

Navigating AND-OR Graph Modifications to Debug Failing Proof Search

Justin Lubin, Marlena Preigh, Max Willsey, Sarah E. Chasins

PLDI 2026

Proof search is ubiquitous and powerful

Proof search is ubiquitous and powerful



Databases
(Datalog)

Proof search is ubiquitous and powerful



Databases
(Datalog)



Program
synthesis

Proof search is ubiquitous and powerful



Databases
(Datalog)

Program
synthesis

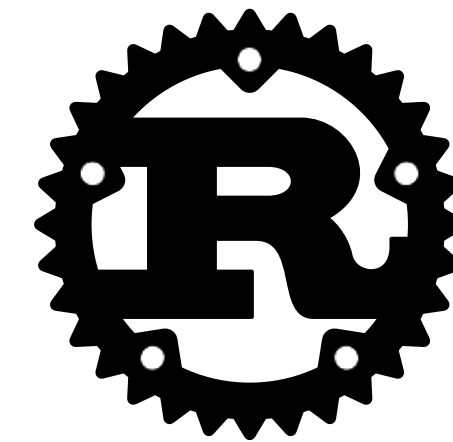
Type
systems

Proof search is ubiquitous and powerful

Databases
(Datalog)

Program
synthesis

Type
systems



Traits

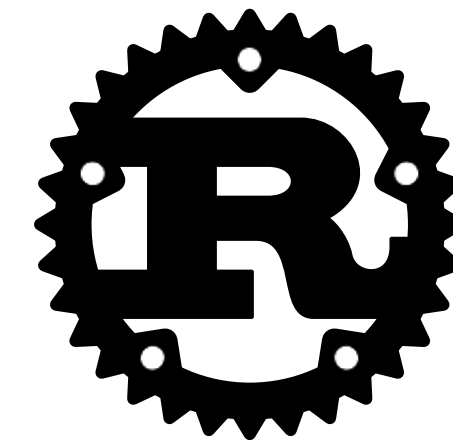
Proof search is ubiquitous and powerful

Databases
(Datalog)

Program
synthesis

Type
systems

Interactive
theorem
provers



Traits

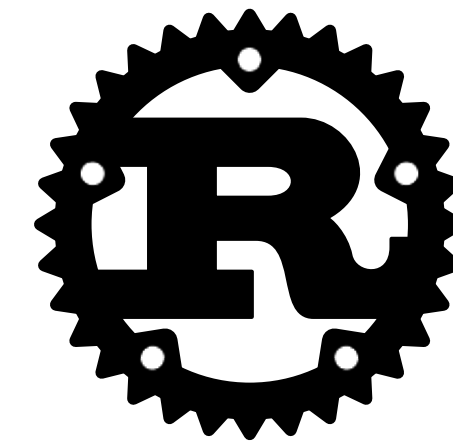
Proof search is ubiquitous and powerful

Databases
(Datalog)

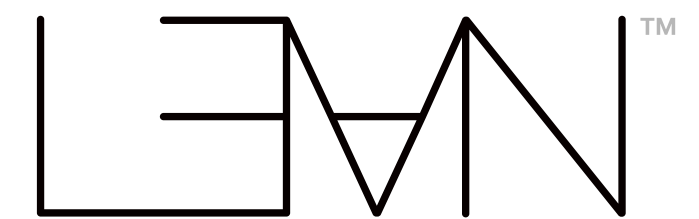
Program
synthesis

Type
systems

Interactive
theorem
provers



Traits



Aesop

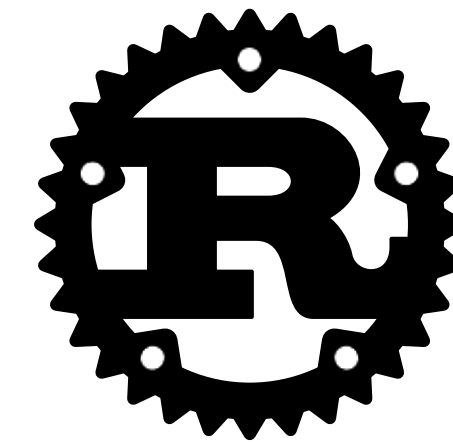
What happens when proof search fails?

Databases
(Datalog)

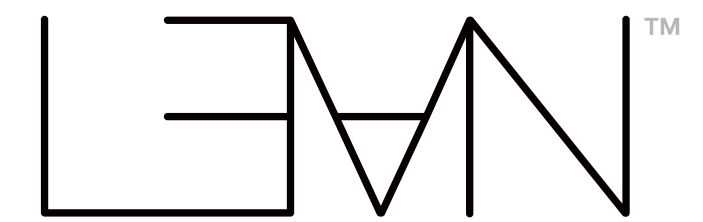
Program
synthesis

Type
systems

Interactive
theorem
provers



Traits



Aesop

What happens when proof search fails?

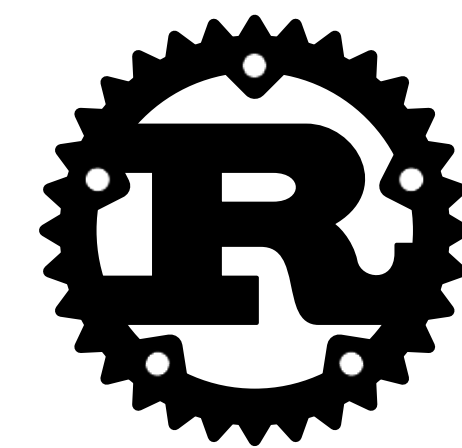
Databases
(Datalog)

Program
synthesis

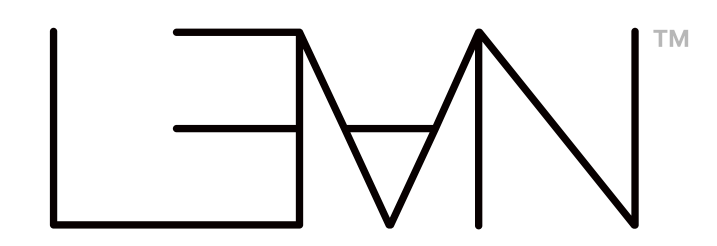
Type
systems

Interactive
theorem
provers

Check failed:
(edge 1 3)



Traits



Aesop

What happens when proof search fails?

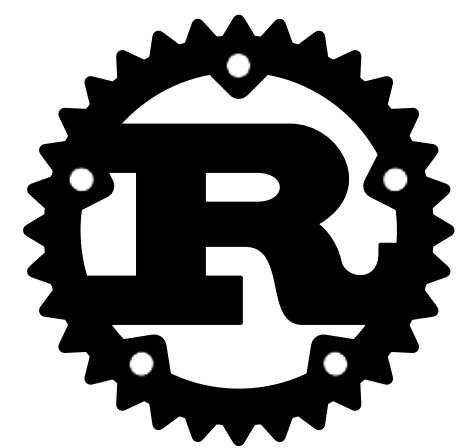
Databases
(Datalog)

Check failed:
(edge 1 3)

Program
synthesis

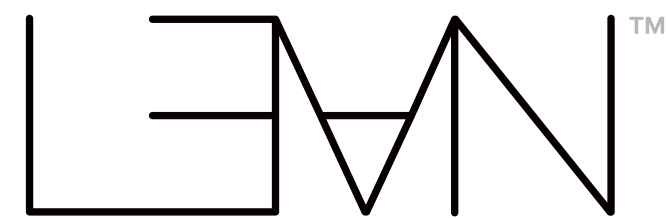
Not possible!

Type
systems



Traits

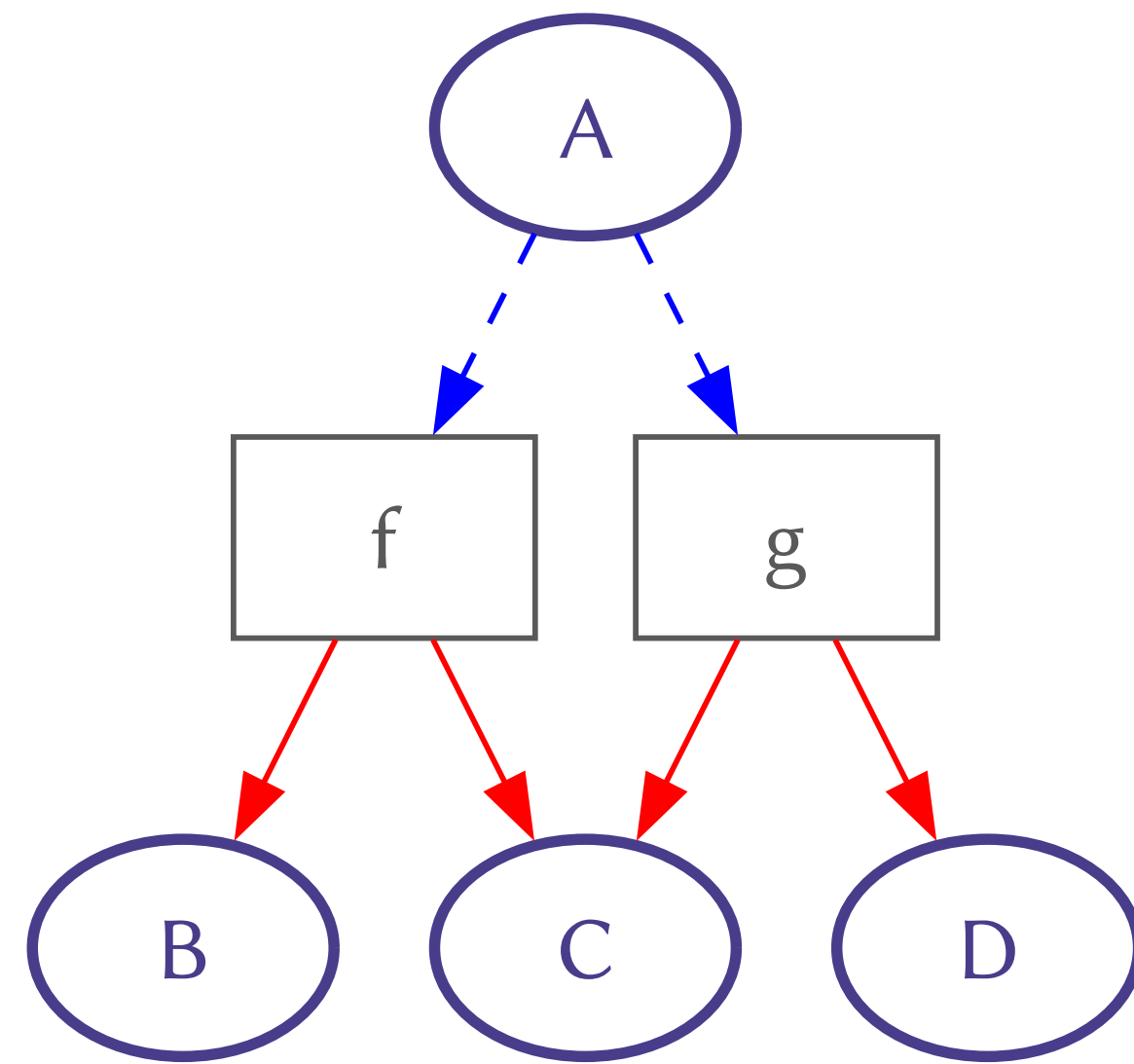
Interactive
theorem
provers



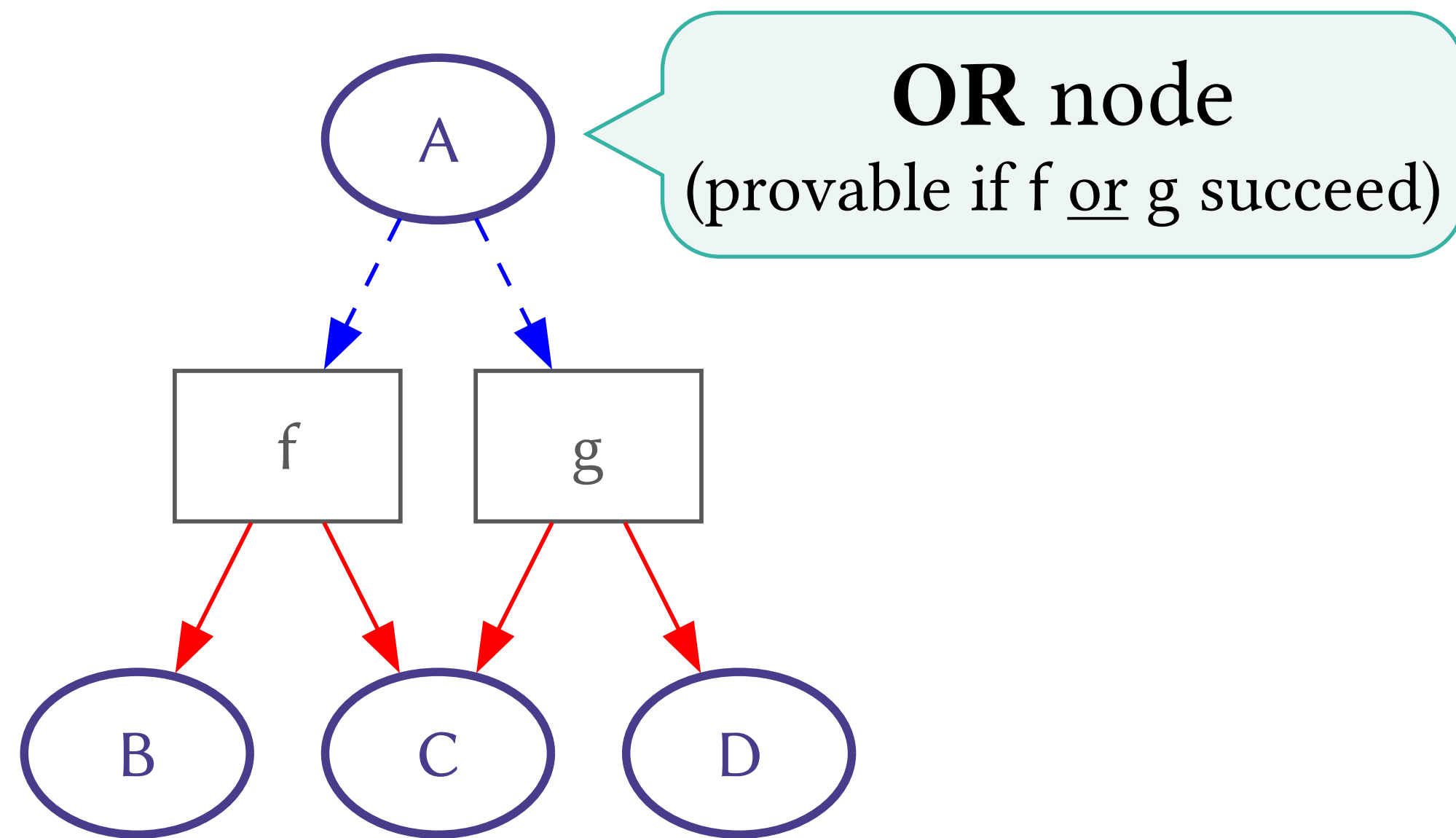
Aesop

We can represent the proof search trace using an AND–OR graph

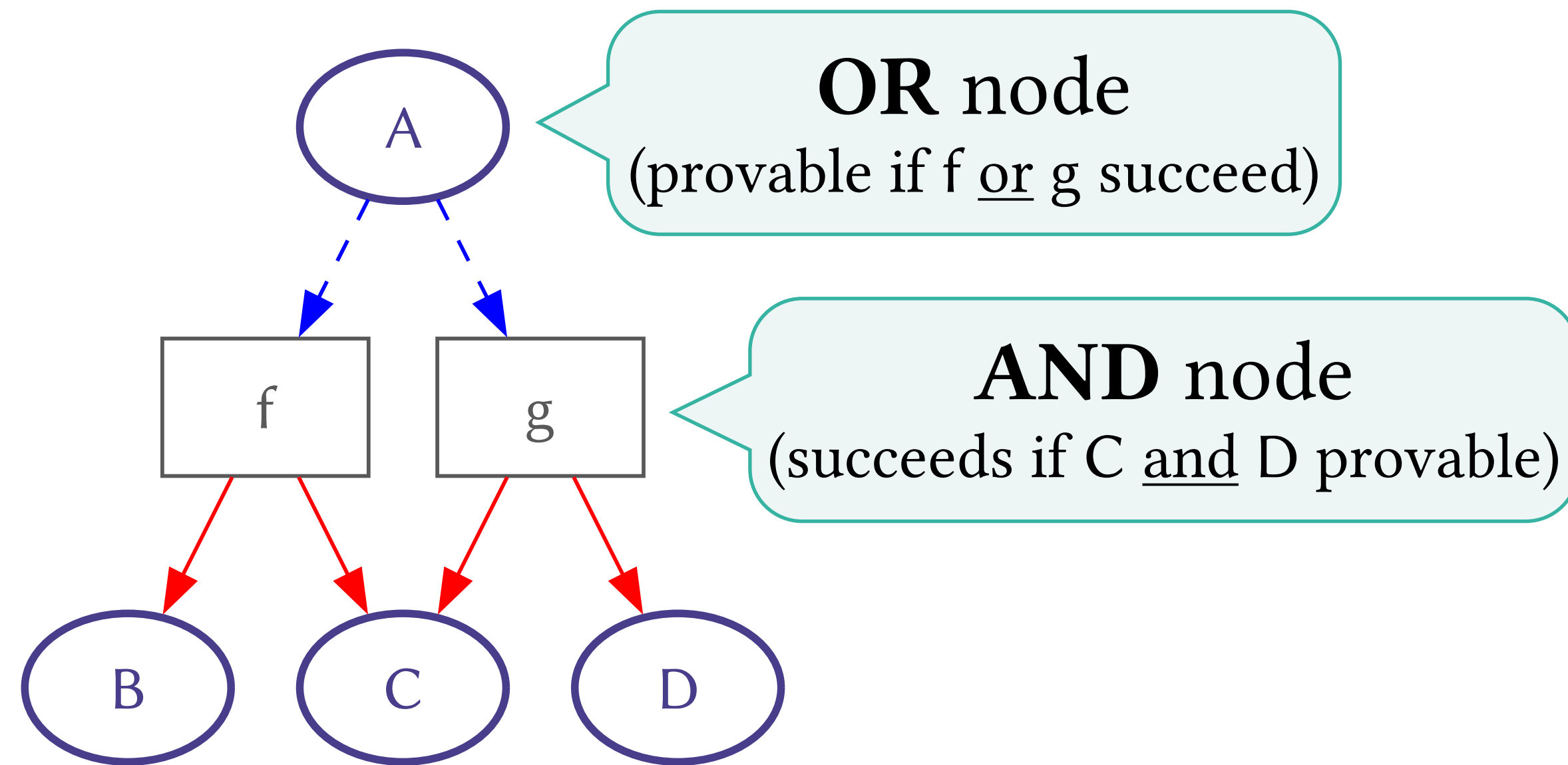
We can represent the proof search trace using an AND-OR graph



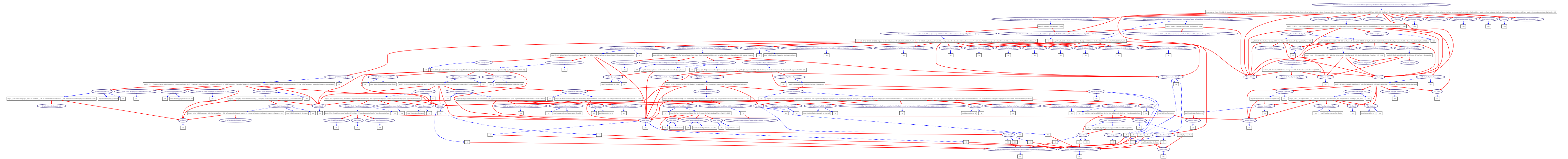
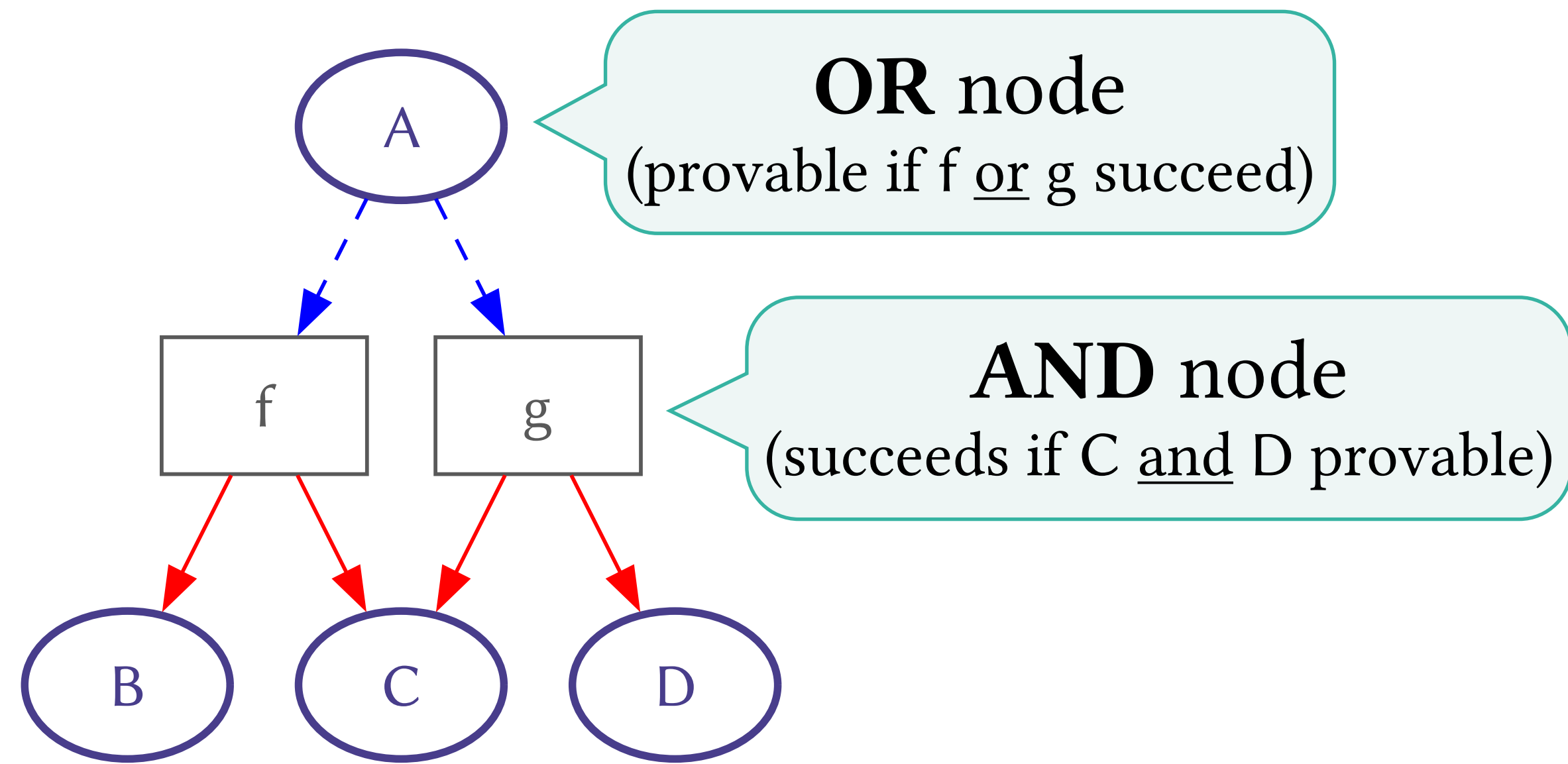
We can represent the proof search trace using an AND-OR graph



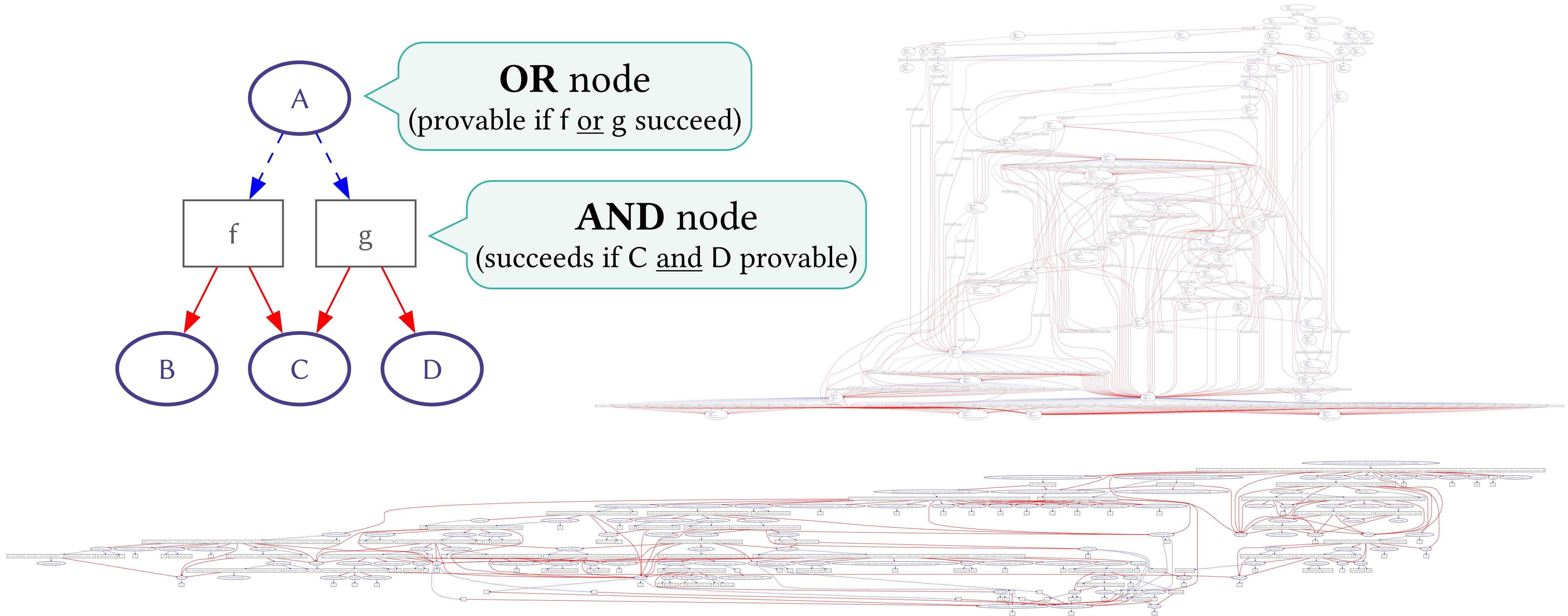
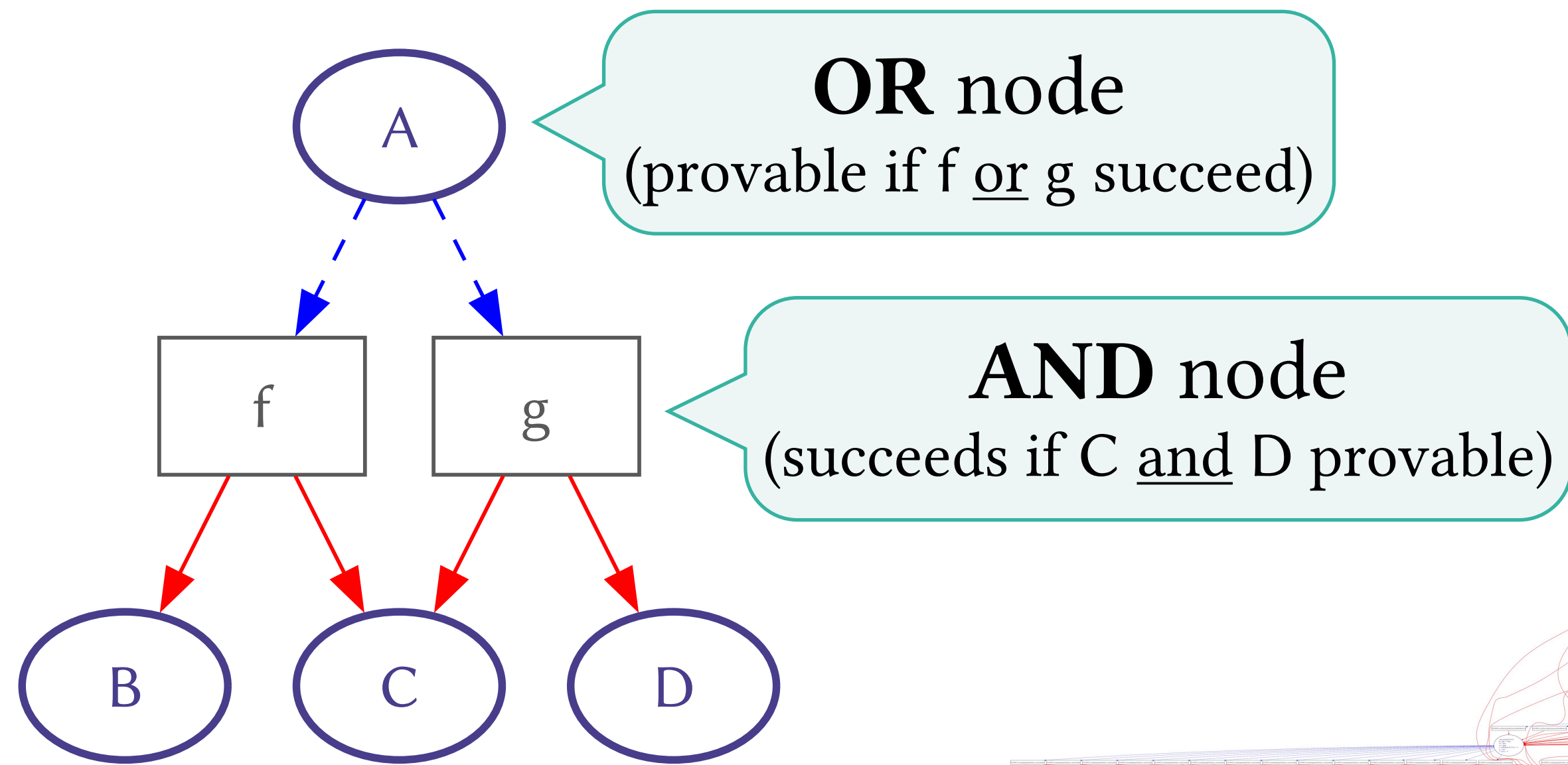
We can represent the proof search trace using an AND-OR graph



We can represent the proof search trace using an AND-OR graph



We can represent the proof search trace using an AND-OR graph



Gray, Crichton, and Krishnamurthi. An Interactive Debugger for Rust Trait Errors (*PLDI 2025*).

