

SEMANTIC CODE SEARCH VIA EQUATIONAL REASONING

Varot Premtoon, James Koppel, Armando Solar-Lezama

Massachusetts Institute of Technology

PLDI 2020

Priya Srikumar

Cornell University

Fall PLDG 2020



```
count = 0
for a in cart:
    if a == item
        count += 1
use (count)
```



```
count = 0
for a in cart:
    if a == item:
        count += 1
use (count)
```



```
use(cart.count(item))
```



```
count = 0
for i in cart:
    if debug:
        print(cart[i])
    if cart[i] == item:
        count += 1
use(count)
```



```
count = 0
for i in range(len(arr)):
    if item != arr[i]:
        continue
    count += 1
use(count)
```



```
count = 0
i = 0
while i < len(cart):
    if cart[i] == k:
        count += 1
    i += 1
use(count)
```



```
for .* in cart
```


Semantic code search

Scope

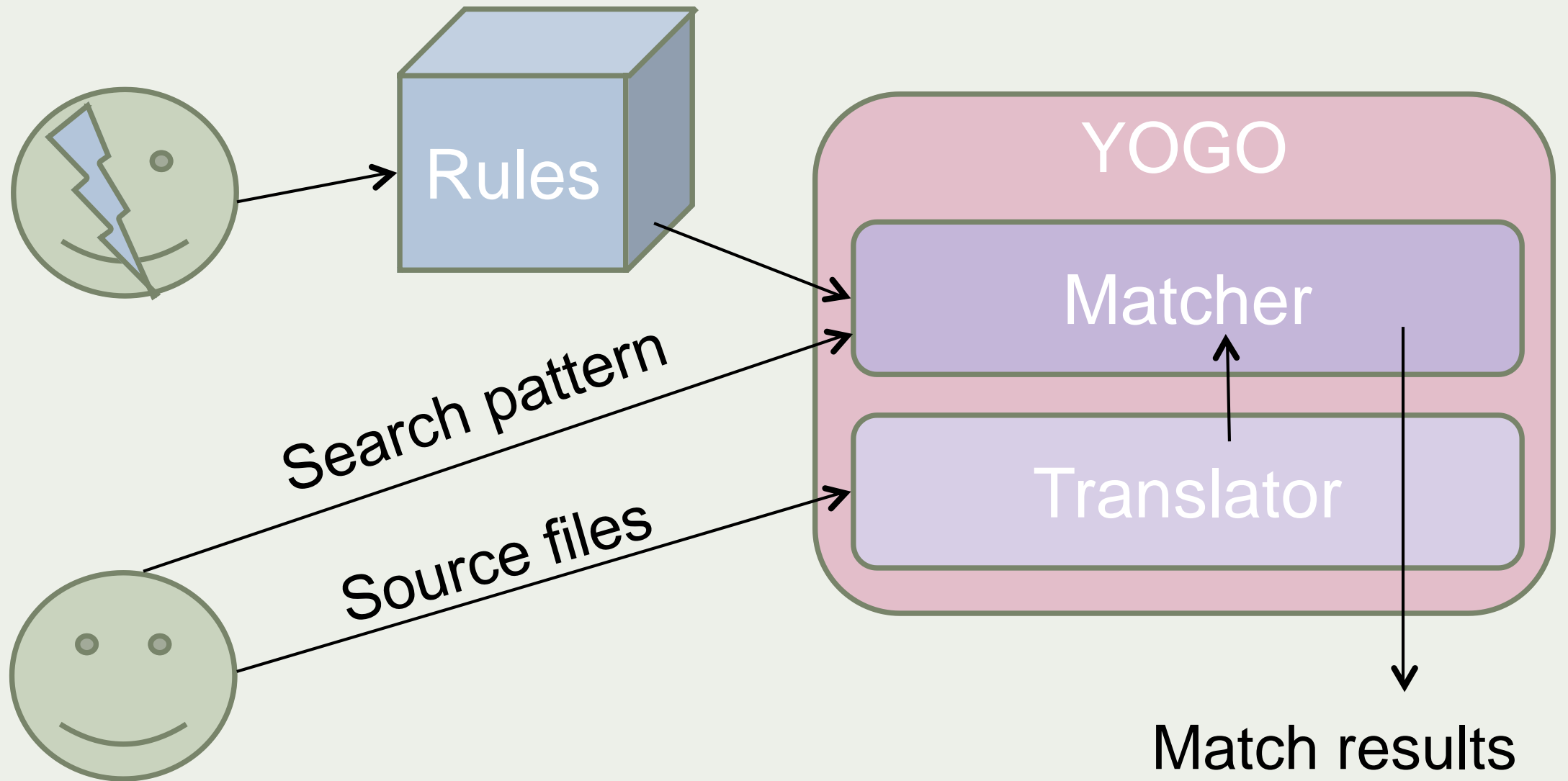
Query type

Performance needs



You Only Grep
Once





Operational evaluation

Search patterns gleaned from common StackOverflow questions

YOGO times out on 1% of the methods searched

False negatives due to tool incompleteness and reasoning

limitations

Graal

13 static analyzers on top of Graal's bytecode analysis

YOGO analyzed 1.2 million LOC in 2.5 hours on 30 AWS instances

Generalizing one of the analyzers' queries revealed a bug it missed!

YOGO query was 60 lines; the analyzer was 330 lines

Usability evaluation

YOGO was straightforward enough to use that the authors

outsourced Graal query coding; it only took 3 days to learn!

