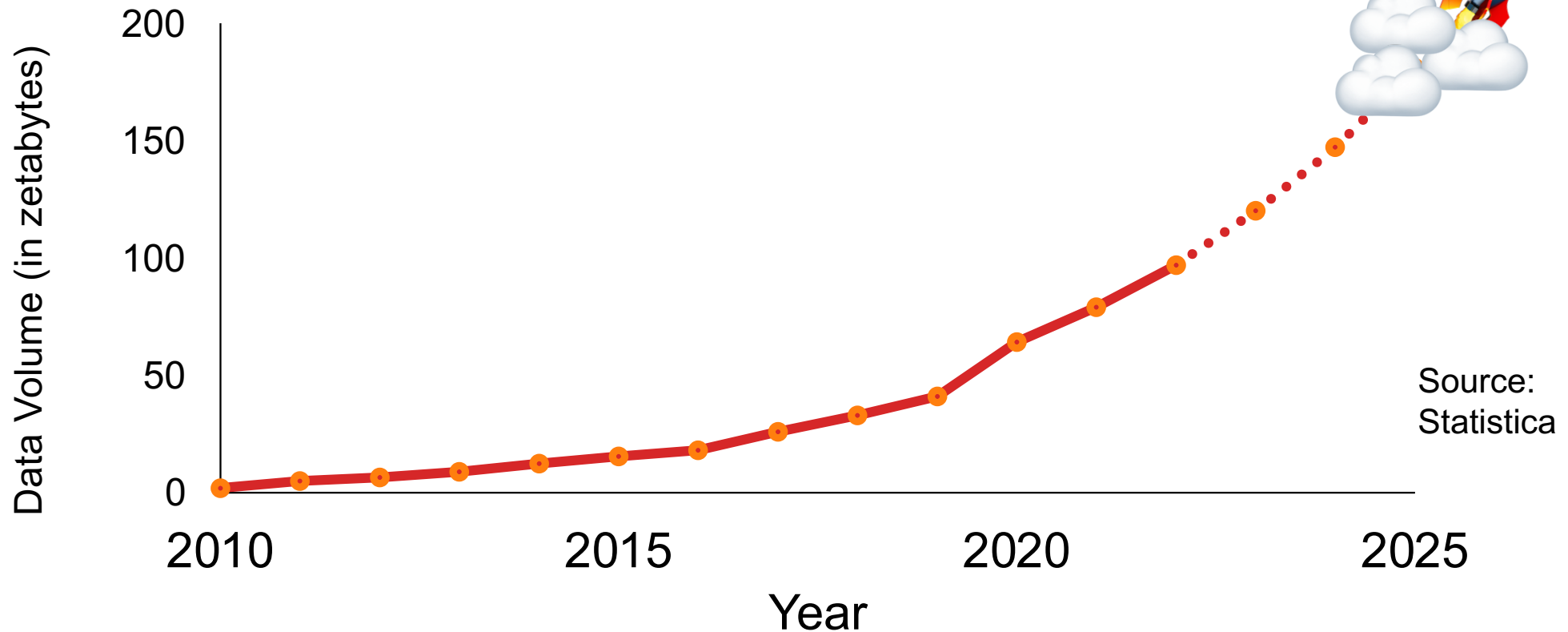


# ✨ A Database of Code ✨ aka Advanced Metaprogramming Queries

Anna Herlihy  
HPTS  
October 10 2022

# *Machine-generated data drives growth*

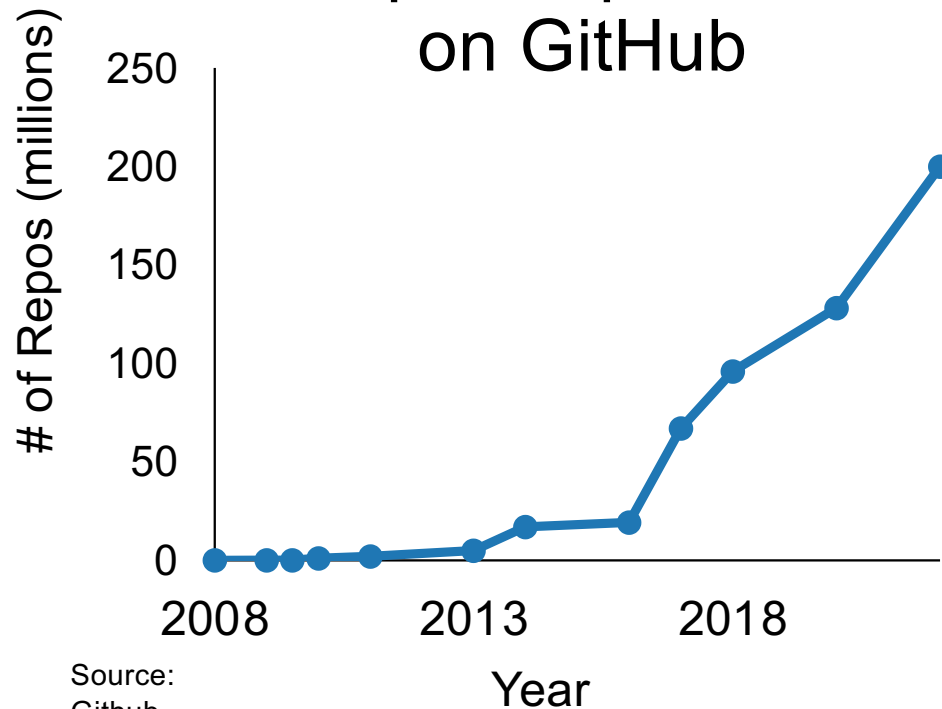
Projected volume of data created and consumed