Limperg and From. Aesop: White-Box Best-First Proof Search for Lean (*CPP 2023*).

Hard and harder: AND graphs vs. AND-OR graphs

Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Senay and Lazzeri, 1991

Explaining Program Execution in Deductive Systems

Arora *et al.*, 1993

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Caballero *et al.*, 2017

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

Pacak and Erdweg, 2023

Datalog is used for complex programming tasks nowadays, consisting of numerous inter-dependent predicates. But Datalog lacks interactive debugging techniques that support the stepwise execution and inspection of the execution state. In this paper, we propose interactive debugging of Datalog programs following a top-down evaluation strategy called recursive query/subquery. While the recursive query/subquery approach is well-known in the literature, we are the first to provide a complete programming-language semantics based on it. Specifically, we develop the first small-step operational semantics for top-down Datalog, where subqueries occur as nested intermediate terms. The small-step semantics forms the basis of step-into interactions in the debugger. Moreover, we show how step-over interactions can be realized efficiently based on a hybrid Datalog semantics that adds a bottom-up database to our top-down operational semantics. We implemented a debugger for core Datalog following these semantics and explain how to adopt it for debugging the frontend languages of Soufflé and InCA. Our evaluation shows that our hybrid Datalog semantics can be used to debug real-world Datalog programs with realistic workloads.

CCS Concepts: • **Theory of computation** → **Constraint and logic programming; Operational semantics; Logic and databases**; • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: debugging, Datalog, small-step operational semantics, QSQR

ACM Reference Format:

André Pacak and Sebastian Erdweg. 2023. Interactive Debugging of Datalog Programs. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 248 (October 2023), 28 pages. <https://doi.org/10.1145/3622824>

1 INTRODUCTION

Datalog is a logic programming language that was invented in the 1980s as a recursive query language for databases [Maier *et al.* 2018]. However, this is not how Datalog is being used nowadays. In the last 20 years or so, it has become increasingly popular to use Datalog as a programming language to solve all kinds of problems, from program analysis [Bravenboer and Smaragdakis 2009; Madsen *et al.* 2016; Szabó *et al.* 2021] to distributed computing [Abiteboul *et al.* 2005] and network monitoring [Alvaro *et al.* 2010, 2011]. Usages of Datalog usually have two things in common. First, they process graph-structured data of considerable size, which is well supported by Datalog's fixpoint semantics. Second, they involve complex Datalog programs, consisting of many inter-dependent rules. For example, even the simplest analysis from the Doop framework for Java bytecode consists of 560 lines of Datalog code that query 32 relations from the database and compute 78 inter-dependent derived relations [Bravenboer and Smaragdakis 2009]. Sophisticated Doop analyses are an order of magnitude larger still.

Unfortunately, development tools for Datalog have not been able to keep up with the growing size and complexity of Datalog programs. Specifically, there is no interactive debugging support available for Datalog to date. Interactive debuggers assist developers in program understanding

This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: the structure of the program, the behavior of the program, and the behavior of the program's execution. This is accomplished by the generation of a separate set of bitmaps, particularly useful in...

1. Introduction

It has long been recognized that the ability to effectively debug complex computer programs is essential to the development of reliable software. This paper provides a graphical representation of program behavior and execution, which is particularly useful in...

In this paper, we describe a framework for graphical representation of logic program behavior and execution. This is accomplished by the generation of a separate set of bitmaps, particularly useful in...

TH0402-8/91/0000/00

Abstract
have a
logical
can be
gram l
possibl
correct
identif
function
a sessi
viewing
'Expla
ductivi
cursor
the in
data a
(3) Su
makes
progra
With t
facility
on the
and th
out du
is the l
nation
approa
relatio

Algorithmic debugging is a technique proposed in 1982 by E. Y. Shapiro in the context of logic programming. This survey shows how it has been adapted to fit many different paradigms and discusses the general framework of algorithmic debuggers. We elaborate a summary of algorithmic debuggers. Categories and Subject Descriptors: D.2.4 [Software Engineering]: Testing and Verification; D.2.7 [Software Engineering]: Debugging. General Terms: Languages, Verification, Theory.

Additional Key Words:
ACM Reference for
Rafael Caballero, Adrián Riesco, Josep Silva, and Ehud Y. Shapiro. Explaining Program Execution in Deductive Systems. *ACM Comput. Surv.* 50, 4, Article 60 (2018). <https://doi.org/10.1145/3181145>

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors' addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; J. Silva, Departament de Sistemes Informàtics i Enginyeria, Universitat Politècnica de València, E-46100, Valencia, Spain; E. Y. Shapiro, IBM Research, Armonk, NY, USA.

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors' addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; J. Silva, Departament de Sistemes Informàtics i Enginyeria, Universitat Politècnica de València, E-46100, Valencia, Spain; E. Y. Shapiro, IBM Research, Armonk, NY, USA.

© 2017 ACM 0360-0300/2017/1145/3181145/31

Hard and harder: AND graphs vs. AND-OR graphs

Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Senay and Lazzeri, 1991

Explaining Program Execution in Deductive Systems

Arora *et al.*, 1993

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Caballero *et al.*, 2017

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

Pacak and Erdweg, 2023

“Even though a proof tree gives a brief account of the solution steps performed during the computation, it only does so for the goals that have been solved in the process. However, it presents neither the goals that have been tried but unsolved nor the alternative clauses that could have been used to find an alternative solution.”

This paper describes a framework for graphical representation of logic program behavior. This is accomplished by the generation of a separate set of bitmaps, particularly useful in...

1. Introduction

It has long been recognized that the ability to effectively debug complex computer programs is essential to the development of reliable software. Graphical debugging programs in the form of graphical user interfaces (GUIs) have been developed so far. However, the graphical user interface for program behavior for program code still remains a mode of operation switch for the programmer. In this paper, we present a graphical user interface for program behavior at a logic program level. The graphical user interface is presented in...

Abstract
have a graphical user interface that can be used to debug logic programs. This is accomplished by the generation of a separate set of bitmaps, particularly useful in...
(3) Summary
makes program debugging easier. With this facility on the graphical user interface, the graphical user interface is the graphical user interface approach...

Algorithmic debugging is a technique proposed in 1982 by E. Y. Shapiro in the context of logic programming. This survey shows how to fit many different paradigms and discuss the general framework of algorithmic debuggers. We elaborate a summary of algorithmic debuggers. Categories and Subject Descriptors: D.2.4 [Software Engineering]: Software Testing and Verification. General Terms: Languages. Additional Key Words: ACM Reference Format: Rafael Caballero, Adrián Risco, and Josep Silva. 2017. A Survey of Algorithmic Debugging. *ACM Comput. Surv.* 50, 4, Article 60. <https://doi.org/10.1145/3122222>

1 INTRODUCTION

In the 1980s, logic programming (Kowalski 2011 and Shapiro 1986) was a paradigm that was used to describe what a program must do, leaving the details of how to do it to the computer. Ehud Y. Shapiro presented his seminal work (Shapiro 1982) as the 1982 ACM Distinguished Achievement Award.

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors' addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, 28040 Madrid, Spain; A. Risco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, 28040 Madrid, Spain; J. Silva, Departamento de Sistemas de Información, Universidad Politécnica de Valencia, E-46100, Valencia, Spain. Permission to make digital or hard copies of part or all of this work is granted by ACM, provided that copies are made without charge and that the full citation on the page is retained. Abstracting with credit is permitted. © 2017 ACM 0360-0300/2017/04-0600-02 <https://doi.org/10.1145/3122222>

Datalog is used for complex programming tasks nowadays, consisting of numerous inter-dependent predicates. But Datalog lacks interactive debugging techniques that support the stepwise execution and inspection of the execution state. In this paper, we propose interactive debugging of Datalog programs following a top-down evaluation strategy called recursive query/subquery. While the recursive query/subquery approach is well-known in the literature, we are the first to provide a complete programming-language semantics based on it. Specifically, we develop the first small-step operational semantics for top-down Datalog, where subqueries occur as nested intermediate terms. The small-step semantics forms the basis of step-into interactions in the debugger. Moreover, we show how step-over interactions can be realized efficiently based on a hybrid Datalog semantics that adds a bottom-up database to our top-down operational semantics. We implemented a debugger for core Datalog following these semantics and explain how to adopt it for debugging the frontend languages of Soufflé and InCA. Our evaluation shows that our hybrid Datalog semantics can be used to debug real-world Datalog programs with realistic workloads.

CCS Concepts: • Theory of computation → Constraint and logic programming; Operational semantics; Logic and databases; • Software and its engineering → Software testing and debugging.

Additional Key Words and Phrases: debugging, Datalog, small-step operational semantics, QSQR

ACM Reference Format:

André Pacak and Sebastian Erdweg. 2023. Interactive Debugging of Datalog Programs. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 248 (October 2023), 28 pages. <https://doi.org/10.1145/3622824>

1 INTRODUCTION

Datalog is a logic programming language that was invented in the 1980s as a recursive query language for databases [Maier *et al.* 2018]. However, this is not how Datalog is being used nowadays. In the last 20 years or so, it has become increasingly popular to use Datalog as a programming language to solve all kinds of problems, from program analysis [Bravenboer and Smaragdakis 2009; Madsen *et al.* 2016; Szabó *et al.* 2021] to distributed computing [Abiteboul *et al.* 2005] and network monitoring [Alvaro *et al.* 2010, 2011]. Usages of Datalog usually have two things in common. First, they process graph-structured data of considerable size, which is well supported by Datalog's fixpoint semantics. Second, they involve complex Datalog programs, consisting of many inter-dependent rules. For example, even the simplest analysis from the Doop framework for Java bytecode consists of 560 lines of Datalog code that query 32 relations from the database and compute 78 inter-dependent derived relations [Bravenboer and Smaragdakis 2009]. Sophisticated Doop analyses are an order of magnitude larger still.

Unfortunately, development tools for Datalog have not been able to keep up with the growing size and complexity of Datalog programs. Specifically, there is no interactive debugging support available for Datalog to date. Interactive debuggers assist developers in program understanding

Authors' addresses: André Pacak, pacak@uni-mainz.de, JGU Mainz, Germany; Sebastian Erdweg, erdweg@uni-mainz.de, JGU Mainz, Germany.

Hard and harder: AND graphs vs. AND-OR graphs

Senay and Lazzeri, 1991

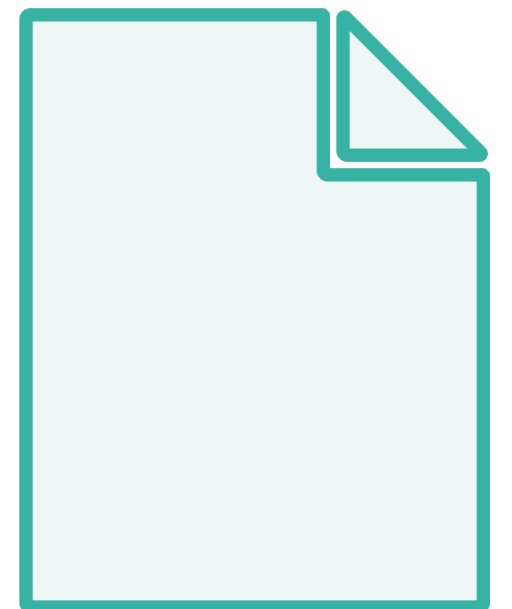
“Even though a proof tree gives a brief account of the solution steps performed during the computation, it only does so for the goals that have been solved in the process. However, it presents neither the goals that have been tried but unsolved nor the alternative clauses that could have been used to find an alternative solution.”

Arora *et al.*, 1993

Caballero *et al.*, 2017

Pacak and Erdweg, 2023

AND graph
(buggy program/proof) \approx



Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Explaining Program Execution in Deductive Systems

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

Datalog is used for complex programming tasks nowadays, consisting of numerous inter-dependent predicates. But Datalog lacks interactive debugging techniques that support the stepwise execution and inspection of the execution state. In this paper, we propose interactive debugging of Datalog programs following a top-down evaluation strategy called recursive query/subquery. While the recursive query/subquery approach is well-known in the literature, we are the first to provide a complete programming-language semantics based on it. Specifically, we develop the first small-step operational semantics for top-down Datalog, where subqueries occur as nested intermediate terms. The small-step semantics forms the basis of step-into interactions in the debugger. Moreover, we show how step-over interactions can be realized efficiently based on a hybrid Datalog semantics that adds a bottom-up database to our top-down operational semantics. We implemented a debugger for core Datalog following these semantics and explain how to adopt it for debugging the frontend languages of Soufflé and Inca. Our evaluation shows that our hybrid Datalog semantics can be used to debug real-world Datalog programs with realistic workloads.

CCS Concepts: • **Theory of computation** → **Constraint and logic programming; Operational semantics; Logic and databases**; • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: debugging, Datalog, small-step operational semantics, QSQR

ACM Reference Format:

André Pacak and Sebastian Erdweg. 2023. Interactive Debugging of Datalog Programs. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 248 (October 2023), 28 pages. <https://doi.org/10.1145/3622824>

1 INTRODUCTION

Datalog is a logic programming language that was invented in the 1980s as a recursive query language for databases [Maier *et al.* 2018]. However, this is not how Datalog is being used nowadays. In the last 20 years or so, it has become increasingly popular to use Datalog as a programming language to solve all kinds of problems, from program analysis [Bravenboer and Smaragdakis 2009; Madsen *et al.* 2016; Szabó *et al.* 2021] to distributed computing [Abiteboul *et al.* 2005] and network monitoring [Alvaro *et al.* 2010, 2011]. Usages of Datalog usually have two things in common. First, they process graph-structured data of considerable size, which is well supported by Datalog’s fixpoint semantics. Second, they involve complex Datalog programs, consisting of many inter-dependent rules. For example, even the simplest analysis from the Doop framework for Java bytecode consists of 560 lines of Datalog code that query 32 relations from the database and compute 78 inter-dependent derived relations [Bravenboer and Smaragdakis 2009]. Sophisticated Doop analyses are an order of magnitude larger still.

Unfortunately, development tools for Datalog have not been able to keep up with the growing size and complexity of Datalog programs. Specifically, there is no interactive debugging support available for Datalog to date. Interactive debuggers assist developers in program understanding

This paper describes a framework for graphical environments that includes three types of representations for these graphs: the most important scheme presented, program behavior, and this is accomplished by the generation of a separate set of bitmaps, particularly useful in...

1. Introduction

It has long been recognized that the ability to effectively debug complex computer programs is a difficult task. The development of graphical debugging environments has made it possible to view the execution of a program in a more intuitive way. This paper describes a framework for graphical environments that includes three types of representations for these graphs: the most important scheme presented, program behavior, and this is accomplished by the generation of a separate set of bitmaps, particularly useful in...

Abstract: This survey shows how to fit many different paradigms and discuss the general framework of algorithmic debuggers. Categories and Subject Descriptors: D.2.4 [Software Engineering]: Testing and Verification. General Terms: Languages. Additional Key Words: ACM Reference Format: Rafael Caballero, Adrián Riesco, and Josep Silva. 2017. Explaining Program Execution in Deductive Systems. *Proc. ACM Program. Lang.* 1, Article 60 (October 2017), 28 pages. <https://doi.org/10.1145/3177224>

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors’ addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, E-28042, Madrid, Spain; J. Silva, Departamento de Ingeniería de Sistemas y Computación, Universidad Politécnica de Valencia, E-46100, Valencia, Spain. Permission to make digital or hard copies of this work for personal or classroom use is granted by ACM, provided that the full citation on the front page is included. Abstracting with credit is permitted. © 2017 ACM 0360-0300/2017/10-1145/31

TH0402-8/91/0000/00

Hard and harder: AND graphs vs. AND-OR graphs

Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Senay and Lazzeri, 1991

Explaining Program Execution in Deductive Systems

Arora *et al.*, 1993

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Caballero *et al.*, 2017

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

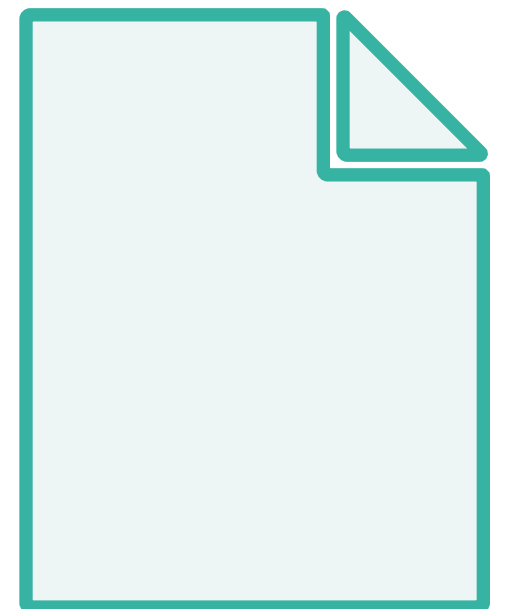
Pacak and Erdweg, 2023

AND graph
(buggy program/proof)

≈

Tools to help: type systems, debuggers
(step-based, algorithmic), *etc.*

“Even though a proof tree gives a brief account of the solution steps performed during the computation, it only does so for the goals that have been solved in the process. However, it presents neither the goals that have been tried but unsolved nor the alternative clauses that could have been used to find an alternative solution.”



This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: a graphical representation of the structure of logic programs, a graphical representation of the behavior of logic programs, and a graphical representation of the execution of logic programs. The most important scheme presented is the graphical representation of the behavior of logic programs. This is accomplished by the generation of a separate set of bitmaps, which are particularly useful in the debugging of logic programs.

1. Introduction

It has long been recognized that the ability to effectively debug complex computer programs is a difficult task. The development of graphical debugging tools has provided a significant improvement in the ability to debug logic programs. Although the development of graphical debugging tools has been a long process, it has now reached a stage where it is possible to provide a graphical representation of the execution of logic programs. This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: a graphical representation of the structure of logic programs, a graphical representation of the behavior of logic programs, and a graphical representation of the execution of logic programs. The most important scheme presented is the graphical representation of the behavior of logic programs. This is accomplished by the generation of a separate set of bitmaps, which are particularly useful in the debugging of logic programs.

Abstract: This survey shows how to fit many different paradigms and discuss the general framework of algorithmic debuggers. We elaborate a summary of algorithmic debuggers. Categories and Subject Descriptors: D.2.4 [Software Engineering]: Testing and Verification. General Terms: Languages. Additional Key Words: Algorithmic debugging, Datalog, small-step operational semantics, QSQR.

1 INTRODUCTION

In the 1980s, logic programming (Kowalski 2011 and Shapiro 1986) was a new paradigm. With the advent of logic programming, the need for a graphical representation of the execution of logic programs became apparent. This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: a graphical representation of the structure of logic programs, a graphical representation of the behavior of logic programs, and a graphical representation of the execution of logic programs. The most important scheme presented is the graphical representation of the behavior of logic programs. This is accomplished by the generation of a separate set of bitmaps, which are particularly useful in the debugging of logic programs.

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors' addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 47, E-28042, Madrid, Spain; J. Silva, Universitat Politècnica de València, E-46100, Valencia, Spain. Permission to make digital or hard copies of this work for personal or classroom use is granted by ACM, provided that the fee of \$15.00 is paid directly to ACM. Abstracting with credit is permitted. © 2017 ACM 0360-0300/2017/08 \$15.00 <https://doi.org/10.1145/3122222>

Hard and harder: AND graphs vs. AND-OR graphs

“Even though a proof tree gives a brief account of the solution steps performed during the computation, it only does so for the goals that have been solved in the process. However, it presents neither the goals that have been tried but unsolved nor the alternative clauses that could have been used to find an alternative solution.”

Senay and Lazzeri, 1991

Arora *et al.*, 1993

Caballero *et al.*, 2017

Pacak and Erdweg, 2023

AND graph
(buggy program/proof)

≈

Tools to help: type systems, debuggers
(step-based, algorithmic), *etc.*

AND-OR graph
(all possible failing programs/proofs)

≈

Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Explaining Program Execution in Deductive Systems

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

Datalog is used for complex programming tasks nowadays, consisting of numerous inter-dependent predicates. But Datalog lacks interactive debugging techniques that support the stepwise execution and inspection of the execution state. In this paper, we propose interactive debugging of Datalog programs following a top-down evaluation strategy called recursive query/subquery. While the recursive query/subquery approach is well-known in the literature, we are the first to provide a complete programming-language semantics based on it. Specifically, we develop the first small-step operational semantics for top-down Datalog, where subqueries occur as nested intermediate terms. The small-step semantics forms the basis of step-into interactions in the debugger. Moreover, we show how step-over interactions can be realized efficiently based on a hybrid Datalog semantics that adds a bottom-up database to our top-down operational semantics. We implemented a debugger for core Datalog following these semantics and explain how to adopt it for debugging the frontend languages of Souflé and InCA. Our evaluation shows that our hybrid Datalog semantics can be used to debug real-world Datalog programs with realistic workloads.

CCS Concepts: • Theory of computation → Constraint and logic programming; Operational semantics; Logic and databases; • Software and its engineering → Software testing and debugging.

Additional Key Words and Phrases: debugging, Datalog, small-step operational semantics, QSQR

ACM Reference Format:

André Pacak and Sebastian Erdweg. 2023. Interactive Debugging of Datalog Programs. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 248 (October 2023), 28 pages. <https://doi.org/10.1145/3622824>

1 INTRODUCTION

Datalog is a logic programming language that was invented in the 1980s as a recursive query language for databases [Maier *et al.* 2018]. However, this is not how Datalog is being used nowadays. In the last 20 years or so, it has become increasingly popular to use Datalog as a programming language to solve all kinds of problems, from program analysis [Bravenboer and Smaragdakis 2009; Madsen *et al.* 2016; Szabó *et al.* 2021] to distributed computing [Abiteboul *et al.* 2005] and network monitoring [Alvaro *et al.* 2010, 2011]. Usages of Datalog usually have two things in common. First, they process graph-structured data of considerable size, which is well supported by Datalog’s fixpoint semantics. Second, they involve complex Datalog programs, consisting of many inter-dependent rules. For example, even the simplest analysis from the Doop framework for Java bytecode consists of 560 lines of Datalog code that query 32 relations from the database and compute 78 inter-dependent derived relations [Bravenboer and Smaragdakis 2009]. Sophisticated Doop analyses are an order of magnitude larger still.

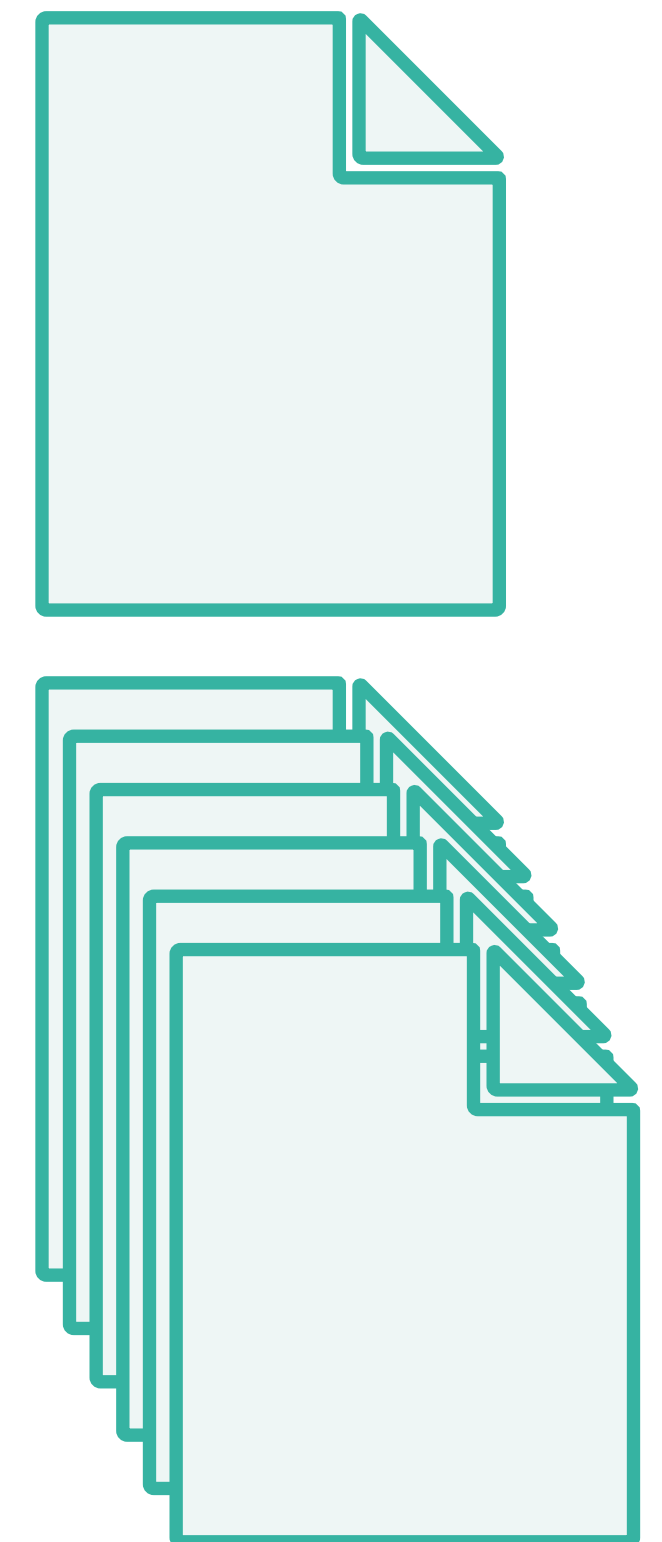
Unfortunately, development tools for Datalog have not been able to keep up with the growing size and complexity of Datalog programs. Specifically, there is no interactive debugging support available for Datalog to date. Interactive debuggers assist developers in program understanding

Abstract
have a
logical
can be
gram l
possibl
correct
identif
function
a sessi
viewing
‘Expla
ductivi
cursor
the in
data a
(3) Su
makes
progra
With t
facility
on the
and th
out du
is the
nation
approa
relatio

Algorithmic debugging is a technique proposed in 1982 by E. V. Shapiro in the context of logic programming. This survey shows how fitting many different paradigms and discuss mic debugger tools that we elaborate a summa algorithmic debuggers. Categories and Subject Specifications and verification definitions and theory. General Terms: Language. Additional Key Words: ACM Reference for Rafael Caballero, Adrián Riesco, Josep Silva, and Josep Silva. *Surveys*, 4, Article 60 (2017). <https://doi.org/10.1145/3122824>

1 INTRODUCTION
In the 1980s, logic programming (Kowalski 2011 and Shapiro 1986) was a program must do, leaving Ehud Y. Shapiro’s seminal work (Shapiro 1982) as the 1982 ACM Distinguished Paper.