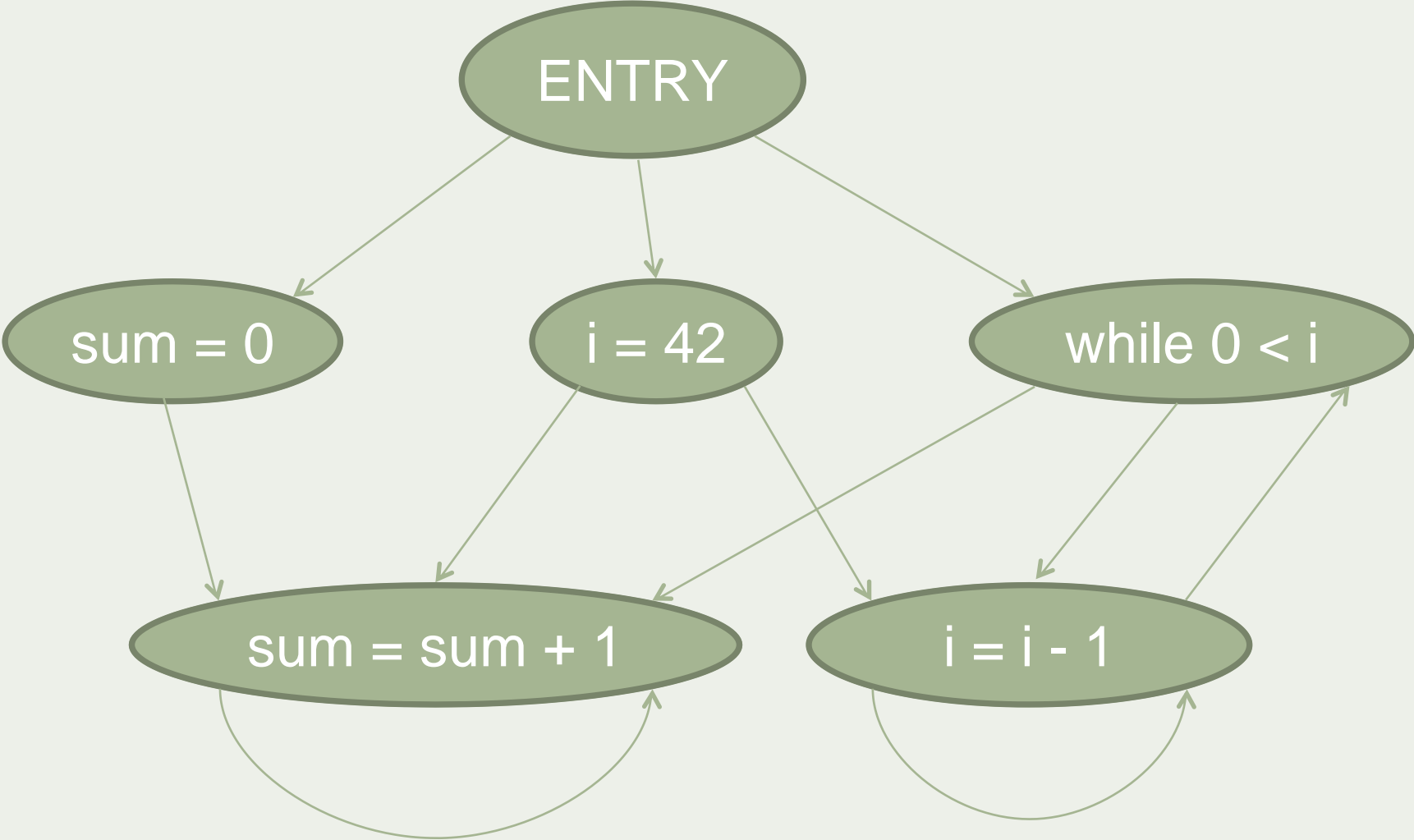


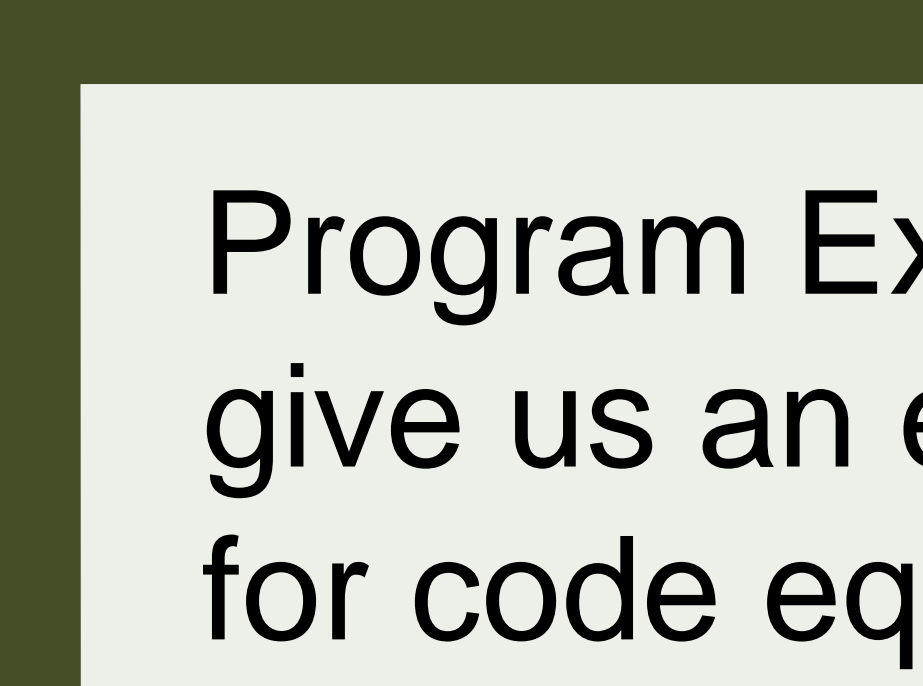
```
i = 42  
while (0 < i)  
    i = i - 1
```