**Prediction: Code Volume & Complexity** 

# More People Are Writing More Code

Unique Repositories  
on GitHub

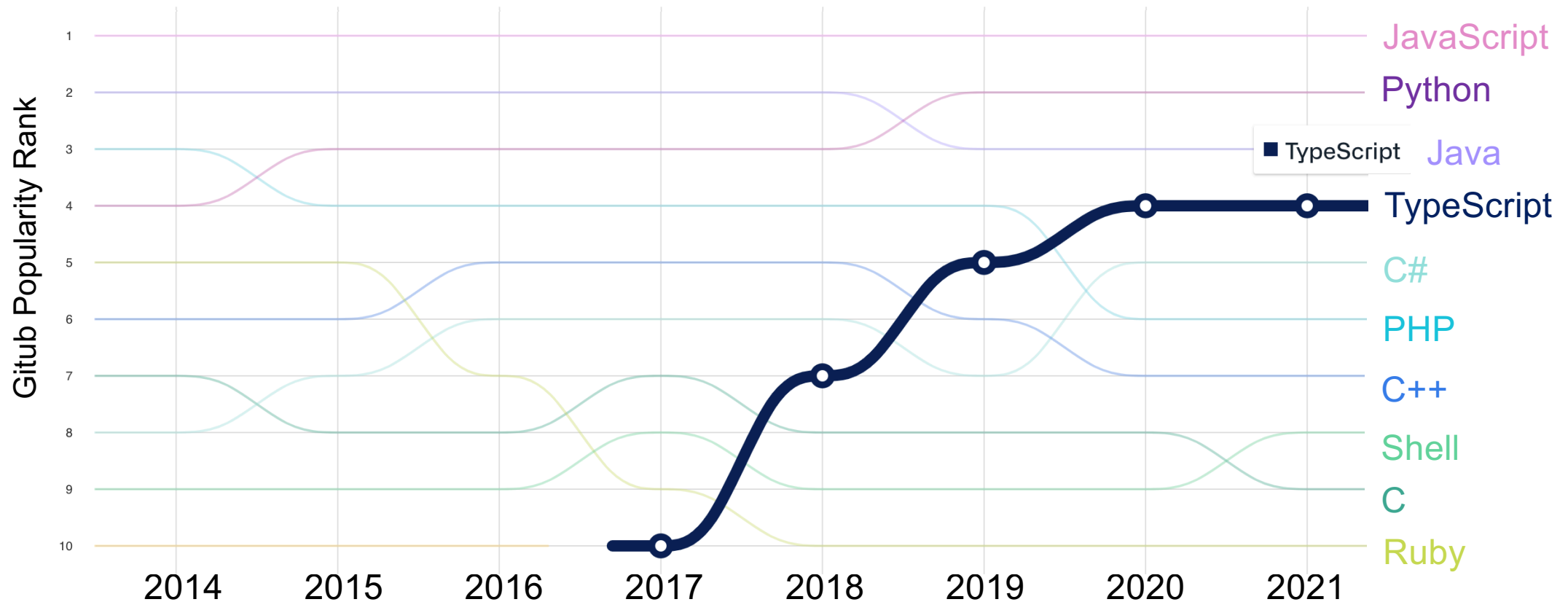


Code complexity is growing

- GitHub Co-Pilot
- Stack Graphs
- Meta Aroma
- Facebook Getafix, Infer
- IBM Project CodeNet

