

# **String Solving Strategies**

**Straight-Line Fragment**

# Simplified Syntax

Rewrite constraints to use simple terms. Conjunctions of

- concatenation:  $X = Y Z$
- regular tests:  $X \in r$

Write a sequences  $\phi, \psi$  instead of conjunctions  $\phi \wedge \psi$

E.g.

$$b X = X c Y \wedge X \in ab^*c$$

becomes

$$Z = B X, \quad B \in b, \quad Z = X Z', \quad Z' = C Y, \quad C \in c, \quad X \in ab^*c$$

# Boolean Combinations

We can generalise to arbitrary boolean combinations of terms

$$\phi, \phi' := X = YZ \mid X \in r \mid \phi \wedge \phi' \mid \phi \vee \phi' \mid \neg \phi$$

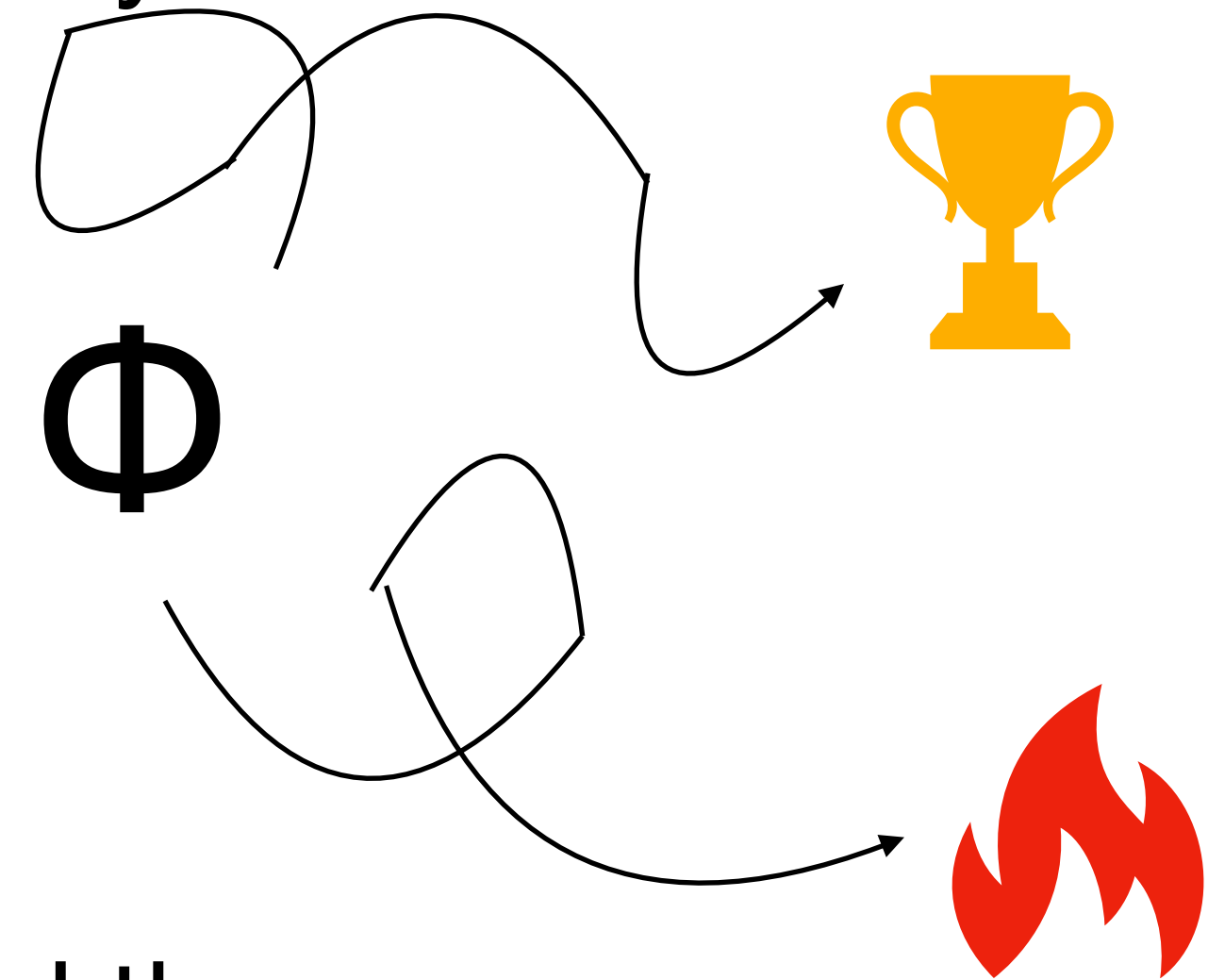
Two quick ways to see this:

- Disjoint Normal Form
  - Rewrite  $\phi$  into a disjunction of conjunctions
  - Consider each conjunction in turn
- DPLL[T] / SMT solving
  - Explore Boolean structure via DPLL algorithm
  - Test feasibility of conjunctive queries

# Proof System

Our approach to solving constraints is to use a proof system

1. start with a constraint
2. apply proof rules
3. find a contradiction or a solution



This gives a search problem: apply rules until you find the answer




Rules applied non-deterministically

We will discuss deterministic strategies at the end

# Proof Rules Syntax

A simple proof rule is


$$\frac{\dots, \Phi, \Psi}{\dots, \Phi \wedge \Psi}$$


(read upwards)

If we have  $\Phi \wedge \Psi$  we can derive both  $\Phi$  and  $\Psi$

A rule that creates two branches:

- find a solution on one of them, or
- find a contradiction on all of them

$$\frac{\dots, \Phi \qquad \dots, \Psi}{\dots, \Phi \vee \Psi}$$


# Constraint Propagation

For each

$$X = Y Z$$

Forwards propagation

- Push constraints on Y and Z into a constraint on X

Backwards propagation

- Pull constraint on X back into constraints on Y and Z

# String Functions

Concatenation is a string function

$$X = Y Z \longrightarrow X = \text{concat}(Y, Z)$$

More generally, any function could appear

$$X = f(X_1, \dots, X_n)$$

where  $f$  is some function with input  $X_1, \dots, X_n$ .

# Forwards Propagation



# Forwards Propagation (concat)

Suppose

$$\dots, X = YZ, Y \in r_1, Z \in r_2, \dots$$

We can infer that also

$$X \in r_1 r_2$$

Or in fact

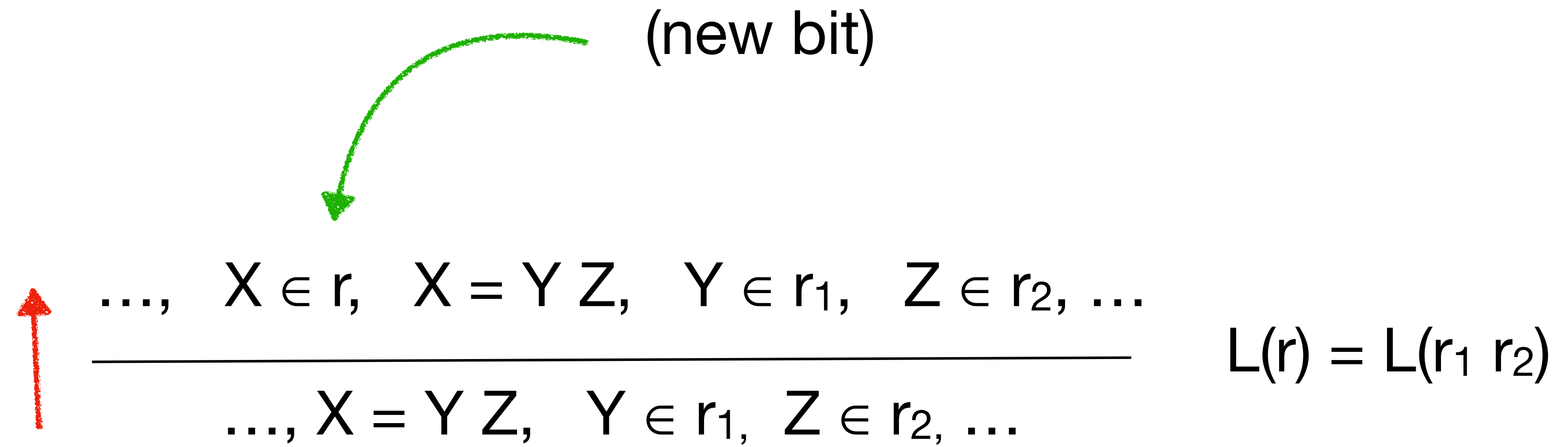
$$X \in r$$

$L(r)$  -- set of words matching  $r$

for any  $r$  with

$$L(r) = \text{concat}(L(r_1), L(r_2))$$

# As a Proof Rule



(new bit)

$$\frac{\dots, X \in r, X = YZ, Y \in r_1, Z \in r_2, \dots}{\dots, X = YZ, Y \in r_1, Z \in r_2, \dots} \quad L(r) = L(r_1 r_2)$$

# Forwards Propagation in General

Suppose

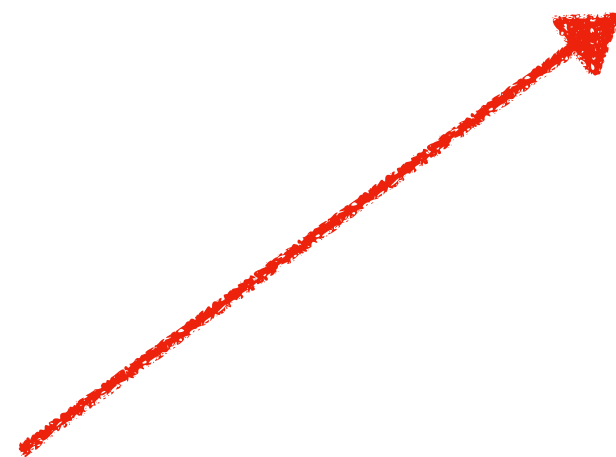
$$\dots, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots$$

and

$$f(L(r_1), \dots, L(r_n)) = L(r) \text{ for some regular expression } r$$


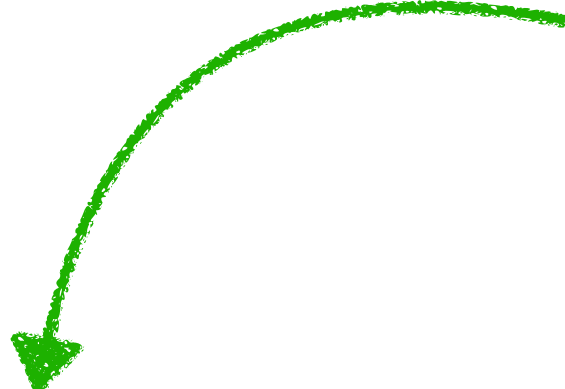
Add constraint  $X \in r$

$$\dots, \underline{X \in r}, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots$$