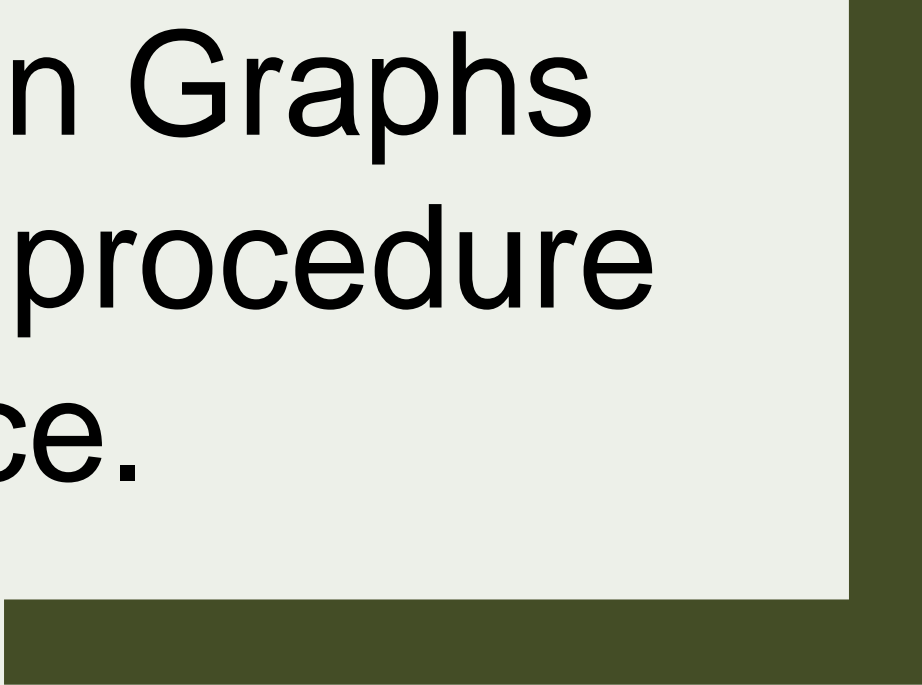
```
sum = 0  
i = 42  
while (0 < i)  
    sum = sum + 1  
    i = i - 1
```

```
i = 42
while (0 < i)
  i = i - 1
```





Program Expression Graphs
give us an efficient procedure
for code equivalence.