**Complexity is outgrowing capacity**

# Programmability Improves Complexity



Source:  
Github

**Trend towards higher level of abstraction**

# Abstraction Means Code Generation

- Advanced type systems
- Metaprogramming (Rust, Scala, etc.)
- Provably correct → rule-based languages
- Low-Code, No-Code systems

***Machine-generated code drives exponential growth*** <sub>5</sub>

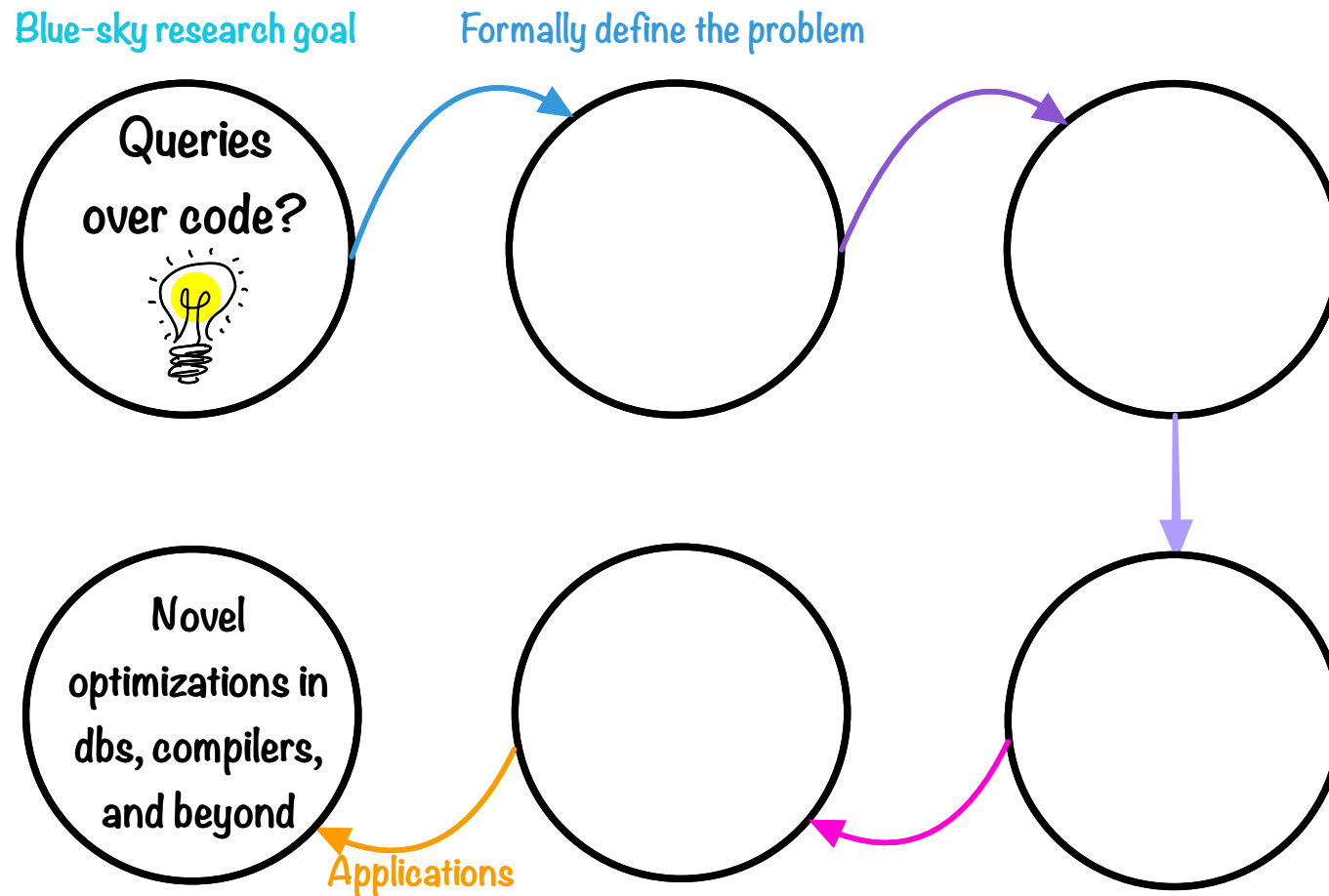
## A Challenge ☺

Database people understand *large volumes of data*

We need to start thinking about systems that store and mine information from ***large volumes of code***

**So what could that look like?**

# The Road to CodeDB



# How to extract meaning from code

Does X have side effects?

Where are the callers of X?

What type is X?

What inherits from X?

Does X undo Y?

Could X be tainted?

Is X in scope at point Y?

