

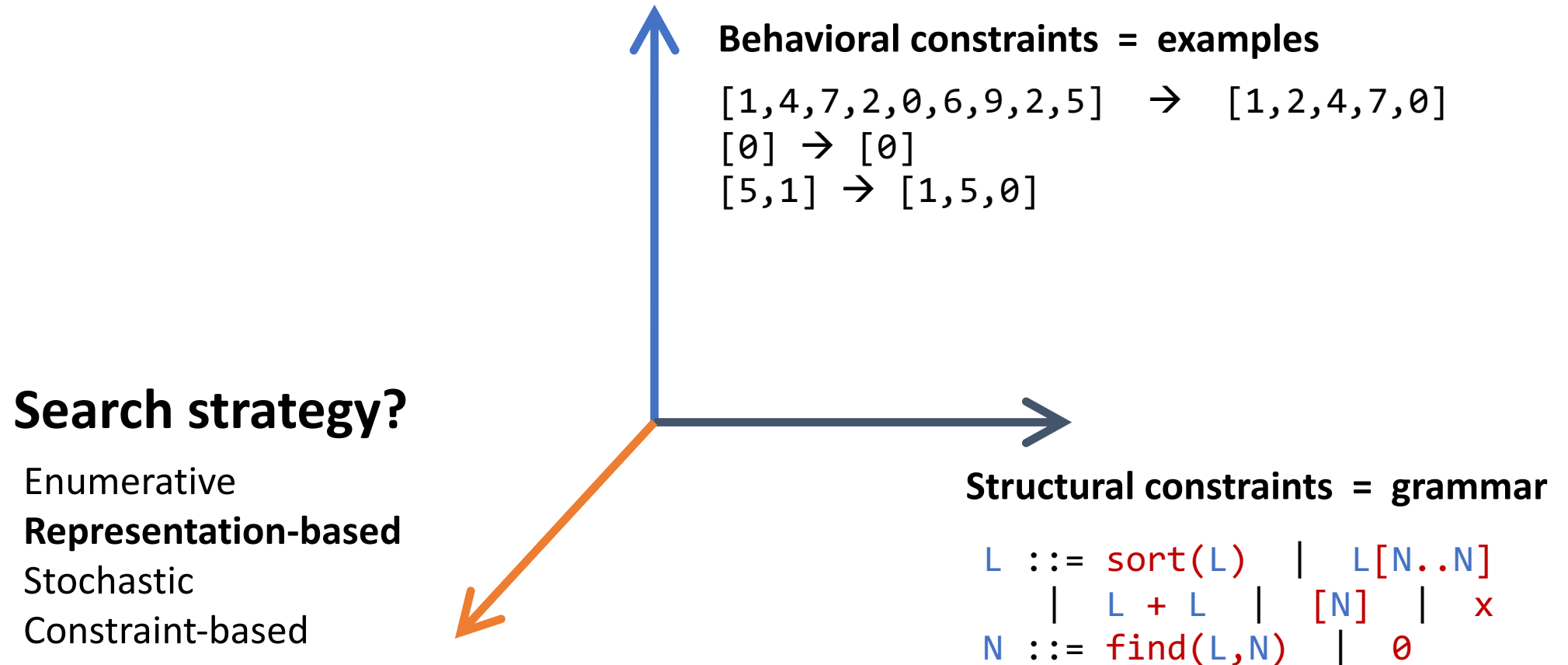
# #9: Representation-based Search

**Sankha Narayan Guria**

EECS 700: Introduction to Program Synthesis



# The problem statement



# Representation-based search

- **Idea:**

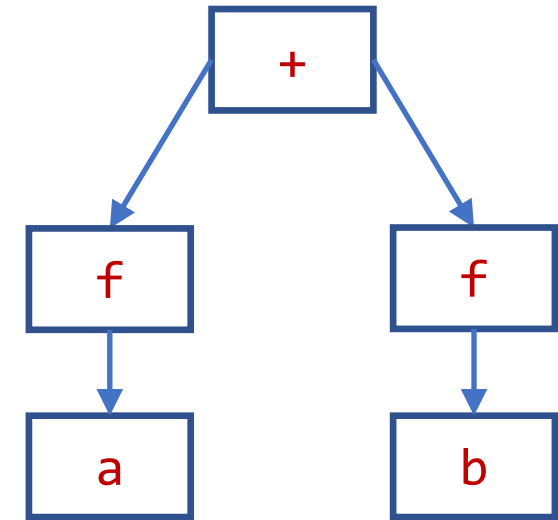
1. build a data structure that compactly represents good parts of the program space
2. extract solution from that data structure

# Compact term representation

- Consider the space of 9 programs:

$f(a) + f(a)$	$f(a) + f(b)$	$f(a) + f(c)$
$f(b) + f(a)$	$f(b) + f(b)$	$f(b) + f(c)$
$f(c) + f(a)$	$f(c) + f(b)$	$f(c) + f(c)$

- Can we represent this compactly?
  - observation 1: same top level structure, independent subterms

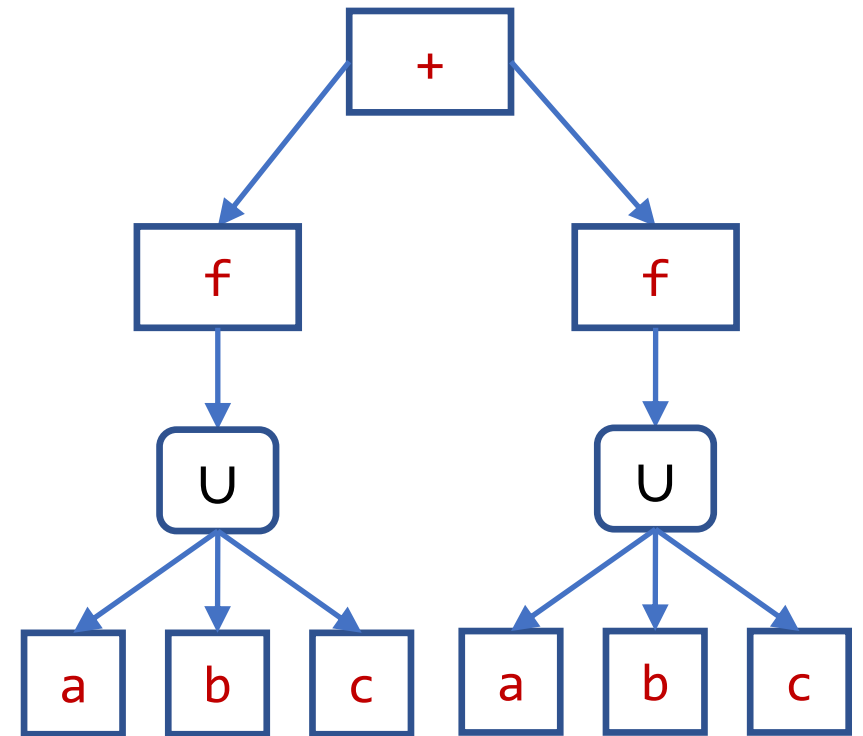


# Compact term representation

- Consider the space of 9 programs:

$f(a) + f(a)$	$f(a) + f(b)$	$f(a) + f(c)$
$f(b) + f(a)$	$f(b) + f(b)$	$f(b) + f(c)$
$f(c) + f(a)$	$f(c) + f(b)$	$f(c) + f(c)$

- Can we represent this compactly?
  - observation 1: same top level structure, independent subterms
  - observation 2: shared sub-spaces

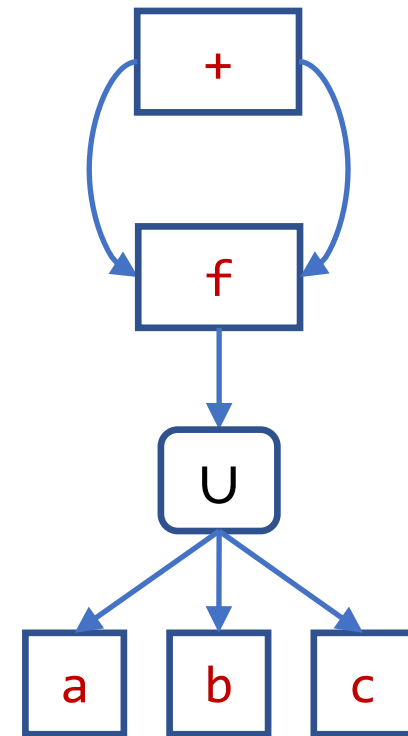


# Compact term representation

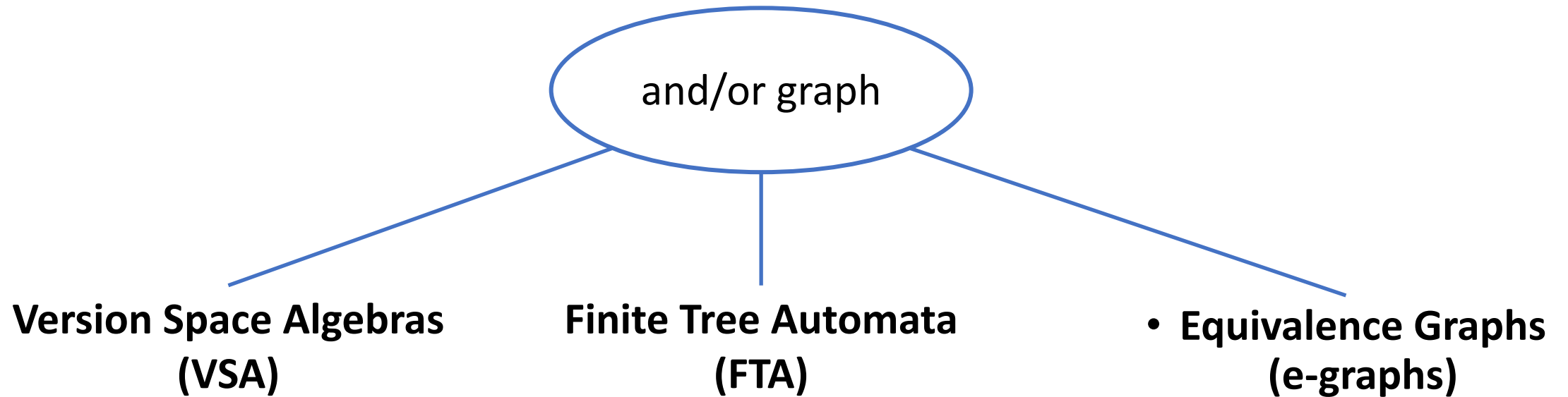
- Consider the space of 9 programs:

$f(a) + f(a)$	$f(a) + f(b)$	$f(a) + f(c)$
$f(b) + f(a)$	$f(b) + f(b)$	$f(b) + f(c)$
$f(c) + f(a)$	$f(c) + f(b)$	$f(c) + f(c)$

- Can we represent this compactly?
  - observation 1: same top level structure, independent subterms
  - observation 2: shared sub-spaces
- Key idea: use and-or graph!