This work has been partially supported by grant TIN2013-44742-C4-3-3-R by the Generalitat Valenciana project N-Greens Software. Authors’ addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 4, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo, 4, E-28042, Madrid, Spain; J. Silva, Departamento de Sistemas Informáticos, Universidad Politécnica de Valencia, E-46100, Valencia, Spain. Permission to make digital or hard copies of this work for personal or classroom use is granted by ACM, provided that copies are made without fee for individuals and small businesses, that they are made for personal or classroom use, and that they bear this notice with the full citation on the first page. Abstracting with credit is permitted. For all other uses, permission should be sought from ACM. © 2017 ACM 0360-0300/2017/11-44742-04 \$15.00. <https://doi.org/10.1145/3122824>



Hard and harder: AND graphs vs. AND-OR graphs

Graphical Representation of Logic Programs and Their Behavior

Hikmet Senay Santos G. Lazzeri

Department of Electrical Engineering and Computer Science
The George Washington University

Senay and Lazzeri, 1991

Explaining Program Execution in Deductive Systems

Arora *et al.*, 1993

A Survey of Algorithmic Debugging

RAFAEL CABALLERO and ADRIÁN RIESCO, Universidad Complutense de Madrid
JOSEP SILVA, Universitat Politècnica de València

Caballero *et al.*, 2017

Interactive Debugging of Datalog Programs

ANDRÉ PACAK, JGU Mainz, Germany
SEBASTIAN ERDWEG, JGU Mainz, Germany

Pacak and Erdweg, 2023

AND graph
(buggy program/proof)

≈

Tools to help: type systems, debuggers
(step-based, algorithmic), *etc.*

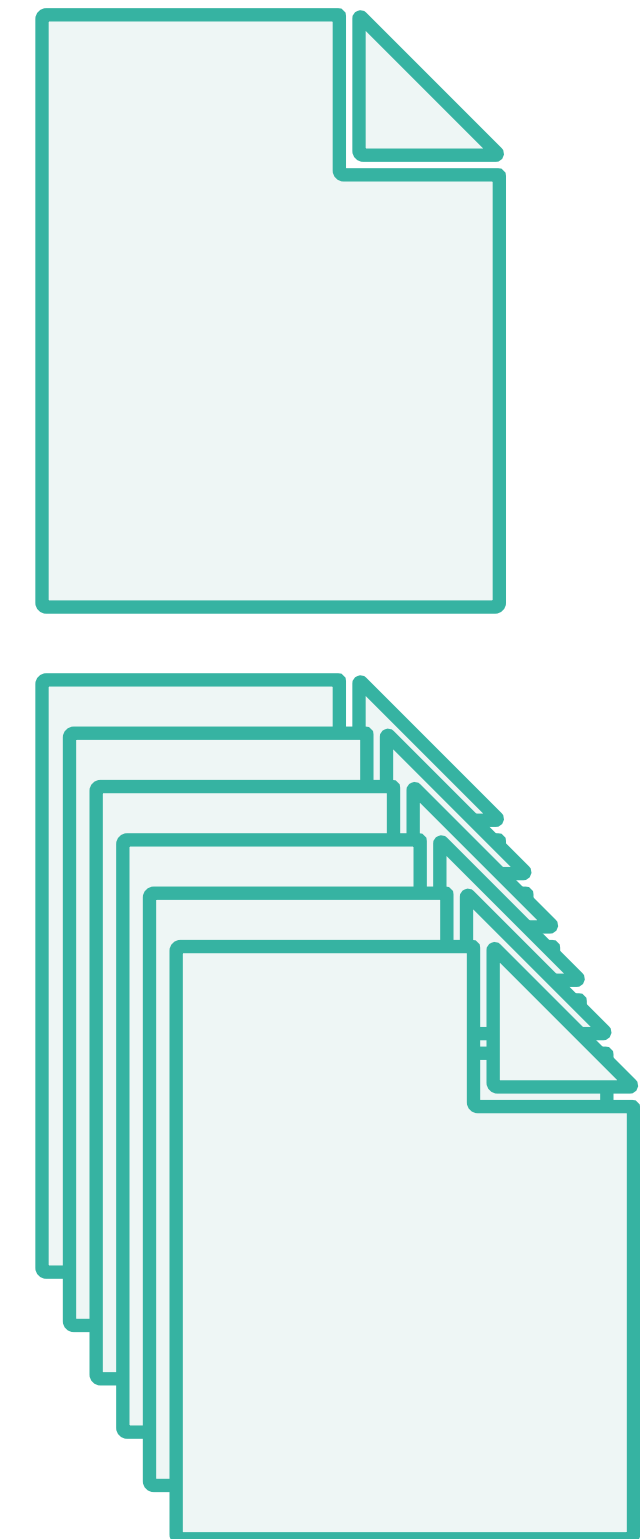
AND-OR graph

(all possible failing programs/proofs)

≈

Tools to help: Graph visualization,
graph traversal

“Even though a proof tree gives a brief account of the solution steps performed during the computation, it only does so for the goals that have been solved in the process. However, it presents neither the goals that have been tried but unsolved nor the alternative clauses that could have been used to find an alternative solution.”



This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: a structural graph, a behavioral graph, and a proof graph. The structural graph shows the structure of the logic program. The behavioral graph shows the behavior of the logic program. The proof graph shows the proof of the logic program. This is accomplished by generating a set of bitmaps, which are particularly useful in the context of logic programming.

1. Introduction

It has long been recognized that the ability to effectively debug complex computer programs is a critical skill. The development of graphical debugging tools for logic programs is a particularly challenging task. Although the development of such tools has been ongoing for many years, the current state of the art is still limited. This paper describes a framework for graphical representation of logic programs. It includes three types of graphs: a structural graph, a behavioral graph, and a proof graph. The structural graph shows the structure of the logic program. The behavioral graph shows the behavior of the logic program. The proof graph shows the proof of the logic program. This is accomplished by generating a set of bitmaps, which are particularly useful in the context of logic programming.

TH0402-8/91/0000/00

Abstract: This survey shows how to fit many different paradigms and discuss the general framework of algorithmic debuggers. We elaborate a summary of algorithmic debuggers. Categories and Subject Descriptors: D.2.4 [Software Engineering]: Testing and Verification. General Terms: Languages. Additional Key Words: Algorithmic debugging. ACM Reference Format: Rafael Caballero, Adrián Riesco, and Josep Silva. 2017. A Survey of Algorithmic Debugging. *ACM Comput. Surv.* 50, 4, Article 60. <https://doi.org/10.1145/3022824>

1 INTRODUCTION

In the 1980s, logic programming (Kowalski 2011 and Shapiro 1986) was a paradigm that was used to solve a wide range of problems. With the advent of logic programming, the need for a facility to debug logic programs became apparent. The development of such tools has been ongoing for many years, but the current state of the art is still limited.

This work has been partially supported by the grant TIN2013-44742-C4-3-R by the Generalitat Valenciana project N-Greens Software. Authors' addresses: R. Caballero, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, E-28042, Madrid, Spain; A. Riesco, Universidad Complutense de Madrid, C/ Arzobispo Morcillo s/n, E-28042, Valencia, Spain; J. Silva, Universitat Politècnica de València, E-46100, Valencia, Spain. Permission to make digital or hard copies of this work for personal or classroom use is granted by ACM, provided that the fee of \$15.00 is paid directly to ACM. Abstracting with credit is permitted. For more information on copying, see the full text of the article. © 2017 ACM 0360-0300/2017/04-0060-08 <https://doi.org/10.1145/3022824>

Datalog is used for complex programming tasks nowadays, consisting of numerous inter-dependent predicates. But Datalog lacks interactive debugging techniques that support the stepwise execution and inspection of the execution state. In this paper, we propose interactive debugging of Datalog programs following a top-down evaluation strategy called recursive query/subquery. While the recursive query/subquery approach is well-known in the literature, we are the first to provide a complete programming-language semantics based on it. Specifically, we develop the first small-step operational semantics for top-down Datalog, where subqueries occur as nested intermediate terms. The small-step semantics forms the basis of step-into interactions in the debugger. Moreover, we show how step-over interactions can be realized efficiently based on a hybrid Datalog semantics that adds a bottom-up database to our top-down operational semantics. We implemented a debugger for core Datalog following these semantics and explain how to adopt it for debugging the frontend languages of Soufflé and InCA. Our evaluation shows that our hybrid Datalog semantics can be used to debug real-world Datalog programs with realistic workloads.

CCS Concepts: • Theory of computation → Constraint and logic programming; Operational semantics; Logic and databases; • Software and its engineering → Software testing and debugging.

Additional Key Words and Phrases: debugging, Datalog, small-step operational semantics, QSQR

ACM Reference Format:

André Pacak and Sebastian Erdweg. 2023. Interactive Debugging of Datalog Programs. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 248 (October 2023), 28 pages. <https://doi.org/10.1145/3622824>

1 INTRODUCTION

Datalog is a logic programming language that was invented in the 1980s as a recursive query language for databases [Maier *et al.* 2018]. However, this is not how Datalog is being used nowadays. In the last 20 years or so, it has become increasingly popular to use Datalog as a programming language to solve all kinds of problems, from program analysis [Bravenboer and Smaragdakis 2009; Madsen *et al.* 2016; Szabó *et al.* 2021] to distributed computing [Abiteboul *et al.* 2005] and network monitoring [Alvaro *et al.* 2010, 2011]. Usages of Datalog usually have two things in common. First, they process graph-structured data of considerable size, which is well supported by Datalog's fixpoint semantics. Second, they involve complex Datalog programs, consisting of many inter-dependent rules. For example, even the simplest analysis from the Doop framework for Java bytecode consists of 560 lines of Datalog code that query 32 relations from the database and compute 78 inter-dependent derived relations [Bravenboer and Smaragdakis 2009]. Sophisticated Doop analyses are an order of magnitude larger still.

Unfortunately, development tools for Datalog have not been able to keep up with the growing size and complexity of Datalog programs. Specifically, there is no interactive debugging support available for Datalog to date. Interactive debuggers assist developers in program understanding

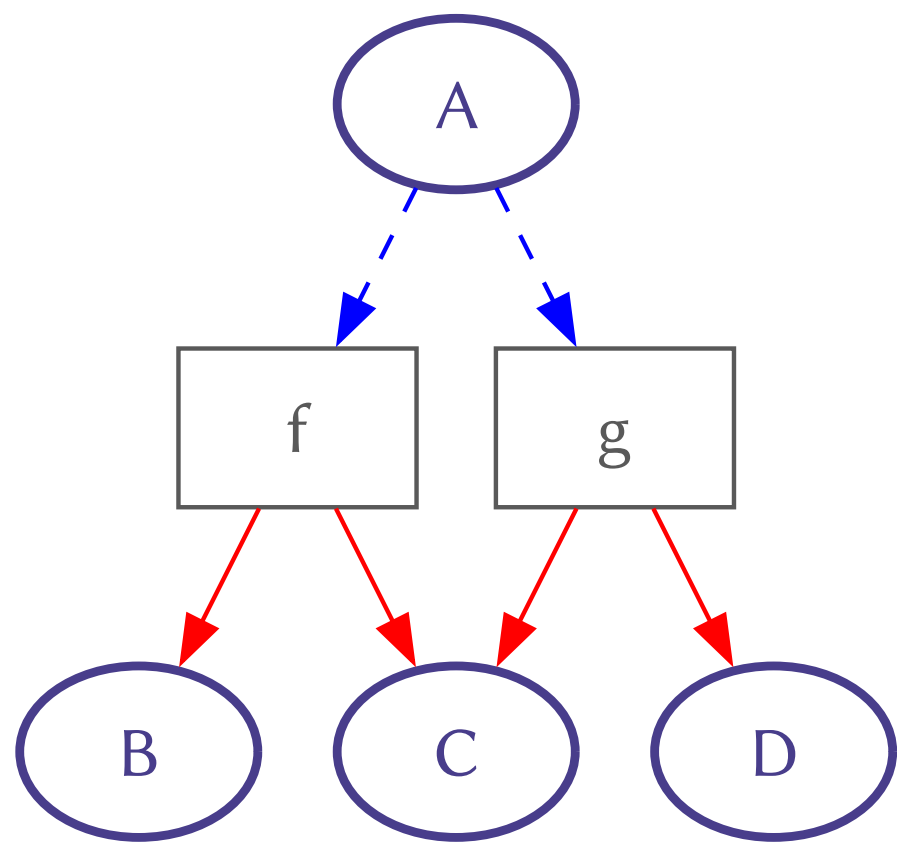
AND-OR graphs are not just graphs

AND-OR graphs are not just graphs

Observation: Hidden inside is a proof system that can be poked and prodded

AND-OR graphs are not just graphs

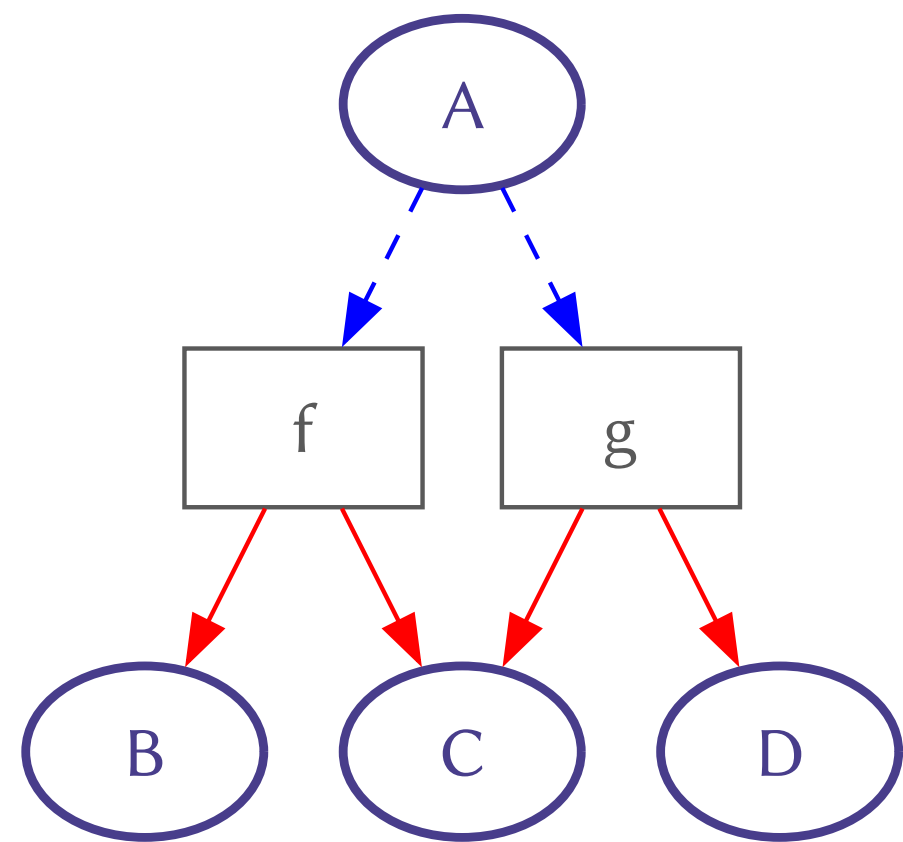
Observation: Hidden inside is a proof system that can be poked and prodded



AND-OR graphs are not just graphs

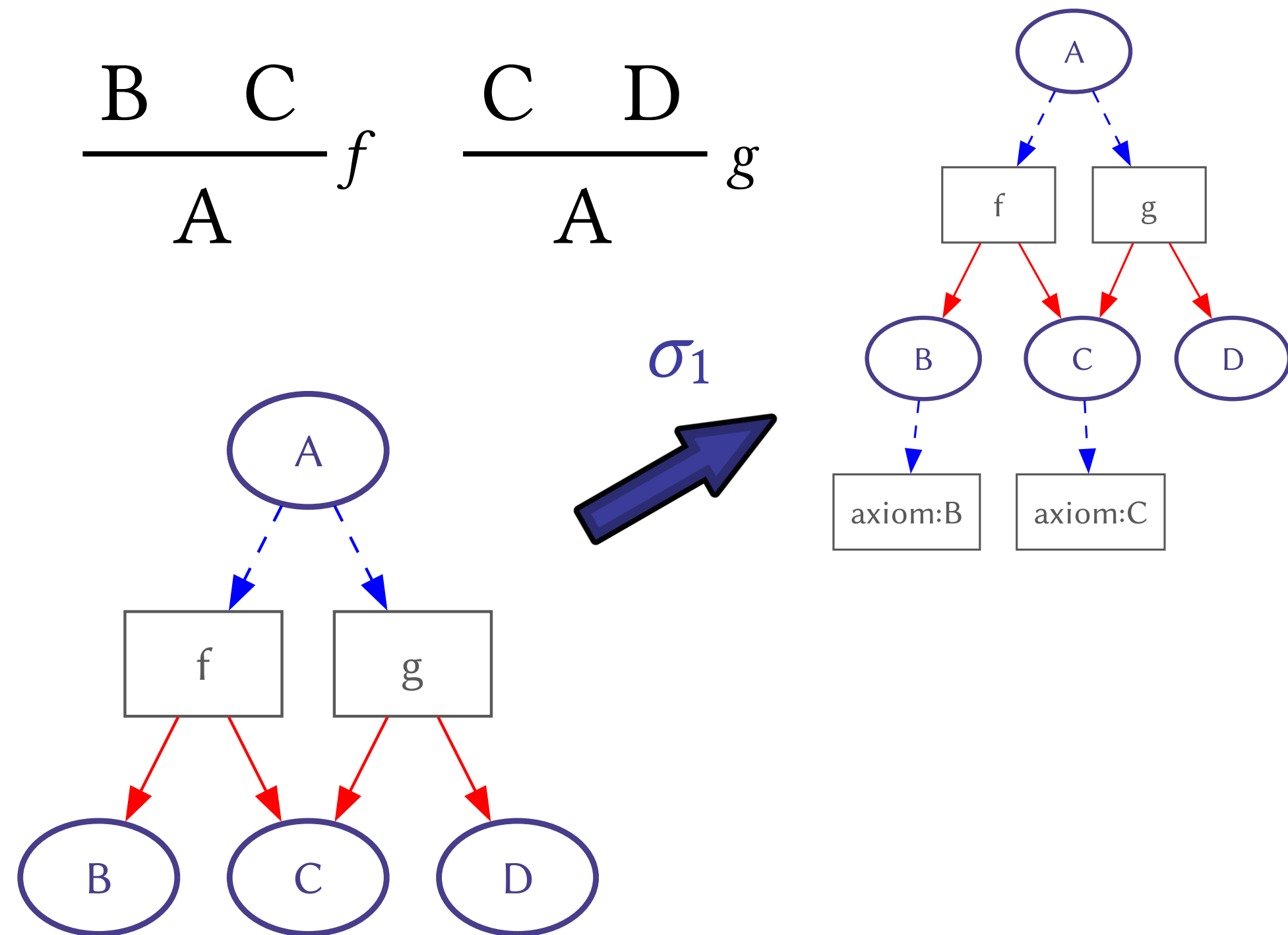
Observation: Hidden inside is a proof system that can be poked and prodded