- AI-based approaches: rigorous enough?
- Formal verification-based approaches: scalable enough?
- Static analysis 😊

**Static analysis can be explainable & efficient**



# Extract guarantees *without running it*

- A program is a function  $= F$ 
  - instructs a machine how to change state at each step
- Static analysis computes an abstract state  $= \check{x}$ 
  - **over-approximates** all possible concrete states

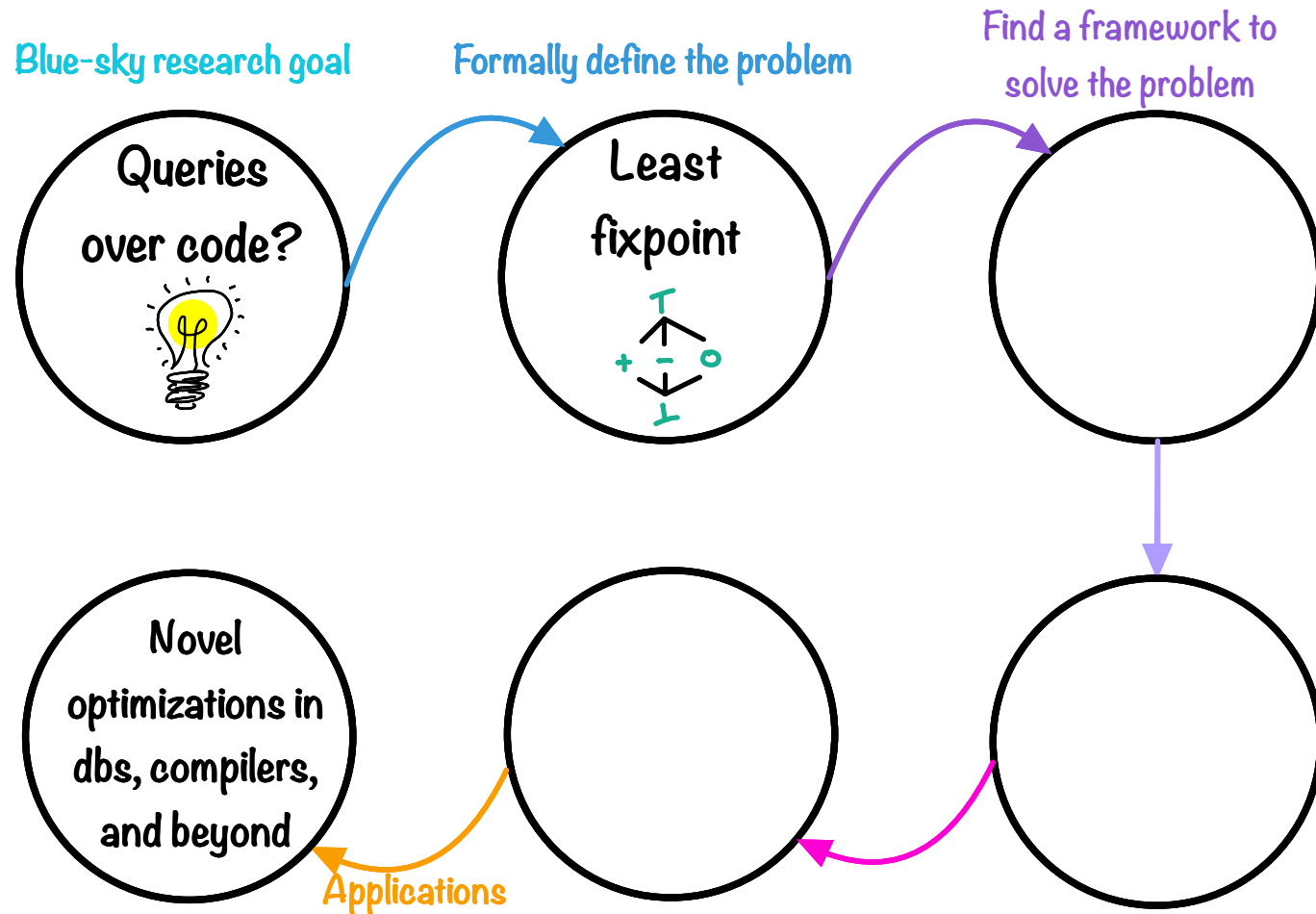
**Need a common framework to express these problems**

## Not magic, math

- $\check{F}$  is an abstraction of the program function  $F$
- Every sound approximation must satisfy  $\check{F}(\check{x}) \sqsubseteq \check{x}$ 
  - Want the smallest  $\check{x}$  that satisfies this property
  - iteratively apply  $\check{F}$  until a fixed point is reached
- This class of problems are called **least fixed-point**