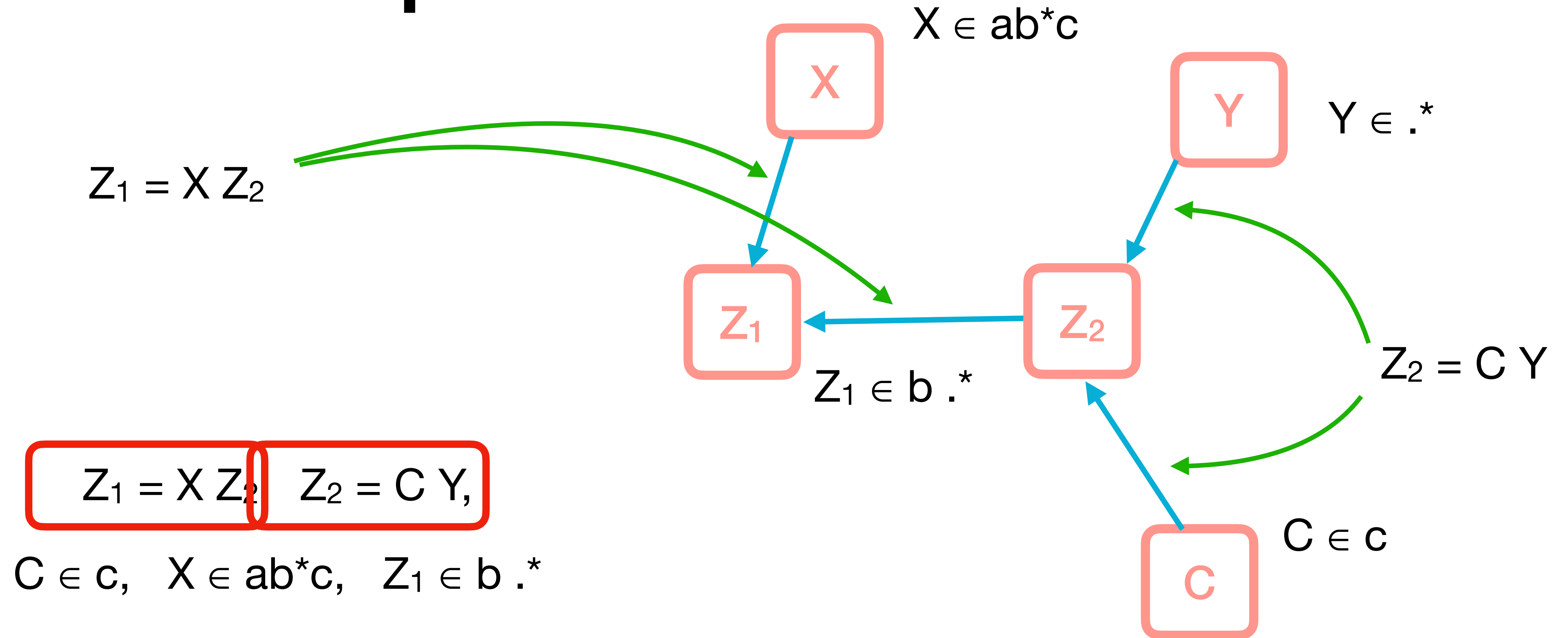


# As a Proof Rule

(new bit)


$$\frac{\dots, X \in r, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots}{\dots, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots} \quad L(r) = f(L(r_1), \dots, L(r_n))$$

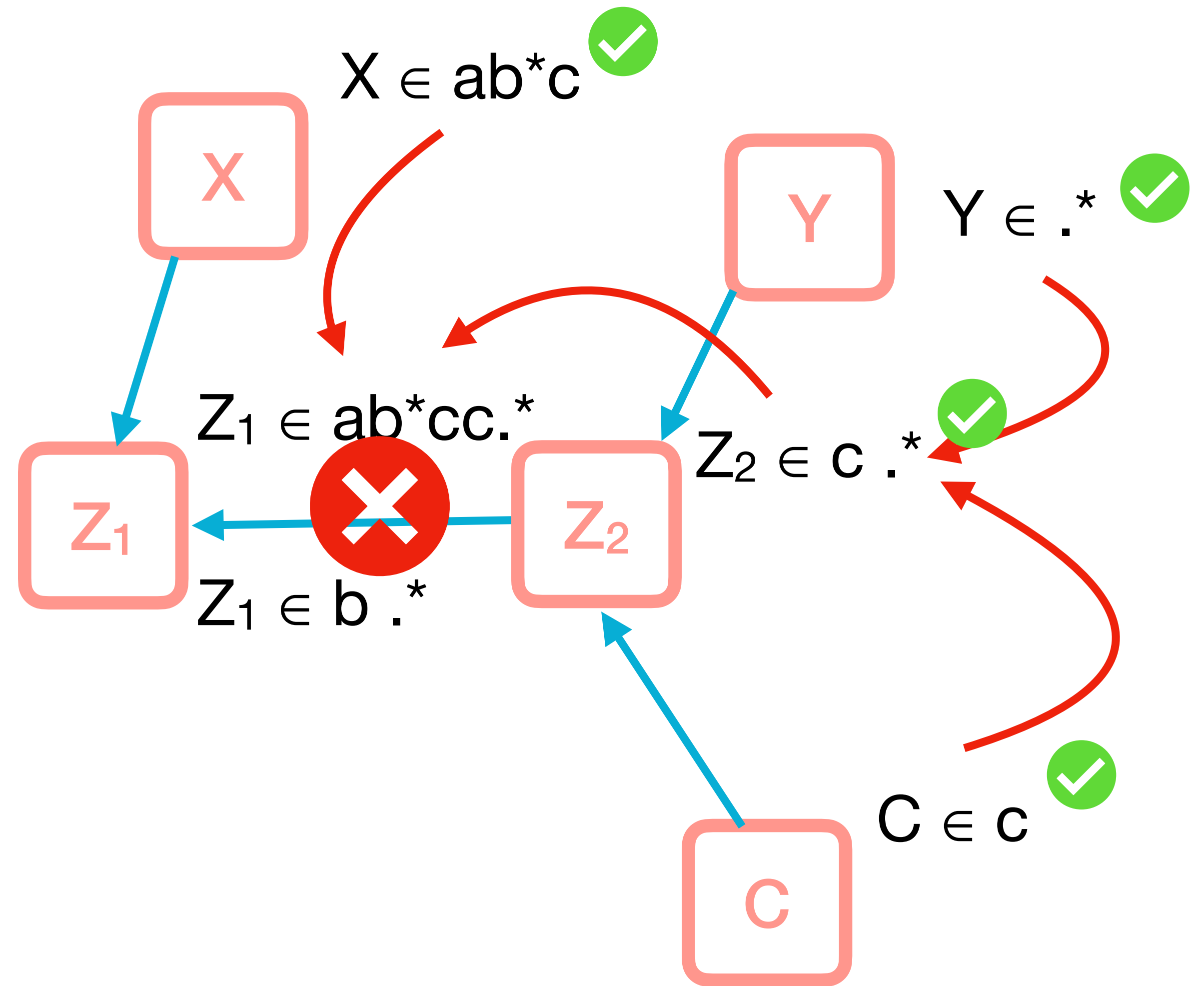
# Visual Example



# Solve via Forwards Propagation

$$Z_1 = X Z_2, \quad Z_1 = C Y,$$

$$C \in c, \quad X \in ab^*c, \quad Z_1 \in b.^*$$



Now check satisfiability for each variable

# Constraint Elimination?

Suppose

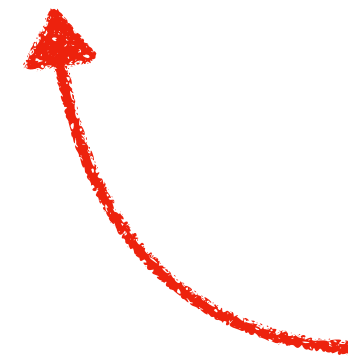
$$X = Y Y, \quad Y \in a^*b^*$$

Notice

$$L(a^*b^*a^*b^*) = \text{concat}(L(a^*b^*), L(a^*b^*))$$

Propagate

$$X \in a^*b^*a^*b^*, \quad X = Y Y, \quad Y \in a^*b^*$$



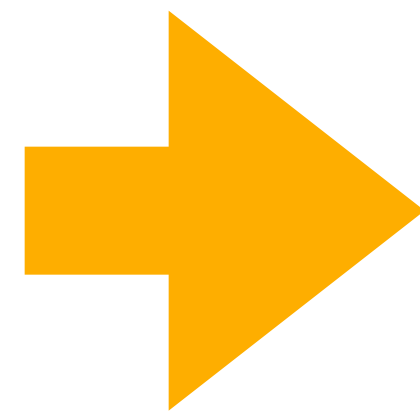
Do we still need  $X = Y Y$ ?

# An Overapproximation

$$X \in a^*b^*a^*b^*, \quad X = Y Y, \quad Y \in a^*b^*$$

Notice

- $X = abaabb$  satisfies  $X \in a^*b^*a^*b^*$
- But not  $X = Y Y$  for any value of  $Y$



We need to keep  $X = Y Y$  too



# Backwards Propagation

# Backwards Propagation (concat)

Suppose

$$\dots, X \in r_1 r_2, X = Y Z, \dots$$

We can infer that also

$$Y \in r_1, Z \in r_2$$

Or in fact

$$Y \in r_3, Z \in r_4$$

for any  $r_3, r_4$  with

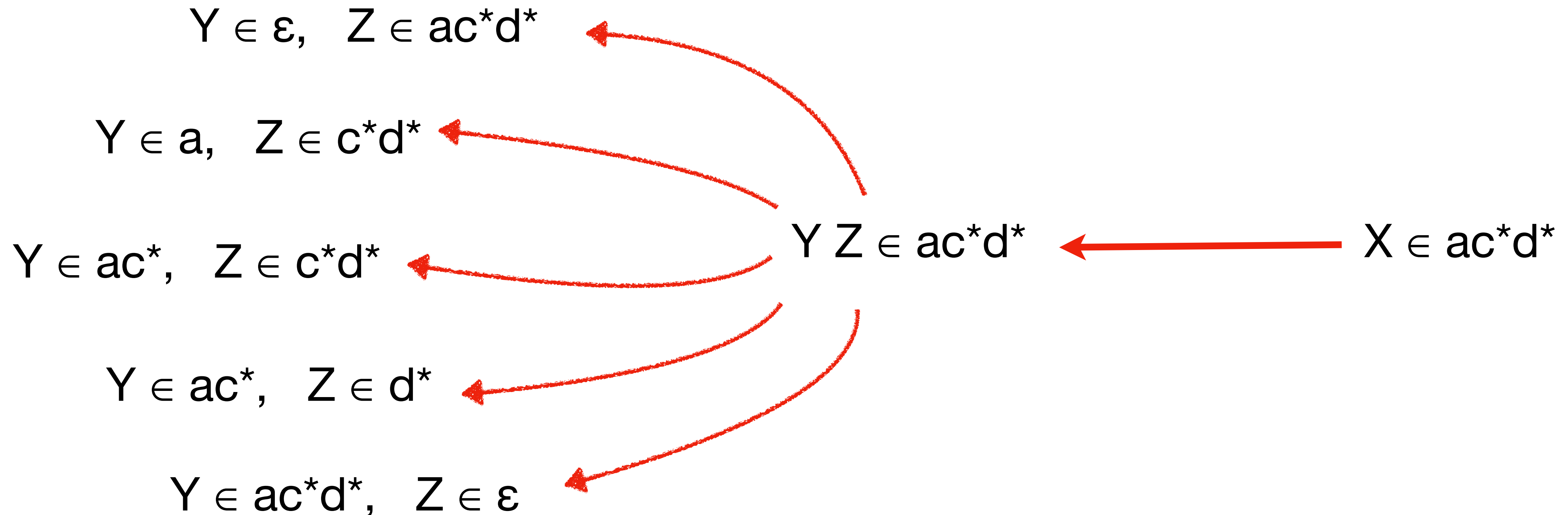
$$\text{concat}(L(r_3), L(r_4)) = L(r_1 r_2)$$