$$\frac{B \quad C}{A} f \quad \frac{C \quad D}{A} g$$



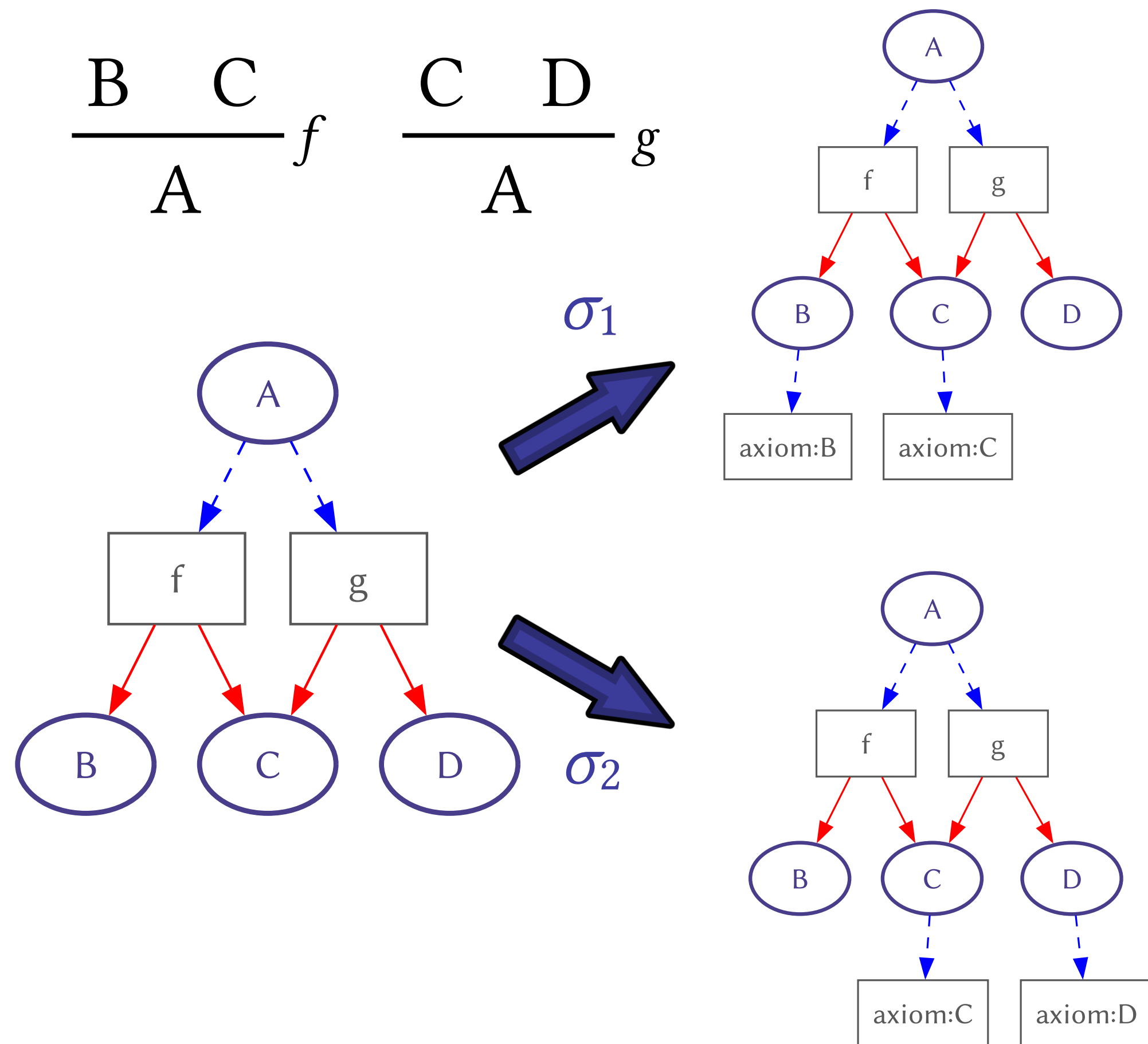
AND-OR graphs are not just graphs

Observation: Hidden inside is a proof system that can be poked and prodded



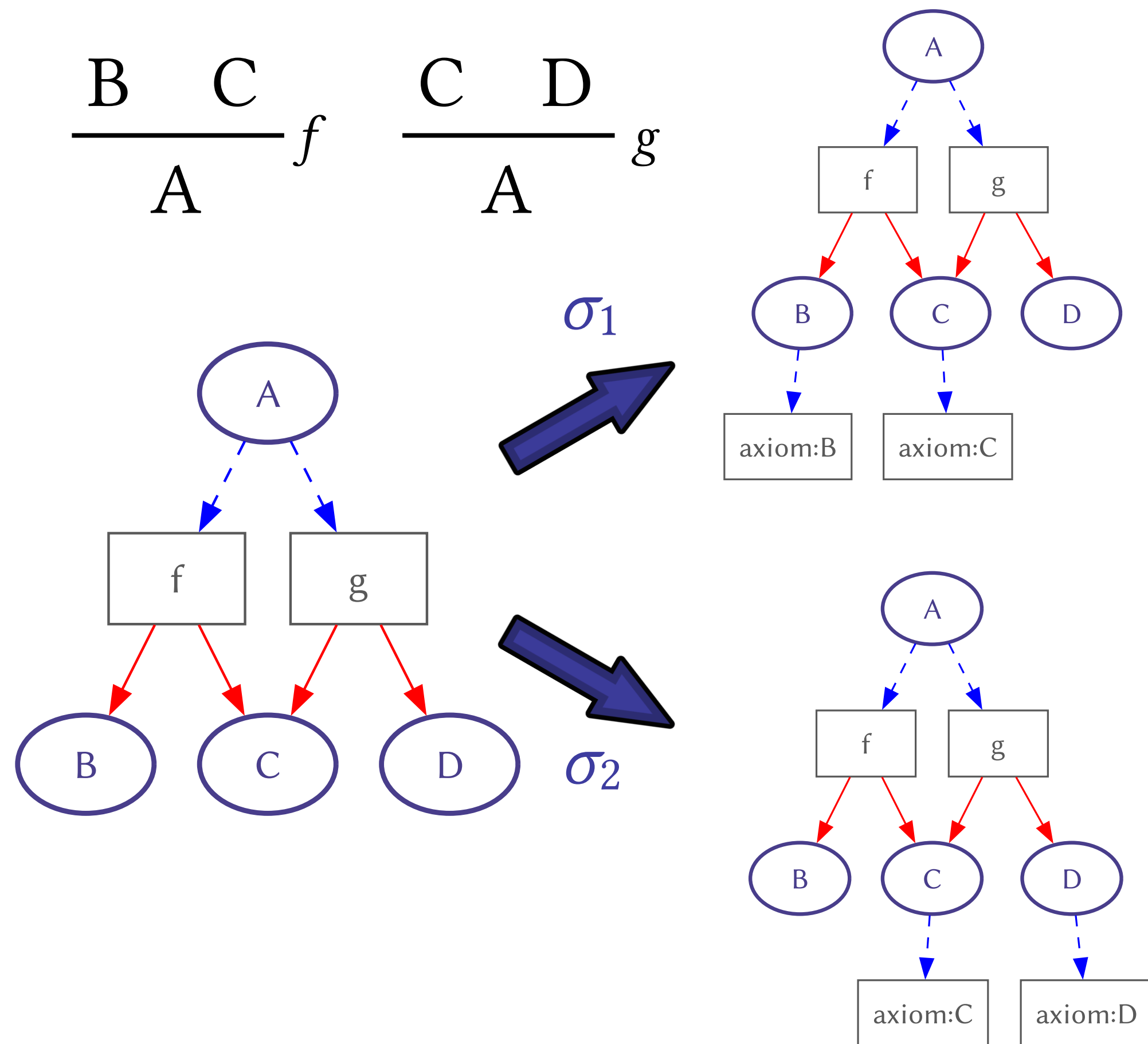
AND-OR graphs are not just graphs

Observation: Hidden inside is a proof system that can be poked and prodded



AND-OR graphs are not just graphs

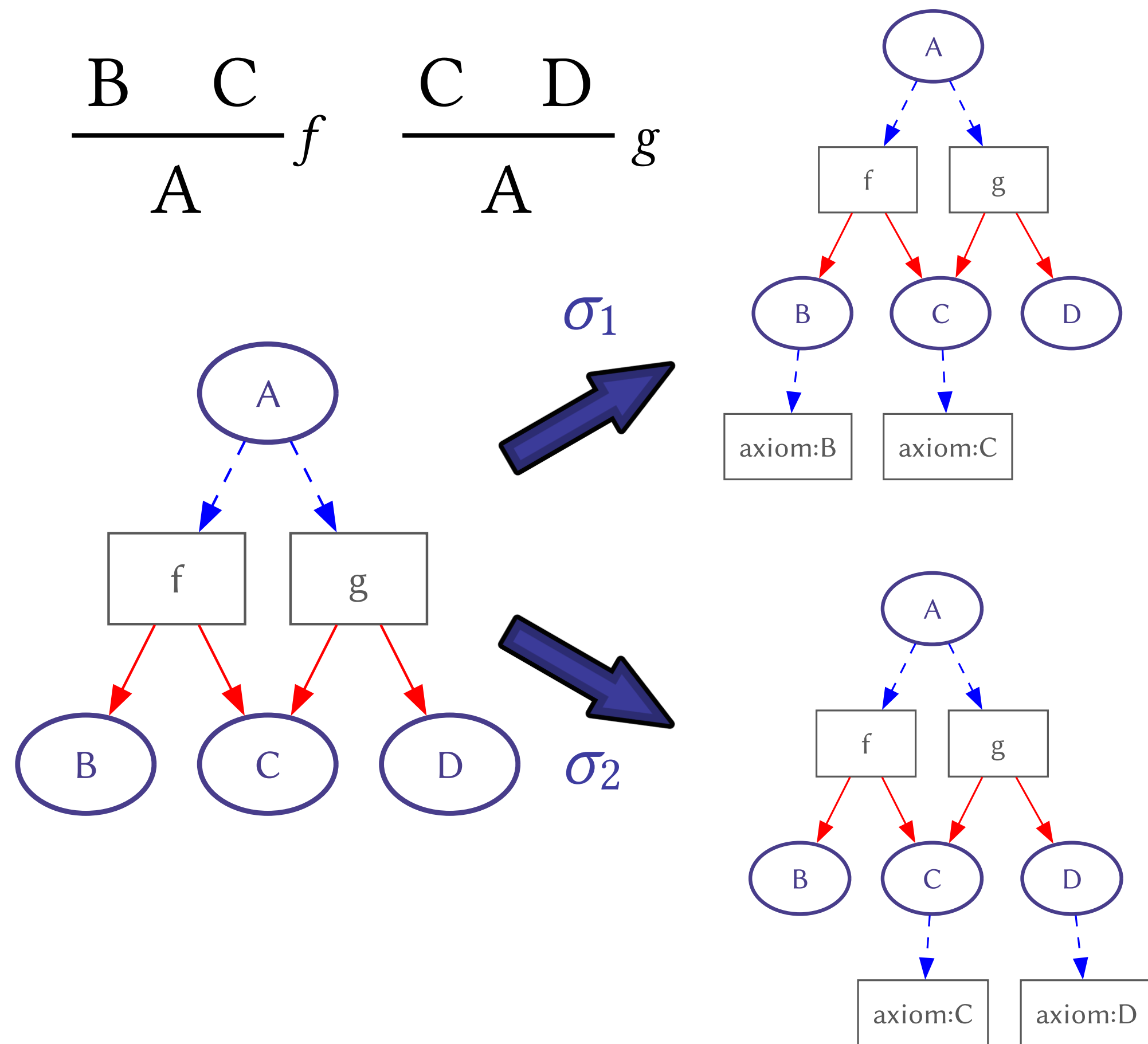
Observation: Hidden inside is a proof system that can be poked and prodded



Key insight:

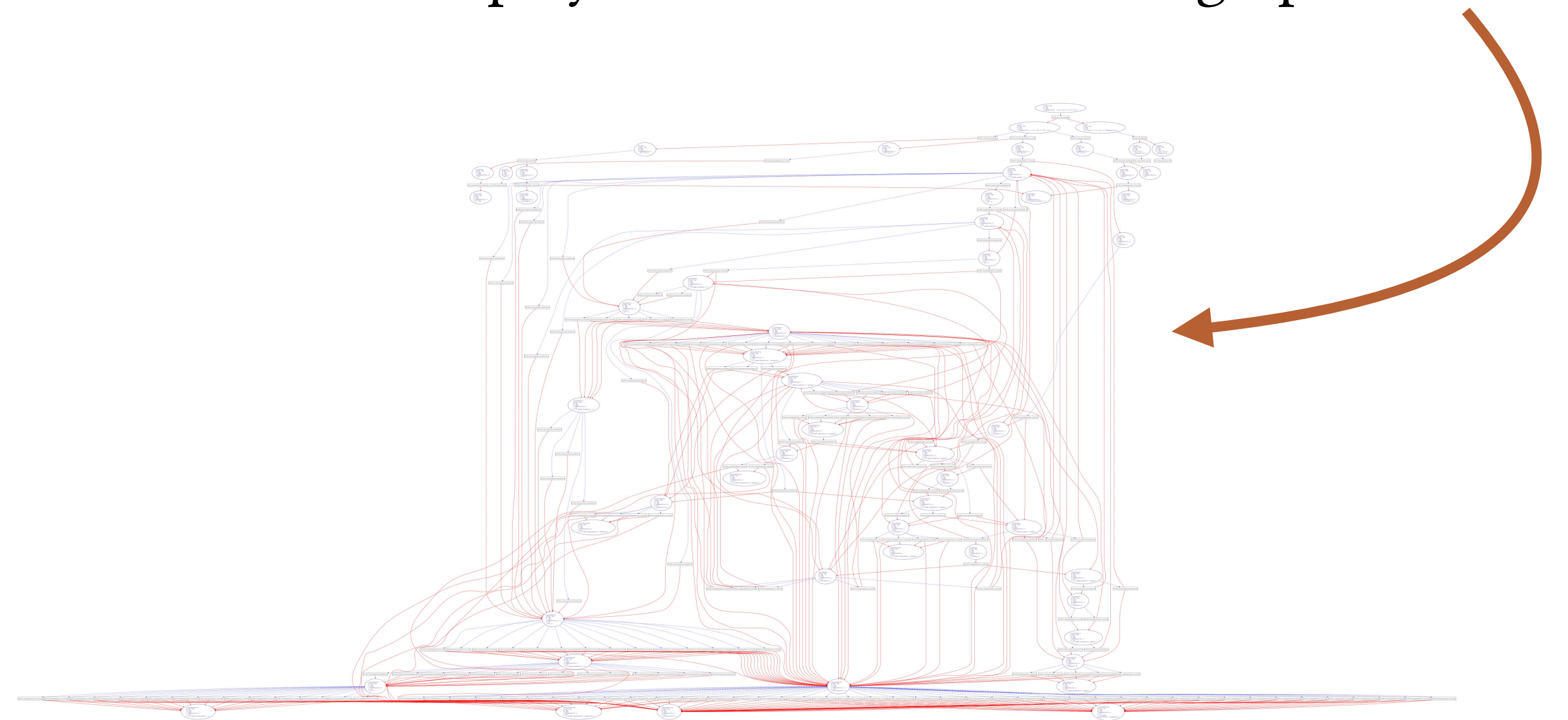
AND-OR graphs are not just graphs

Observation: Hidden inside is a proof system that can be poked and prodded



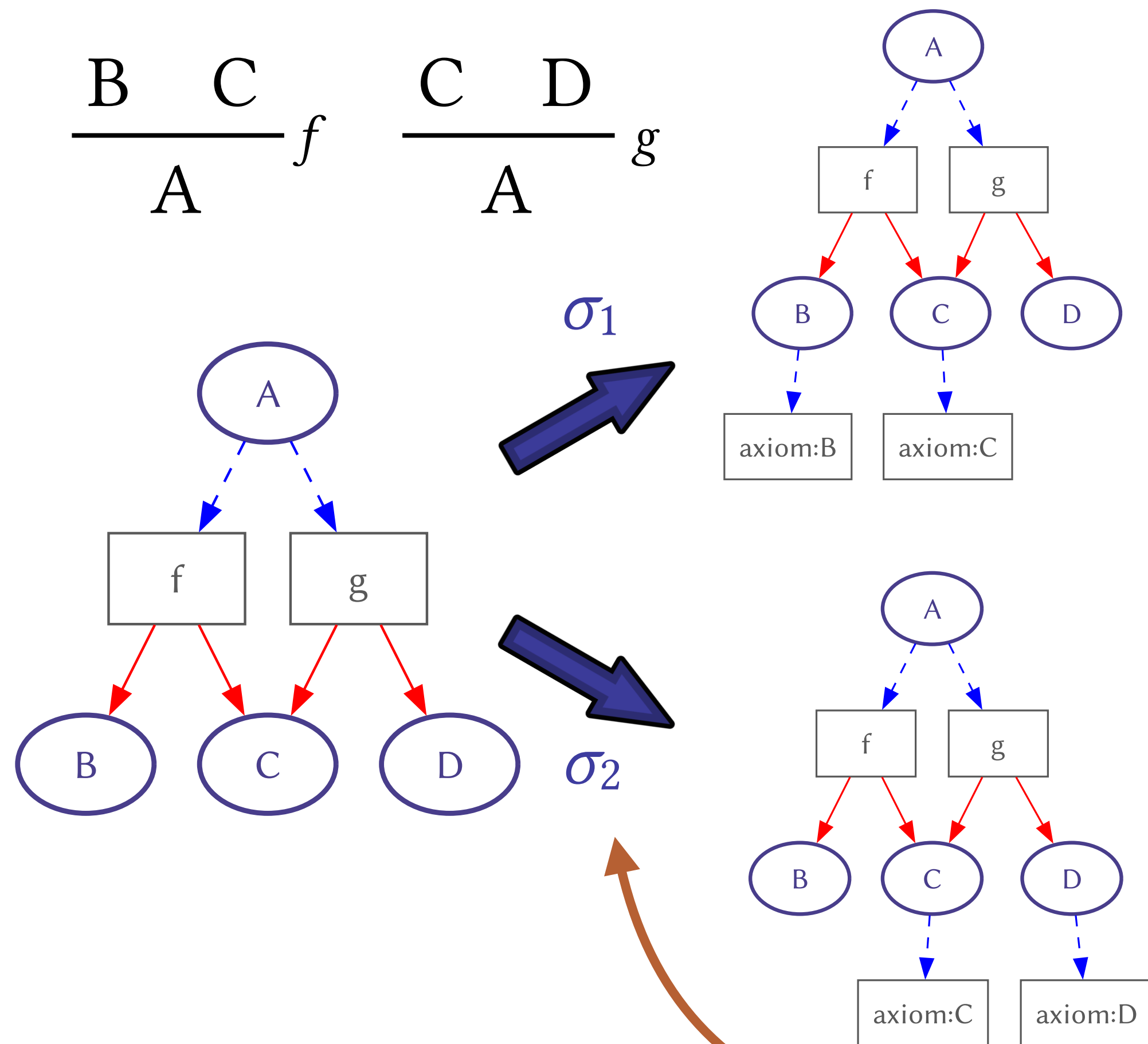
Key insight:

We don't want to display/traverse the AND-OR graph...



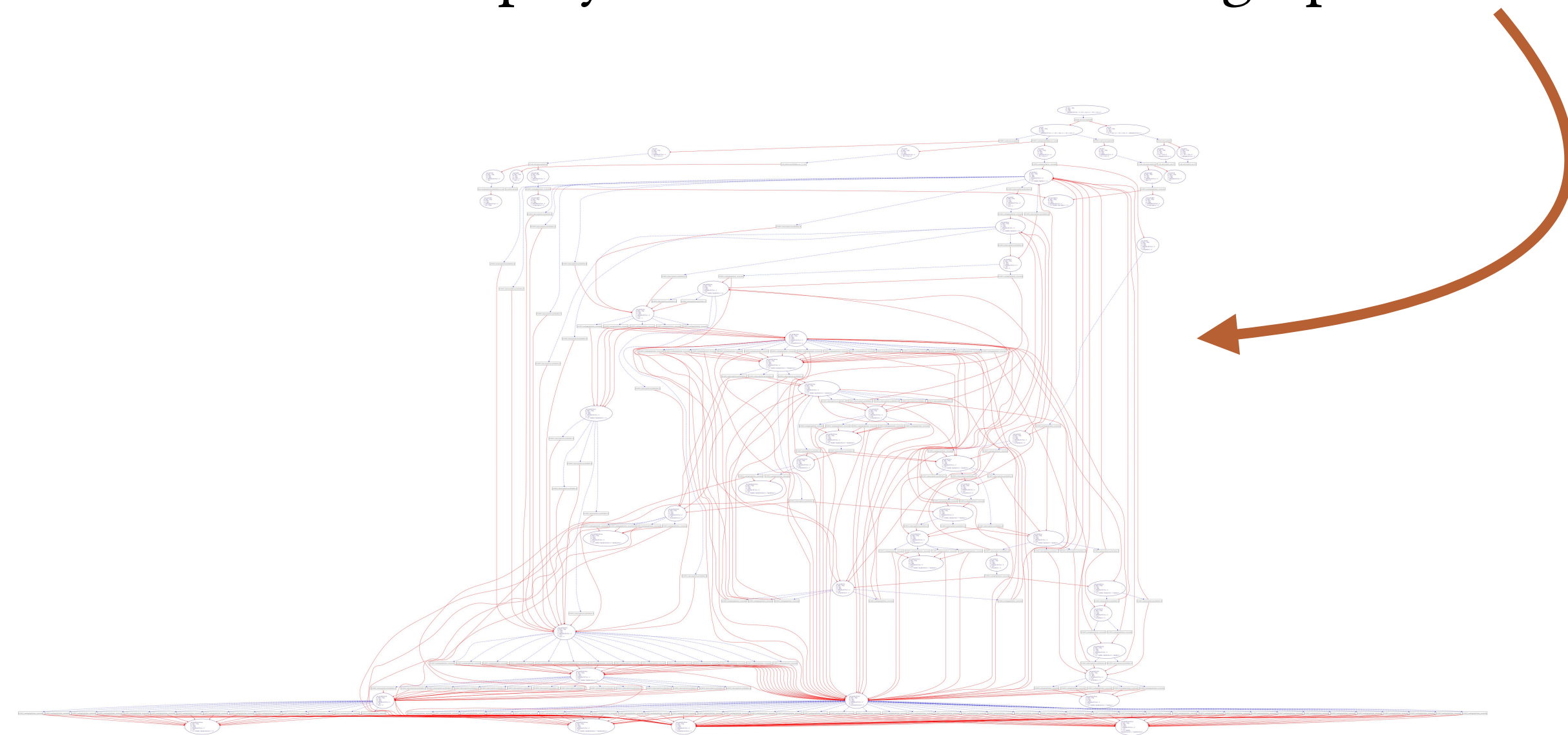
AND-OR graphs are not just graphs

Observation: Hidden inside is a proof system that can be poked and prodded



Key insight:

We don't want to display/traverse the AND-OR graph...



...we want to navigate the **space of modifications** to the graph

Concrete example: Debugging Aesop for Lean

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Concrete example: Debugging Aesop for Lean

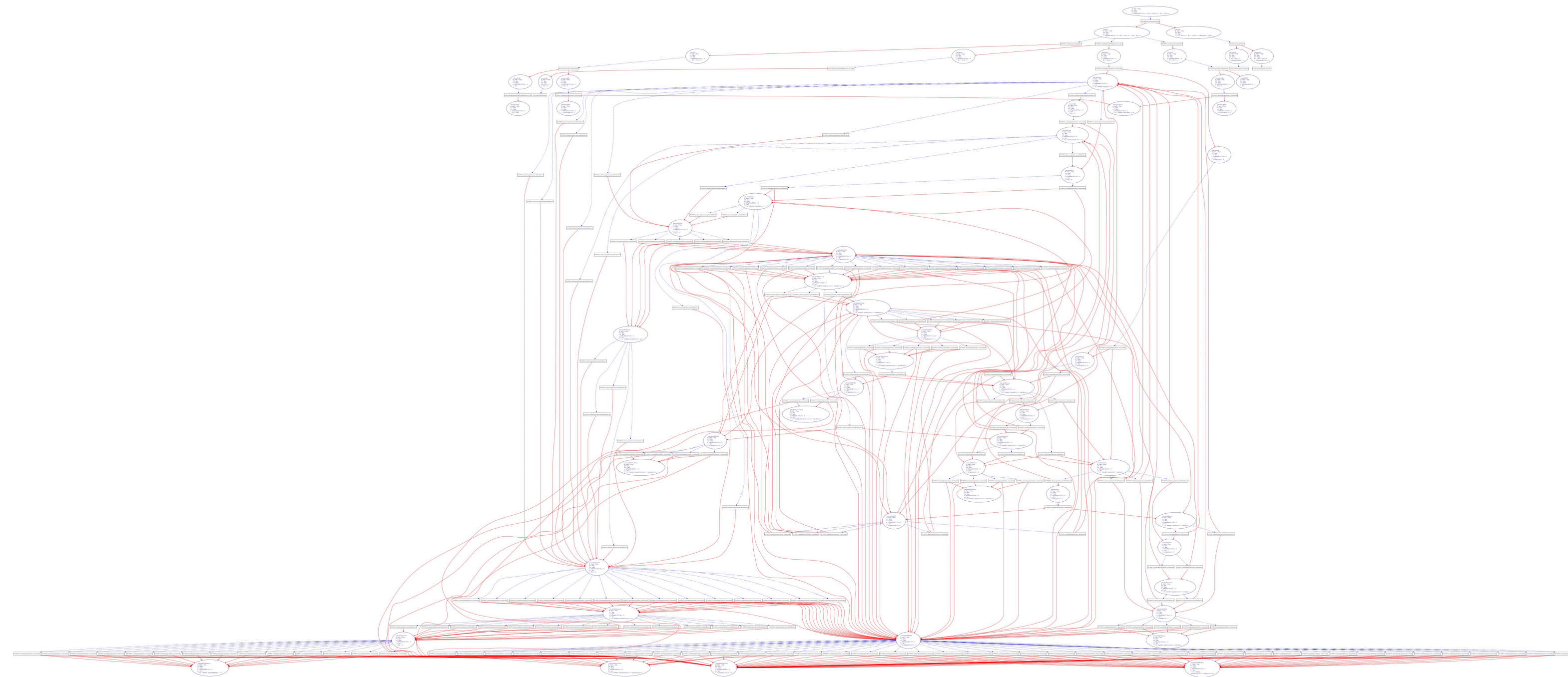
Modifications allowed: labeling/assuming possible lemmas

Goal: $\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

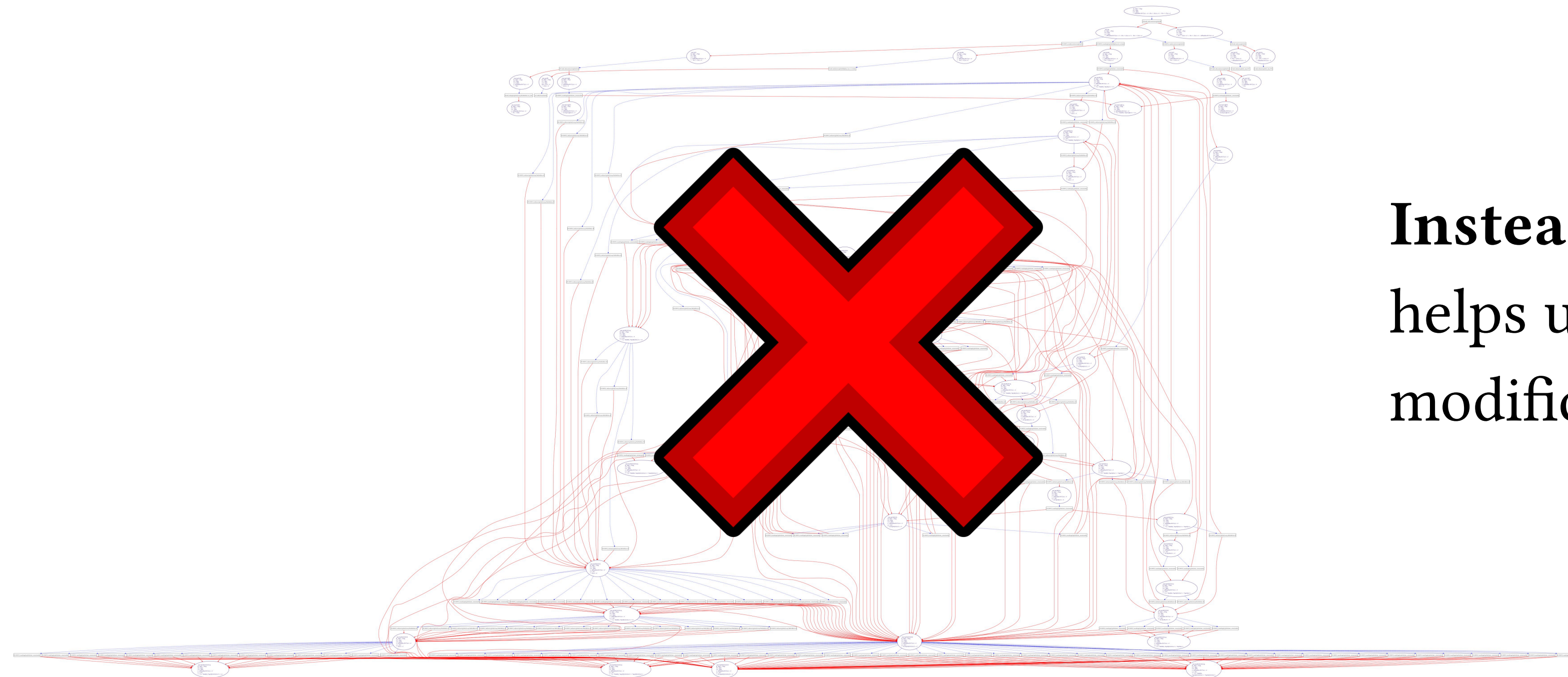
Goal: $\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$



Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal: $\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$



Instead: System helps user navigate modifications

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal: $\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal: $\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \leftrightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$$

**System
selects:**

$$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \rightarrow B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$$

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove
manually”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove
manually”

“Not willing to prove
manually, but should be true”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove
manually”

“Not willing to prove
manually, but should be true”

above + “must be used”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove
manually”

“Not willing to prove
manually, but should be true”

above + “must be used”

above + “must be used”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove
manually”

“Not willing to prove
manually, but should be true”

above + “must be used”

above + “must be used”

“I don’t know”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\leftrightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \boxed{\rightarrow} B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be false”

above + “must be used”

above + “must be used”

“I don’t know”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

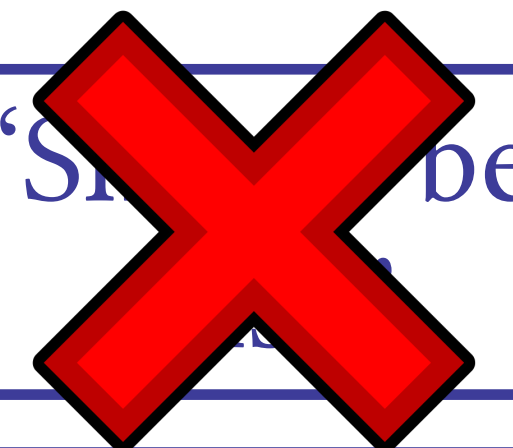
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be true”



above + “must be used”

above + “must be used”

“I don't know”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

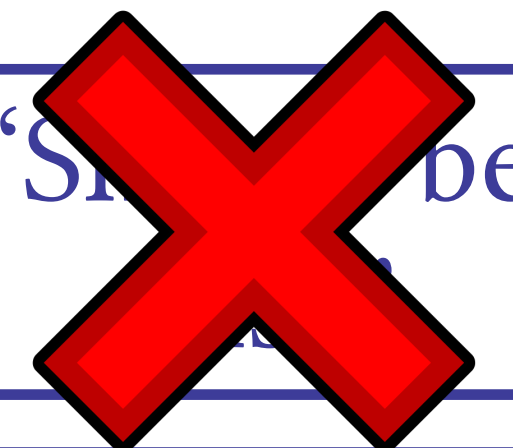
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be true”



above + “must be used”

above + “must be used”

“I don't know”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

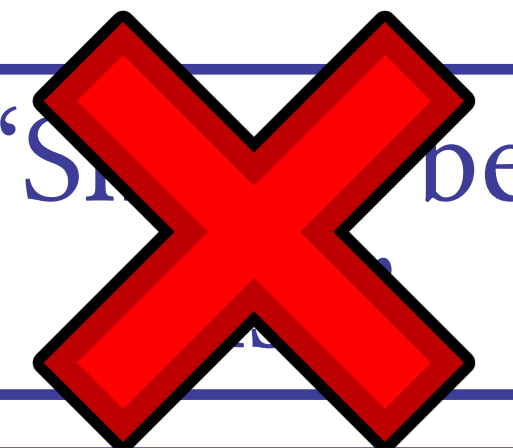
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be true”



above + “must be used”

above + “must be used”

“I don’t know”

**System
selects:**

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

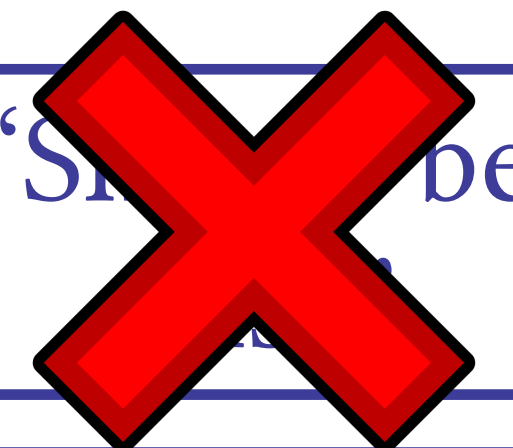
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be true”



above + “must be used”

above + “must be used”

“I don’t know”

**System
selects:**

$\dots, (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \vdash B(\sigma) \rightarrow (S, \sigma) \Downarrow \sigma'$