**Abstract interpretation is a mathematical foundation<sup>0</sup>**

# The Road to CodeDB

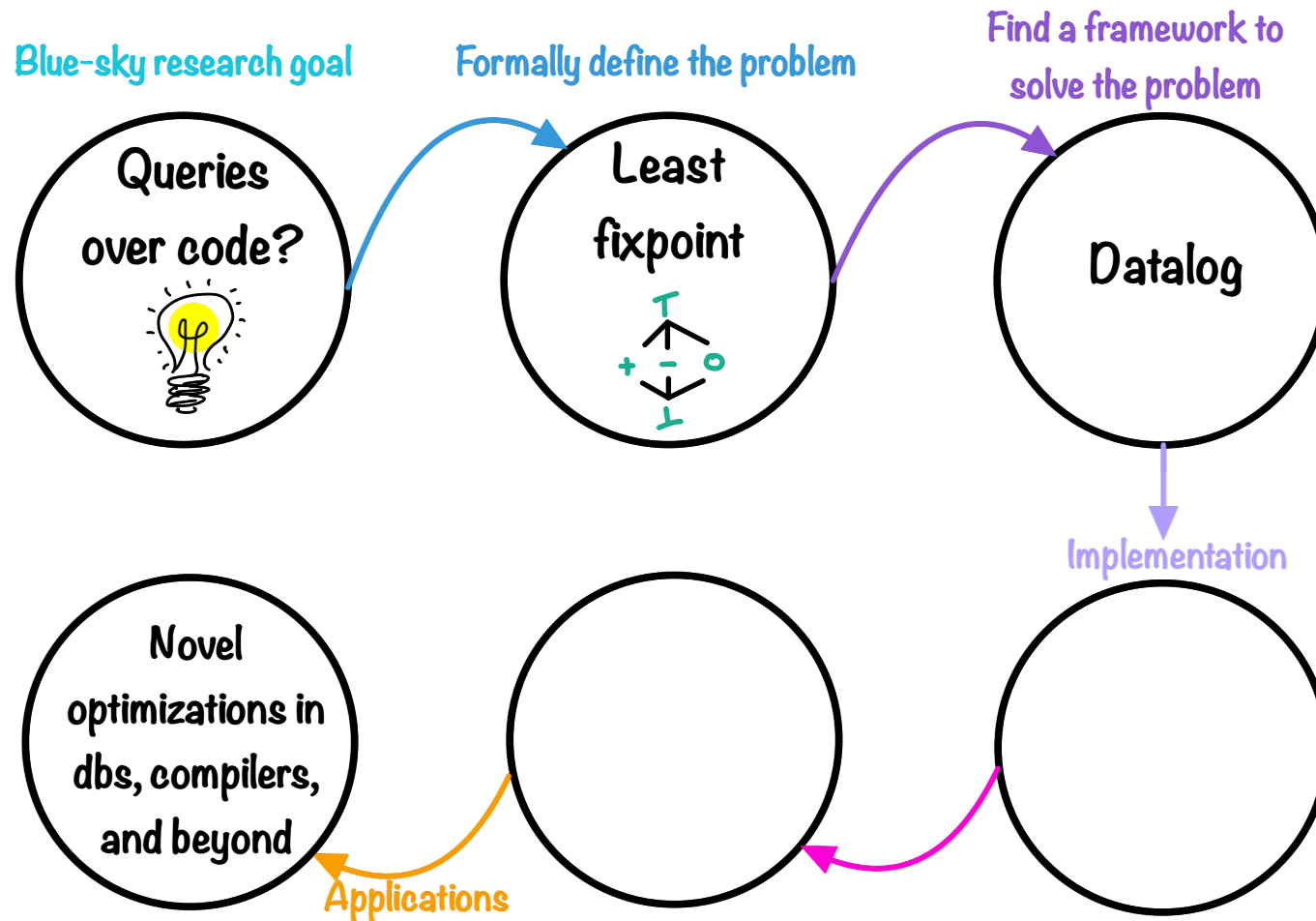


# Datalog 🤝 Fixpoint problems

- Datalog has a fixed-point semantics:
  - a Datalog program is a set of rules
  - its solution is the **minimal model**
- Systems like Souffle use Datalog for static analysis
- Datalog rules are tedious to write and easy to automate

**We can build our query engine on Datalog!**

# The Road to CodeDB



# JIT Data Virtualizing, over Code

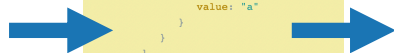
- *Virtualize* [Karpathiotakis et al] over typed IR to extract Datalog facts
  - Avoid ingesting data, JIT generation of operators
- TASTy = Typed Abstract Syntax Trees
  - Includes source positions, types, etc. + TastyQuery
- Scala metaprogramming generates + stores TASTy
  - Runtime inspection, operator specialization, incremental compilation, etc.
  - Compile-time or runtime