# Backwards Propagation

Take constraint

$$X \in ac^*d^* \wedge X = Y Z$$

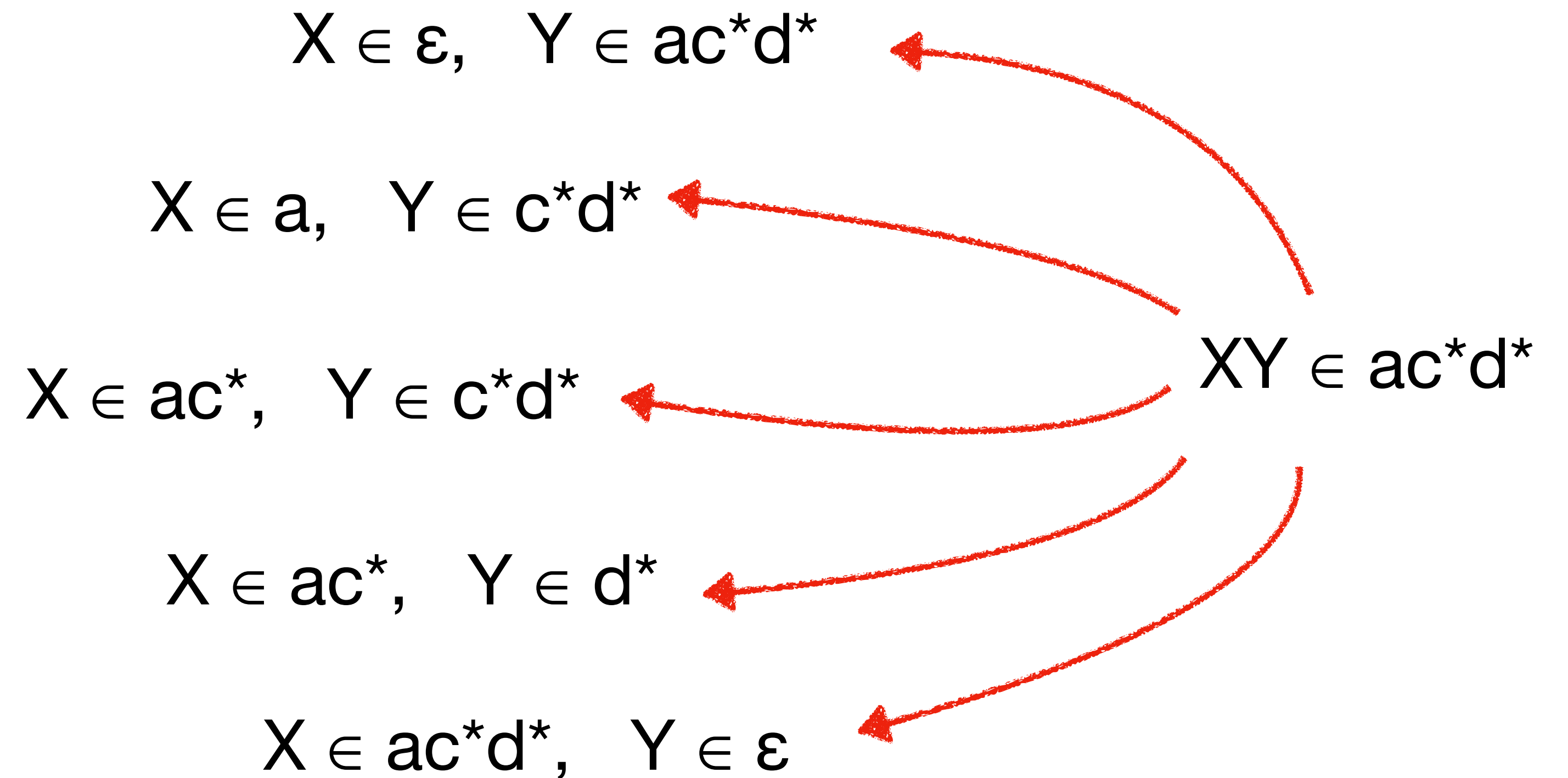
Pull the constraint on  $X$  back through  $X = Y Z$



# Other String Functions

Backwards propagation creates branching for concat

What other string functions can be supported?



# Conditions on String Functions

Consider

$$Z = \text{concat}(X, Y)$$

Given a regular language  $L$ , what is  $\text{concat}^{-1}(L)$ ?

$$X \in \varepsilon, Y \in ac^*d^*$$

$$X \in a, Y \in c^*d^*$$

$$X \in ac^*, Y \in c^*d^*$$

$$X \in ac^*, Y \in d^*$$

$$X \in ac^*d^*, Y \in \varepsilon$$

A finite union of regular alternatives

$$\bigcup_i L_i^X \times L_i^Y$$

(A recognisable relation)

$$XY \in ac^*d^*$$


# Conditions on String Functions

Consider

$$X = f(X_1, \dots, X_n)$$

Given a regular language  $L$ , require  $f^{-1}(L)$  is recognisable

$$X_1 \in r^1_1, \dots, X_n \in r^1_n$$

$$X_1 \in r^2_1, \dots, X_n \in r^2_n$$

$$X_1 \in r^3_1, \dots, X_n \in r^3_n$$

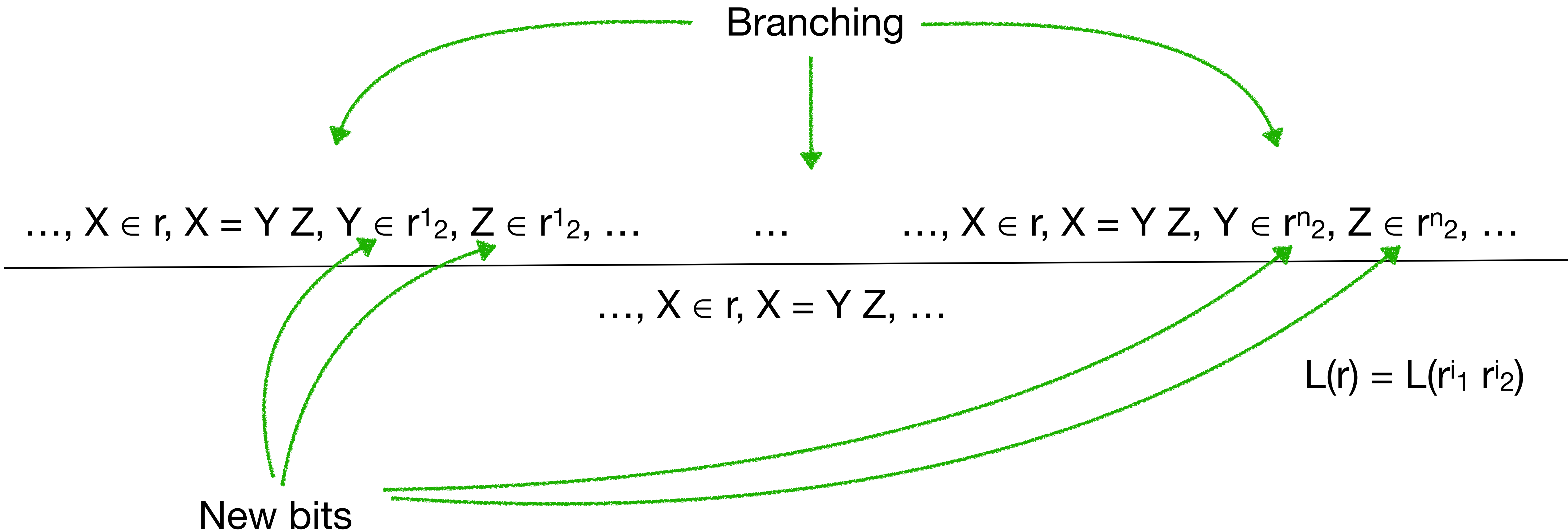
$$X_1 \in r^4_1, \dots, X_n \in r^4_n$$

$$X_1 \in r^5_1, \dots, X_n \in r^5_n$$

$$f^{-1}(L) = \bigcup_i L^{i_1} \times \dots \times L^{i_n}$$

$$f(X_1, \dots, X_n) \in r$$

# Proof Rule



# Example

$\phi := W \in b, Y \in b^*c^*, Z \in a^*bc^*,$

$$X = W, Z = XY$$

$\phi$   
↓

$\phi,$   
 $X \in a^*b,$   
 $Y \in c^*$

↓

$\phi,$   
 $X \in a^*b,$   
 $Y \in c^*,$   
 $W \in a^*b$

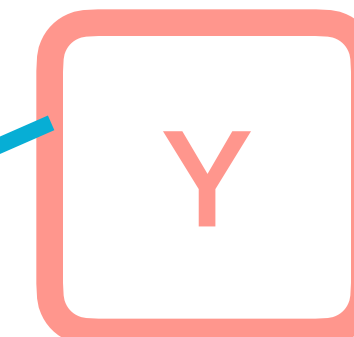


?

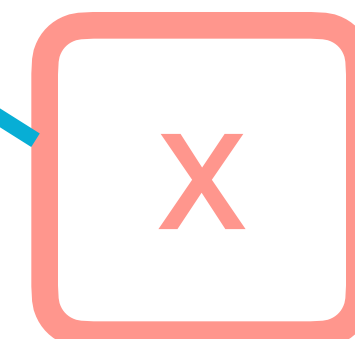
We have finished propagating

There are no contradictions -- are we done?