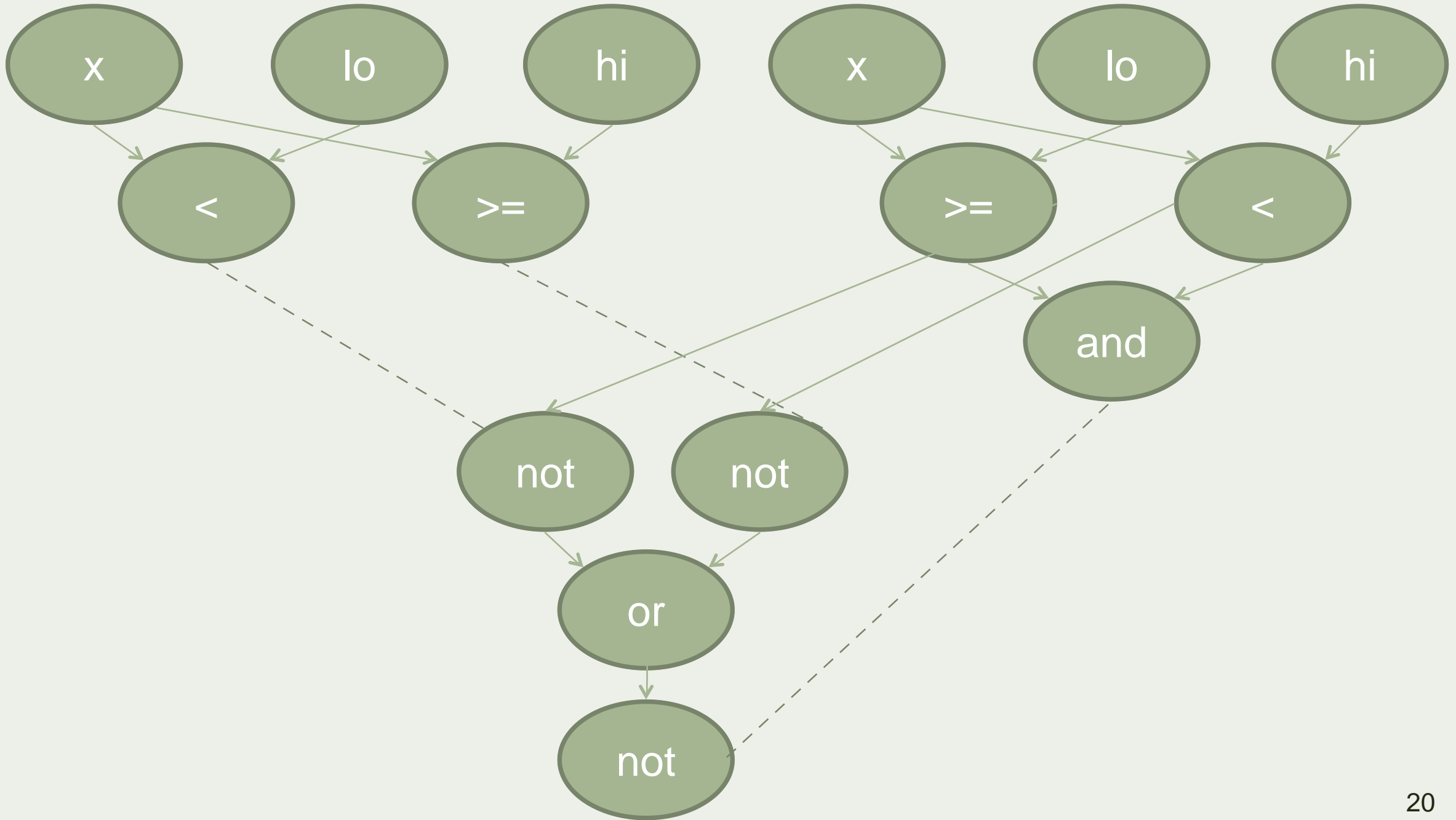




```
x >= lo && x < hi
```

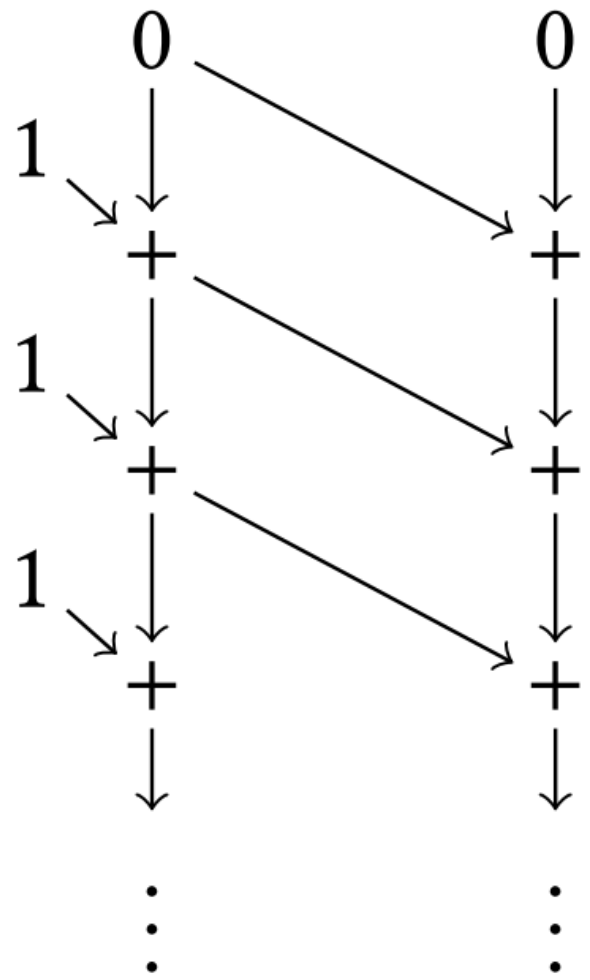


```
(defsearch bound-checking
  (root <- (generic/binop :or (generic/binop :< x lo)
                             (generic/binop :>= x hi))))
```