# Representation-based search

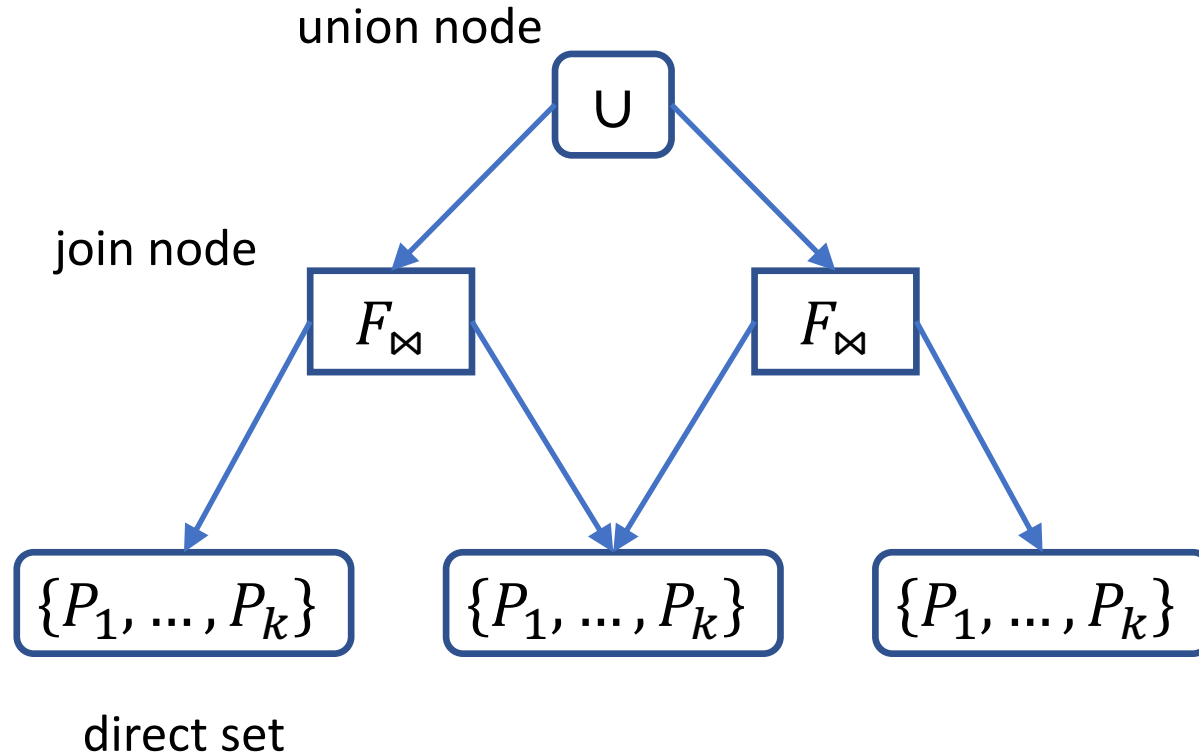


# Version Space Algebra

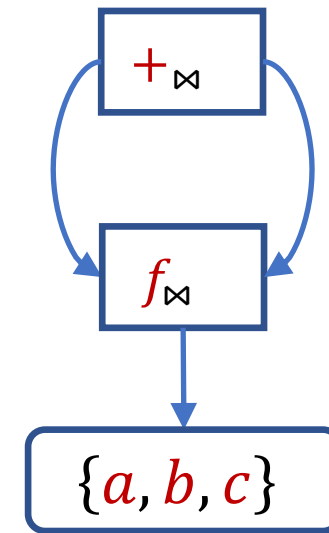
- **Idea:** build a graph that succinctly represents the space of *all* programs consistent with examples
  - called a **version space**
- Operations on version spaces:
  - learn  $\langle \textcolor{brown}{i}, \textcolor{blue}{o} \rangle \rightarrow \text{VS}$
  - $\text{VS}_1 \cap \text{VS}_2 \rightarrow \text{VS}$
  - extract  $\text{VS} \rightarrow \text{program}$
- Algorithm:
  1. learn a VS for each example
  2. intersect them all
  3. extract any (or best) program



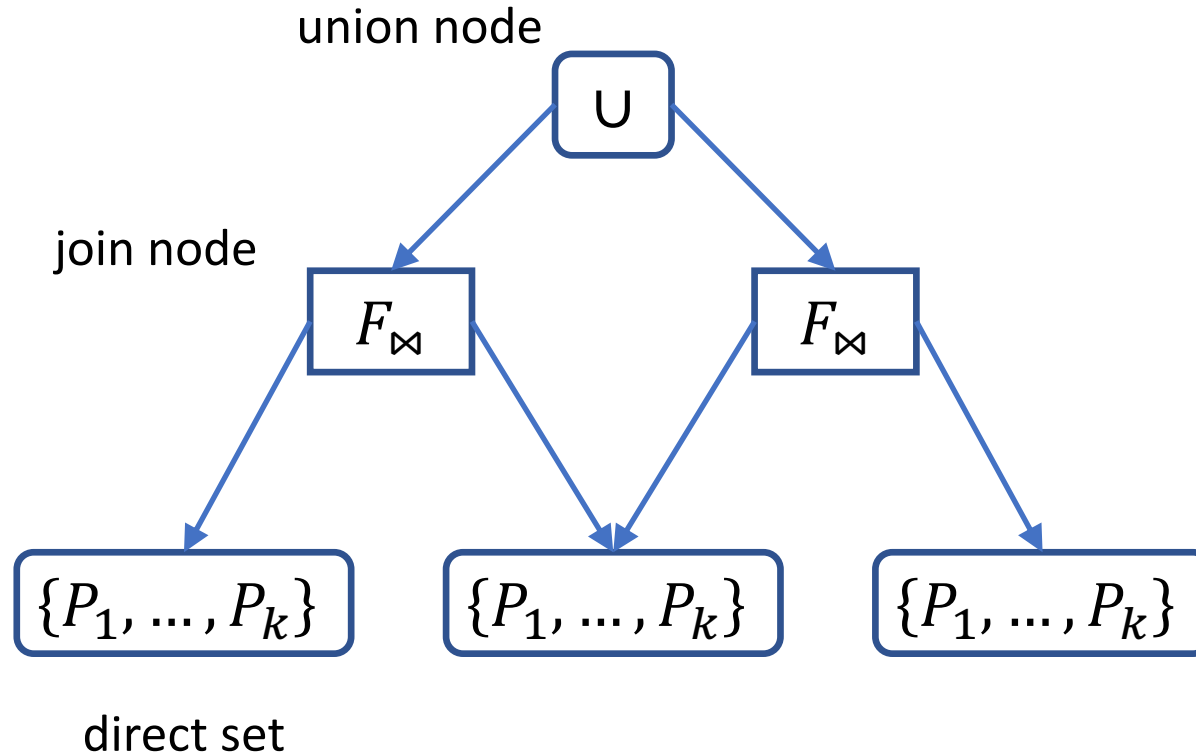
# Version Space Algebra



example:



# Version Space Algebra



Volume of a VSA  
(the number of nodes)  $V(VSA)$

Size of a VSA  
(the number of programs)  $|VSA|$

$$V(VSA) = O(\log|VSA|)$$

# VSA-based search

- Mitchell: *Generalization as search*. AI 1982
- Lau, Domingos, Weld. *Version space algebra and its application to programming by example*. ICML 2000
- Gulwani: *Automating string processing in spreadsheets using input-output examples*. POPL 2011.
  - Follow-up work: BlinkFill, FlashExtract, FlashRelate, ...
  - generalized in the PROSE framework

# FlashFill

[Gulwani '11]

Simplified grammar:

$E ::= F \mid \text{concat}(F, E)$	“Trace” expression
$F ::= \text{cstr}(str) \mid \text{sub}(P, P)$	Atomic expression
$P ::= \text{cpos}(num) \mid \text{pos}(R, R)$	Position expression
$R ::= \text{tokens}(T_1, \dots, T_n)$	Regular expression
$T ::= C \mid C^+$	Token expression
$C ::= ws \mid digit \mid alpha \mid Alpha \mid \$ \mid ^ \mid \dots$	

# FlashFill: example

0 1 2 3 4 5 6 7 8 9 ...

“Hello POPL 2023” → “POPL’2023”

“Goodbye PLDI 2021” → “PLDI’2021”

```
concat(  
  sub(pos(ws, Alpha), pos(Alpha, ws)),  
  concat(  
    cstr(“’”),  
    sub(pos(ws, digit), pos(digit, $))))
```

$E ::= F \mid \text{concat}(F, E)$

$F ::= \text{cstr}(str) \mid \text{sub}(P, P)$

$P ::= \text{cpos}(num) \mid \text{pos}(R, R)$

$R ::= \text{tokens}(T_1, \dots, T_n)$

$T ::= C \mid C+$

# VSAs for Flashfill

- Recall operations on version spaces:
  - $\text{learn } \langle i, o \rangle \rightarrow \text{VS}$
  - $\text{VS}_1 \cap \text{VS}_2 \rightarrow \text{VS}$
  - $\text{extract VS} \rightarrow \text{program}$
- How do we implement learn?
  - define  $\text{learn}_N \langle i, o \rangle$   
for every non-terminal N
  - build VS top-down,  
propagating  $\langle i, o \rangle$  the example

$E ::= F \mid \text{concat}(F, E)$

$F ::= \text{cstr}(str) \mid \text{sub}(P, P)$

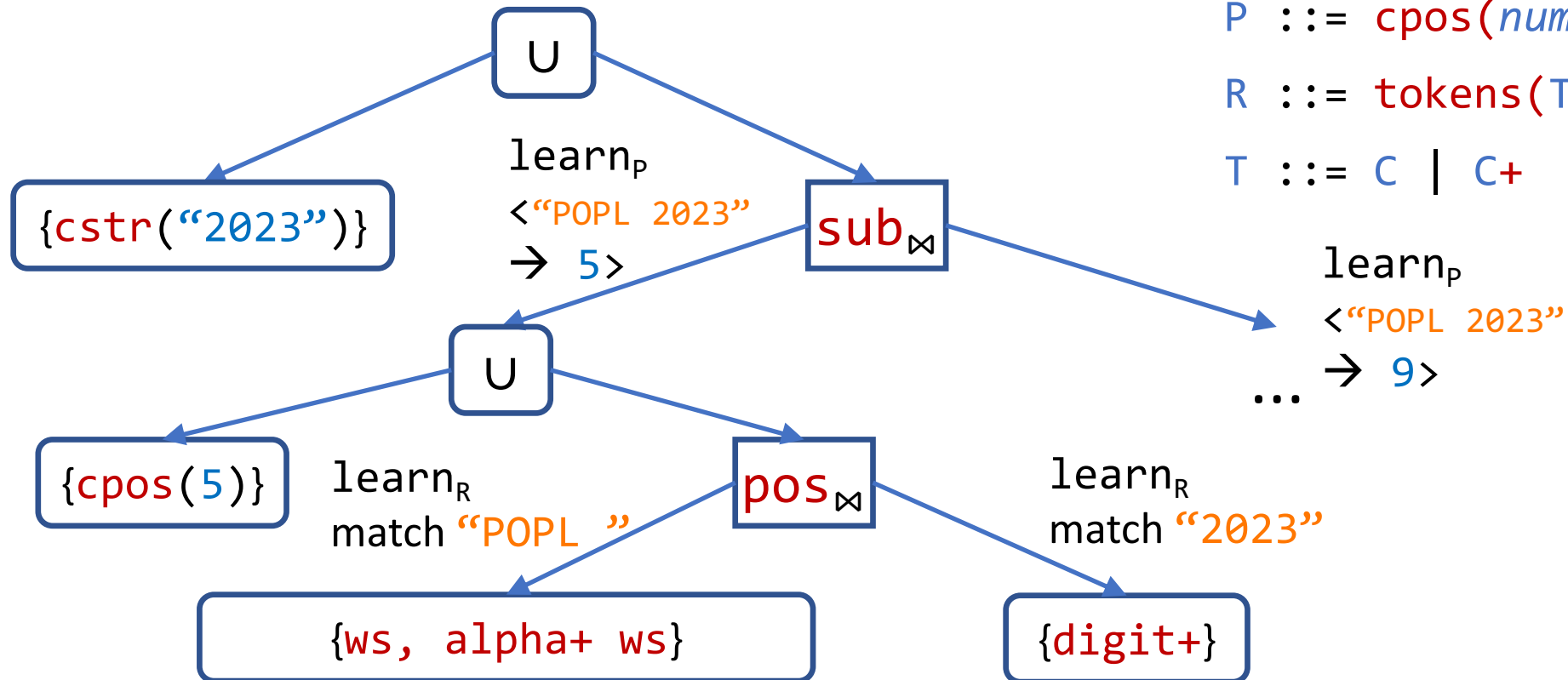
$P ::= \text{cpos}(num) \mid \text{pos}(R, R)$