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

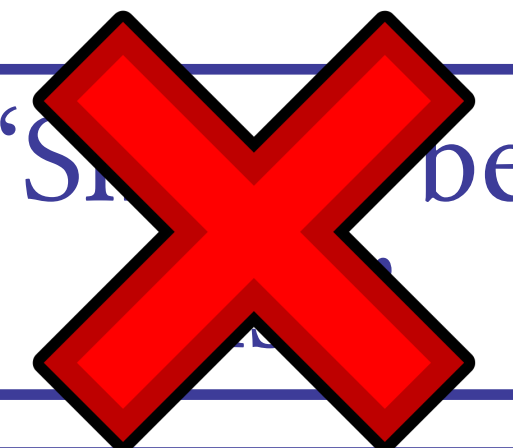
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be true”



above + “must be used”

above + “must be used”

“I don't know”

**System
selects:**

$\dots, (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \vdash B(\sigma) \rightarrow (S, \sigma) \Downarrow \sigma'$

**User can
respond:**

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**System
selects:**

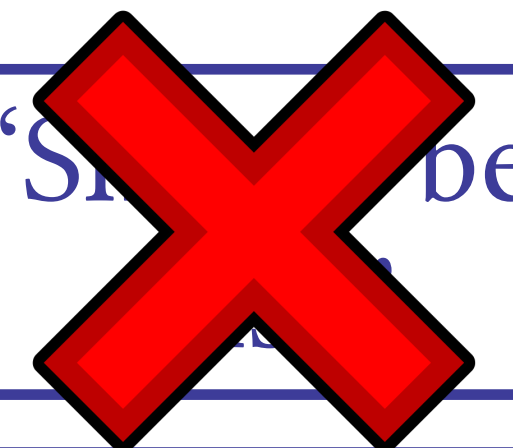
$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be false”



above + “must be used”

above + “must be used”

“I don't know”

**System
selects:**

$\dots, (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \vdash B(\sigma) \rightarrow (S, \sigma) \Downarrow \sigma'$

**User can
respond:**

above options

+

“Should be false”

Concrete example: Debugging Aesop for Lean

Modifications allowed: labeling/assuming possible lemmas

Goal:

$$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \iff B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$$

**System
selects:**

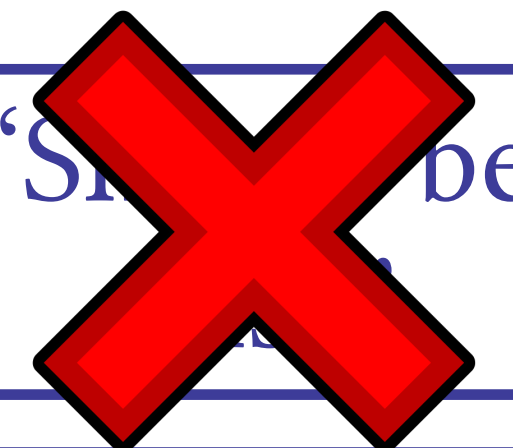
$$\dots \vdash (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \implies B(\sigma) \wedge (S, \sigma) \Downarrow \sigma' \vee \neg B(\sigma) \wedge (T, \sigma) \Downarrow \sigma'$$

**User can
respond:**

“Willing to prove manually”

“Not willing to prove manually, but should be true”

“Should be false”



above + “must be used”

above + “must be used”

“I don’t know”

**System
selects:**

$$\dots, (\text{ifThenElse } B \ S \ T, \sigma) \Downarrow \sigma' \vdash B(\sigma) \rightarrow (S, \sigma) \Downarrow \sigma'$$

**User can
respond:**

above options

+

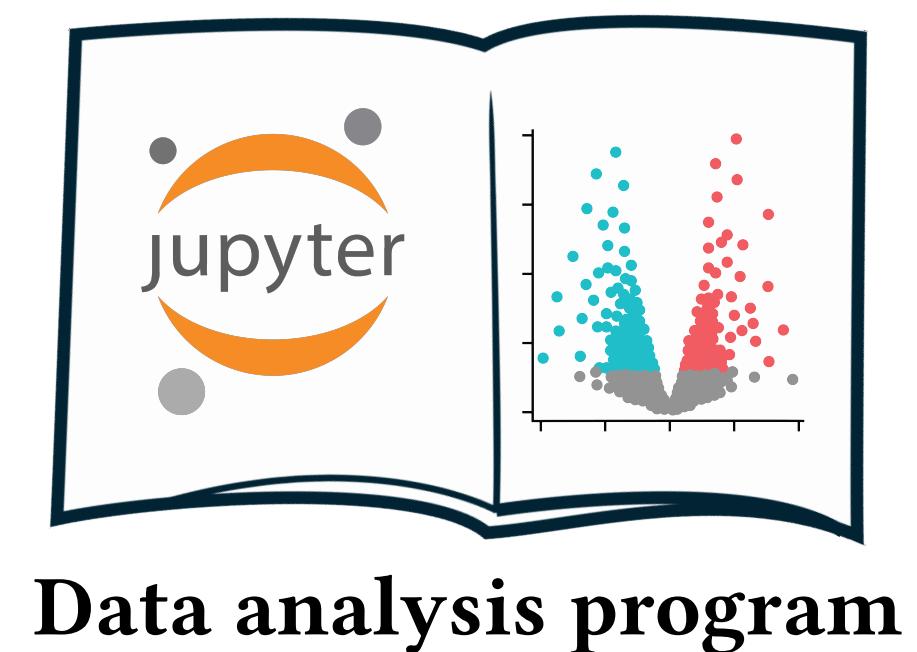
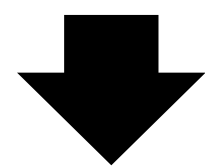
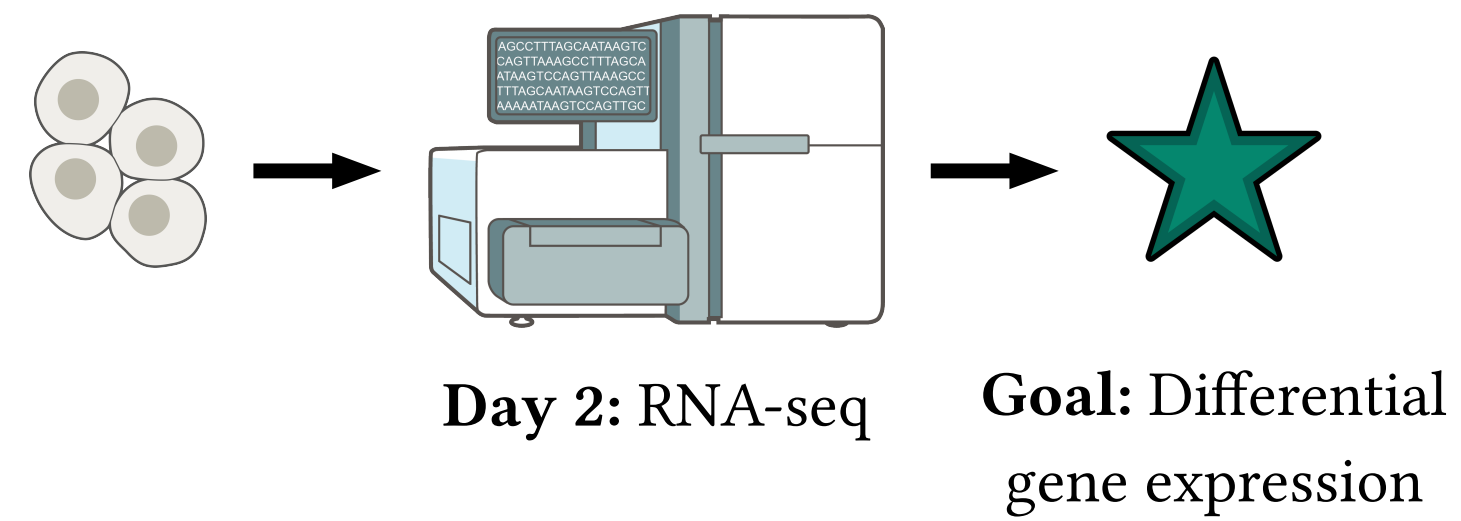
“Should be false”

“Willing to prove manually”

Programming by Navigation used for...

Programming by Navigation used for...

**Synthesizing bioinformatics analyses
from wet lab experimental workflows**

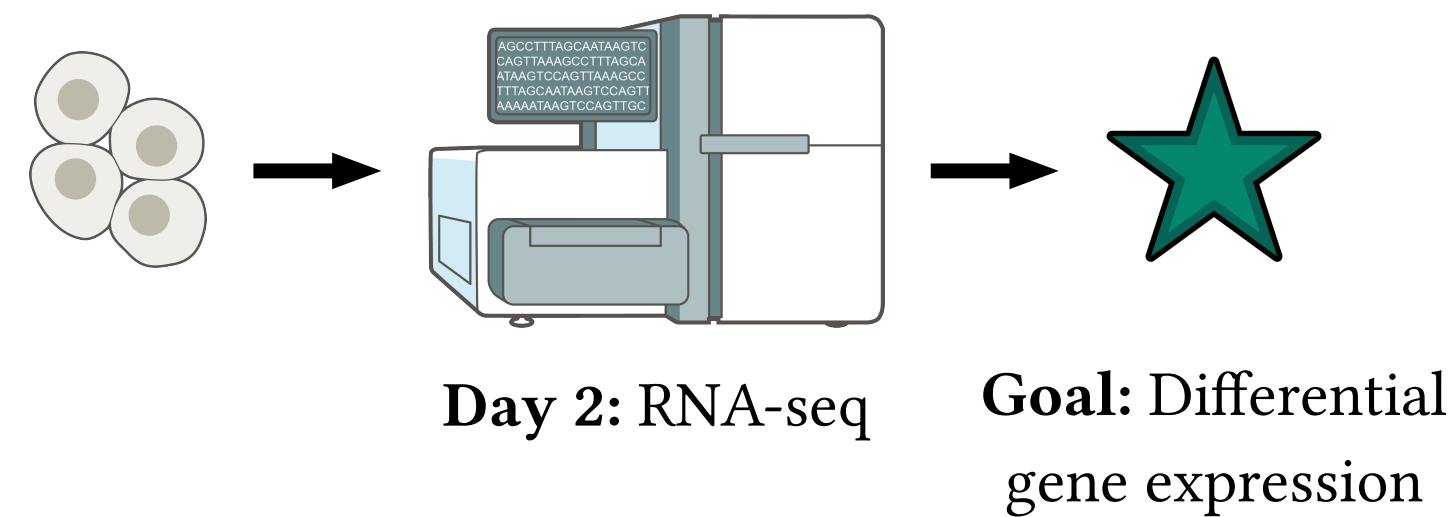


Programming by Navigation used for...

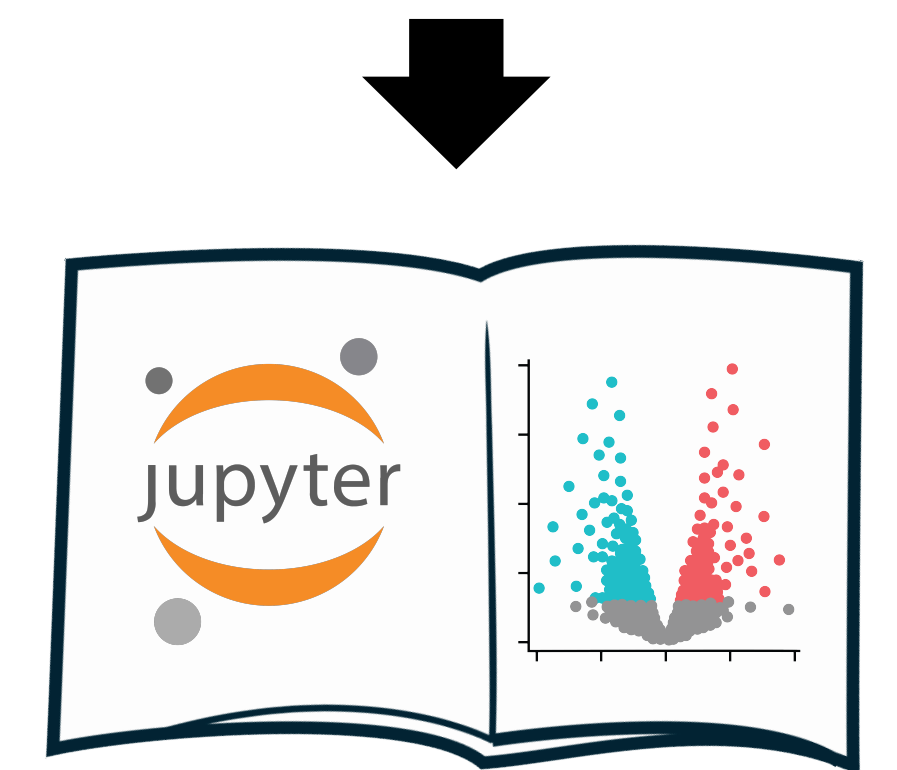
Synthesizing bioinformatics analyses from wet lab experimental workflows

Creating user schedules for a sparse tensor algebra compiler

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \quad \begin{array}{l} \text{Parallelization} \\ \text{factor} \end{array} \begin{array}{|c|} \hline ? \\ \hline \end{array} \begin{array}{|c|} \hline ? \\ \hline \end{array} \begin{array}{l} \text{Memory} \\ \text{layout} \end{array} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \quad \begin{array}{|c|} \hline ? \\ \hline \end{array} \begin{array}{l} \text{Vector} \\ \text{width} \end{array} \\ \begin{array}{|c|} \hline ? \\ \hline \end{array} \begin{array}{l} \text{Dataflow} \\ \text{ordering} \end{array} \quad \begin{array}{|c|} \hline ? \\ \hline \end{array} \begin{array}{l} \text{Memory} \\ \text{layout} \end{array} \end{array} \right\}$$

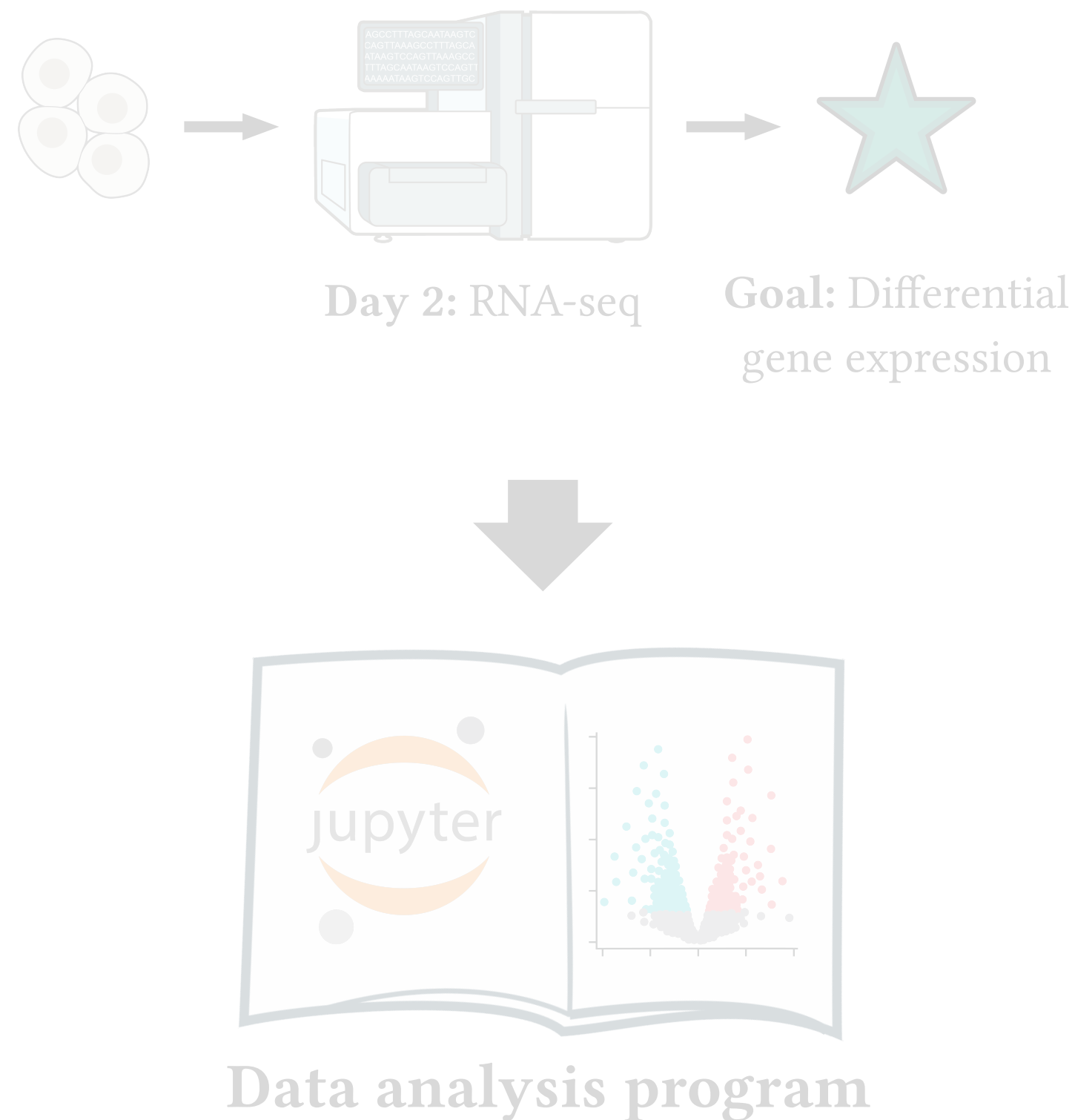


Data analysis program

```
{ "dataflow-ordering": "i0 i1 i2 i4 i8 i5 i9 i6 i7 i3", "sparse-annotations": { "t0": { "format": "CSR", "indices": { "i": true, "j": false } }, "t1": { "format": "CSR", "indices": { "i": true, "j": false } }, "t2": { "format": "CSR", "indices": { "i": true, "j": false } }, "t3": { "format": "CSR", "indices": { "i": true, "j": false } }, "t4": { "format": "custom", "indices": { "i": false, "j": . . .
```

Programming by Navigation used for...

Synthesizing bioinformatics analyses from wet lab experimental workflows



Creating user schedules for a sparse tensor algebra compiler

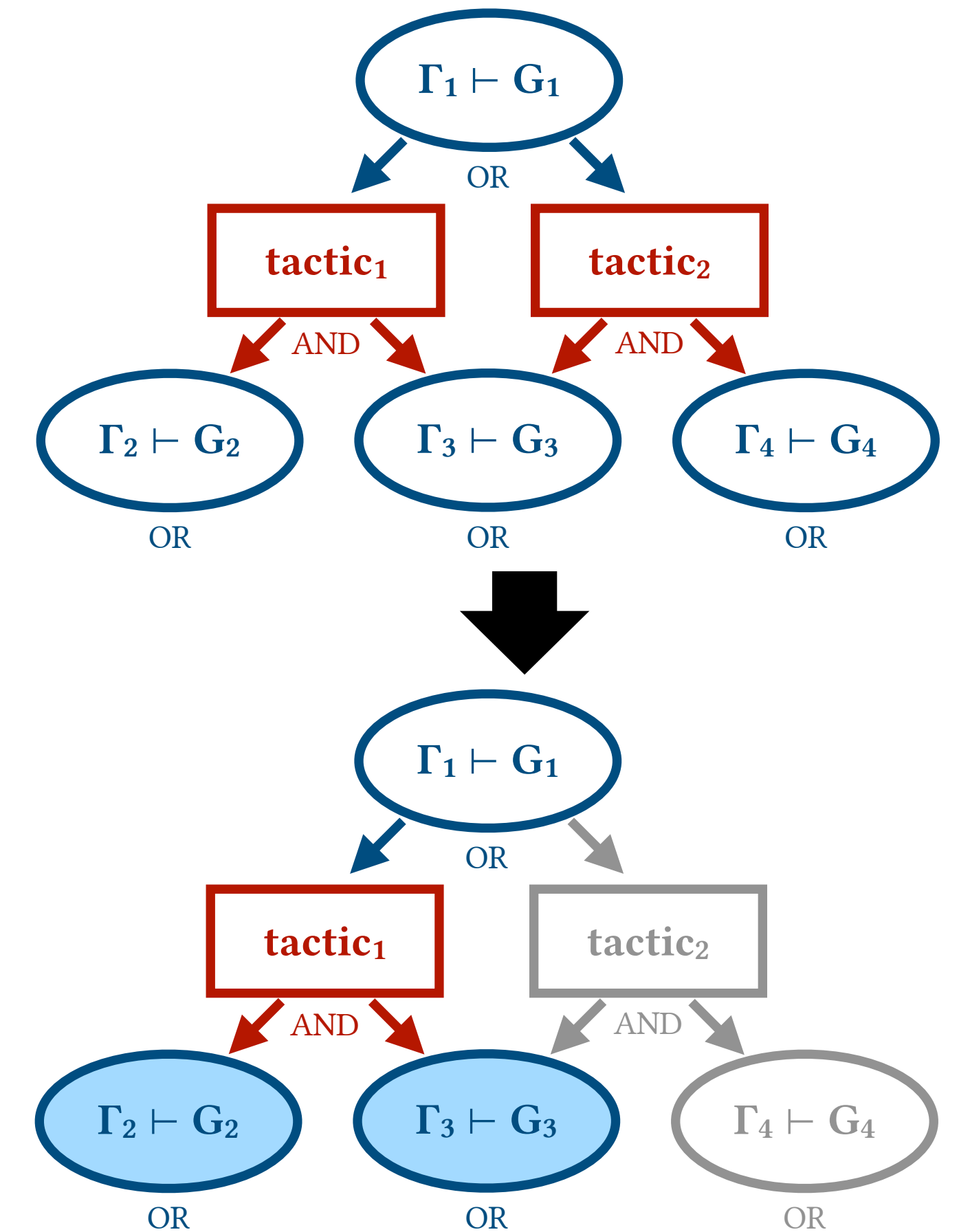
Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).

$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right.$$

Annotations: Parallelization factor, Memory layout, Vector width, Dataflow ordering.

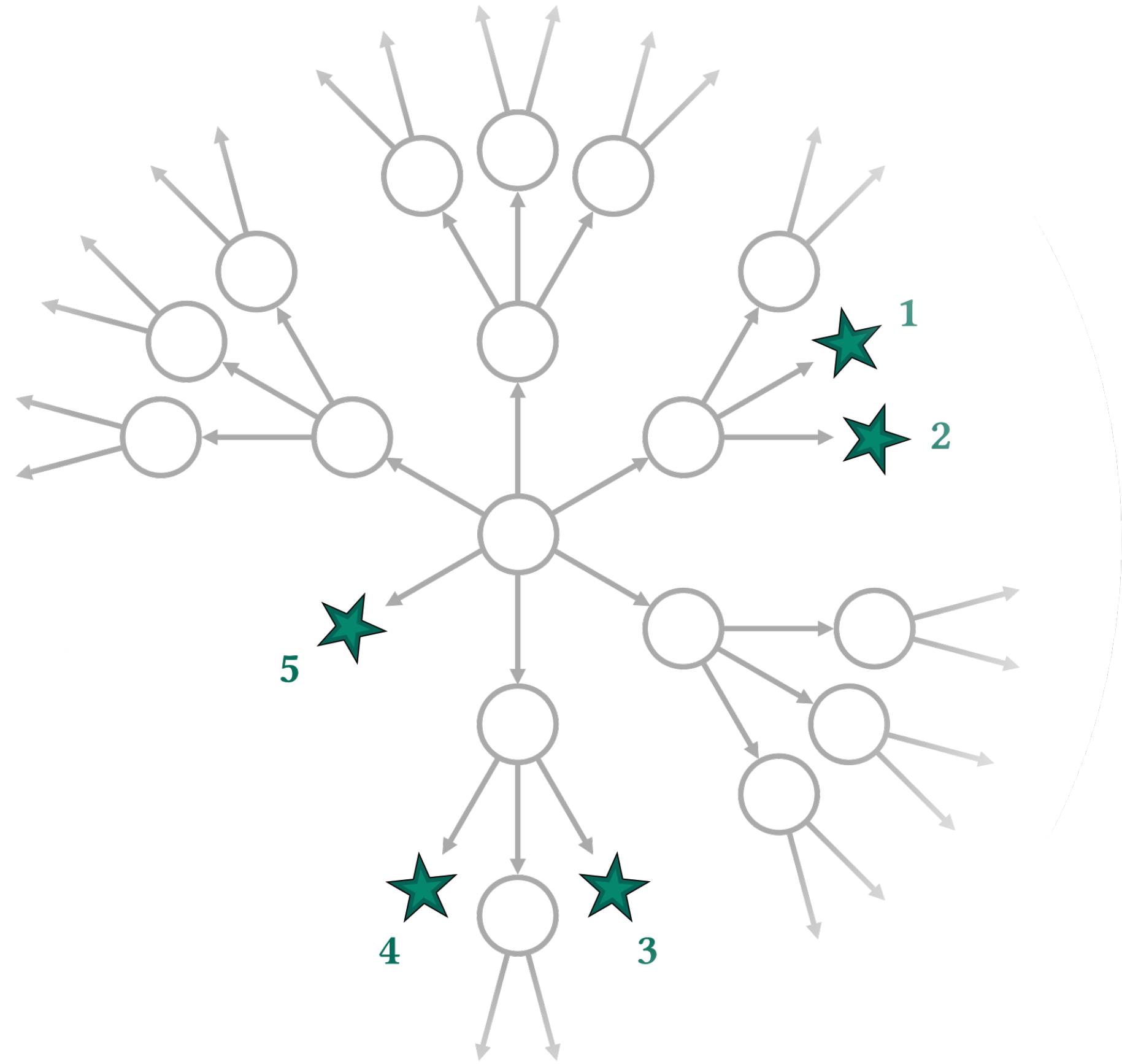
```
{ "dataflow-ordering": "i0 i1 i2 i4 i8 i5 i9 i6 i7 i3", "sparse-annotations": { "t0": { "format": "CSR", "indices": { "i": true, "j": false } }, "t1": { "format": "CSR", "indices": { "i": true, "j": false } }, "t2": { "format": "CSR", "indices": { "i": true, "j": false } }, "t3": { "format": "CSR", "indices": { "i": true, "j": false } }, "t4": { "format": "custom", "indices": { "i": false, "j": ...
```

This talk: Debugging failing proof search (e.g., Aesop for Lean, Rust traits, Datalog, ...)