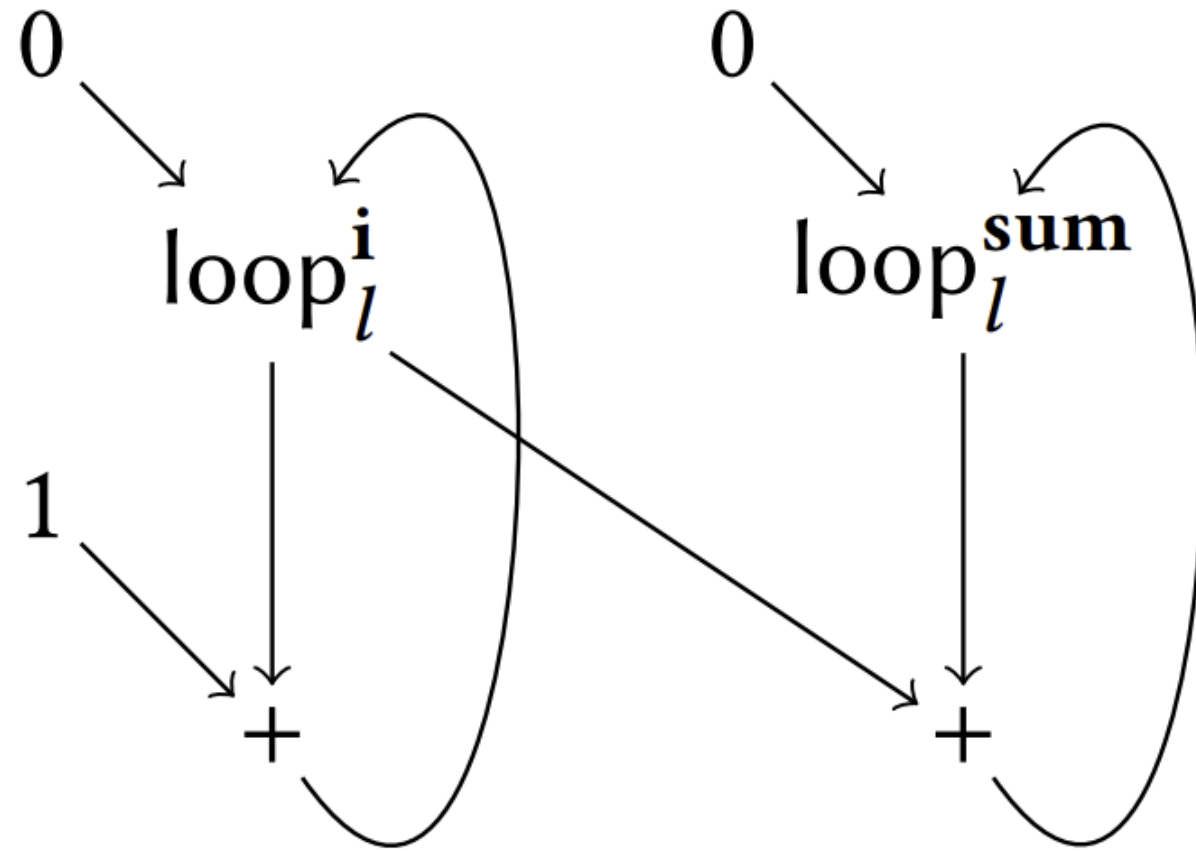




```
sum = 0;  
for(i = 0; ...; i++) {  
    sum += i;  
}
```