$R ::= \text{tokens}(T_1, \dots, T_n)$

$T ::= C \mid C+$

# Learning atomic expressions

- $\text{learn}_F \langle \text{"POPL 2023"} \rangle \rightarrow \text{"2023"} \rangle$



$F ::= \text{cstr}(str) \mid \text{sub}(P_1, P_2)$

$P ::= \text{cpos}(num) \mid \text{pos}(R_1, R_2)$

$R ::= \text{tokens}(T_1, \dots, T_n)$

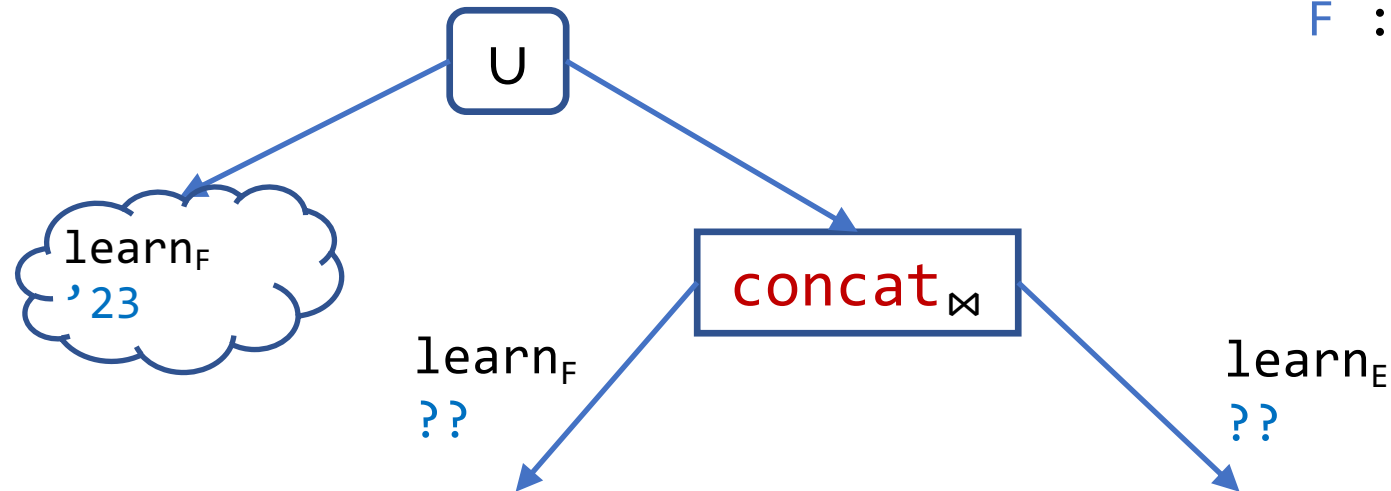
$T ::= C \mid C+$

# Learning trace expressions

- $\text{learn}_E \langle \text{"POPL 2023"} \rightarrow \text{"'23"} \rangle$

$E ::= F \mid \text{concat}(F, E)$

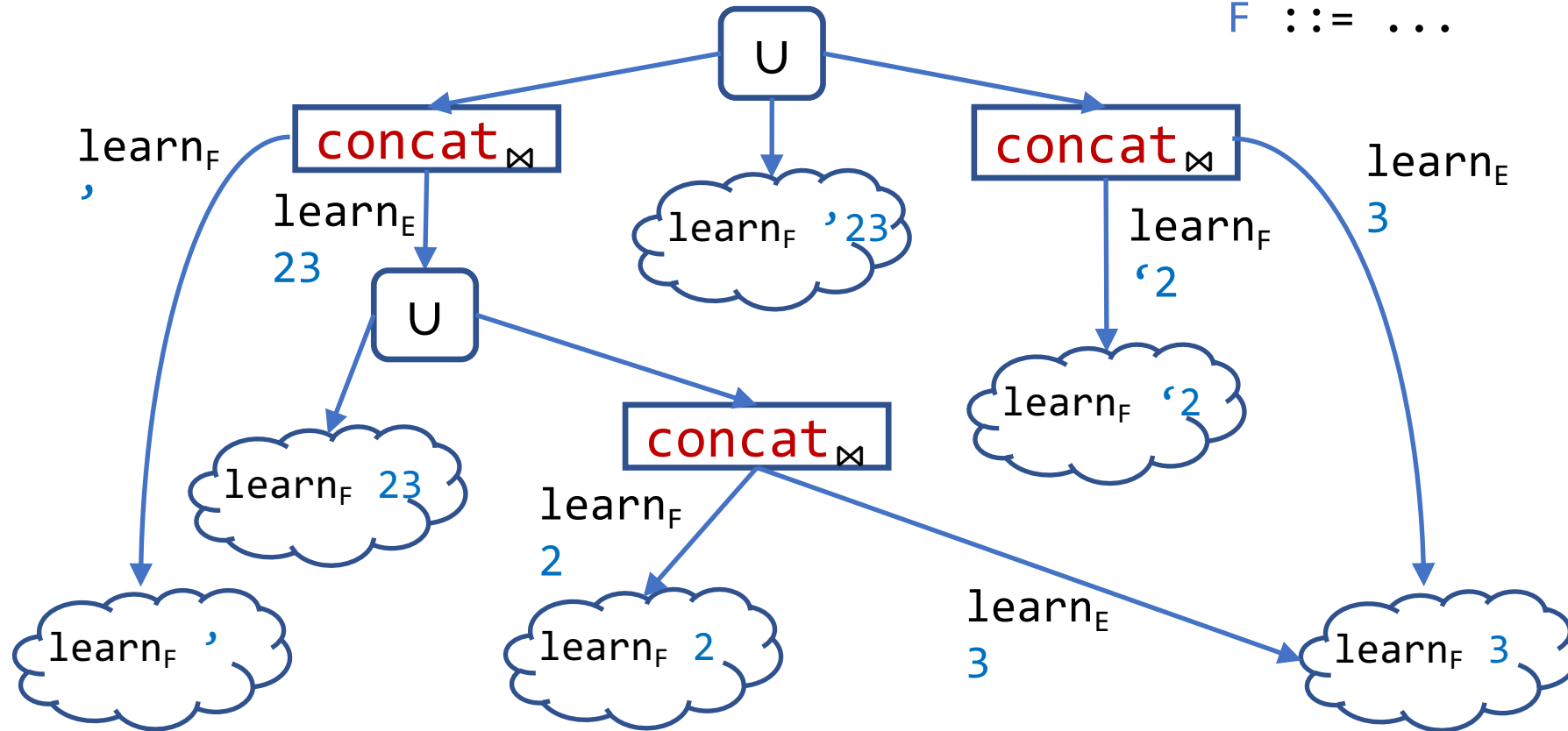
$F ::= \dots$





# Learning trace expressions

- $\text{learn}_E \langle \text{"POPL 2023"} \rightarrow \text{"'23"} \rangle$        $E ::= F \mid \text{concat}(F, E)$   
 $F ::= \dots$



# VSAs for Flashfill

- Recall operations on version spaces:
  - learn  $\langle i, o \rangle \rightarrow VS$
  - $VS_1 \cap VS_2 \rightarrow VS$
  - extract  $VS \rightarrow \text{program}$
- How do we implement intersection?
  - top-down
  - union: intersect all pairs of children
  - join: intersect children pairwise

$E ::= F \mid \text{concat}(F, E)$

$F ::= \text{cstr}(str) \mid \text{sub}(P, P)$

$P ::= \text{cpos}(num) \mid \text{pos}(R, R)$

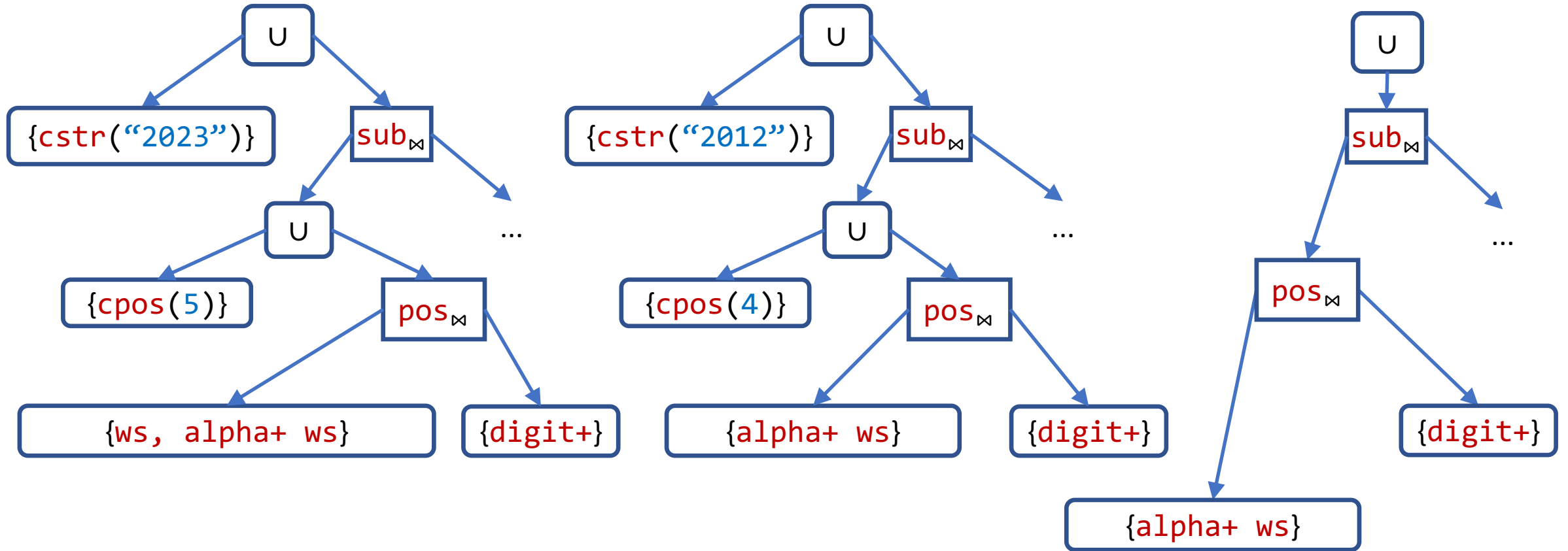
$R ::= \text{tokens}(T_1, \dots, T_n)$

$T ::= C \mid C+$

# Intersection

“POPL 2023” → “2023”

“ 3M 2012” → “2012”



# VSAs for Flashfill

- Recall operations on version spaces:
  - learn  $\langle i, o \rangle \rightarrow VS$
  - $VS_1 \cap VS_2 \rightarrow VS$
  - extract  $VS \rightarrow \text{program}$
- How do we implement extract?
  - any program: just pick one child from every union
  - best program: shortest path in a DAG

$E ::= F \mid \text{concat}(F, E)$

$F ::= \text{cstr}(str) \mid \text{sub}(P, P)$

$P ::= \text{cpos}(num) \mid \text{pos}(R, R)$

$R ::= \text{tokens}(T_1, \dots, T_n)$

$T ::= C \mid C+$

# Discussion

- What do VSAs remind you of in the enumerative world?
  - VSA learning ~ top-down search with top-down propagation
- How are they different?
  - Caching of sub-problems (DAG!)
  - Can construct one per example and intersect
  - This allows it to guess arbitrary constants!