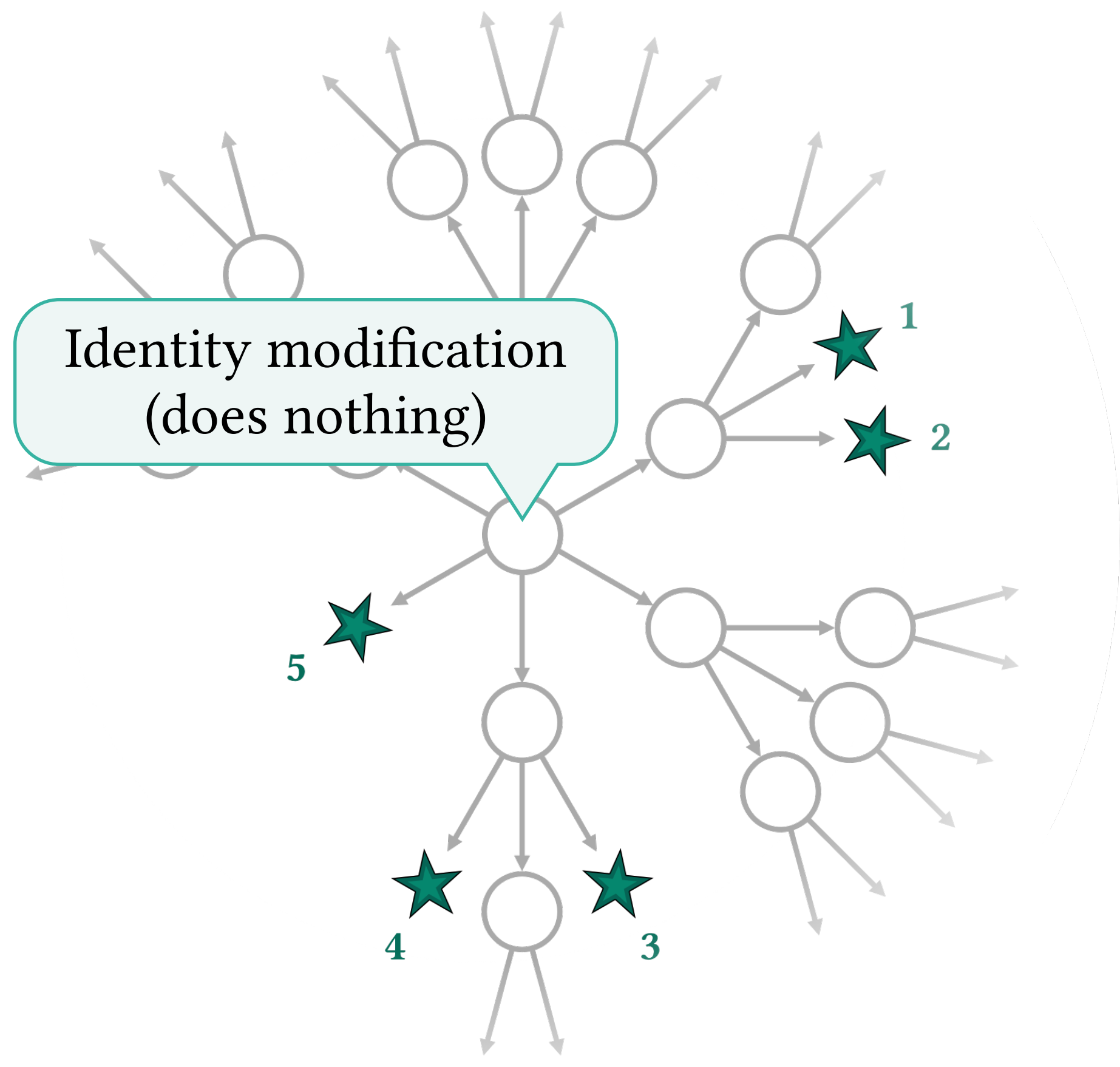


Programming by Navigation

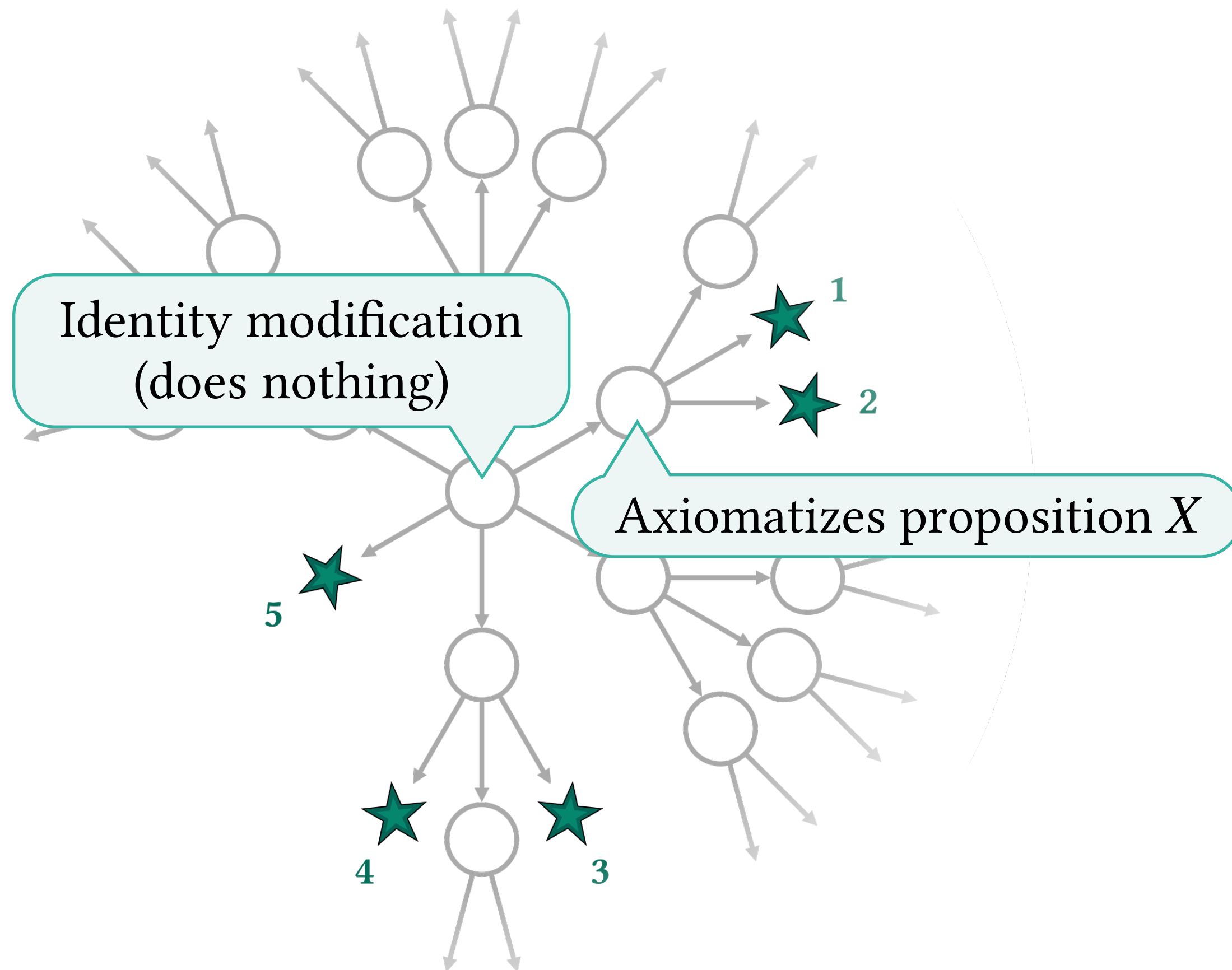
Programming by Navigation



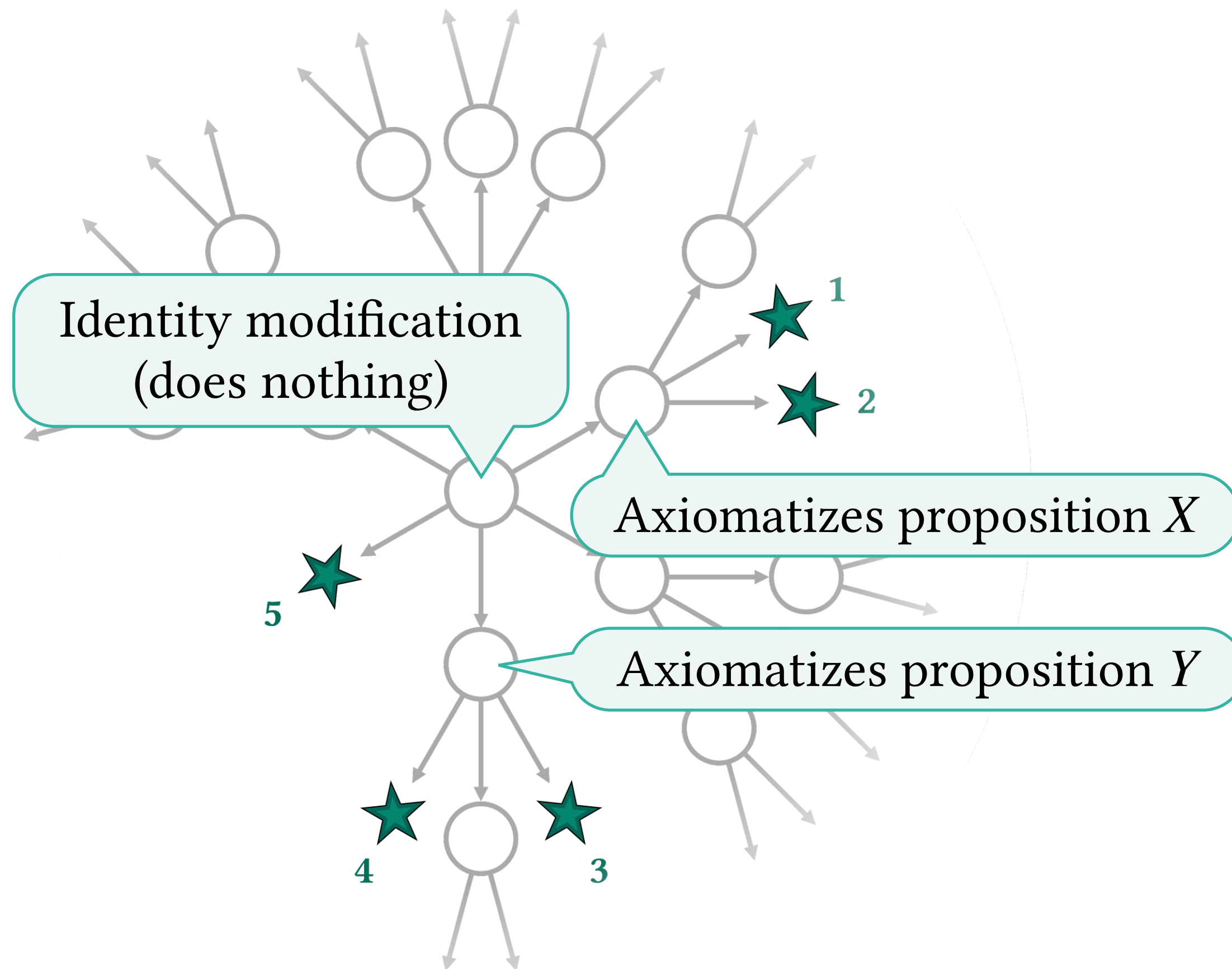
Programming by Navigation



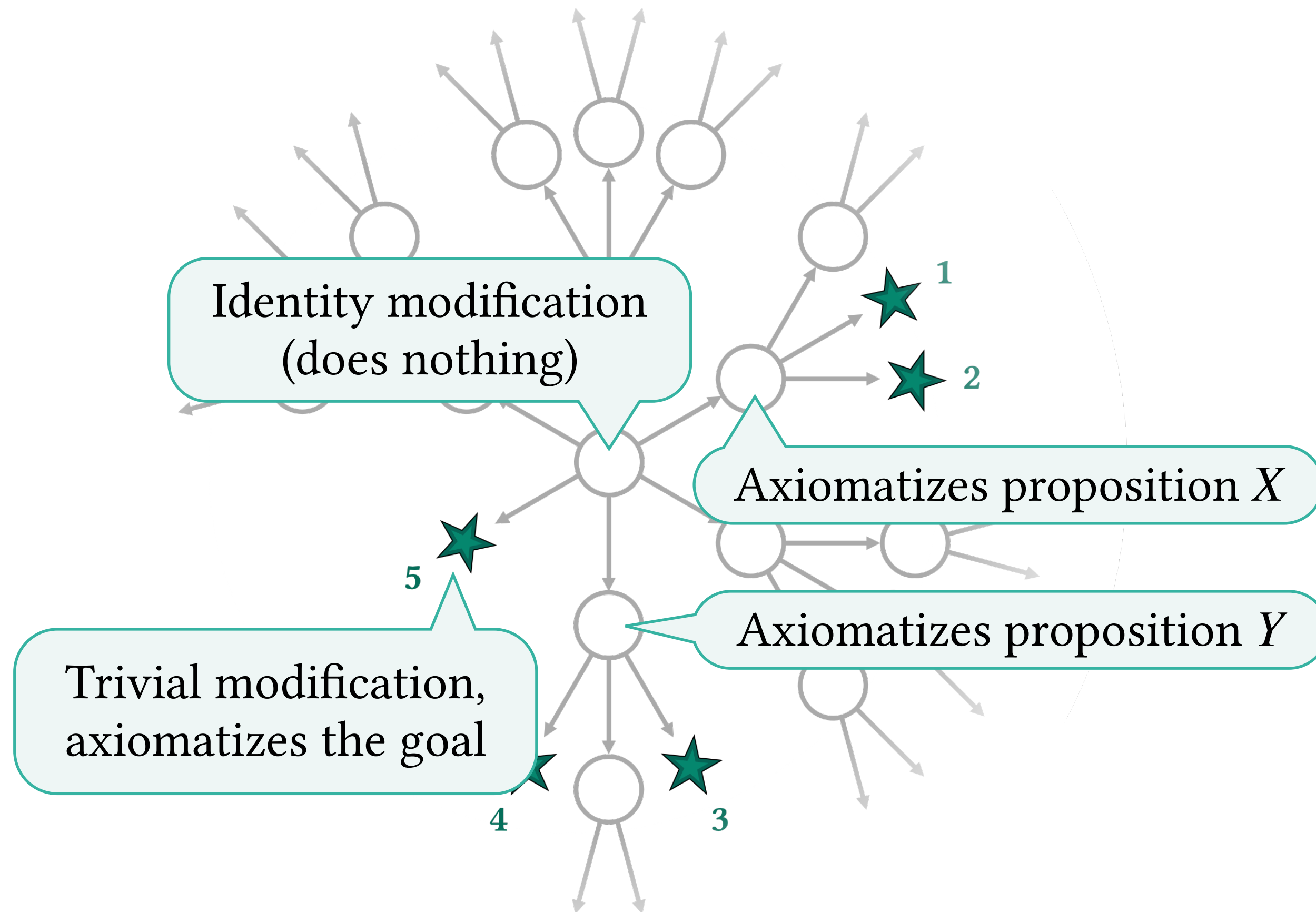
Programming by Navigation



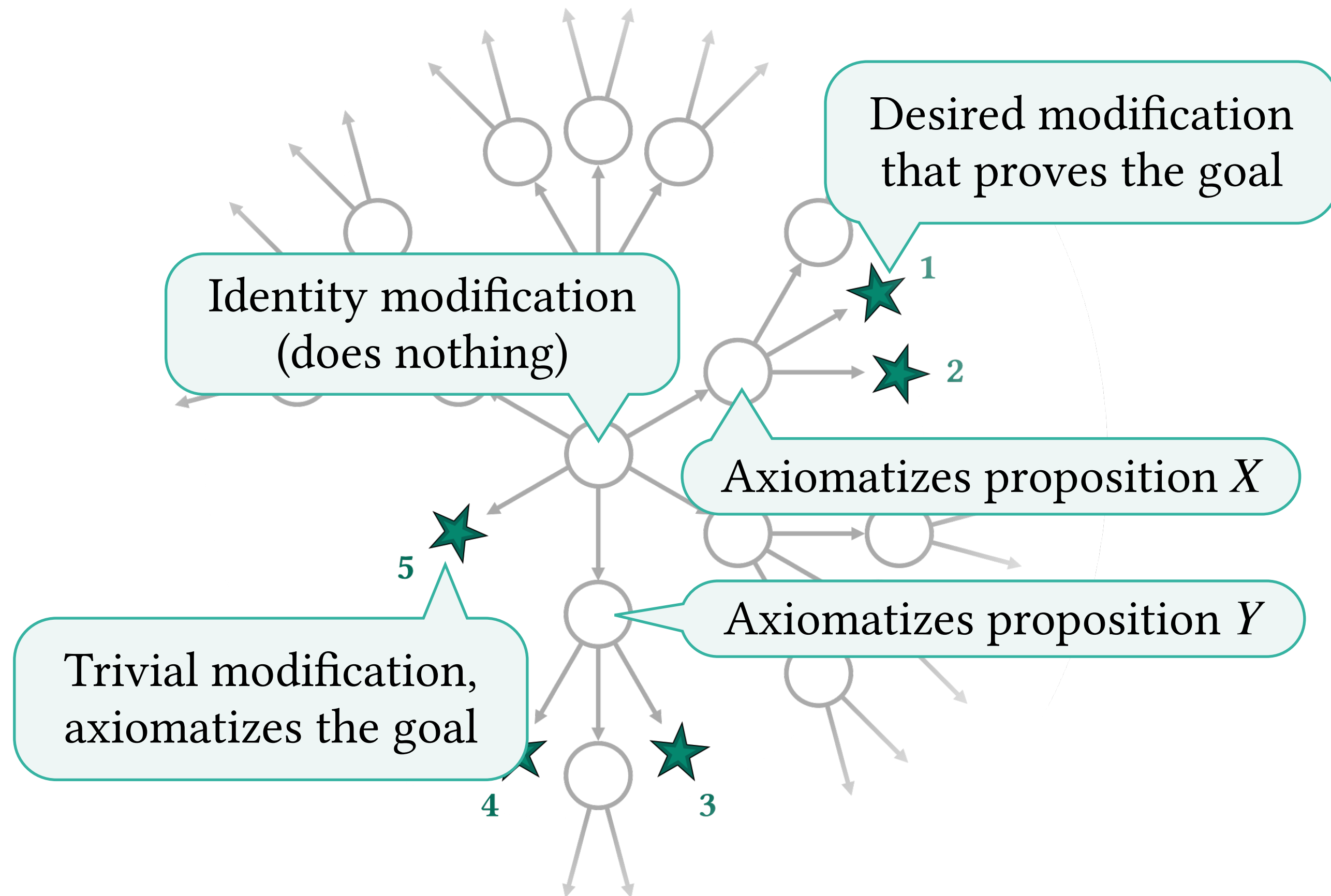
Programming by Navigation



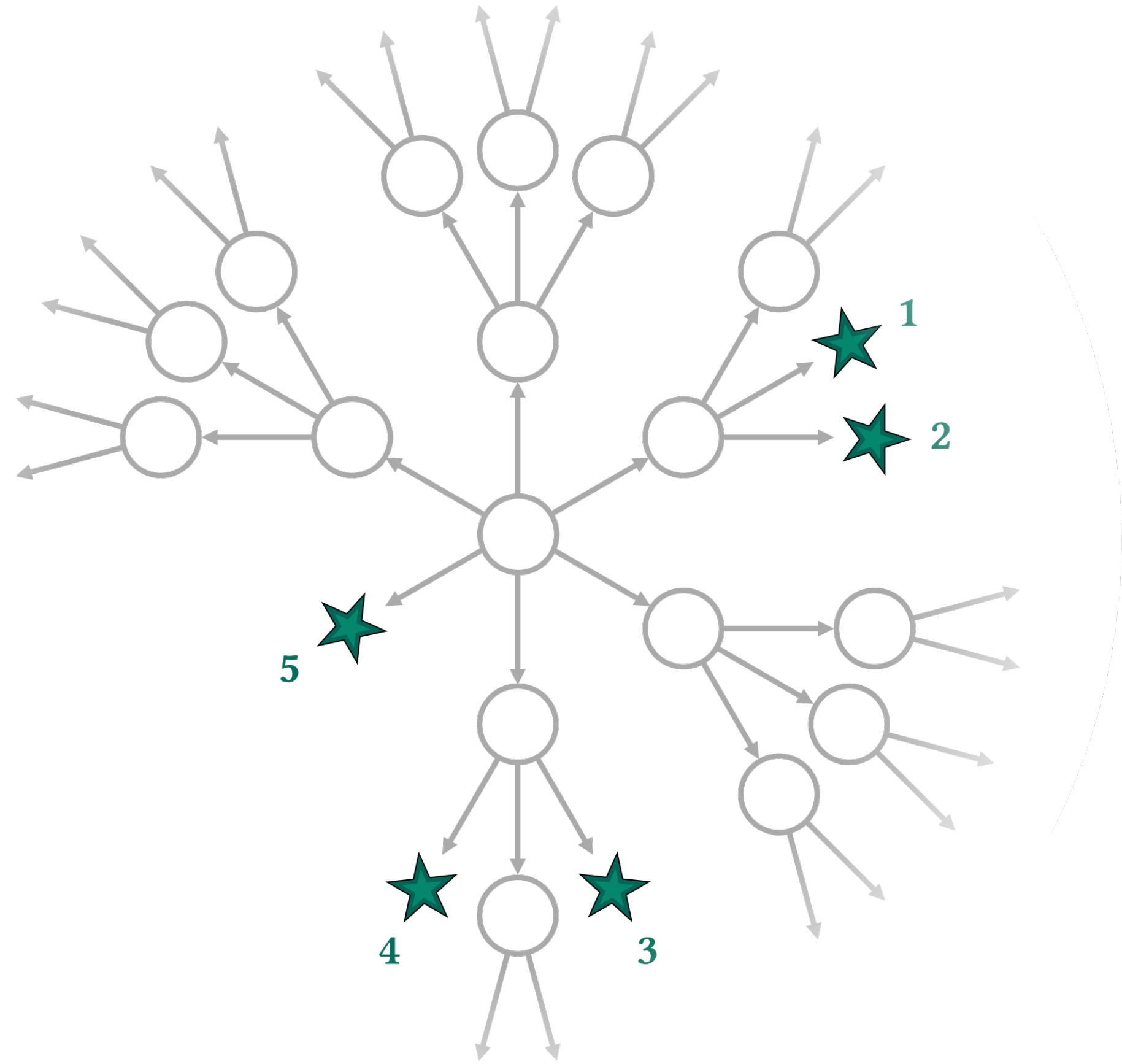
Programming by Navigation



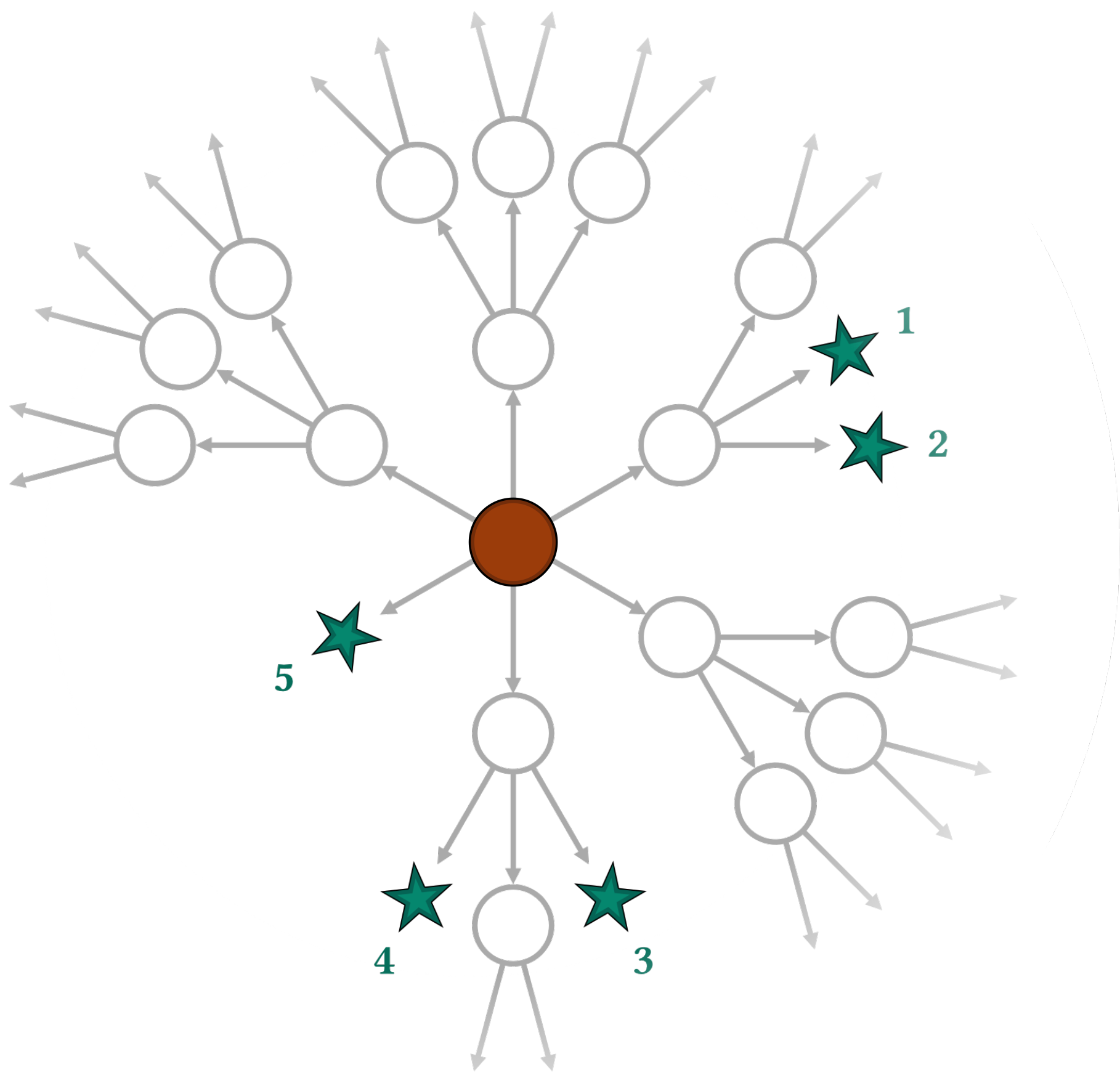
Programming by Navigation



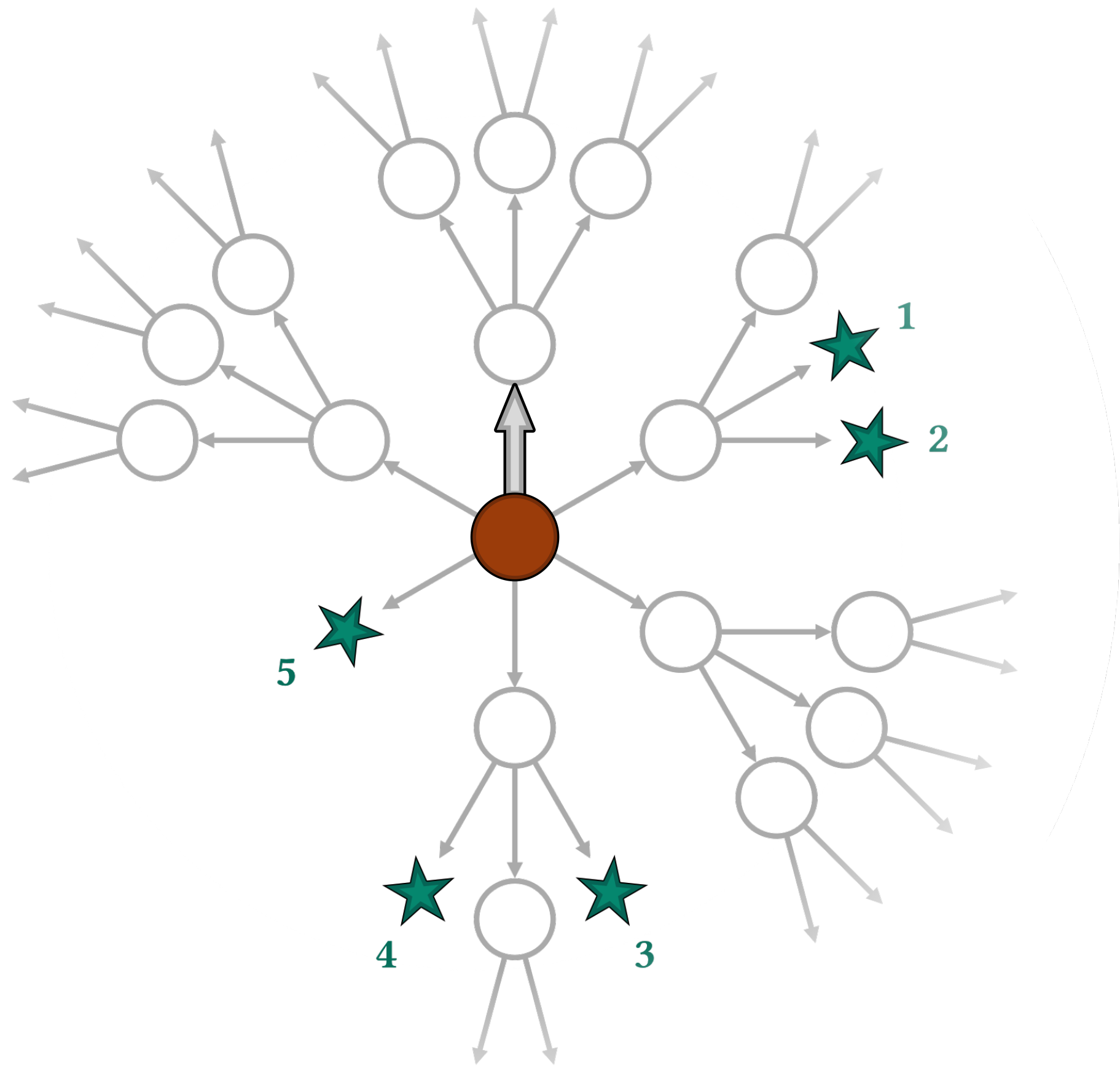
Programming by Navigation