**This is my main research focus 😊**

# Datalog Queries == Datalog Data

- Our data layer is essentially facts derived from TASTy
- **Predefine Datalog rules as out-of-the-box queries**

```
val a = F
val b = a
val c = b
if(c) e
```



```
body: Term.Block {
  - state: {
    - Defn.Val {
      mods: [ ]
      - pats: [
        - Pat.Var {
          - name: Term.Name ~ Node {
            value: "a"
          }
        }
      ]
      + rhs: Lit.Boolean {value, syntax}
    }
    - Defn.Val {
      mods: [ ]
      + pats: [1 element]
      + rhs: Term.Name {value}
    }
  }
}
```

Assign(a, F)

Assign(b, a)

Assign(c, b)



PointsTo(v<sub>1</sub>, v<sub>2</sub>)

isFalse(v)

Reachable(f<sub>1</sub>, f<sub>2</sub>)

deadCode(expr)

deadCode?(e)

**Data  
code**

**TASTy**

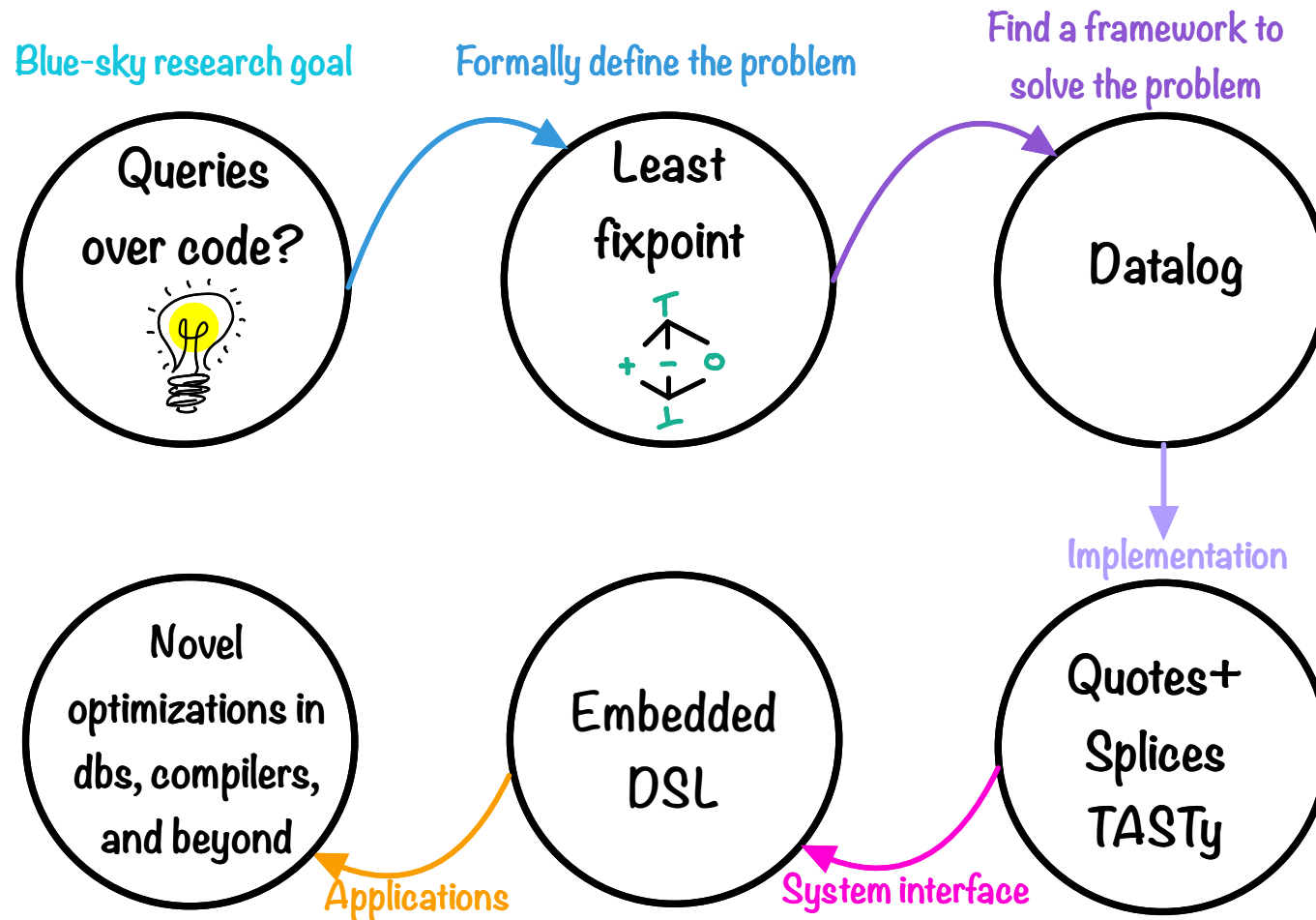
**Datalog  
Facts**

**Datalog  
Rules**

**Query  
code**

**Query language + execution engine + storage layer** <sup>15</sup>

# The Road to CodeDB



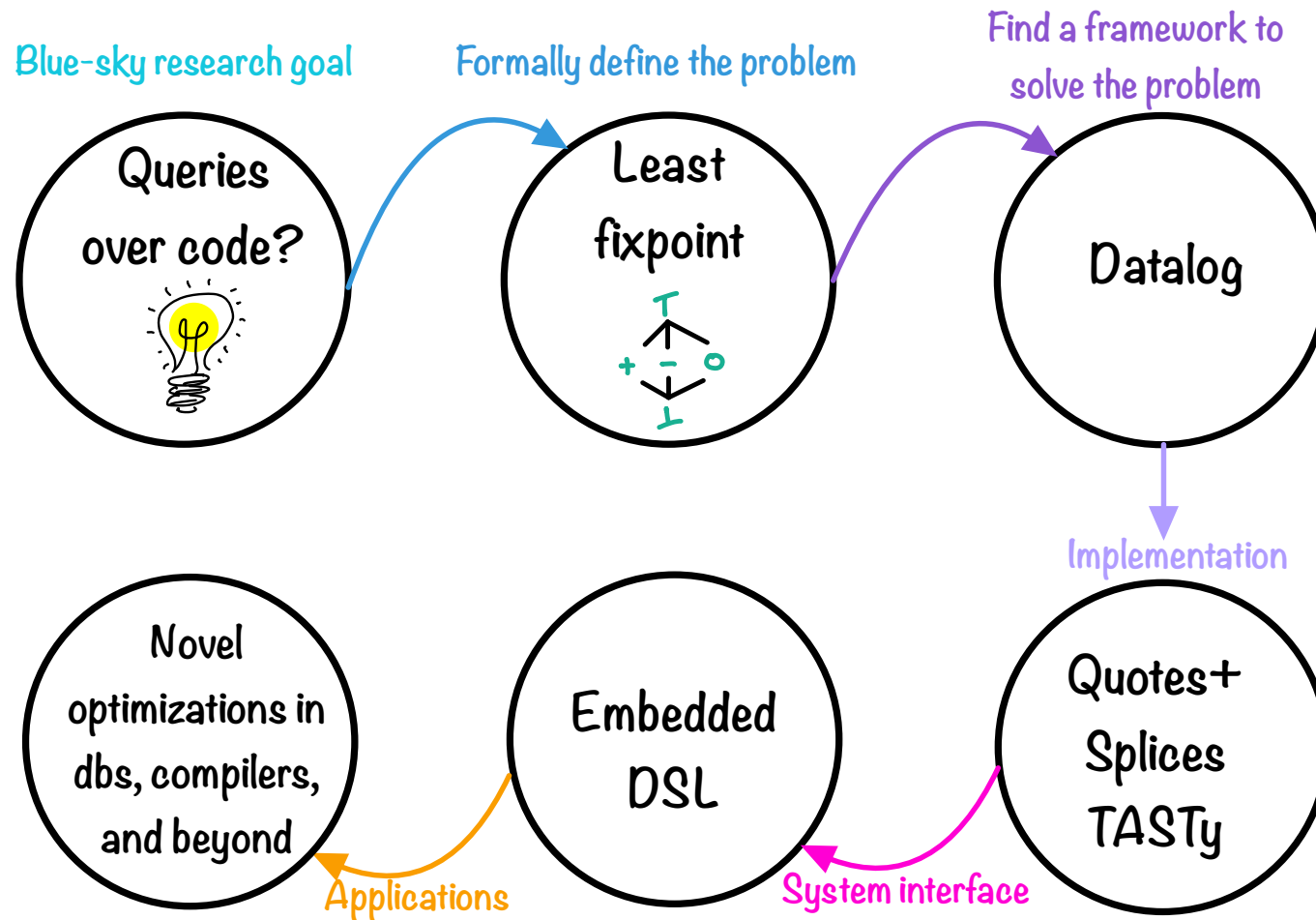


# Composing Queries

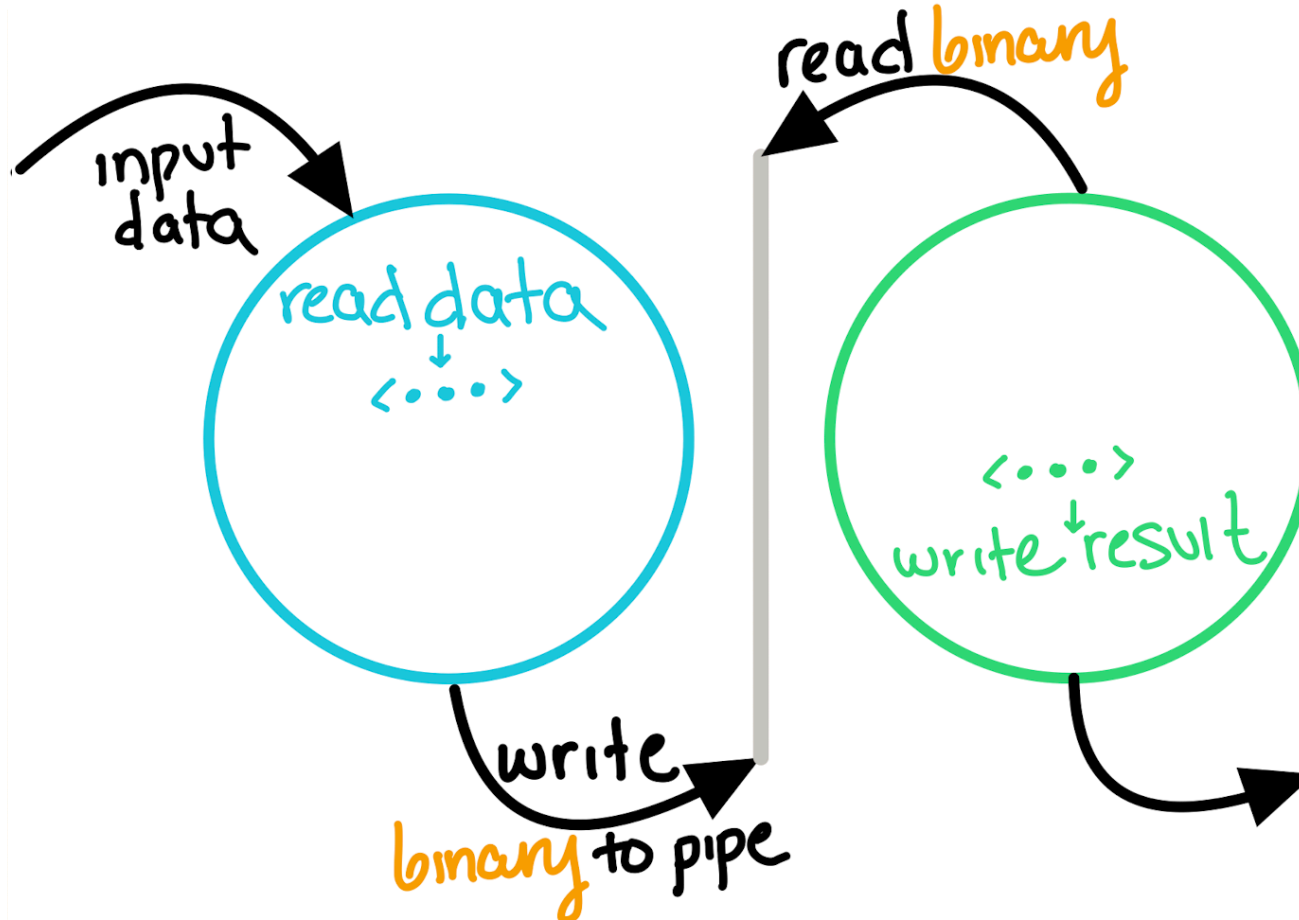
- Constant propagation analysis: `isTrue(s)`; `isFalse(s)`
- Reachability analysis: `isReachable(s)`
- Combination = a conditional constant propagation analysis

**Can compose queries to construct new static analyses**

# The Road to CodeDB



# Identifying Wasted Work



[CIDR2022]

**CodeDB opens new optimization opportunities**

## Arriving to CodeDB

- We need systems that improve code performance & correctness
- We get there with powerful data management techniques

WIP @ [github.com/aherlihy/datalog](https://github.com/aherlihy/datalog)  
[herlihyap@gmail.com](mailto:herlihyap@gmail.com)

**Thank you!!**