# Discussion

- Why could we build a finite representation of all solutions?
  - Could we do it for this language?

$E ::= F + F$

$F ::= k \mid x$

$k \in \mathbb{Z}$      $+$  is integer addition

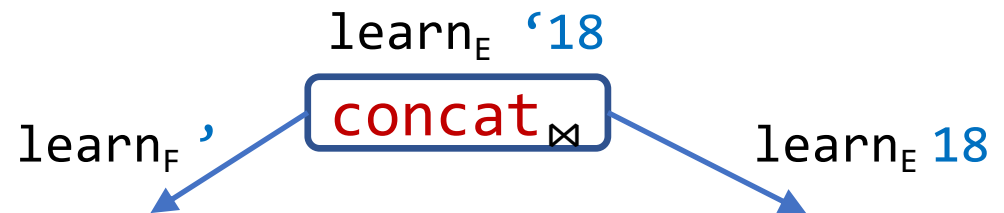
- What about this language?

$E ::= E + 1 \mid x$

# DSL restrictions: efficiently invertible

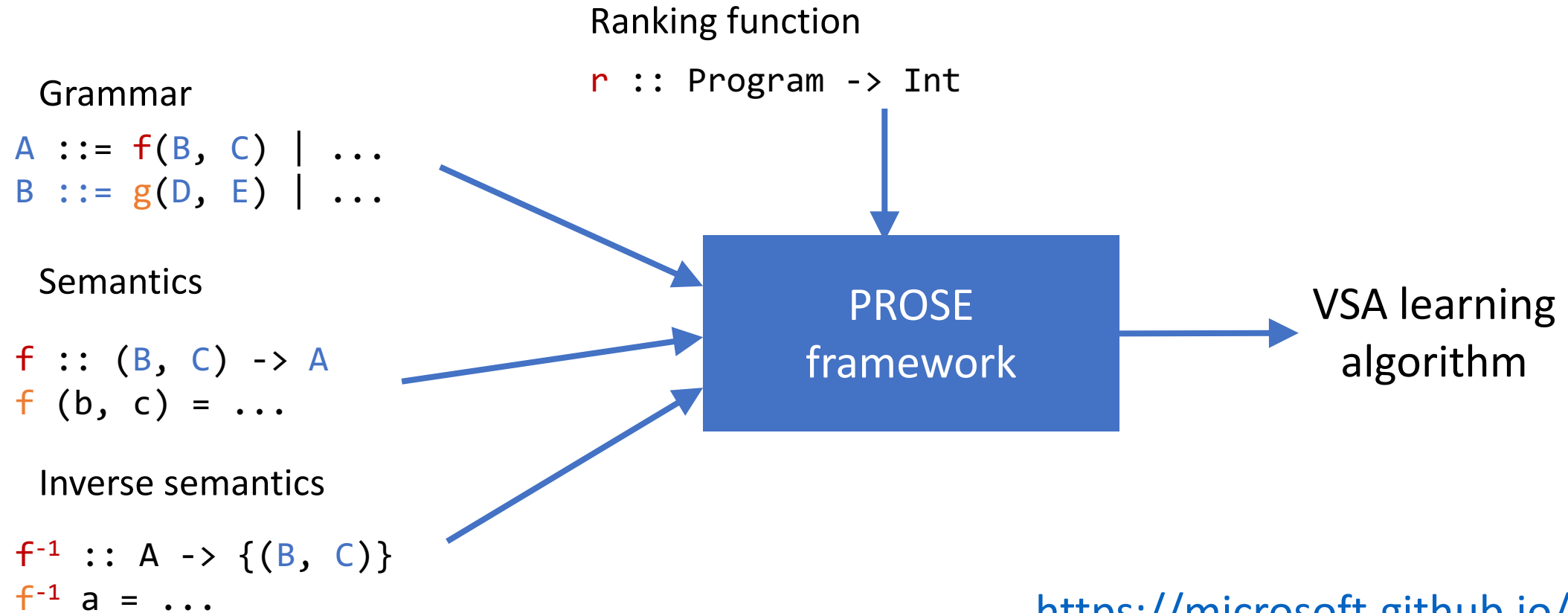
- Every operator has a small, easily computable inverse
  - Example when an inverse is small but hard to compute?
- The space of sub-specs is finite
  - either non-recursive grammar
  - or finite space of values for the recursive non-terminal (e.g. bit-vectors)
  - or every recursive production generates a strictly smaller spec

$E ::= F \mid \text{concat}(F, E)$



# PROSE

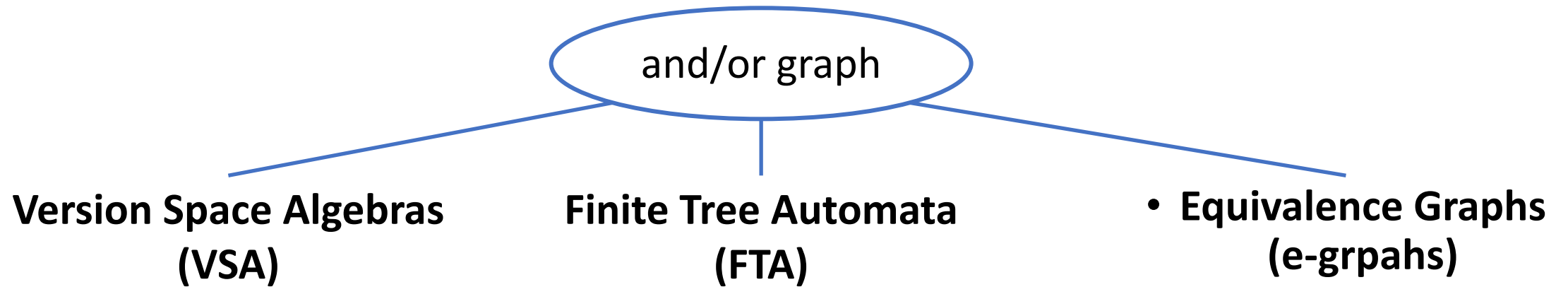
[Polozov, Gulwani '15]



<https://microsoft.github.io/prose/>



# Representation-based search

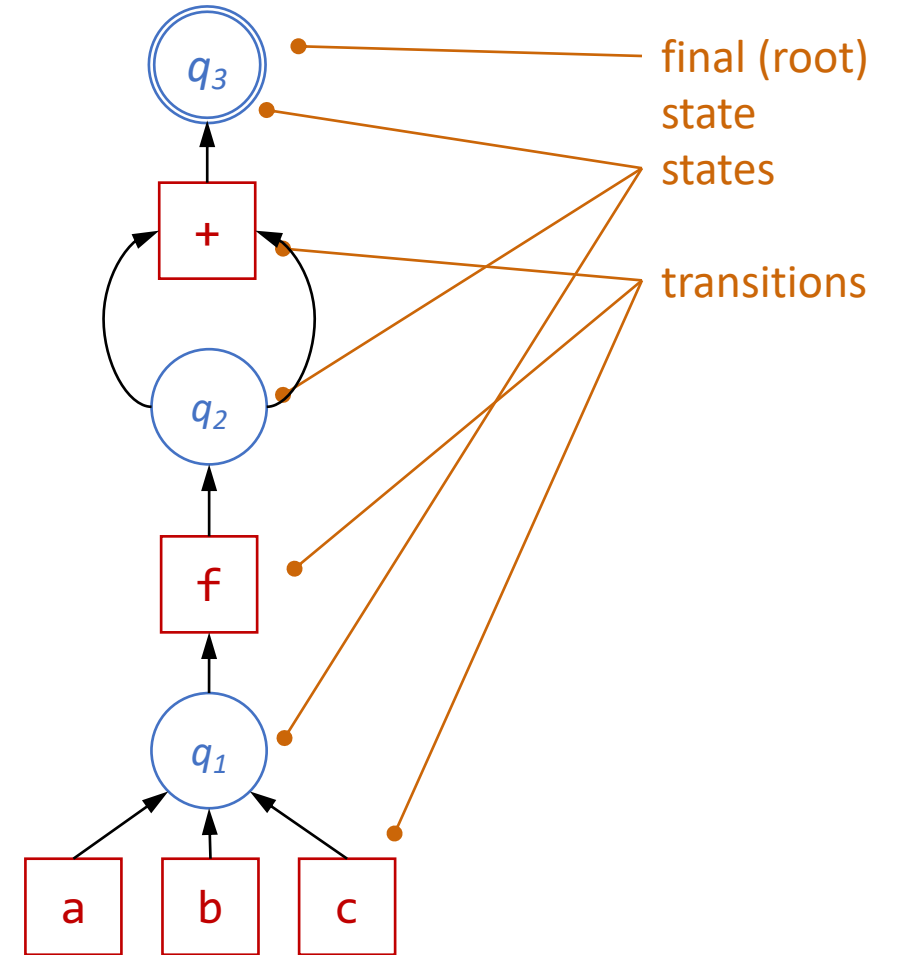
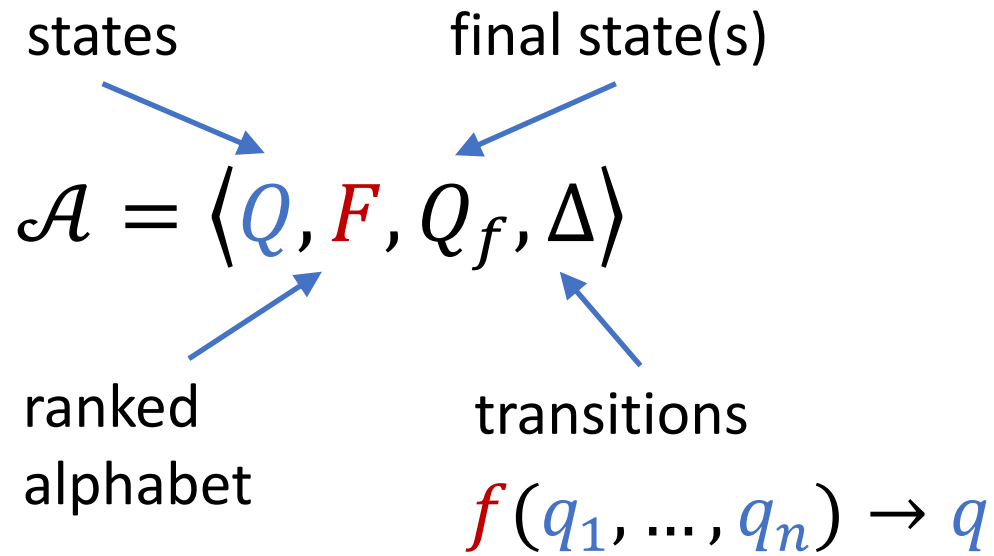


