorial explosion! Naïve won't do.

What are we really trying to compute?

- ▶ Orderings, but also a summary of binding requirements
- ▶ Mode inference for derived predicates
- ▶ What should the ideal mode of check is?

$\text{check}_{bb}(U, P) \text{ :- hash}_{bf^{+?}}(P, H), \text{valid}_{bb}(U, H).$

- ▶ $+?$ or $??? ??$ is strictly better
- ▶ We have submoding! Only want the most general mode
- ▶ Each mode is coupled to orderings that lead to it

Multiple modes due to ordering

- ▶ `clientCheck(P) :- weak(P,H).`
`serverCheck(H) :- weak(P,H).`
`weak(P,H) :- hash+?(P,H), rainbow+?(H,P).`
- ▶ The mode of weak is `+?, ?+`

Mode inference algorithm (sketch)

1. Start with input modes, assume others are pure
2. For each clause, construct a greedy scheduling graph
3. Find the most relaxed modes and their orderings
4. Combine modes of clauses sharing head predicates
5. Go to (2) repeat until convergence

Two heuristics against naïveté

- ▶ Pure atoms can always be scheduled first in any order
- ▶ Scheduling might make other atoms pure
- ▶ ?- ● $p^+(Z)$, $q^{+?}(X, Z)$, $r^{??}(X, Y)$.
?- $r^{??}(X, Y)$, ● $p^+(Z)$, $q^{??}(X, Z)$.
?- $r^{??}(X, Y)$, $q^{??}(X, Z)$, ● $p^?(Z)$.
?- $r^{??}(X, Y)$, $q^{??}(X, Z)$, $p^?(Z)$ ●.

Properties

- ▶ Sound and complete
- ▶ Incremental
- ▶ Highly parallelisable
- ▶ Closely mimics properties of Datalog evaluation
- ▶ Still exponential in degenerate case but only locally
?– $p^+(X)$, $q^+(Y)$, $r^+(Z)$.

Conclusion

- ▶ Well-modedness and range restriction are sources of "is not bound to a value" errors
- ▶ They can be fixed via program transformations
- ▶ Dataflow restricted atoms can be used as if purely logical
- ▶ Modes (or binding sets) can always be inferred
- ▶ Datalog analysis is nice because Datalog is nice!

Thanks. Questions?

Mistral.Contrastin@cl.cam.ac.uk
dodisturb.me
@madgen_