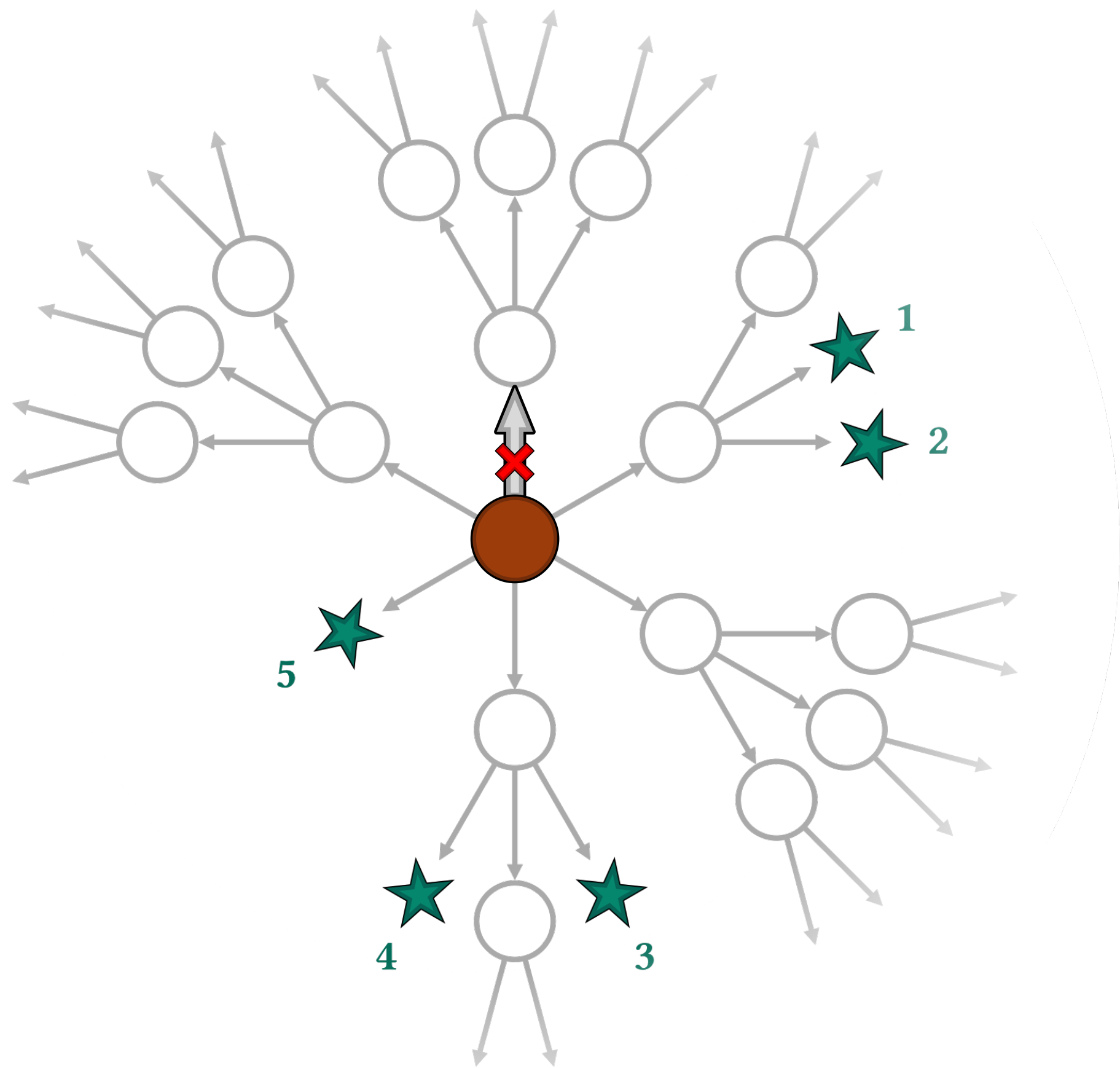
Programming by Navigation



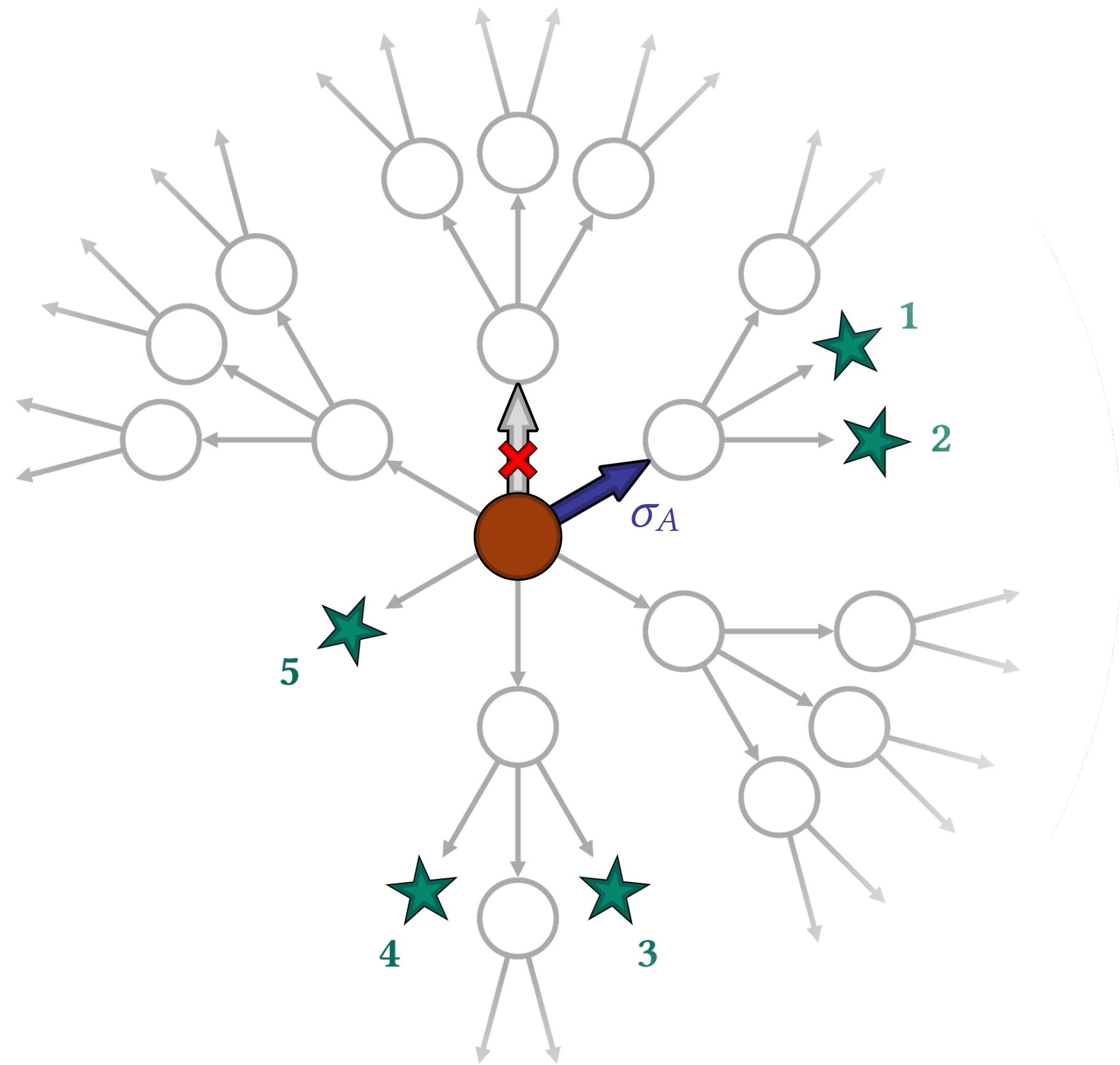
Programming by Navigation



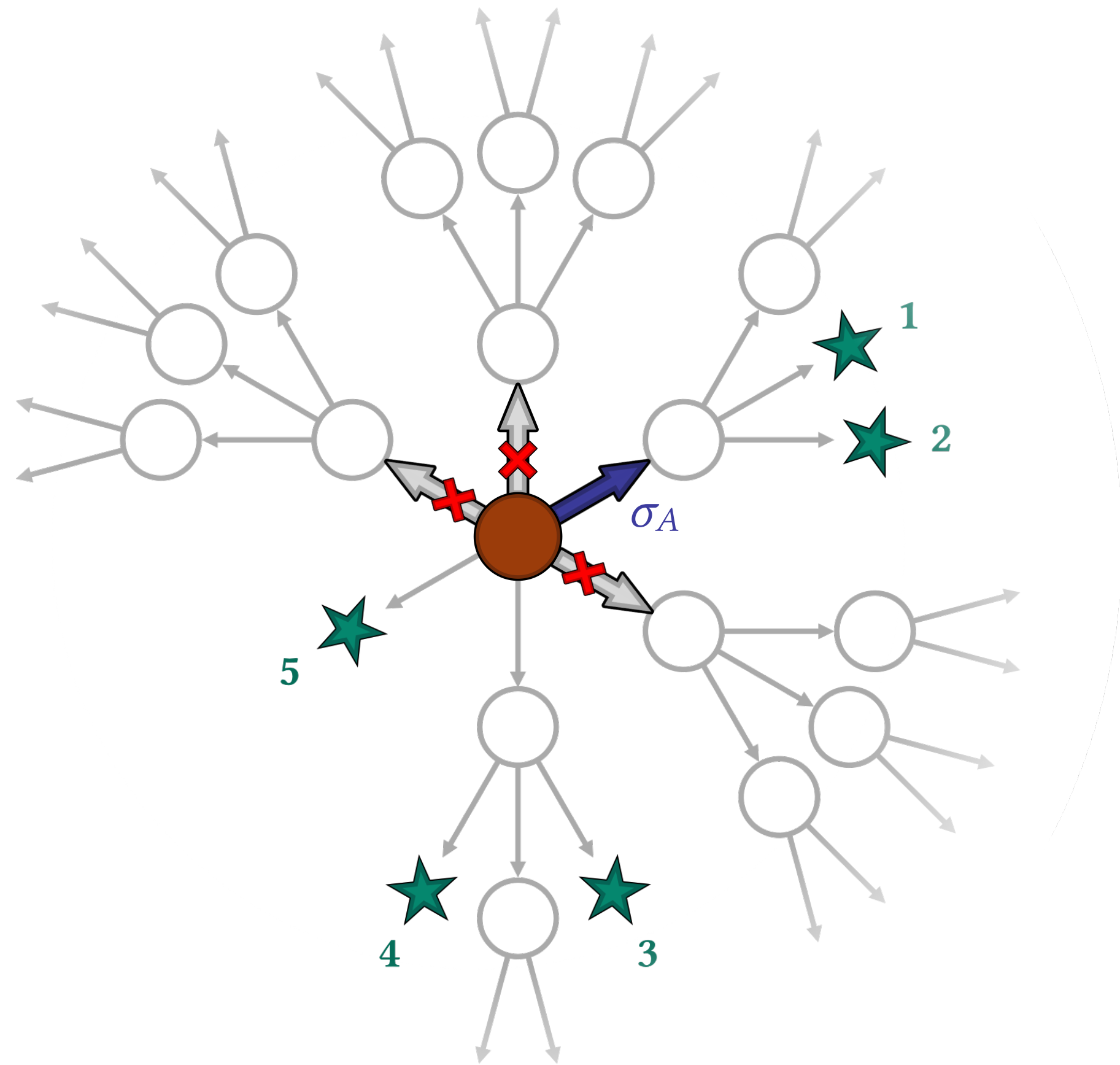
Programming by Navigation



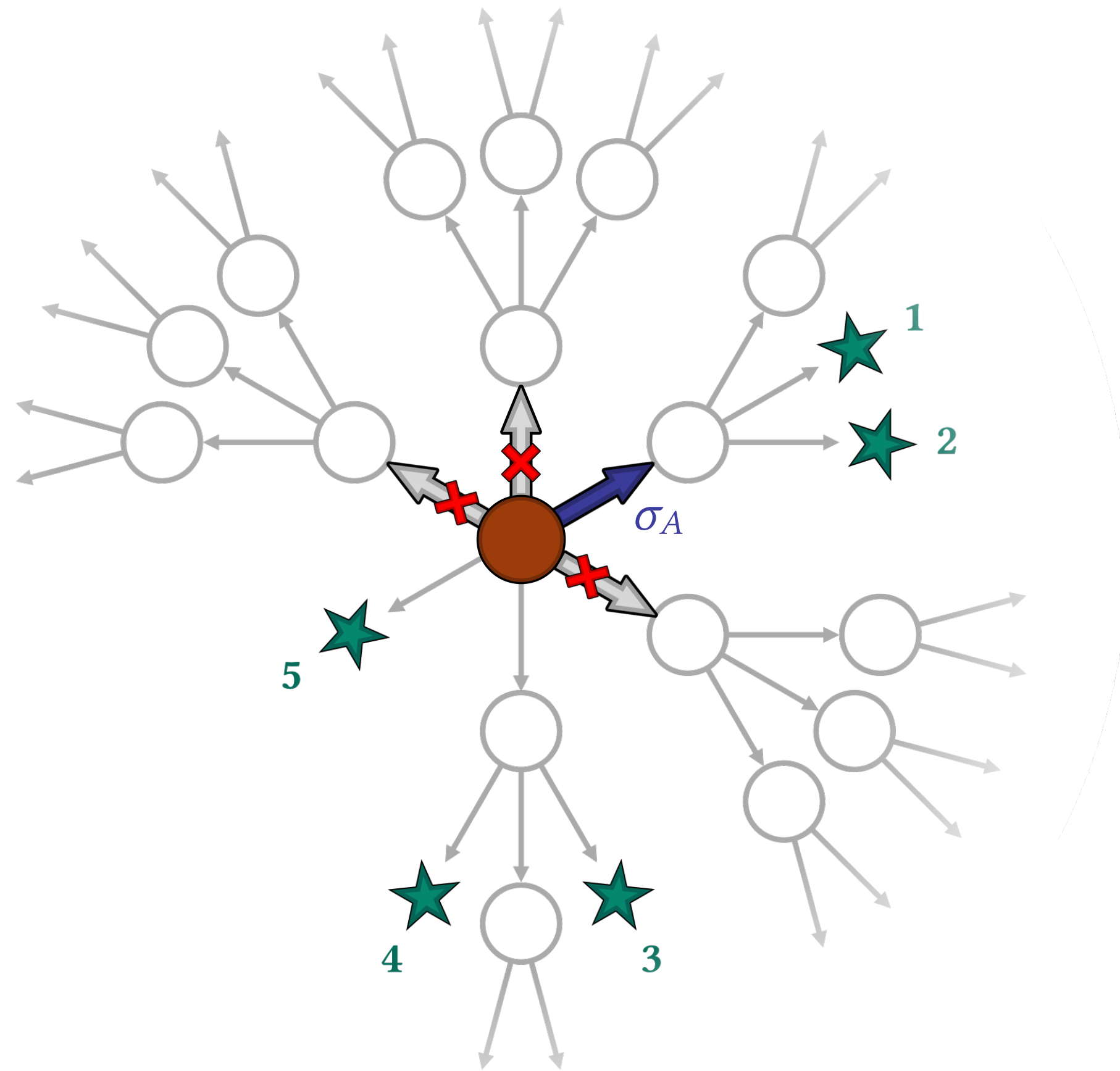
Programming by Navigation



Programming by Navigation

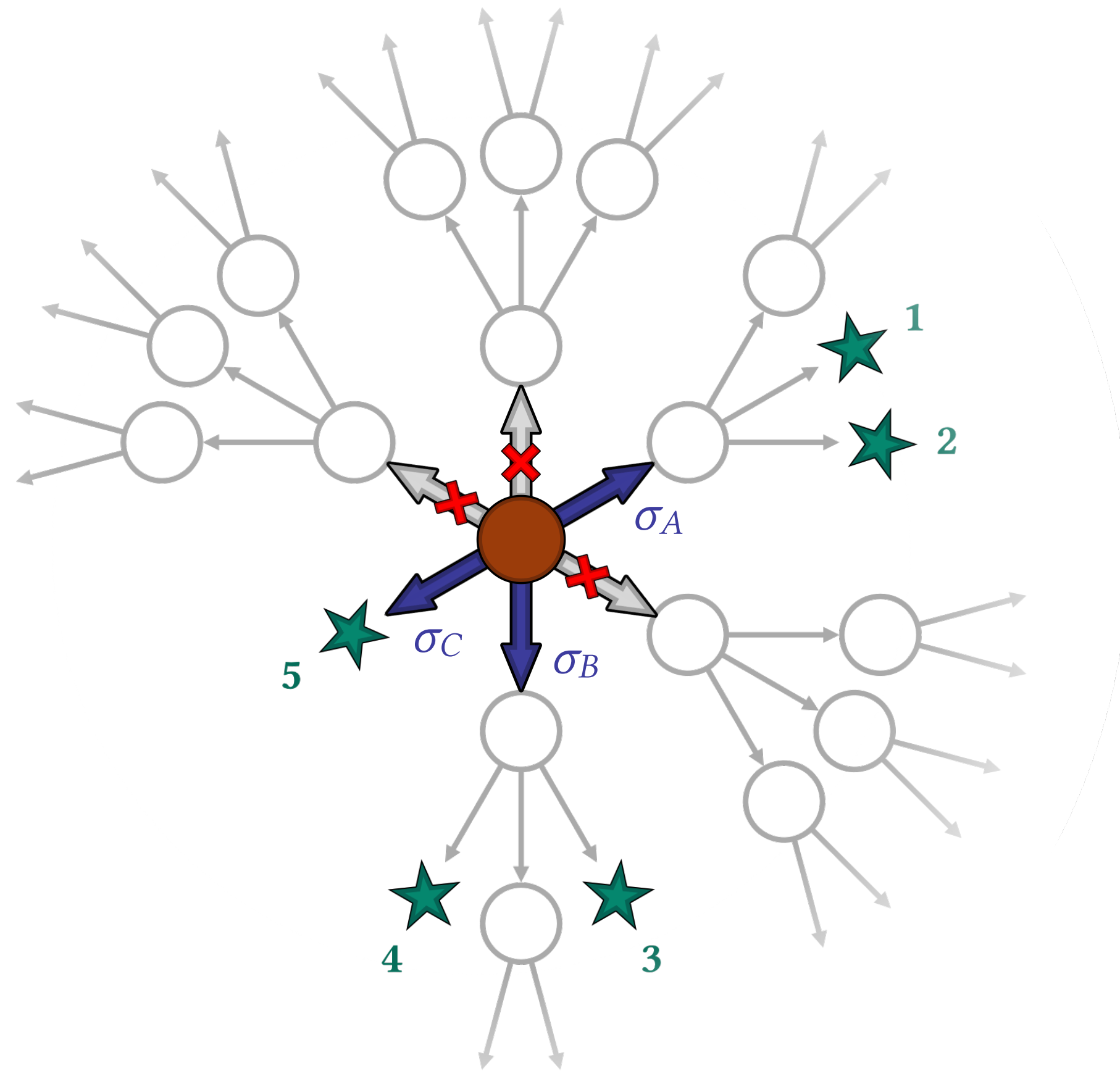


Programming by Navigation



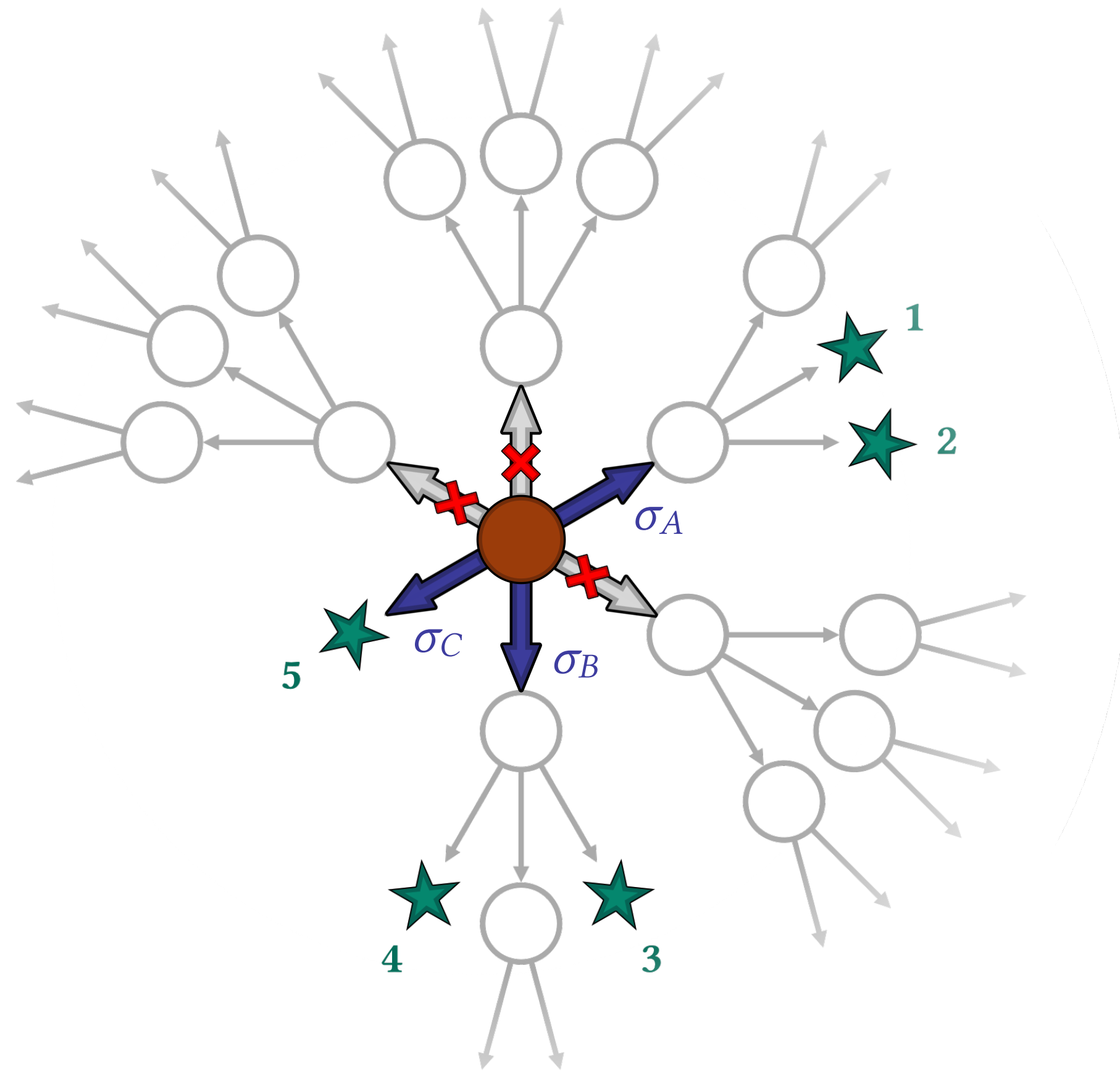
- **STRONG SOUNDNESS:**
Only valid steps are shown

Programming by Navigation



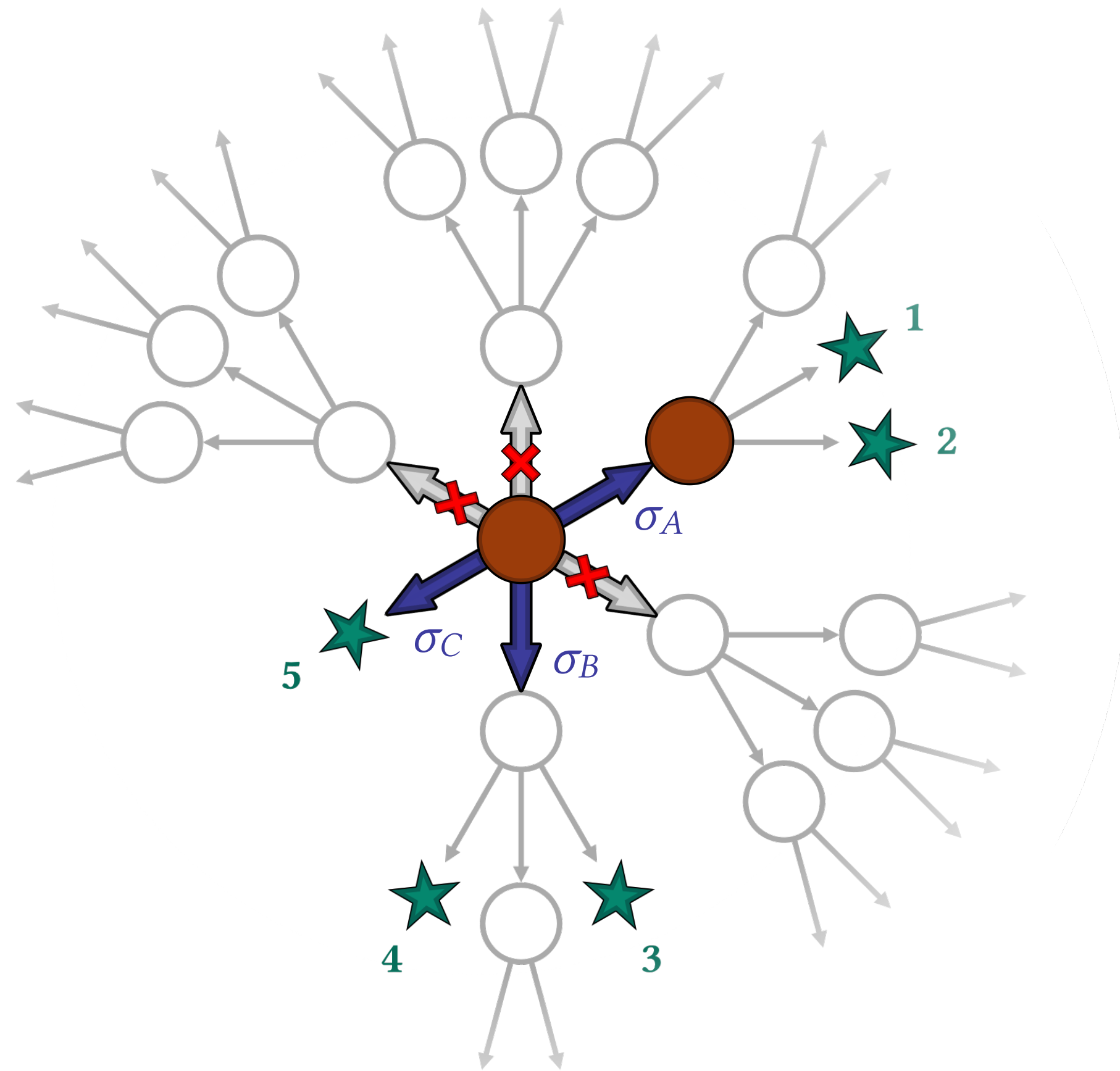
- **STRONG SOUNDNESS:**
Only valid steps are shown