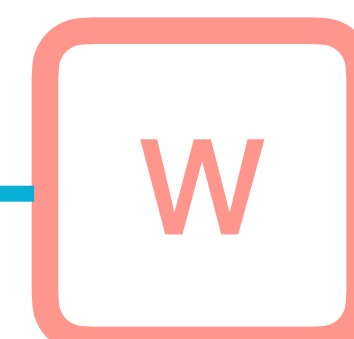
$Z \in a^*bc^*$



$Y \in b^*c^*$   
 $Y \in c^*$



$X \in a^*b$



$W \in b$   
 $W \in a^*b$





# With Proof Rules

$W \in b, Y \in b^*c^*, Z \in a^*bc^*, X = W, Z = X Y, X \in a^*b, Y \in c, X \in a^*b, Y \in c^*, W \in a^*b$

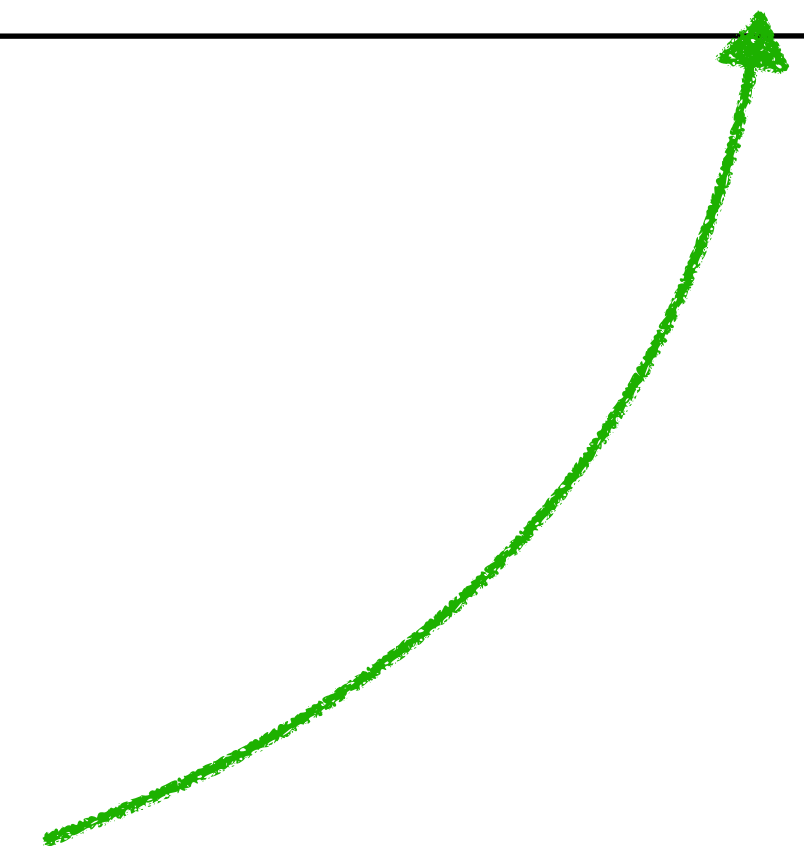
---

$W \in b, Y \in b^*c^*, Z \in a^*bc^*, X = W, Z = X Y, X \in a^*b, Y \in c$  ...

---

$W \in b, Y \in b^*c^*, Z \in a^*bc^*, X = W, Z = X Y$

Other branches



# **Extracting Solutions**

# Extracting Solutions

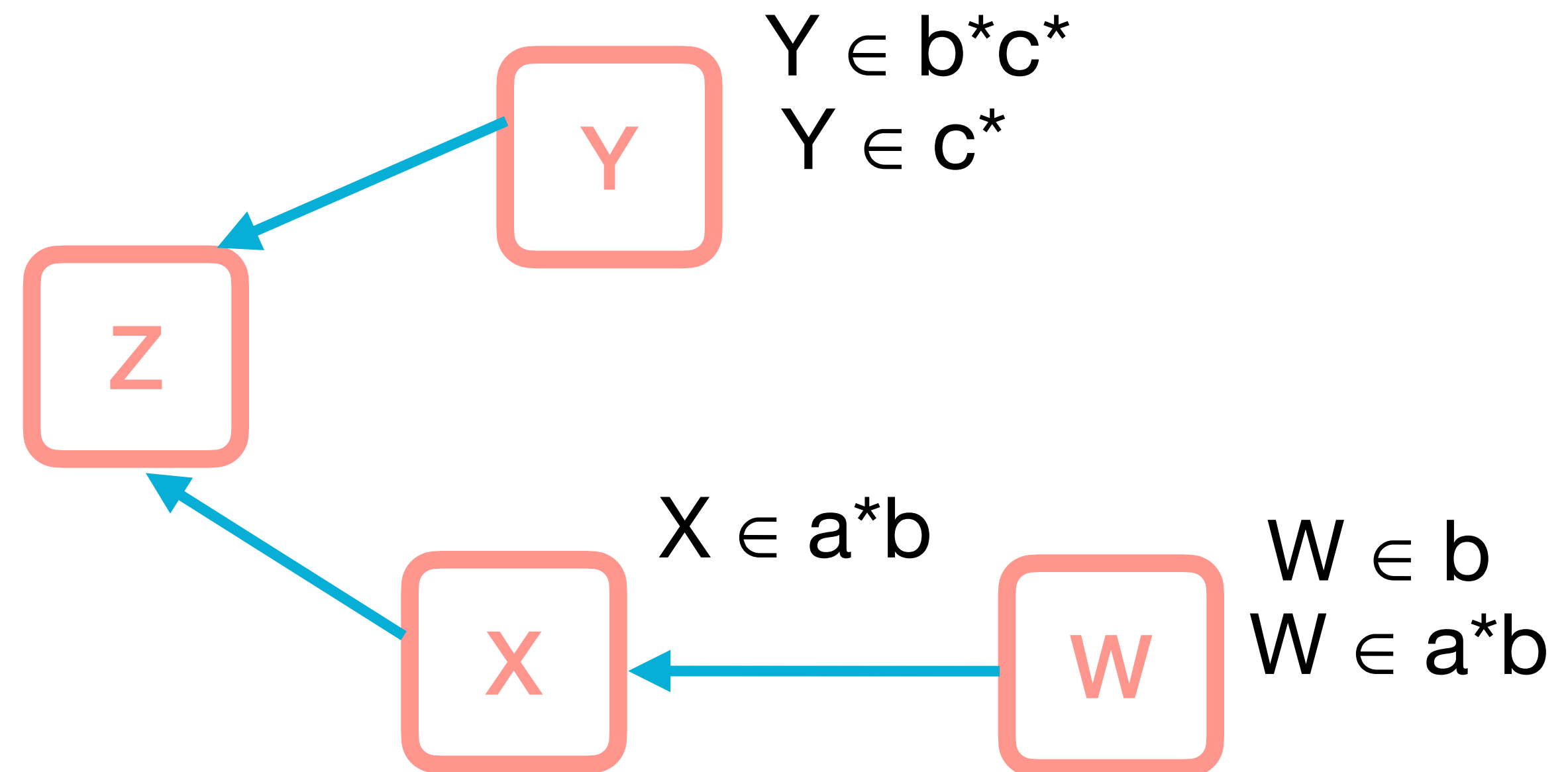
Not all solutions to the regular constraints are correct

$$W \in b, Y \in b^*c^*, Z \in a^*bc^*,$$
$$X = W, Z = XY$$

E.g.

$$X = ab$$
$$Y = c$$
$$Z = b$$
$$W = b$$

$$Z \in a^*bc^*$$



Satisfies each node but not constraint

$$(Z \neq XY)$$

# A Cut Rule

As part of the proof search we can introduce two branches

$$X \in r \qquad X \notin r$$

for some regular expression  $r$

Can be used to guess a solution

- Pick some  $w$  consistent with constraints on  $X$
- Add  $X \in w$  on one branch,  $X \notin w$  on the other

Then propagate to other variables

# As a Proof Rule

...,  $X \in r$

...,  $X \notin r$

---

...

# Closing Branches

Constraints introduced via cut may be consistent or not

Close a branch if an inconsistency is found

$$\dots, X \in r_1, \dots, X \in r_n, \dots \text{ with } L(r_1) \cap \dots \cap L(r_n) = \emptyset$$

Solution found if we have a "pure" solution

$$X_1 \in W_1, \dots, X_n \in W_n$$

But how do we eliminate constraints?

# As Proof Rules

Contradiction:

$$\frac{\dots, X \in r_1, \dots, X \in r_n}{L(r_1) \cap \dots \cap L(r_n) = \emptyset}$$

Solution:

$$\frac{}{X_1 \in W_1, \dots, X_n \in W_n}$$

# Eliminating Constraints

To end up with a pure solution

$$X_1 \in W_1, \dots, X_n \in W_n$$

We need to eliminate all other constraints.

Two methods

- Eliminate multiple regular constraints:  $X \in r_1, \dots, X \in r_n$
- Eliminate assignments:  $X = f(X_1, \dots, X_n)$



# Eliminate Multiple Constraints

$$X \in r_1, \dots, X \in r_n$$

1. Replace with its intersection

$$X \in r_1 \cap \dots \cap r_n$$

2. Remove a subsumed constraint

$$X \in r_2, \dots, X \in r_n \quad \text{when} \quad L(r_2) \cap \dots \cap L(r_n) \subseteq L(r_1)$$

# Proof Rules

Intersection:

$$\frac{\dots, X \in r_1 \cap \dots \cap r_n}{\dots, X \in r_1, \dots, X \in r_n}$$

Subsumed:

$$\frac{\dots, X \in r_2, \dots, X \in r_n}{\dots, X \in r_1, \dots, X \in r_n} \quad L(r_2) \cap \dots \cap L(r_n) \subseteq L(r_1)$$

# Eliminate Function Assignments

$$X = f(X_1, \dots, X_n)$$

Eliminating version of forwards propagation

Suppose

$$\dots, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots$$

and

$$f(L(r_1), \dots, L(r_n)) = L(r) \text{ for some regular expression } r$$

If  $L(r)$  only has one solution, replace  $X = f(X_1, \dots, X_n)$  with  $X \in r$

$$\dots, X \in r, X_1 \in r_1, \dots, X_n \in r_n, \dots$$

# Proof Rules

(we dropped  $X = f(X_1, \dots, X_n)$ )

$\dots, X \in r, X_1 \in r_1, \dots, X_n \in r_n, \dots$

$L(r) = f(L(r_1), \dots, L(r_n))$

---

$\dots, X = f(X_1, \dots, X_n), X_1 \in r_1, \dots, X_n \in r_n, \dots$

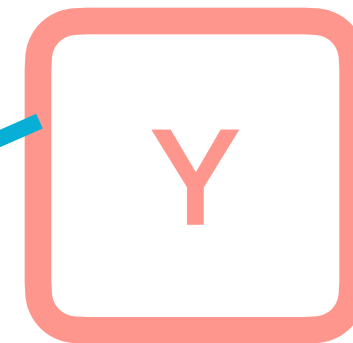
$|L(r)| = 1$

# Deriving Solutions

Eliminate via intersection  
 Pick value for node  
 $c \cap c^* \cap b^*c^* = c$