**ops:** learn-1, intersect, extract

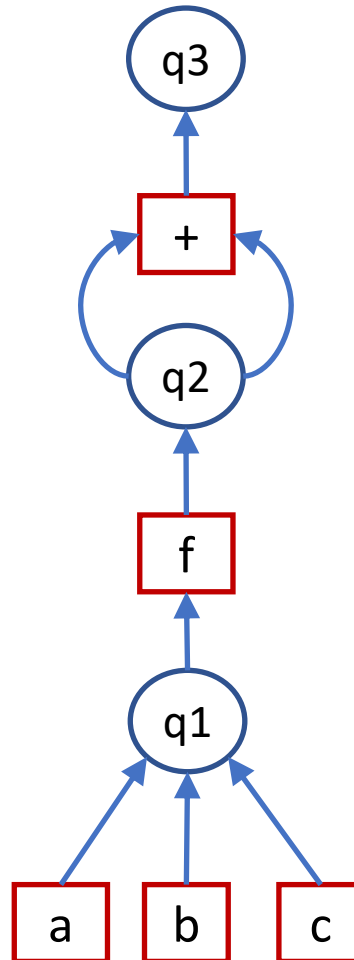
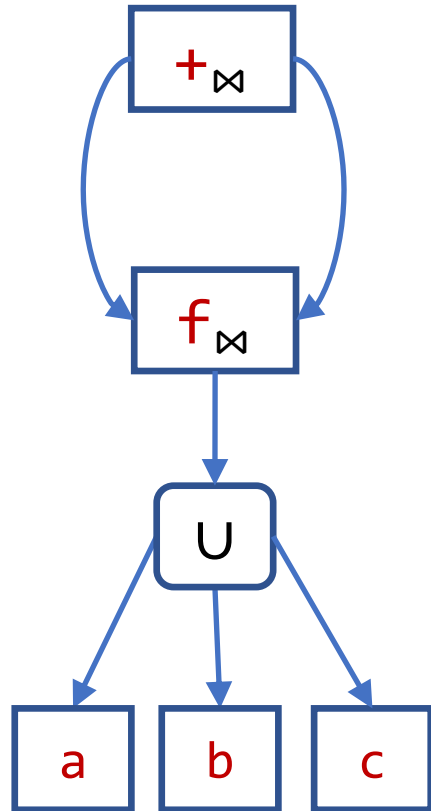
**DSL:** efficiently invertible

**similar to:** top-down prop,  
but can infer constants

# Finite Tree Automata



# VSA vs FTA



- Both are and-or graphs
- FTA state = VSA union node
  - in VSAs singleton joins are omitted
- FTA transition = VSA join node

# FTA-based search

- **Synthesis of Data Completion Scripts using Finite Tree Automata**

Xinyu Wang, Isil Dillig, Rishabh Singh, *OOPSLA'17*

- **Program Synthesis using Abstraction Refinement**

Xinyu Wang, Isil Dillig, Rishabh Singh, *POPL'18*

- **Searching Entangled Program Spaces**

James Koppel, Zheng Guo, Edsko de Vries, Armando Solar-Lezama, Nadia Polikarpova. *ICFP'22*

# FTA-based search

- **Synthesis of Data Completion Scripts using Finite Tree Automata**  
Xinyu Wang, Isil Dillig, Rishabh Singh, *OOPSLA'17*
- **Program Synthesis using Abstraction Refinement**  
Xinyu Wang, Isil Dillig, Rishabh Singh, *POPL'18*
- **Searching Entangled Program Spaces**  
James Koppel, Zheng Guo, Edsko de Vries, Armando Solar-Lezama, Nadia Polikarpova. *ICFP'22*

# Example

Grammar

$N ::= \text{id}(V) \mid N + T \mid N * T$

$T ::= 2 \mid 3$

$V ::= x$

Spec

$1 \rightarrow 9$

# PBE with Finite Tree Automata

$\langle A, \mathbb{Z} \rangle$   
 $A \in \{\text{N}, \text{T}, \text{X}\}$

$\{\langle \text{N}, 9 \rangle\}$

states

final states

$$\mathcal{A} = \langle Q, F, Q_f, \Delta \rangle$$

alphabet

transitions

$\text{id}, +, *$

$$f(q_1, \dots, q_n) \rightarrow q$$

$$+(\langle \text{N}, 1 \rangle, \langle \text{T}, 2 \rangle) \rightarrow \langle \text{N}, 3 \rangle$$

...

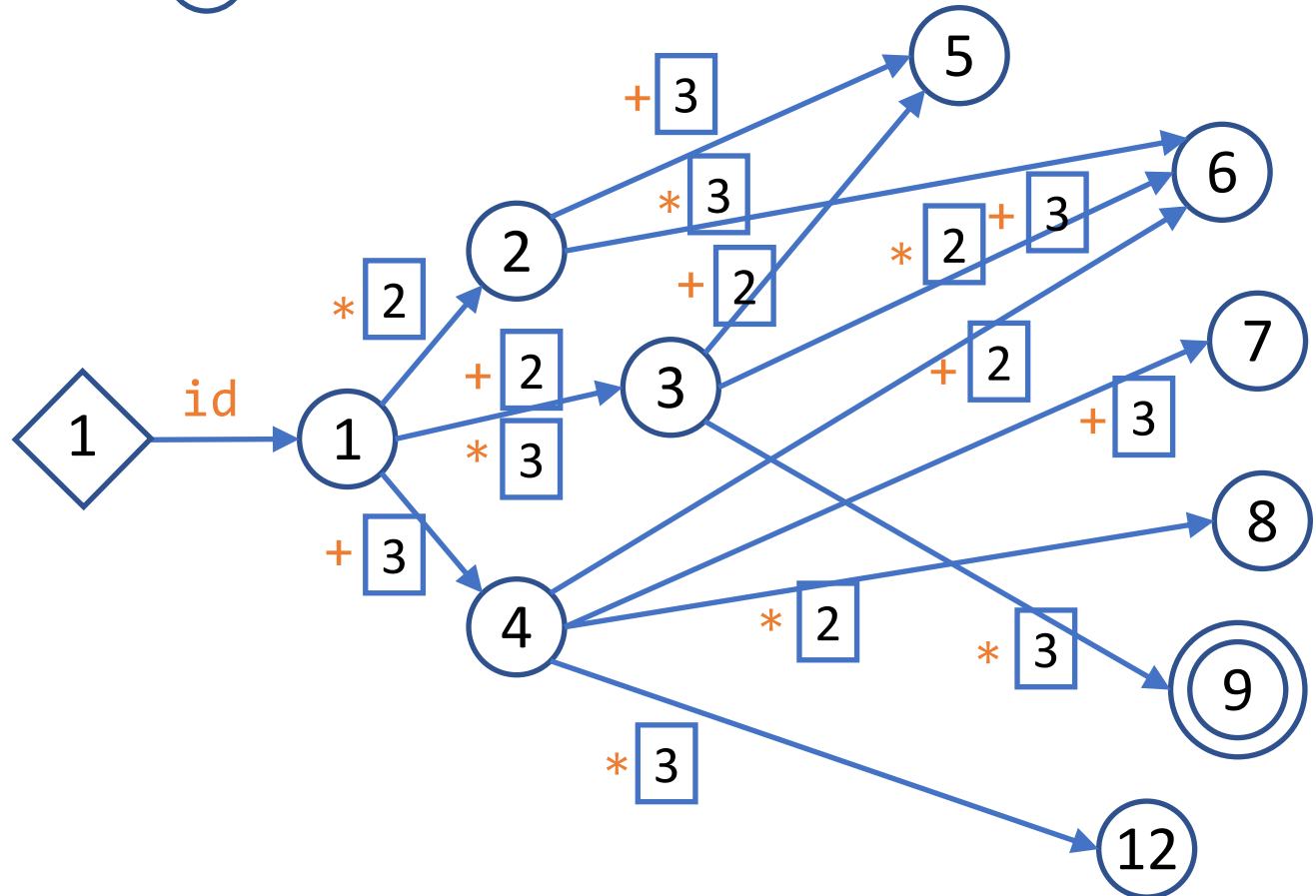
# PBE with Finite Tree Automata

$N ::= \text{id}(V) \mid N + T \mid N * T \quad \bigcirc$

$T ::= 2 \mid 3 \quad \square$

$V ::= x \quad \diamond$

$1 \rightarrow 9$





# Discussion

- What do FTAs remind you of in the enumerative world?
  - FTA ~ bottom-up search with OE
- How are they different?
  - More size-efficient: sub-terms in the bank are replicated, while in the FTA they are shared
  - Hence, can store all terms, not just one representative per class
  - Can construct one FTA per example and intersect
  - More incremental in the CEGIS context!

# FTA-based search

- **Synthesis of Data Completion Scripts using Finite Tree Automata**  
Xinyu Wang, Isil Dillig, Rishabh Singh, *OOPSLA'17*
- **Program Synthesis using Abstraction Refinement**  
Xinyu Wang, Isil Dillig, Rishabh Singh, *POPL'18*
- **Searching Entangled Program Spaces**  
James Koppel, Zheng Guo, Edsko de Vries, Armando Solar-Lezama, Nadia Polikarpova. *ICFP'22*

# Abstract FTA

- **Challenge:** FTA still has too many states
- Idea:
  - instead of one state = one value
  - we can do one state = set of values (= abstract value)

# Abstract FTA

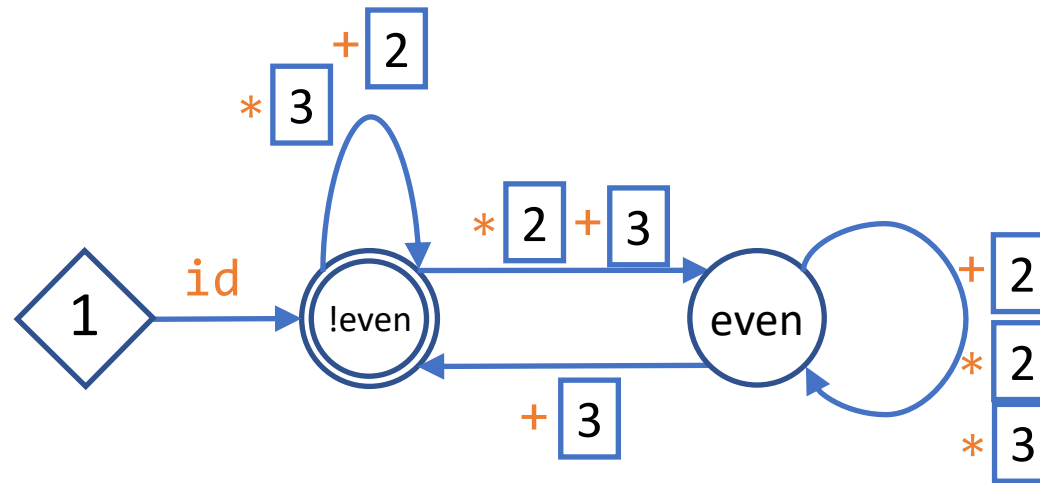
[Wang, Dillig, Singh POPL'18]

$N ::= \text{id}(V) \mid N + T \mid N * T \quad \bigcirc$

$T ::= 2 \mid 3 \quad \square$

$V ::= x \quad \diamond$

$1 \rightarrow 9$



## • What now?

- idea 1: enumerate from reduced space
- idea 2: refine abstraction!