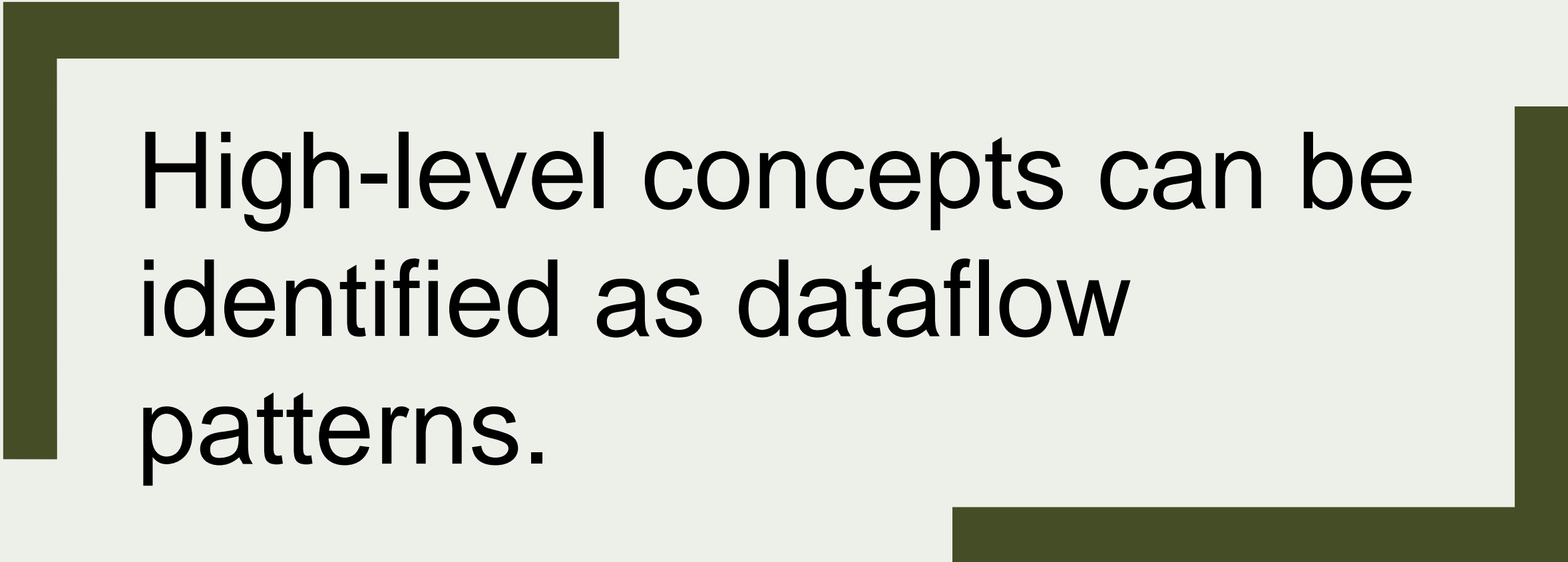
loop_{*l*} sum



Node	Denotation
loop(e_0, e)	$\lambda i. \begin{cases} e_0 & i = 0 \\ e(i-1) & \text{otherwise} \end{cases}$
final(e, l)	$\lambda i. l(\min_{j \in \mathbb{N}} e(j) = \mathbf{false})$



```
i = 0  
while ...:  
    i = i + 1
```



High-level concepts can be identified as dataflow patterns.

Node	Denotation
$\text{seq}(i, k)$	$(i, i + K, i + 2k, \dots)$



```
i = 0
while ...:
    i = i + 1
```

Node	Denotation
iterV((t_0, t_1, \dots))	$\lambda i. \begin{cases} t_i & t_i \neq \perp \\ \perp & \text{otherwise} \end{cases}$
iterP((t_0, t_1, \dots))	$\lambda i. \begin{cases} \mathbf{true} & t_i \neq \perp \\ \mathbf{false} & \text{otherwise} \end{cases}$

INVARIANT(l, e)

PURE(f)

INDEPENDENT(λ_1, λ_2)



```
count = 0
for a in cart:
    if a == item:
        count += 1
use (count)
```

counter \leftarrow loop_l(0, next)

counter \leftarrow loop_{*l*}(0, next)

next \leftarrow cond(iterV_{*l*}(coll) = *k*,
counter + 1,
counter)

answer \leftarrow final_{*l*}(iterP_{*l*}(coll), counter)

INVARIANT_{*l*}(*k*)

Conceptual rewrite soundness

WTS: Code refines behavior of concept specification

Rewrites must be sound wrt their concept's state specification

Rule writers prove soundness; users enjoy the benefits!

Key contributions

Coarser notion of equality supports hierarchical abstraction

Higher-level denotations in the DSL = more abstract
representations

Not a fixed set of equational rules! User-defined rewrite rules

Limitations

Rewriting is hard: syntactic differences, interleaving, overlapping

DSL cannot handle arbitrarily-nested expressions i.e. $f(g(x))$

Impure function calls assume no relation between memory states

Interprocedural patterns not supported; cause graph explosion

Comparison to other code search methods

Static analysis abstractions are too coarse

SMT methods struggle with loops and function calls

Other searches do not scale well

Reflections

Code searcher + bug finder

Generality makes abilities hard to characterize

Presents a technique with a lot of reusability