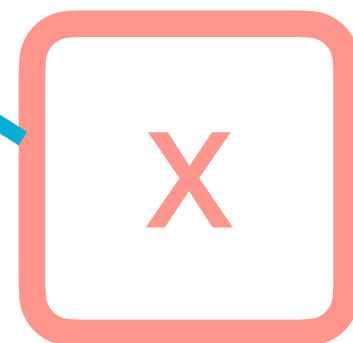
Eliminate ~~Z~~ with intersection  $Z \in b c$

$Z \in a^*bc^*$   
 $Z \in b c$

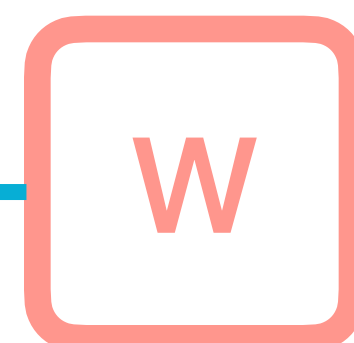


$Y \in b^*c^*$   
 $Y \in c^*$   
 $Y \in c$

Do the same here

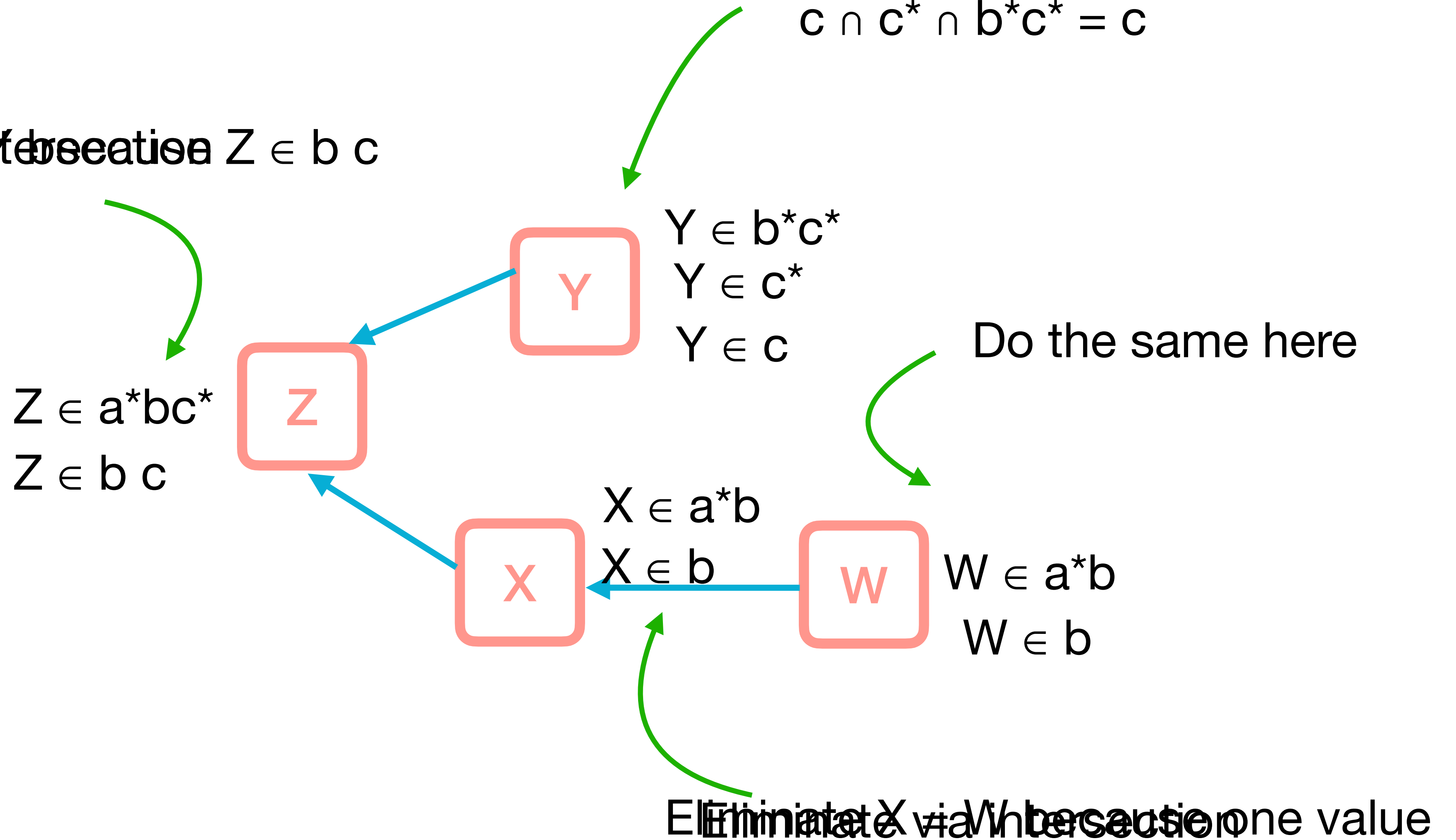


$X \in a^*b$   
 $X \in b$



$W \in a^*b$   
 $W \in b$

Eliminate X via intersection one value



# With Proof Rules

$$\frac{}{W \in b, X \in b, Z \in bc, Y \in c}$$
 $nZ$ 
$$\frac{W \in b, Z \in a^*bc^*, X \in b, Z \in bc, Y \in c, X \in a^*b}{}$$
 $F\text{-Elim } Z = X Y$ 
$$W \in b, Z \in a^*bc^*, X \in b, Z = X Y, Y \in c, X \in a^*b$$
 $nX$ 
$$\frac{W \in b, Z \in a^*bc^*, X \in b, Z = X Y, X \in a^*b, Y \in c, X \in a^*b}{}$$
 $F\text{-Elim } X=W$ 
$$W \in b, Z \in a^*bc^*, X = W, Z = X Y, X \in a^*b, Y \in c, X \in a^*b$$
 $nW$ 
$$W \in b, Z \in a^*bc^*, X = W, Z = X Y, X \in a^*b, Y \in c, X \in a^*b, W \in a^*b$$
 $nY$ 
$$W \in b, Y \in b^*c^*, Z \in a^*bc^*, X = W, Z = X Y, X \in a^*b, Y \in c, X \in a^*b, Y \in c^*, W \in a^*b$$

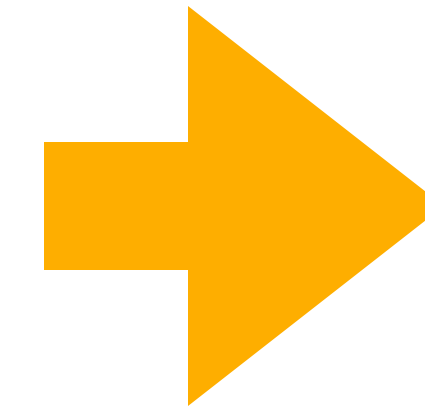
**Forwards or Backwards?**

# Comparing Propagation Direction

We've seen forwards and backwards propagation.

Forwards:

- Regular inputs lead to regular outputs
- But had to overapproximate e.g. concatenation
- No branching required



Backwards:

- Regular outputs lead exactly to recognisable inputs
- No approximation
- Branching required





# Functions Supporting Backwards

Some common string functions support backwards propagation

- Concatenation
- Reverse
- Replace(All)
  - $X = \text{ReplaceAll}(Y, e_1, e_2)$
  - $e_1$  -- a string or regular expression (with capture groups)
  - $e_2$  -- a string, or variable, or replacement pattern (using capture groups)
- Transductions (e.g. HTML escape)

**Completeness**

# Proof Search

Our proof rules are sound

- We find only correct answers

Not necessarily complete

- Propagation may not terminate
- Cut rule introduces infinite alternatives

Currently we've shown proof rules but no strategy

That is, apply the rules non-deterministically to solve a constraint.



Can we provide an algorithm with guarantees?

# Complete Fragments

Implementations needs to apply rules algorithmically

1. Apply forwards propagation until fixed point
  - Complete for constraints with the "tree property" [Kan et al, CPP 2022]
2. Apply backwards propagation once (then extract solutions)
  - Complete for "straight-line" (see later) [Chen et al, POPL 2022]
3. "Chain-free" fragment
  - More complex mix of forward and backwards [Abdulla et al, ATVA 2019]

# **Straight-Line Constraints**

# Straight Line Conjunctions

A constraint is straight-line if all of its conjunctions are straight-line

Conjunctions are "SSA" (Single Static Assignment): each variable assigned once.

A conjunction is straight-line if

- Each variable appears on the LHS at most once
- Variables can be ordered  $X_1, \dots, X_n$ 
  - If  $X_i = \dots X_j \dots$  then  $j < i$

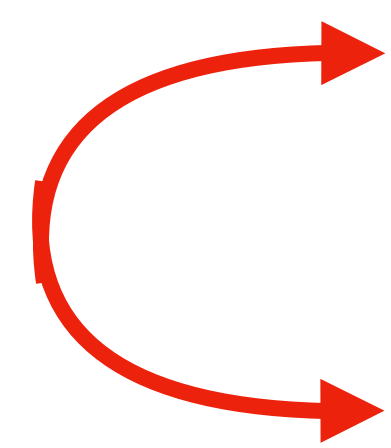
# Examples

Straight-line:

$X \in a^*b^*$ ,  
 $Y = X X$ ,  
 $Y \in abab^*$ ,  
 $Z = X Y$ ,  
 $Z \in ababab$

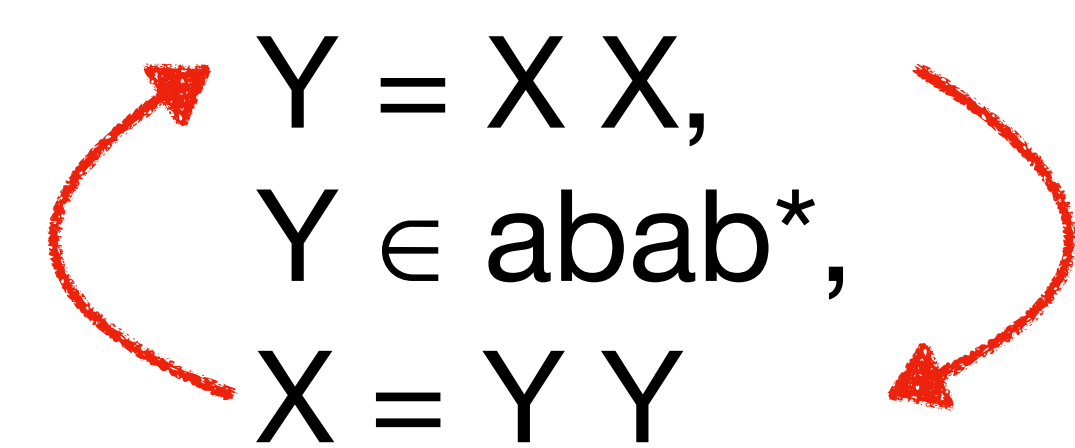
Not straight-line:

$X \in a^*b^*$ ,  
 $Y = X X$ ,  
 $Y \in abab^*$ ,  
 $Y = X Y$



Not straight-line:

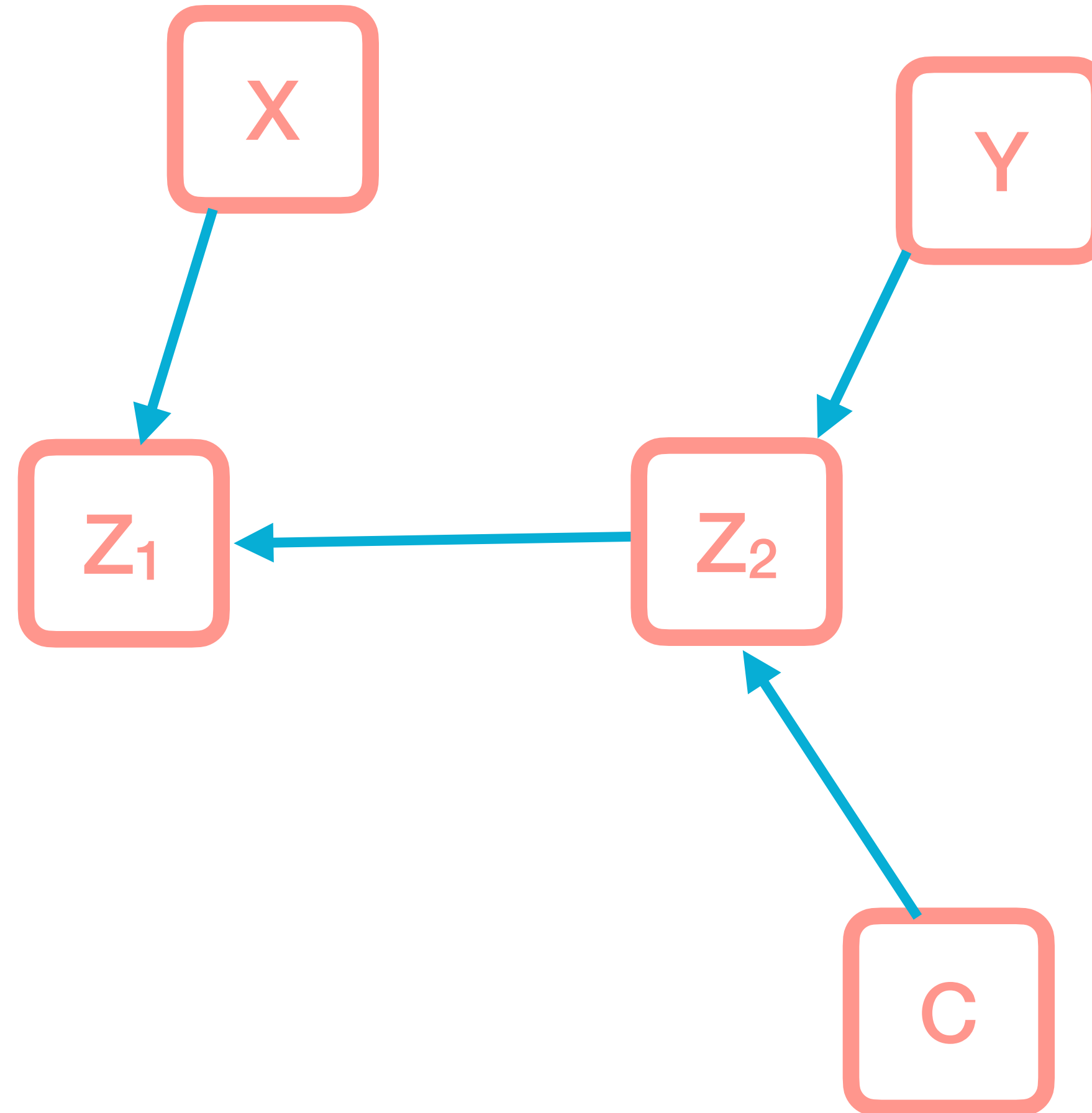
$X \in a^*b^*$ ,  
 $Y = X X$ ,  
 $Y \in abab^*$ ,  
 $X = Y Y$



# Visual Example

$$Z_1 = X Z_2, \quad Z_2 = C Y,$$

$$C \in c, \quad X \in ab^*c, \quad Z_1 \in b.^*$$



Straight-line:

- No cycles
- Each node has incoming arrows from one constraint

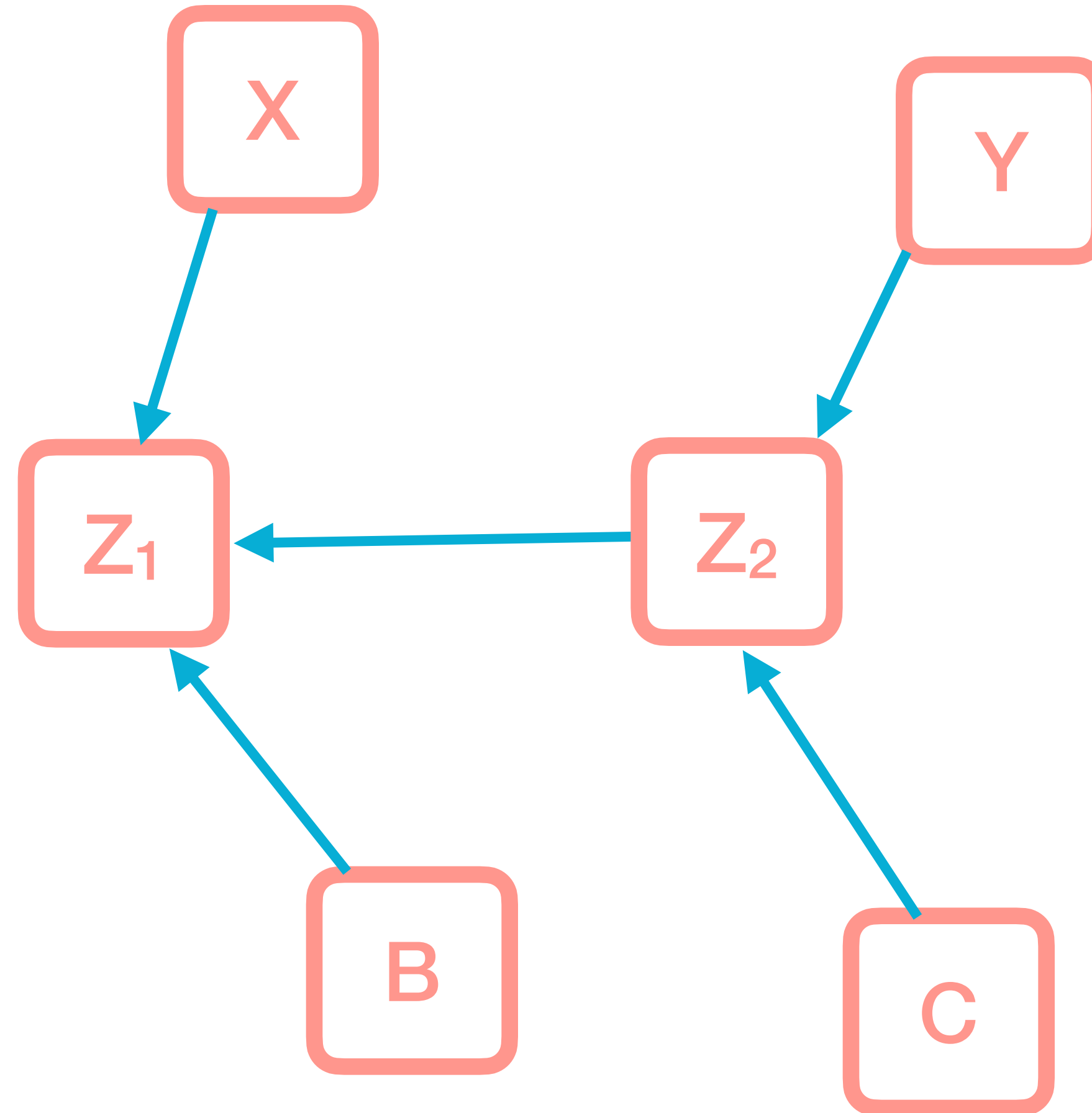


# Visual Non-Example 1

$$Z_1 = X Z_2, \quad Z_2 = C Y,$$

$$C \in c, \quad X \in ab^*c, \quad Z_1 \in b.^*,$$

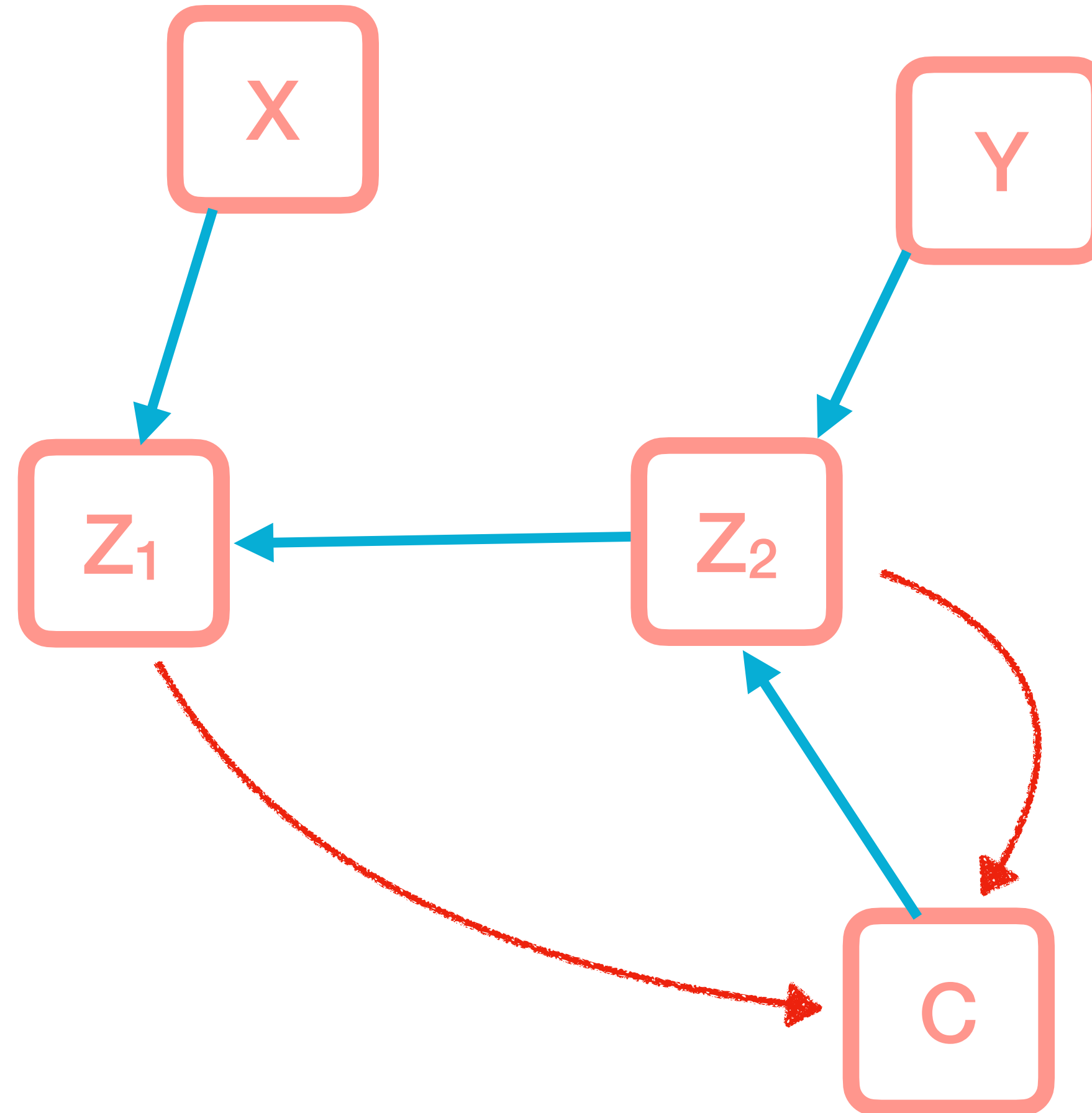
$$Z_1 = Z_2 B$$



Z1 has two incoming constraints

# Visual Non-Example 2

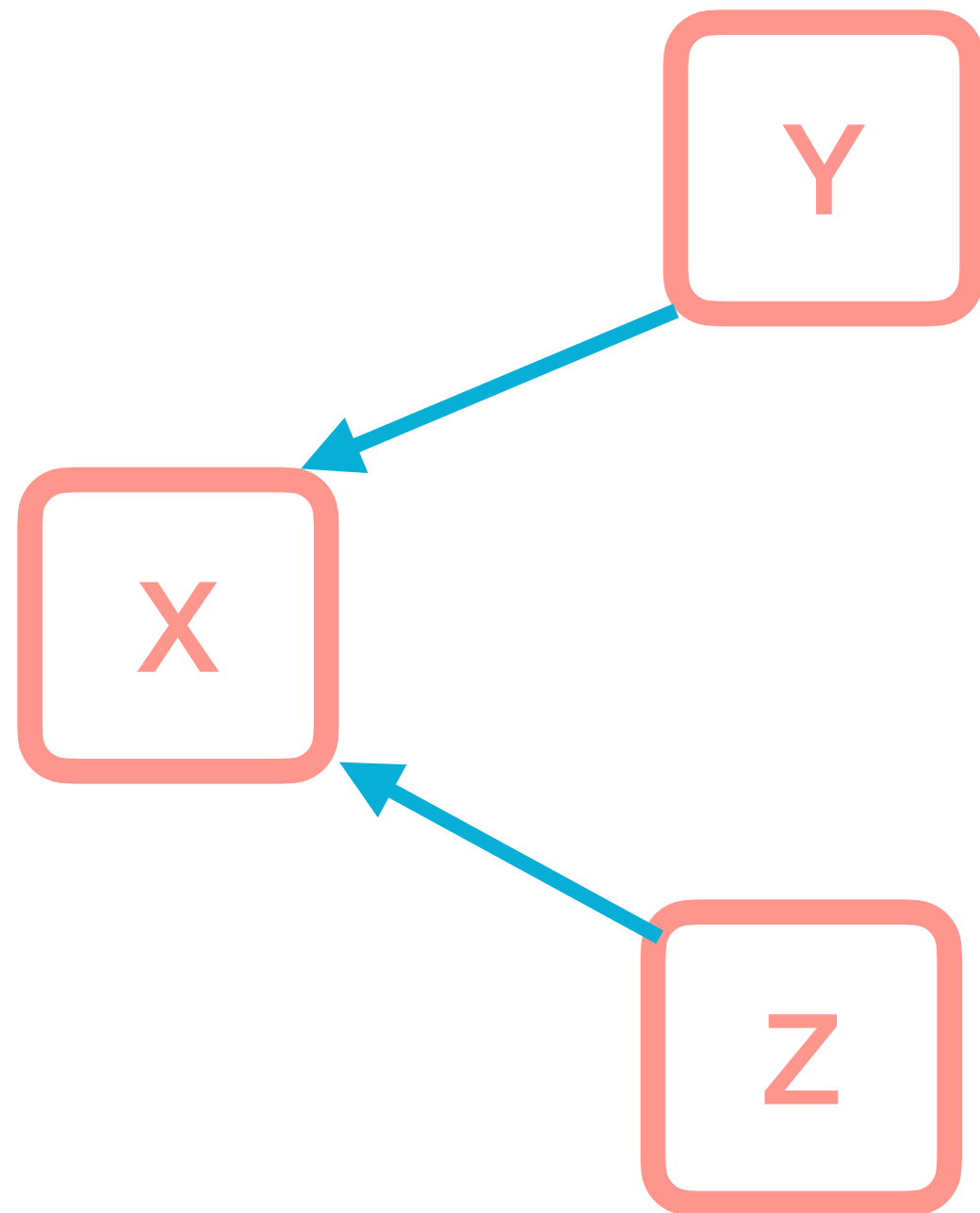
$$\begin{aligned} Z_1 &= X Z_2, & Z_2 &= C Y, \\ C &\in c, & X &\in ab^*c, & Z_1 &\in b.^*, \\ C &= Z_1 Z_2 \end{aligned}$$



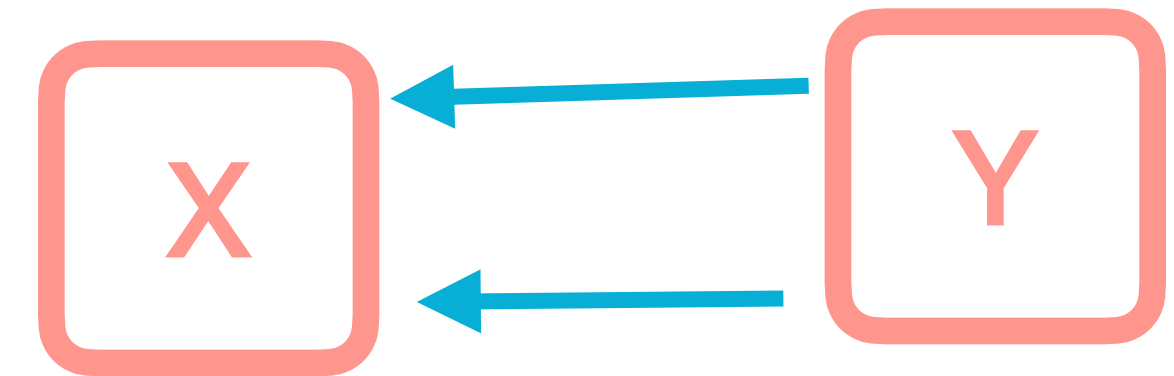
Graph has cycles

# Comparison with Tree Fragment

Straight-line is more general than the tree property



$X = Y Z$  (Straight-line and tree)

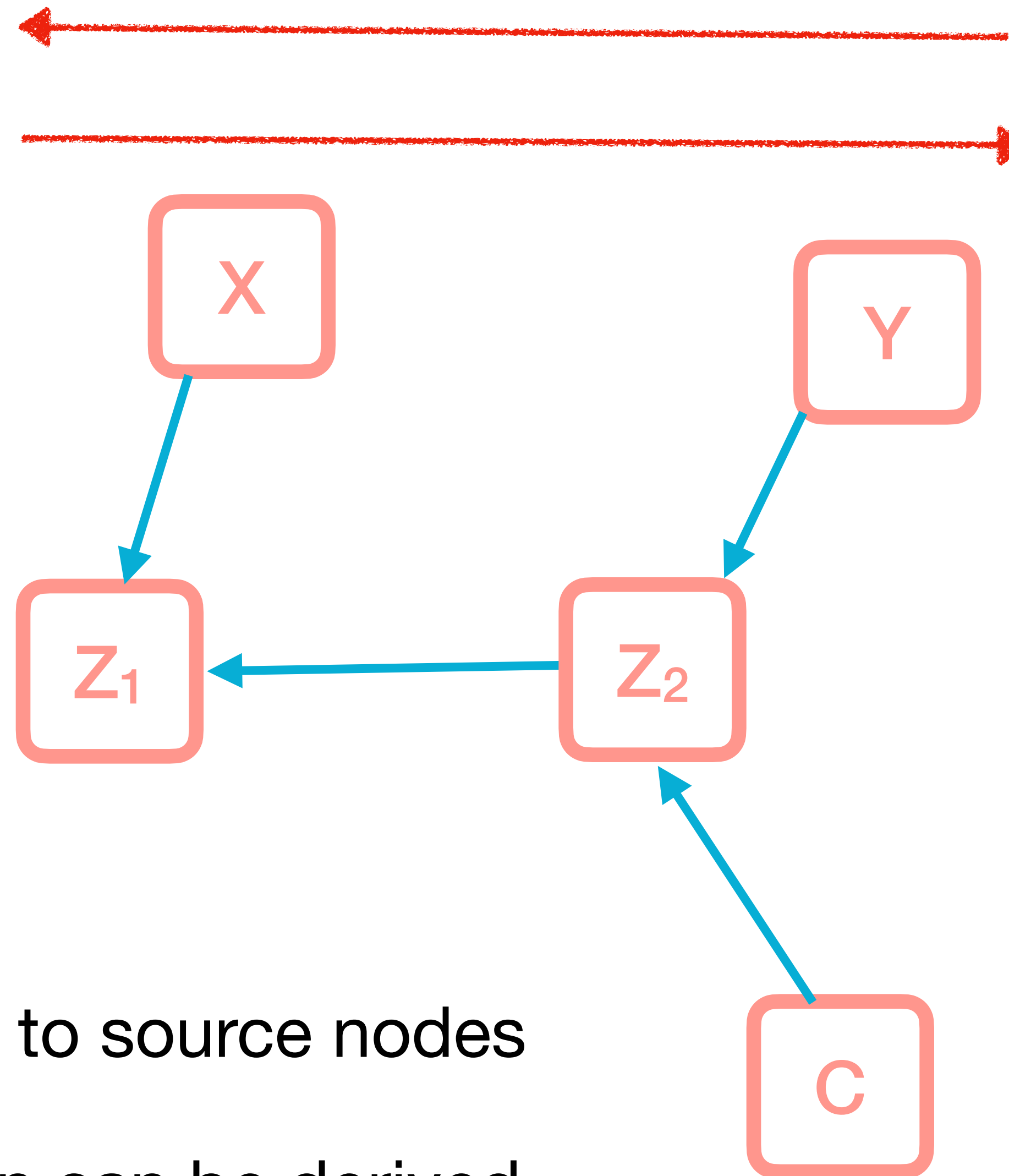


$X = Y Y$  (straight-line, not tree)