# Abstract FTA

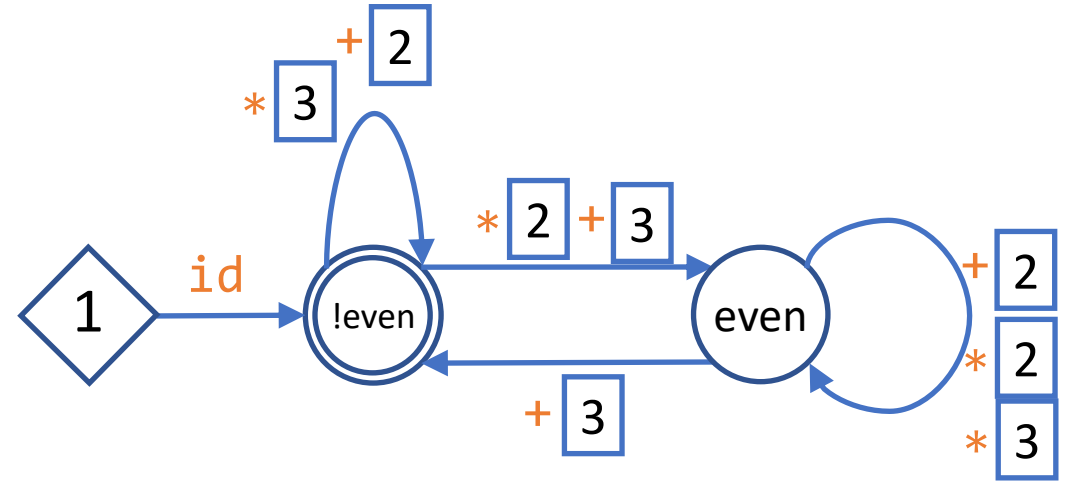
$N ::= \text{id}(V) \mid N + T \mid N * T \quad \bigcirc$

$T ::= 2 \mid 3 \quad \square$

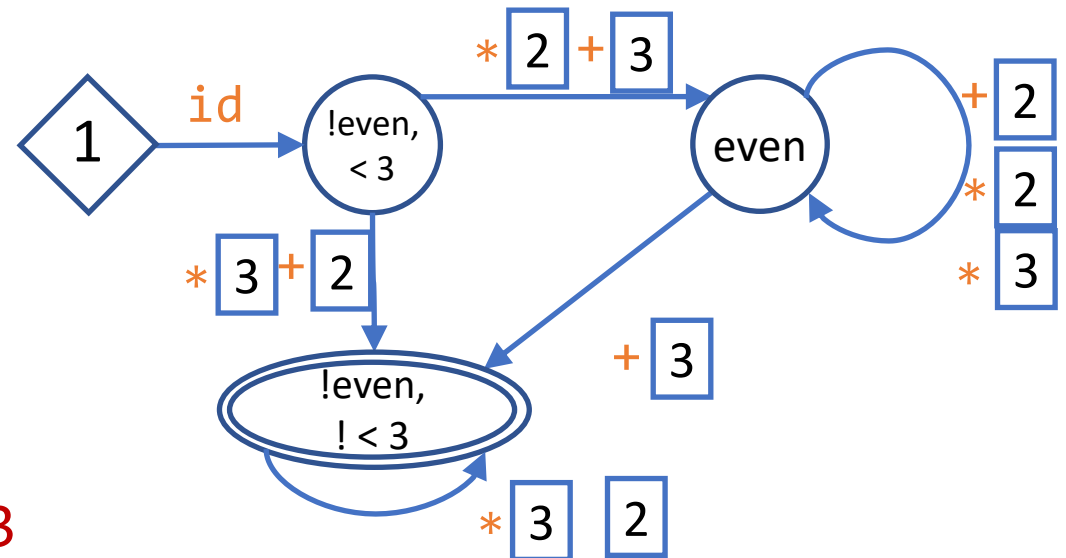
$V ::= x \quad \diamond$

$1 \rightarrow 9$

Predicates: {even, < 3, ...}

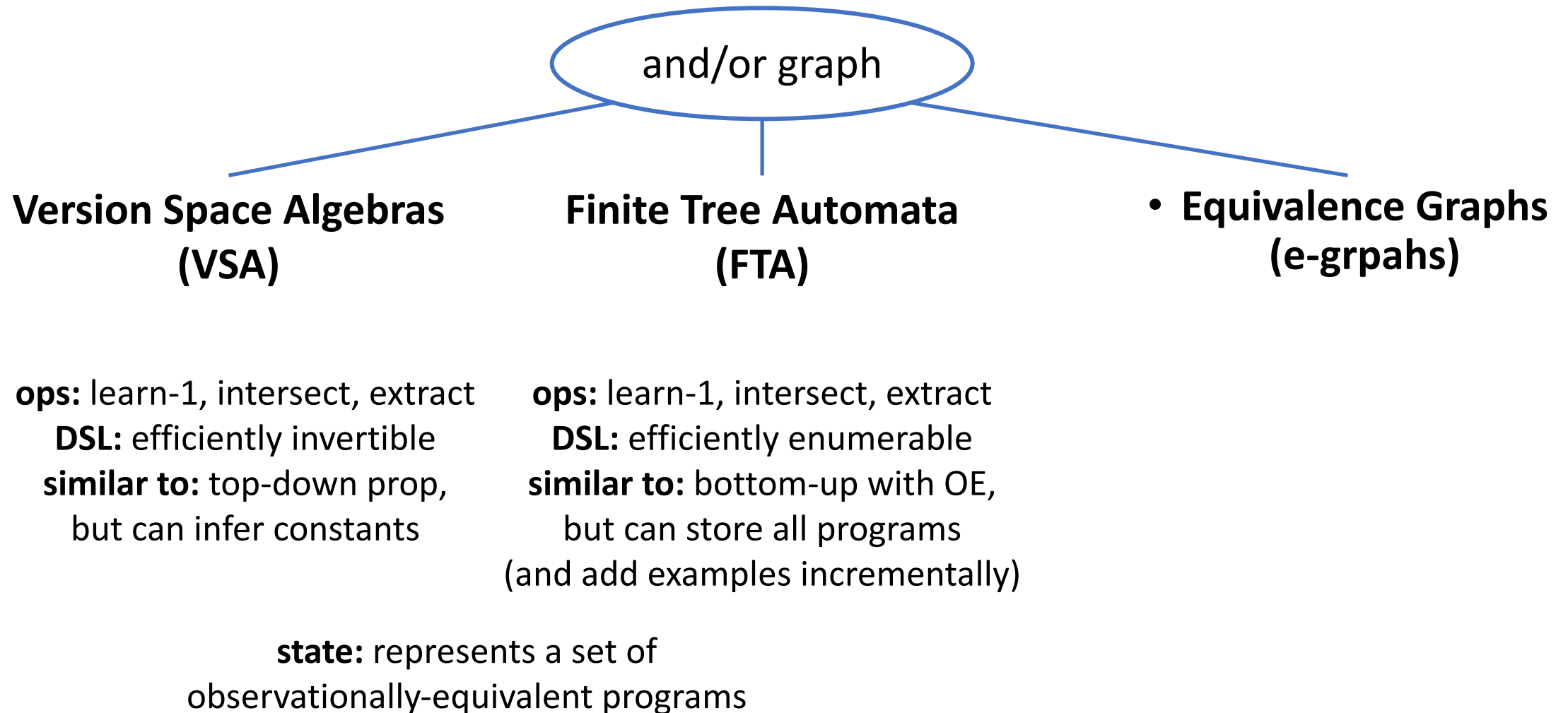


solution:  $\text{id}(x)$

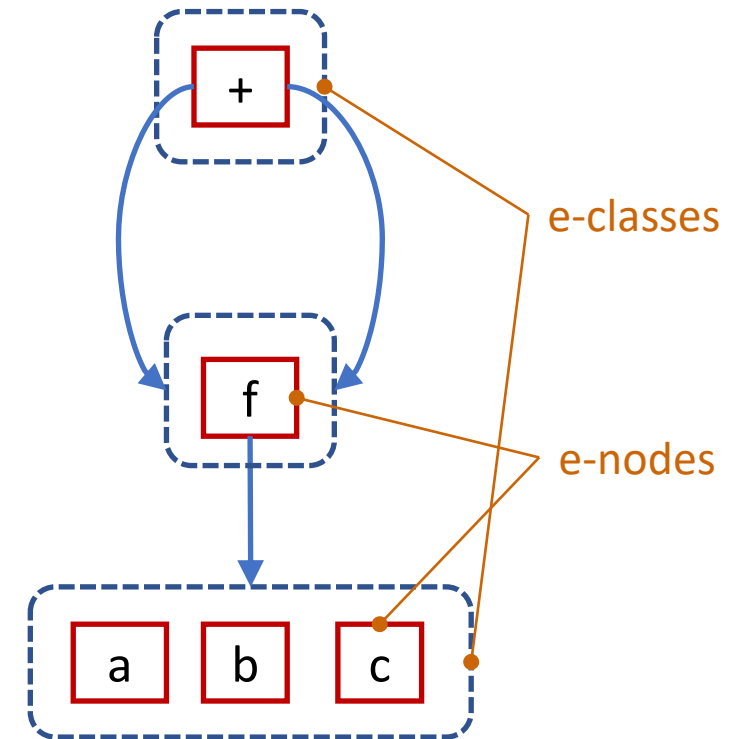
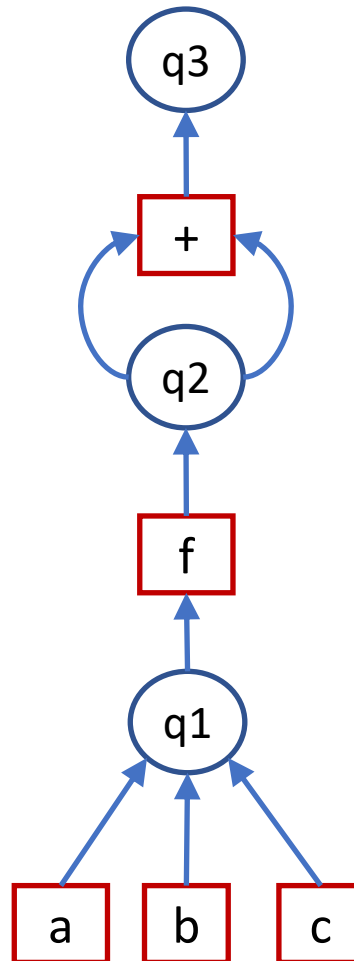
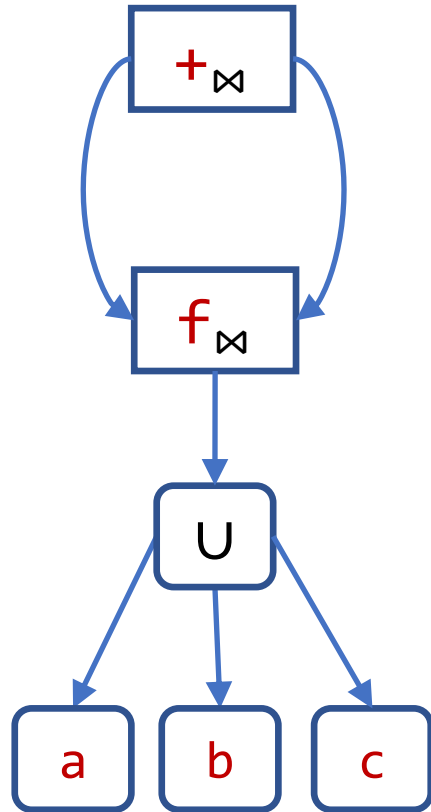


solution:  $\text{id}(x) * 3$

# Representation-based search



# VSA vs FTA vs E-Graphs



# Program search with e-grpahs

- **Equality saturation: a new approach to optimization**  
Ross Tate, Michael Stepp, Zachary Tatlock, Sorin Lerner, *POPL'09*
- **egg: Fast and Extensible Equality Saturation**  
Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, Pavel Panchekha, *POPL'21*
- **Semantic Code Search via Equational Reasoning**  
Varot Premtoon, James Koppel, Armando Solar-Lezama. *PLDI'20*



# Equality saturation

Program optimization via rewriting:

$$\begin{aligned} & (a * 2) / 2 \\ \Rightarrow & a * (2 / 2) \\ \Rightarrow & a * 1 \\ \Rightarrow & a \end{aligned}$$

useful rules:

$$\begin{aligned} (x * y) / z &= x * (y / z) \\ x / x &= 1 \\ x * 1 &= x \end{aligned}$$

not so useful:

$$\begin{aligned} x * 2 &= x \ll 1 \\ x * y &= y * x \end{aligned}$$

**Challenge:** which ones to apply and in what order?

**Idea:** all of them all the time!

# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

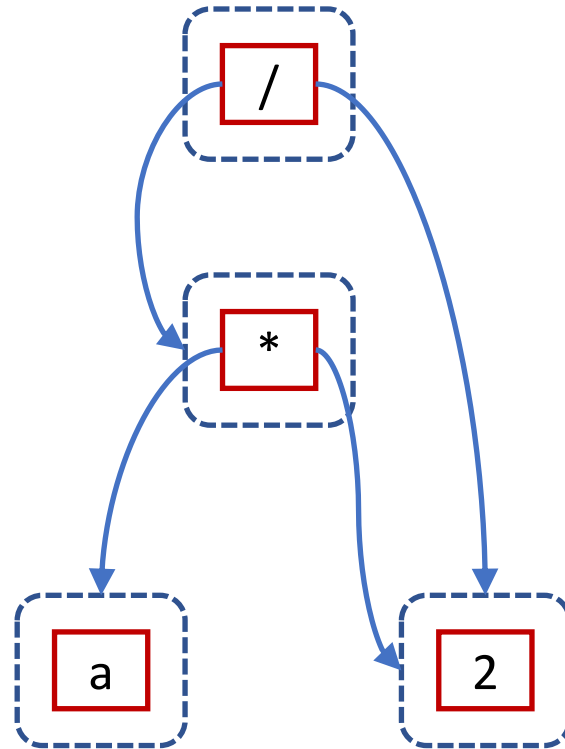
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

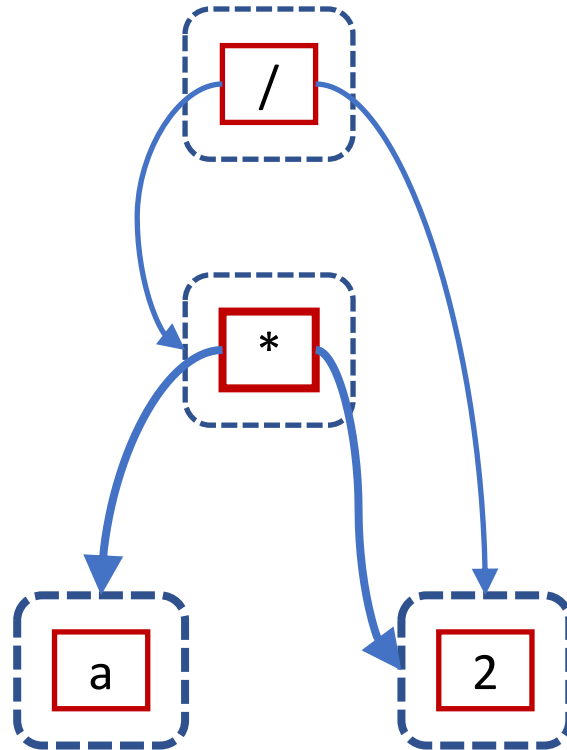
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

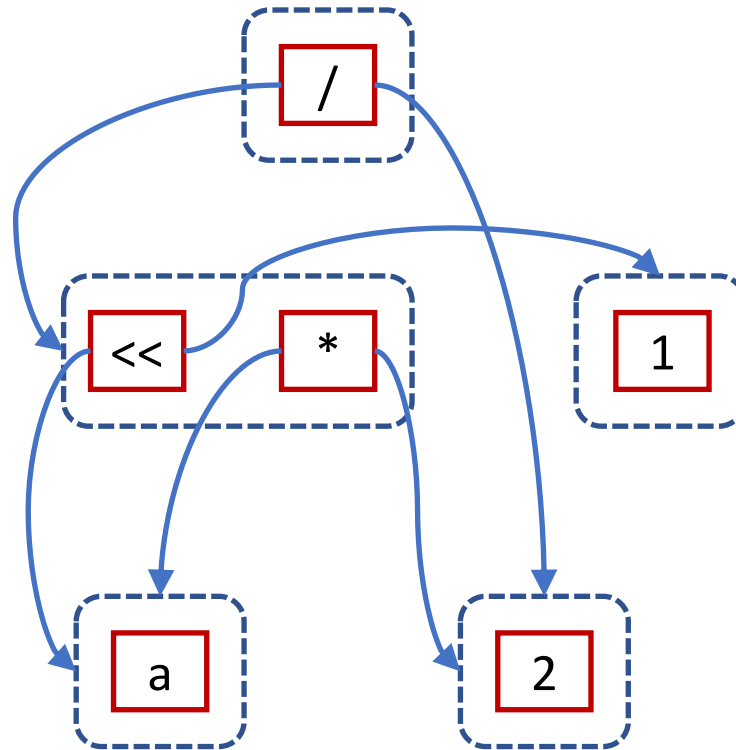
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

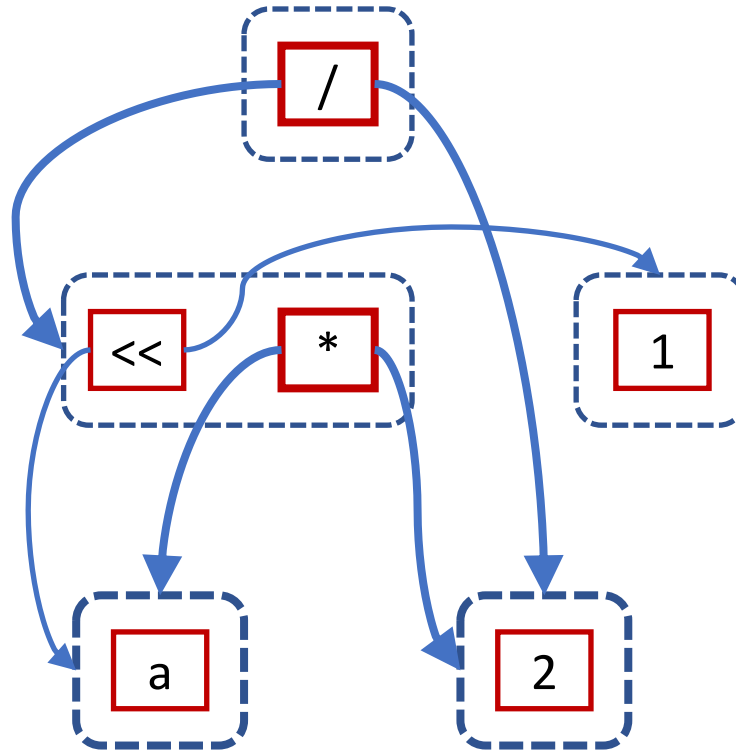
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

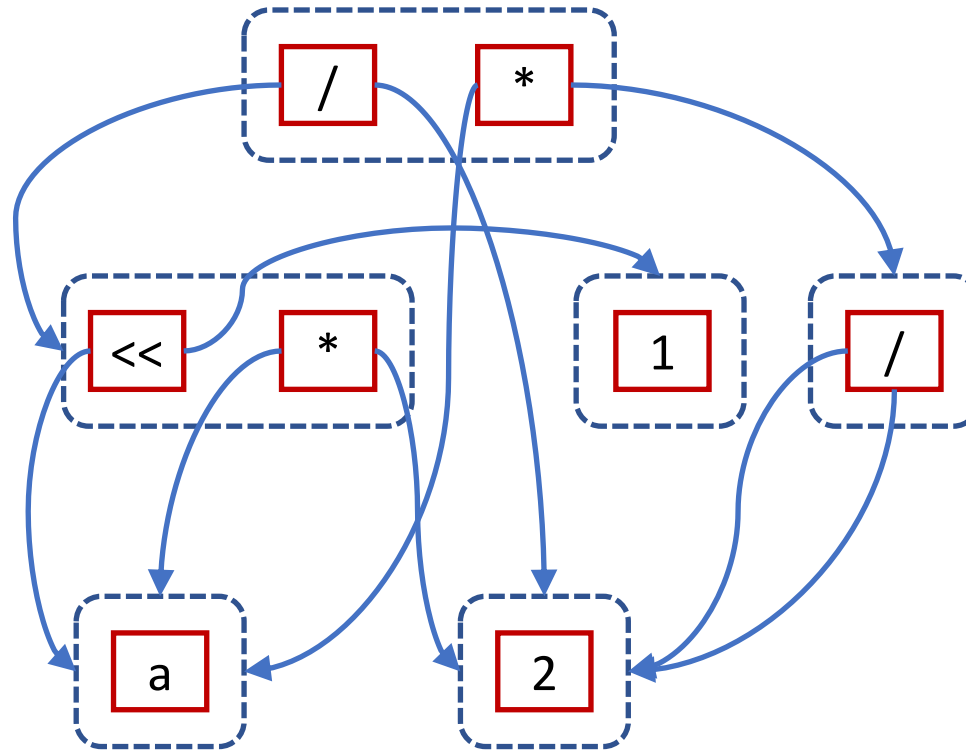
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

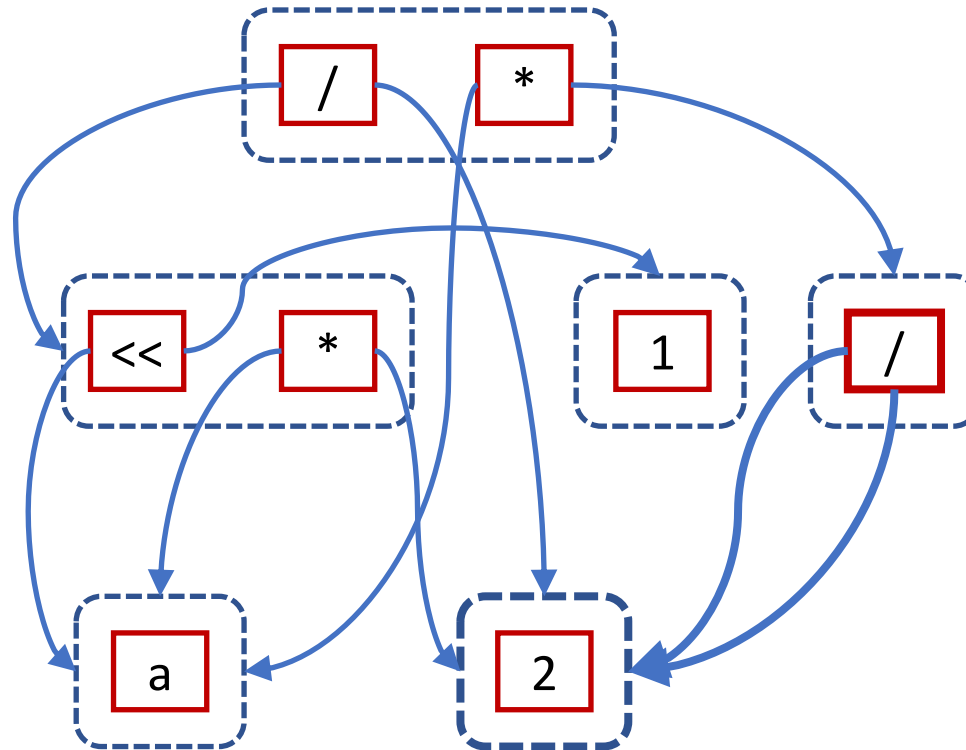
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