# Proof Strategy

$$Z_1 = X Z_2, \quad Z_2 = C Y,$$

$$C \in c, \quad X \in ab^*c, \quad Z_1 \in b.^*$$



Start at sink nodes

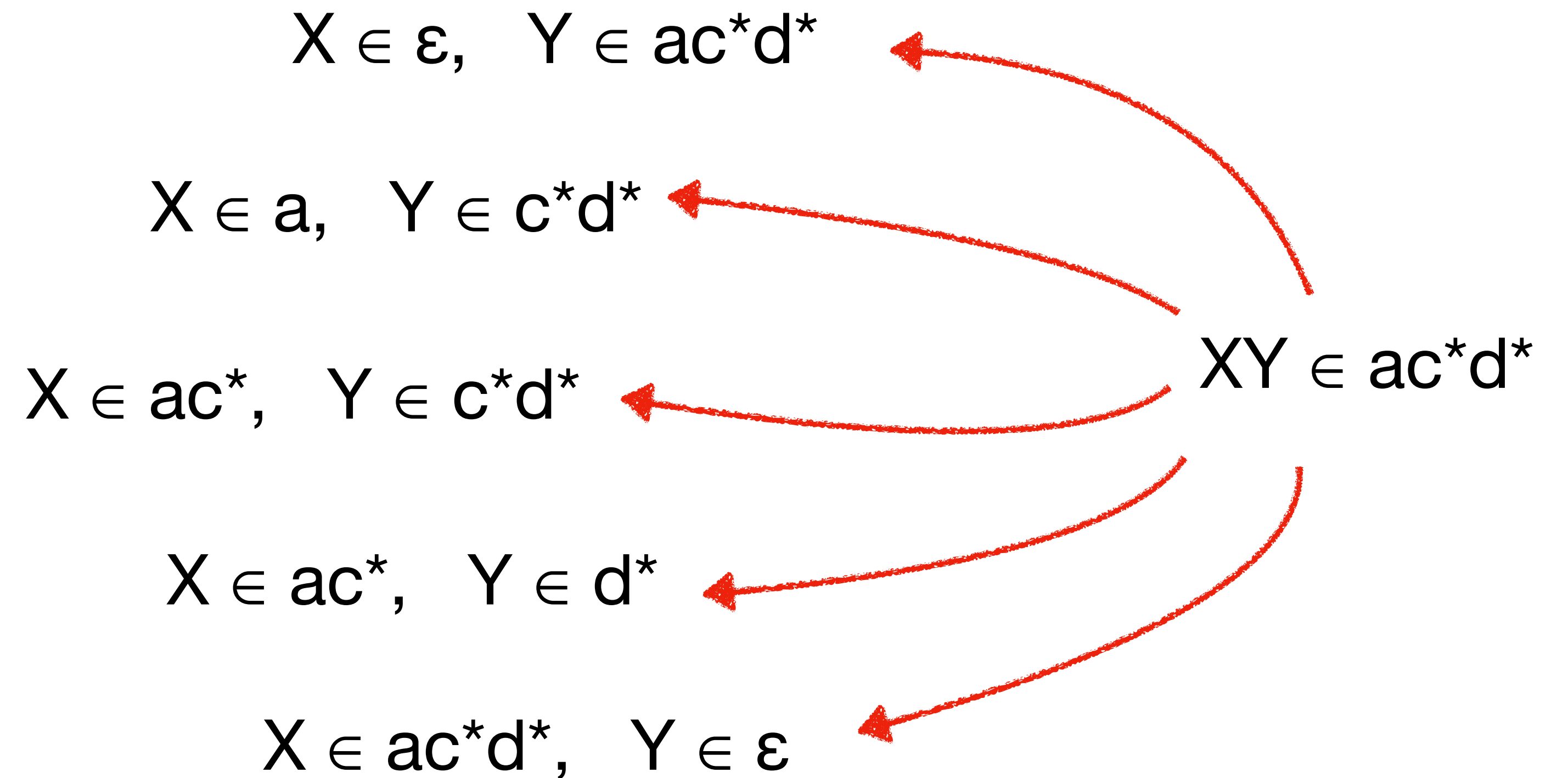
- Propagate constraints backwards to source nodes
- If no contradictions found, solution can be derived
- Terminates because no cycles

# Branching

Backwards propagation creates branching

Each branch has to be explored for unsat

For sat, only one branch needs to succeed

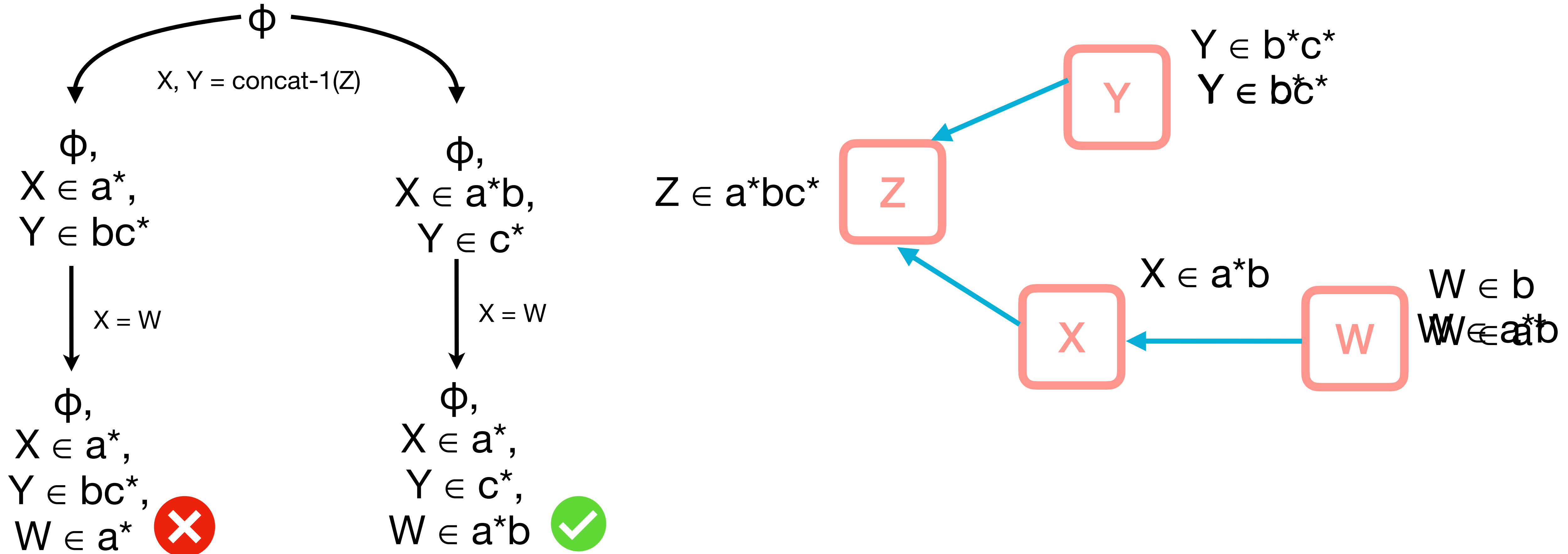


# Example

$$\phi := W \in b, \quad Y \in b^*c^*, \quad Z \in a^*bc^*,$$

$$X = W, \quad Z = XY$$

Search tree:

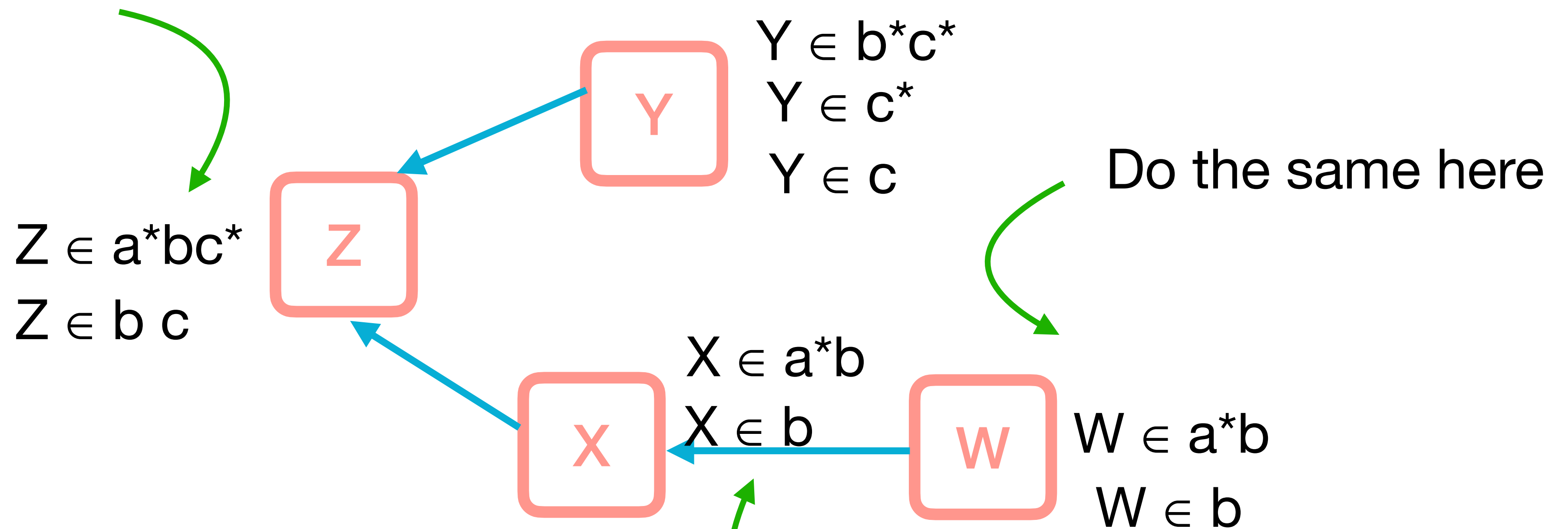


# Deriving Solutions

To derive solutions, we start at the source nodes

Eliminate via intersection  
Pick value for source node  
 $c \cap c^* \cap b^*c^* = c$

**Eliminate  $Z$  via Intersection**  $Z \in b c$



Eliminate  $X$  via Intersection one value

# Experimental Evaluation

Initial experiments on propagation strategies on QF\_S from SMT-COMP

1. Only backward propagation: T/O on 20 (out of ~7k benchmarks)
2. Only forward propagation: T/O on 33
3. Forward then backward: T/O on 18

OSTRICH: portfolio approach supporting richer constraints



# Summary

- Backwards and forwards propagation of constraints
- Solution extraction via cut
- Can support a range of string functions
- Sound, but not always complete
- Completeness for fragments:
  - Tree property
  - Straight-Line
  - Chain-free