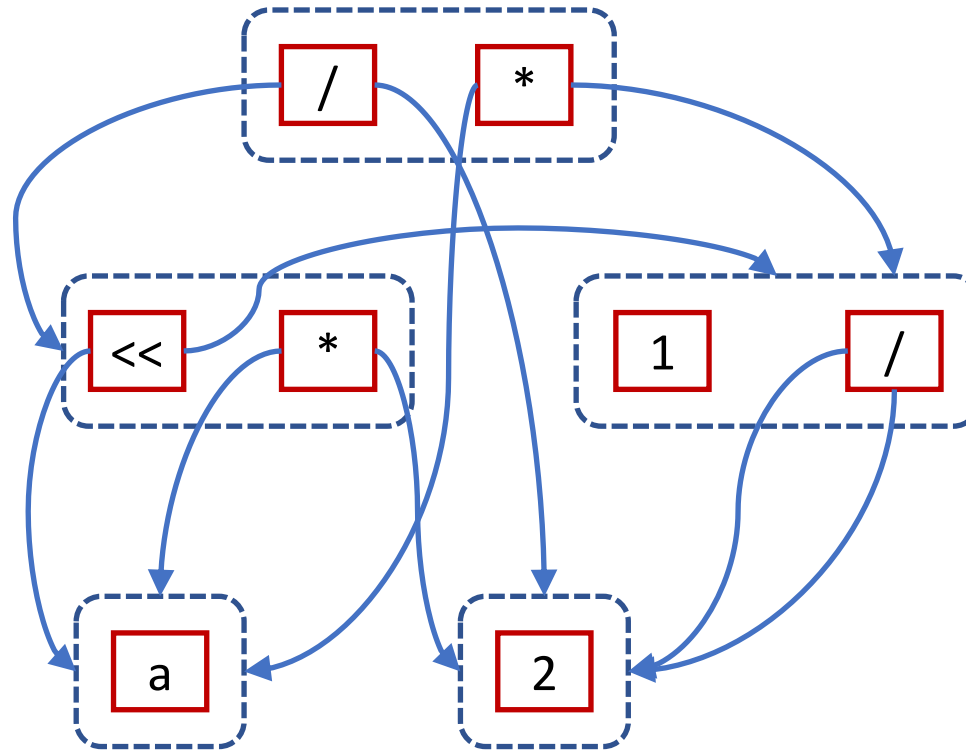
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x << 1$$

$$x * y = y * x$$





# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

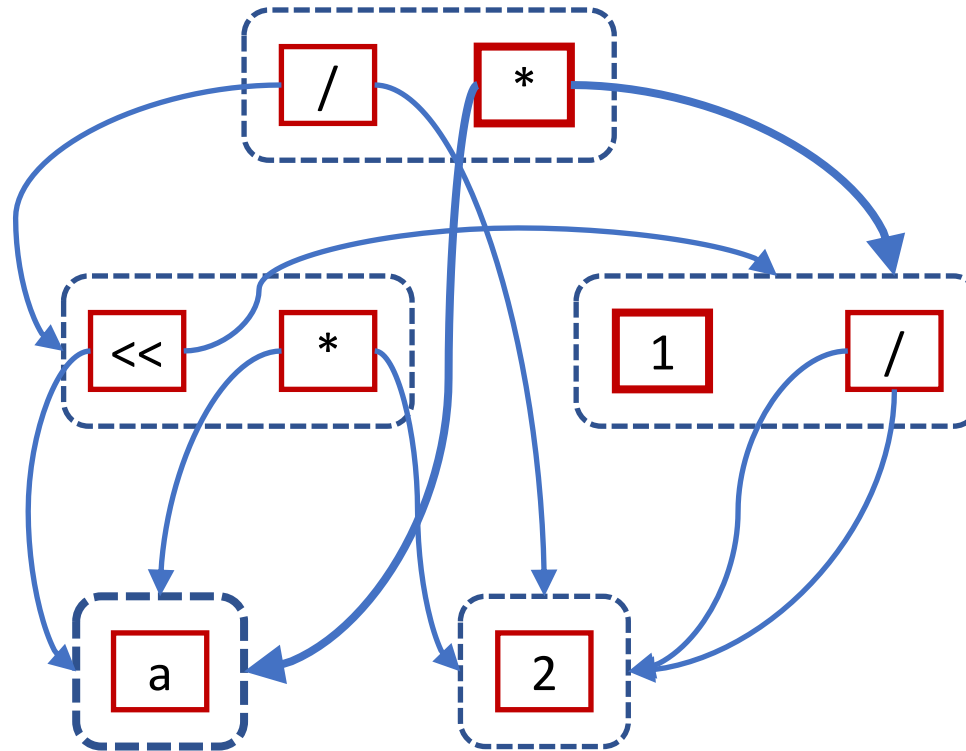
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Equality saturation

Initial term:  $(a * 2) / 2$

Rewrite rules:

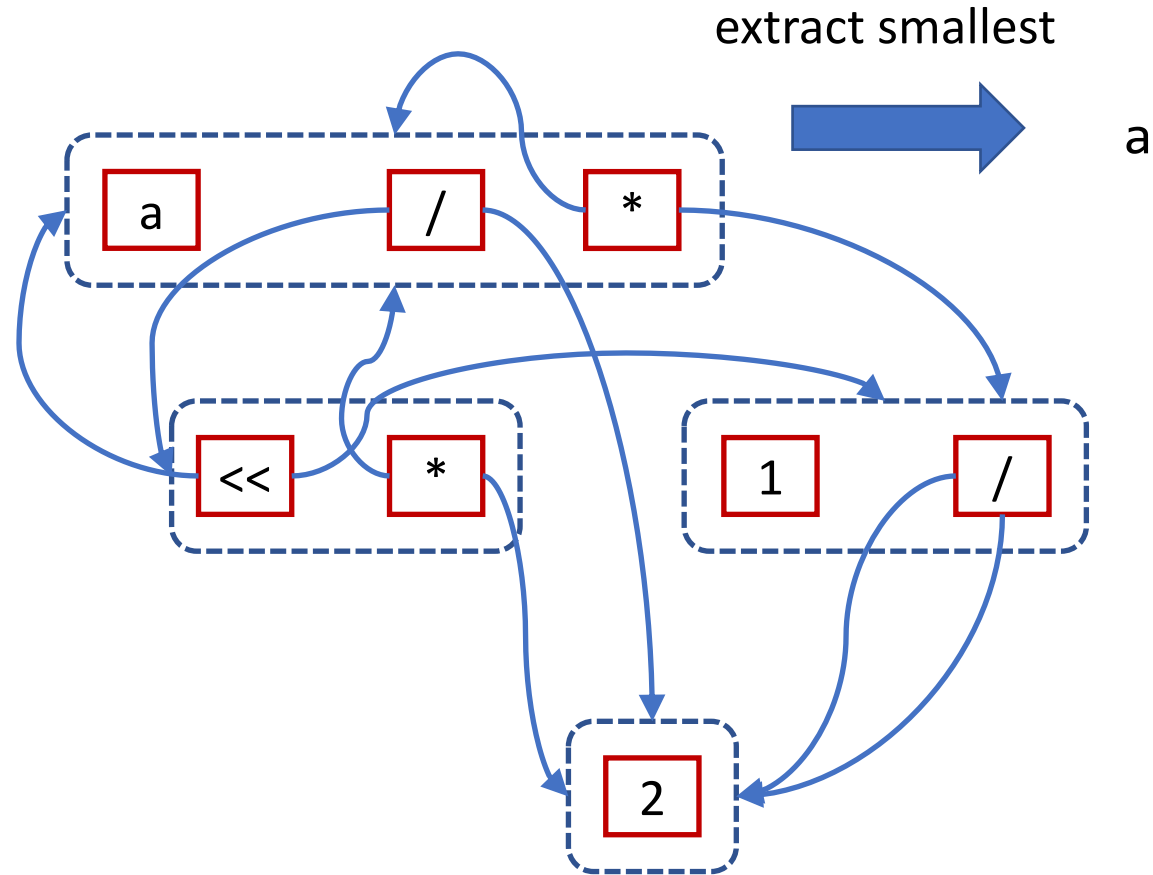
$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

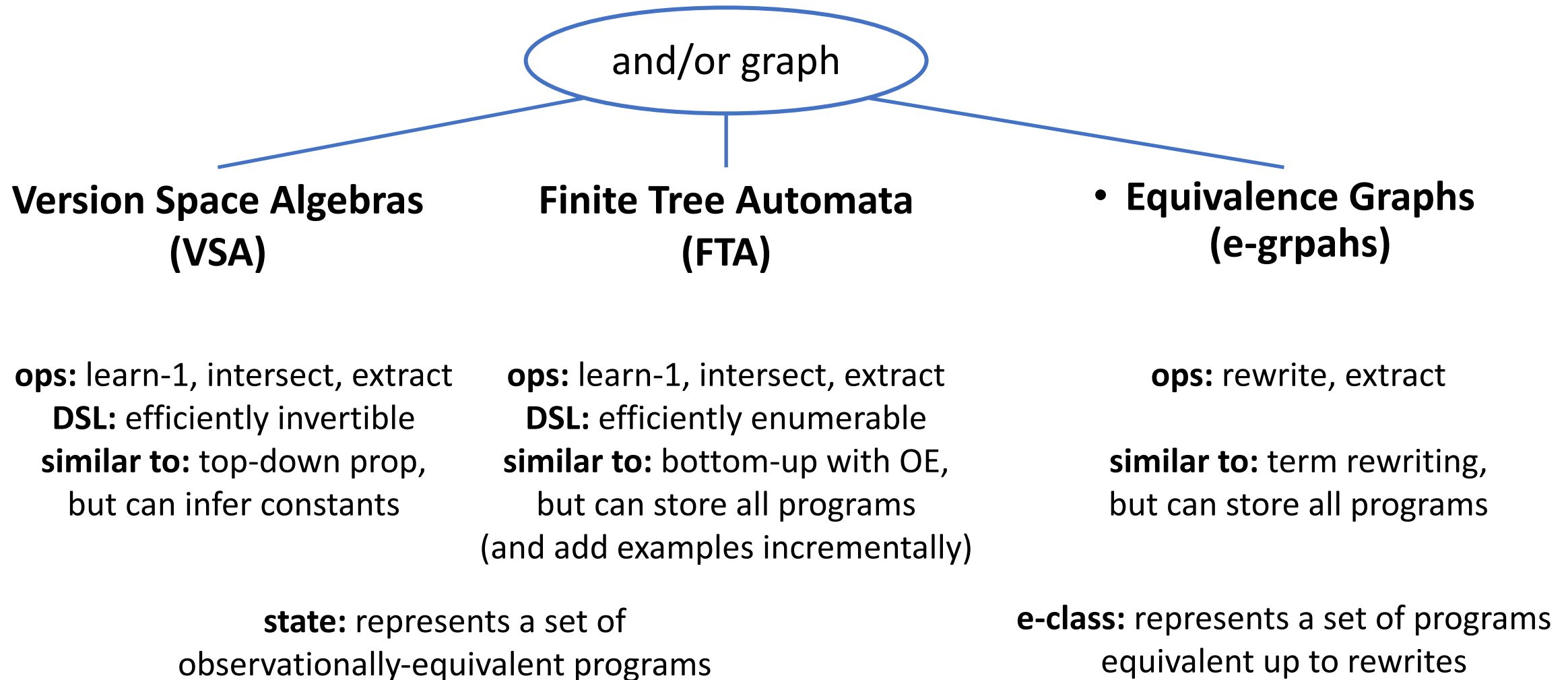
$$x * 1 = x$$

$$x * 2 = x \ll 1$$

$$x * y = y * x$$



# Representation-based search



# BlinkFill

- What does BlinkFill use as behavioral constraints? Structural constraints? Search strategy?
  - input-output examples (+ input examples); custom string DSL; VSA
- What is the main technical insight of BlinkFill wrt FlashFill?
  - BlinkFill uses the available inputs (with no outputs) to infer structure (segmentation) common to all inputs
  - it uses this structure to shrink the DAG and to rank substring expressions

# Example

- $\text{learn}_F \langle \text{"Mumbai, India"} \rightarrow \text{"India"} \rangle$

"Los Angeles, United States"