Programming by Navigation



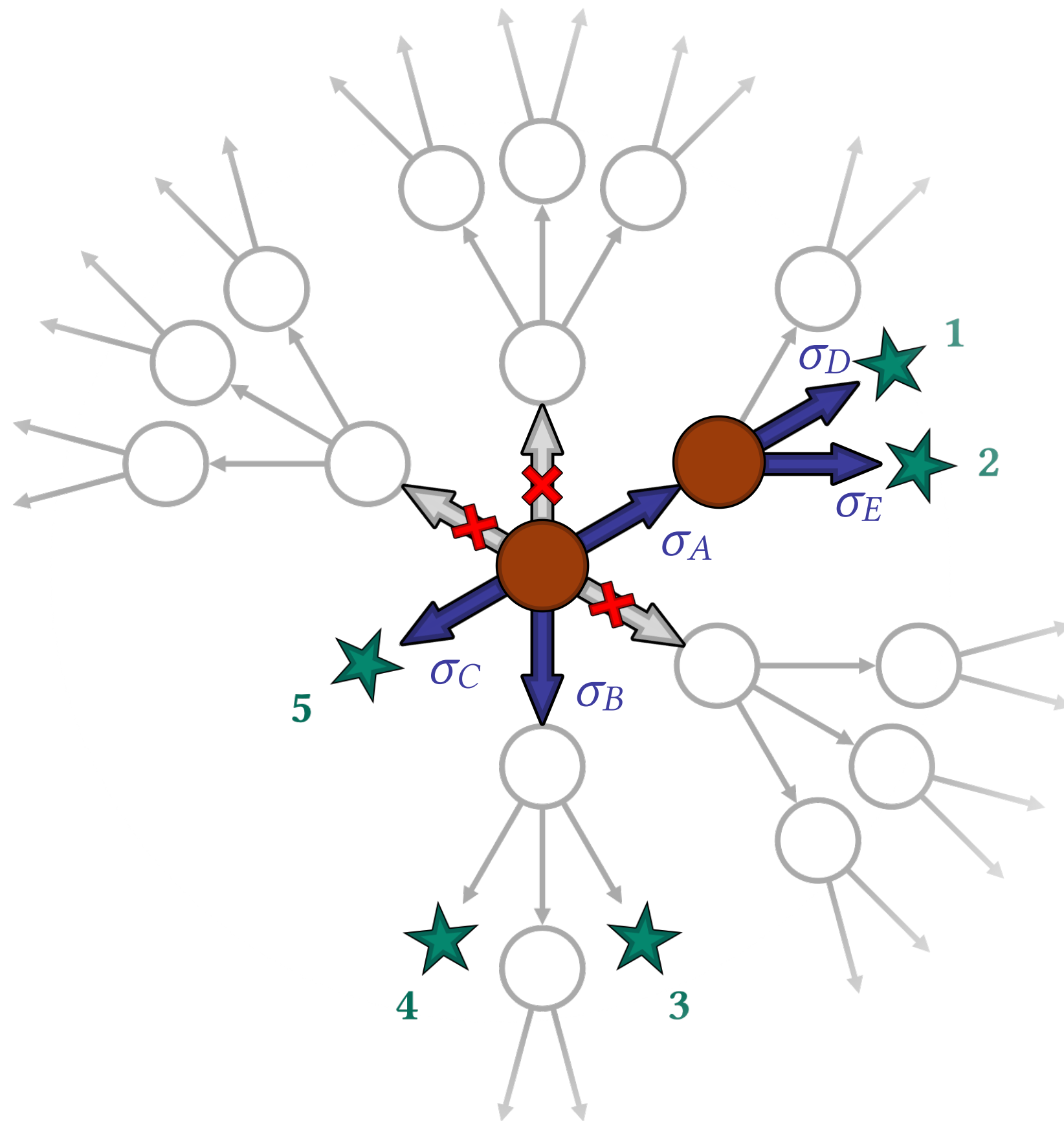
- **STRONG SOUNDNESS:**
Only valid steps are shown
- **STRONG COMPLETENESS:**
Shown steps can reach *all* valid solutions

Programming by Navigation



- **STRONG SOUNDNESS:**
Only valid steps are shown
- **STRONG COMPLETENESS:**
Shown steps can reach *all* valid solutions

Programming by Navigation



- **STRONG SOUNDNESS:**
Only valid steps are shown
- **STRONG COMPLETENESS:**
Shown steps can reach *all* valid solutions

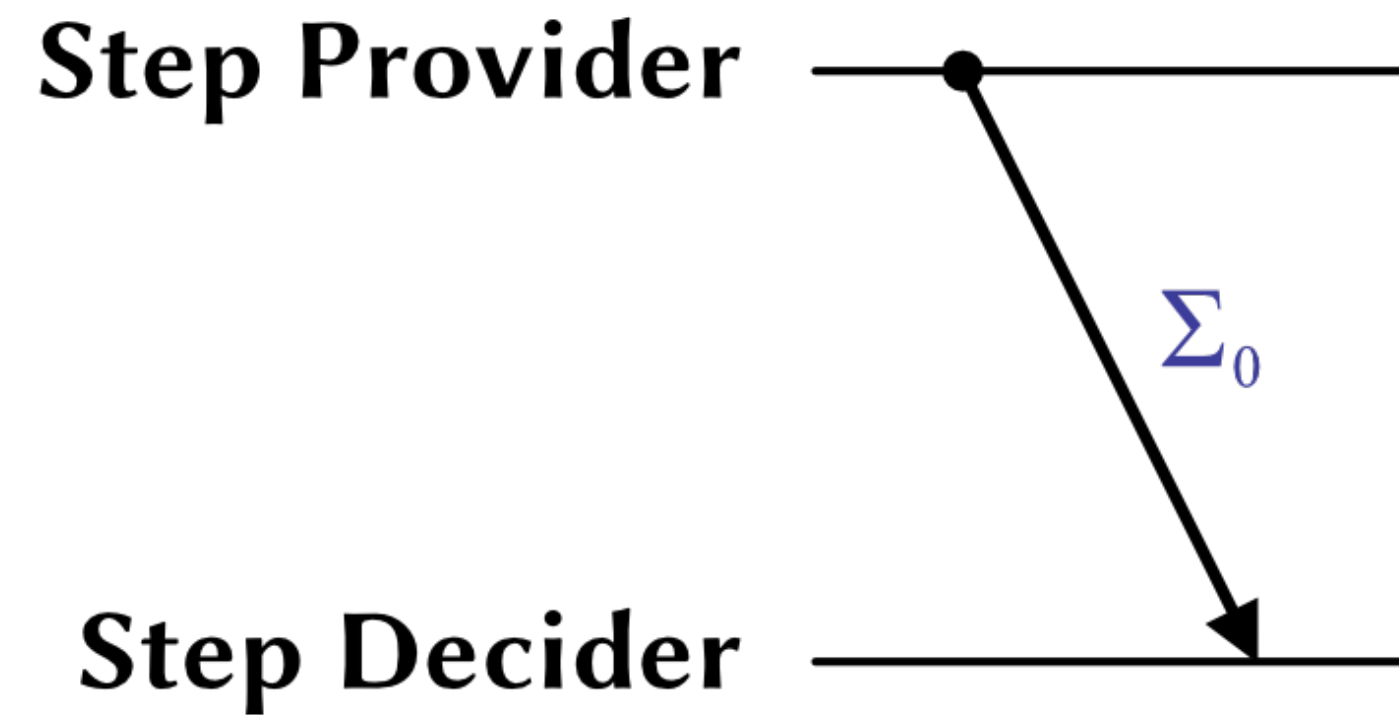
Users select between steps in Programming by Navigation

Users select between steps in Programming by Navigation

Step Provider

Step Decider

Users select between steps in Programming by Navigation



e.g. $\Sigma_0 =$

“Willing to prove manually”

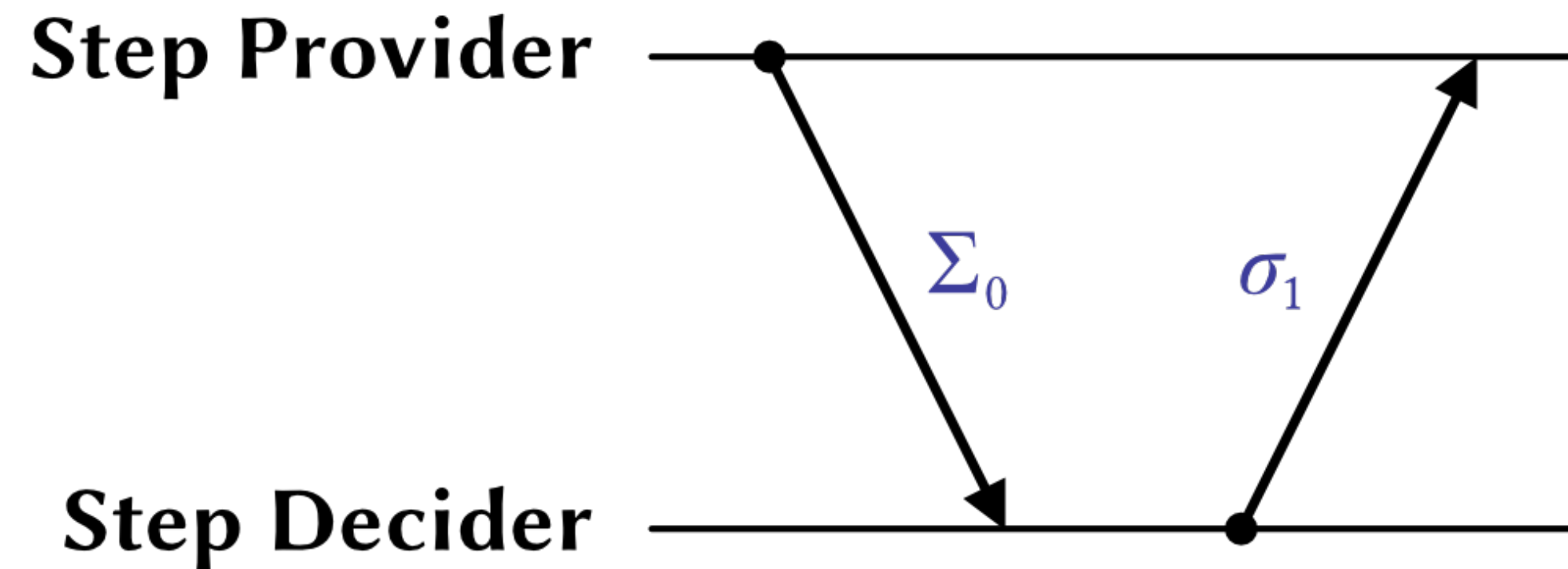
“Not willing to prove manually, but should be true”

above + “must be used”

above + “must be used”

“I don’t know”

Users select between steps in Programming by Navigation



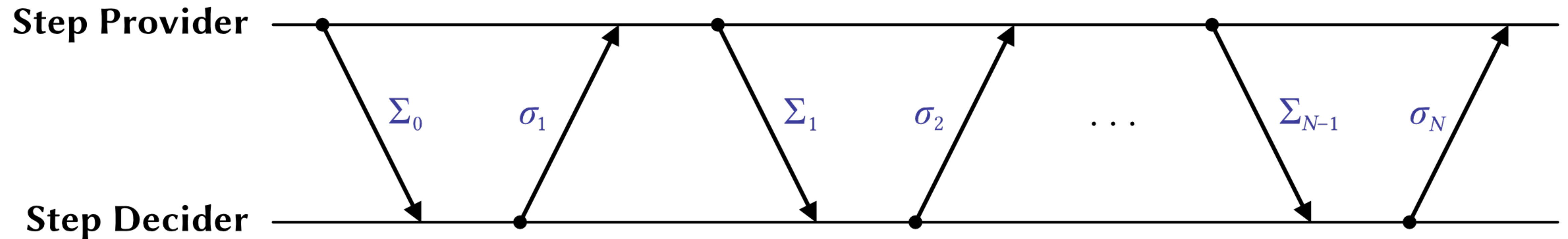
e.g. Σ_0 =

“Willing to prove manually”	“Not willing to prove manually, but should be true”	
<i>above</i> + “must be used”	<i>above</i> + “must be used”	“I don’t know”

e.g. σ_1 =

“Not willing to prove manually, but should be true”

Users select between steps in Programming by Navigation



e.g. $\Sigma_0 =$

“Willing to prove manually”	“Not willing to prove manually, but should be true”	
<i>above</i> + “must be used”	<i>above</i> + “must be used”	“I don’t know”

e.g. $\sigma_1 =$

“Not willing to prove manually, but should be true”

Programming by Navigation Ingredients

Programming by Navigation Ingredients

Expressions

Programming by Navigation Ingredients

Expressions

Validity

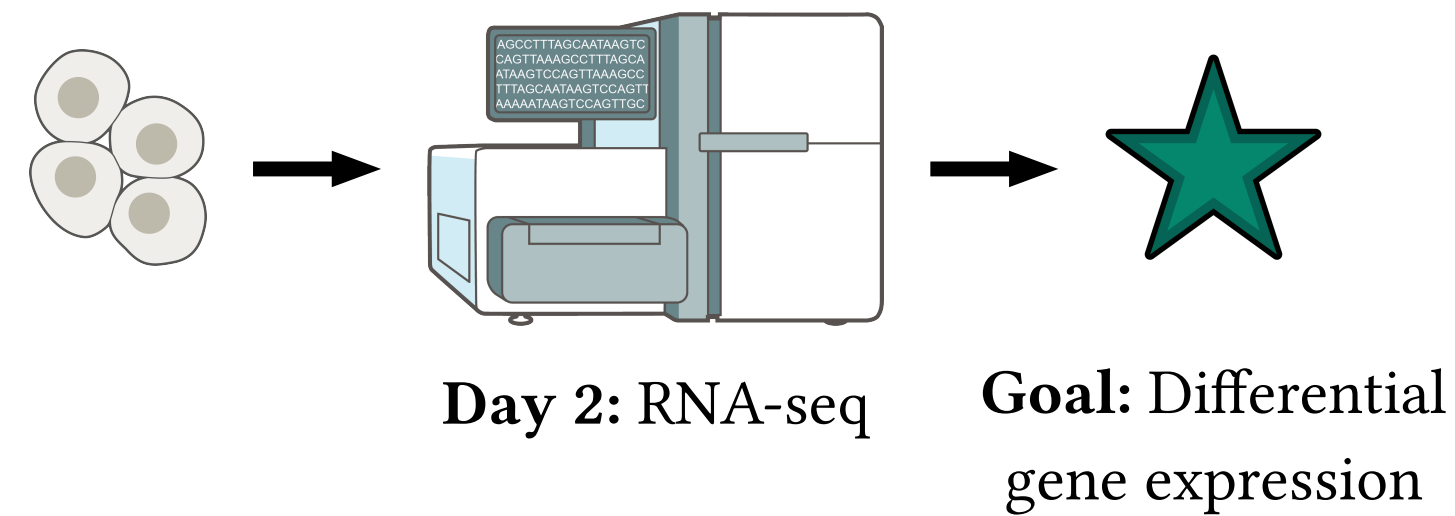
Programming by Navigation Ingredients

Expressions

Validity

Step provider

Programming by Navigation Ingredients



Expressions

Bioinformatics pipelines

Validity

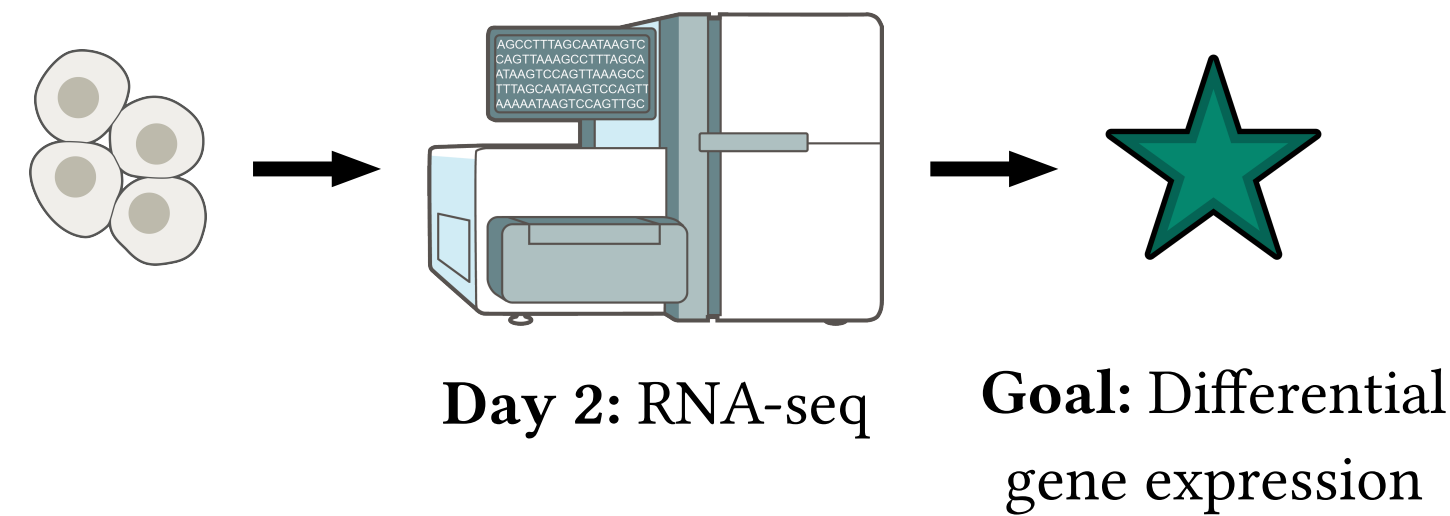
Meets user-defined validity conditions

Step provider

Top-down classical-constructive synthesis

Programming by Navigation Ingredients

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{ij} E_{ij} = A_{ik} D_{kj} \end{array} \right.$$

Annotations for the equations:

- For $\forall_{ikj} A_{ij} = B_{ik} C_{kj}$:
 - Parallelization factor: $\boxed{?}$ $\boxed{?}$
 - Memory layout: $\boxed{?}$
 - Vector width: $\boxed{?}$
- For $\forall_{ij} E_{ij} = A_{ik} D_{kj}$:
 - Dataflow ordering: $\boxed{?}$
 - Memory layout: $\boxed{?}$

Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

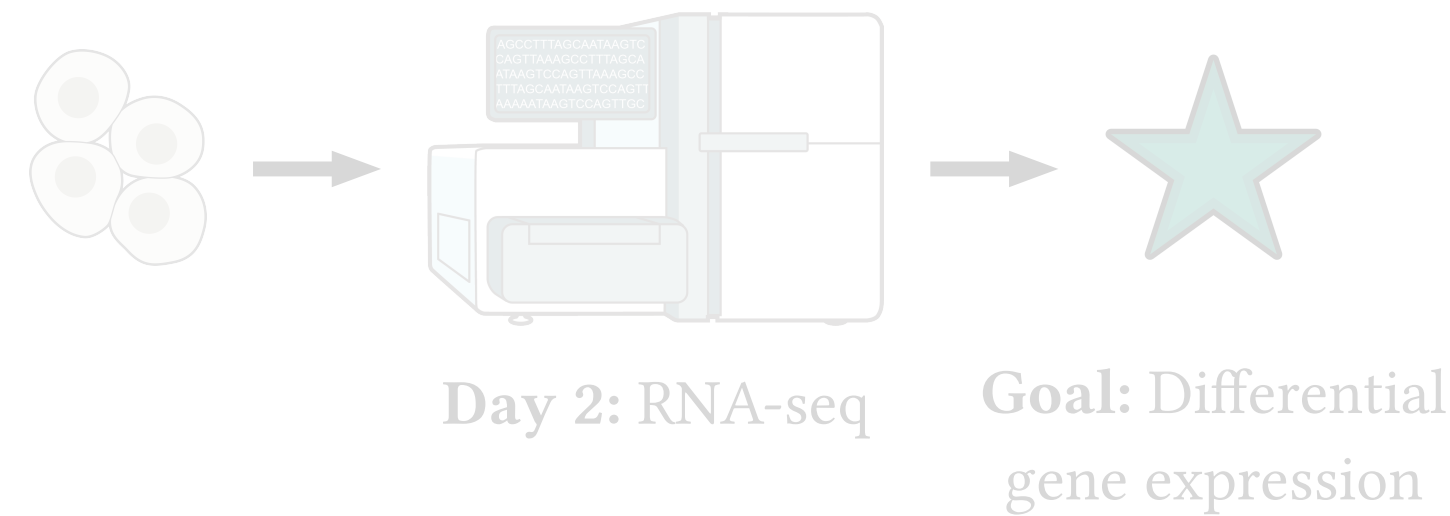
Step provider

Top-down classical-constructive synthesis

Big-step synthesis

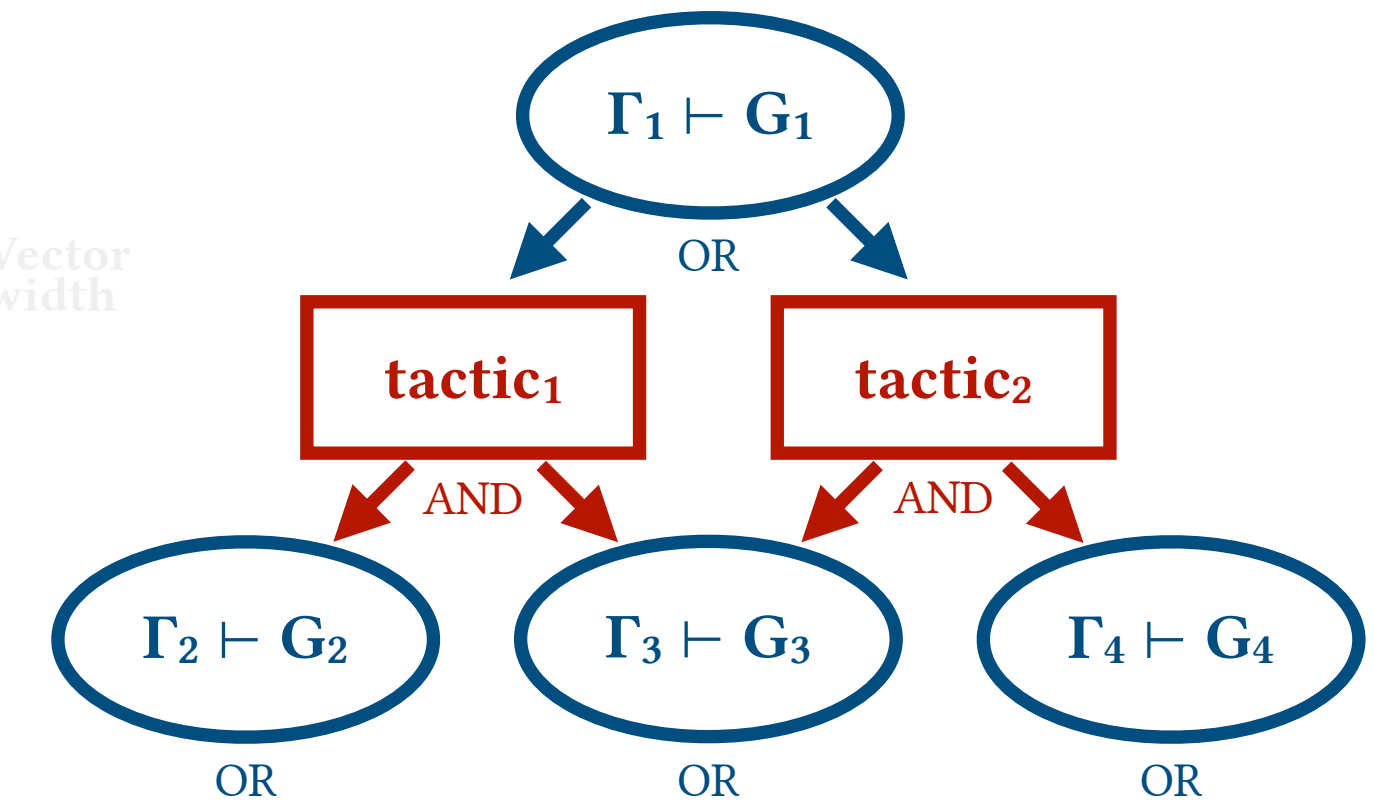
Programming by Navigation Ingredients

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor [?], Memory layout [?], Vector width [?], Dataflow ordering [?], Memory layout [?]



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

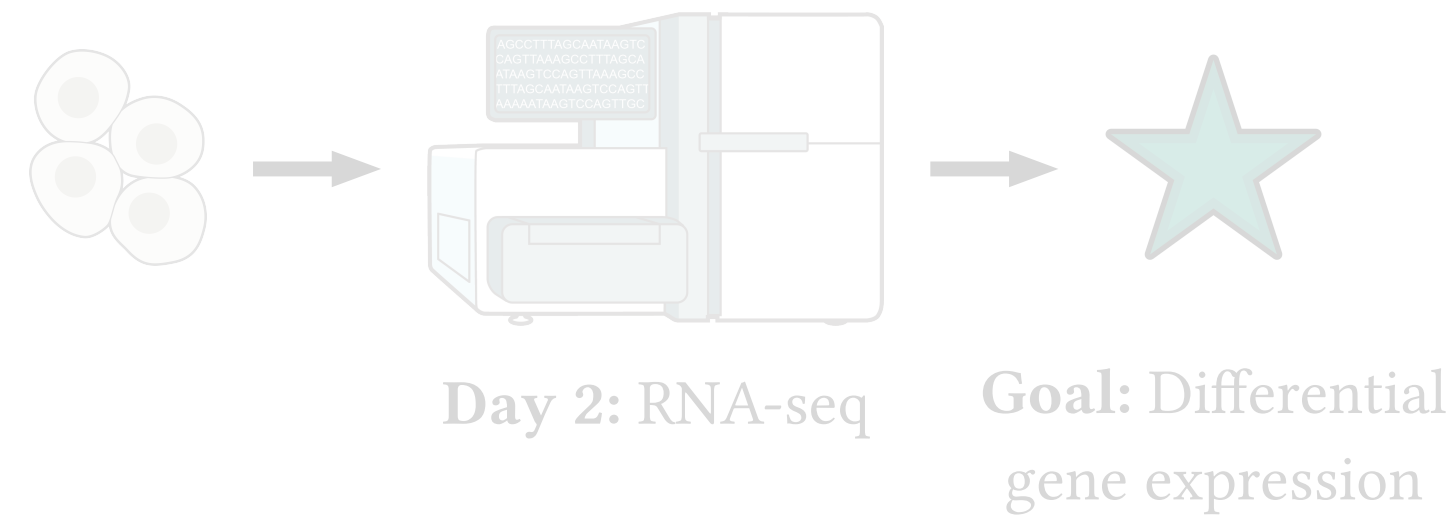
Step provider

Top-down classical-constructive synthesis

Big-step synthesis

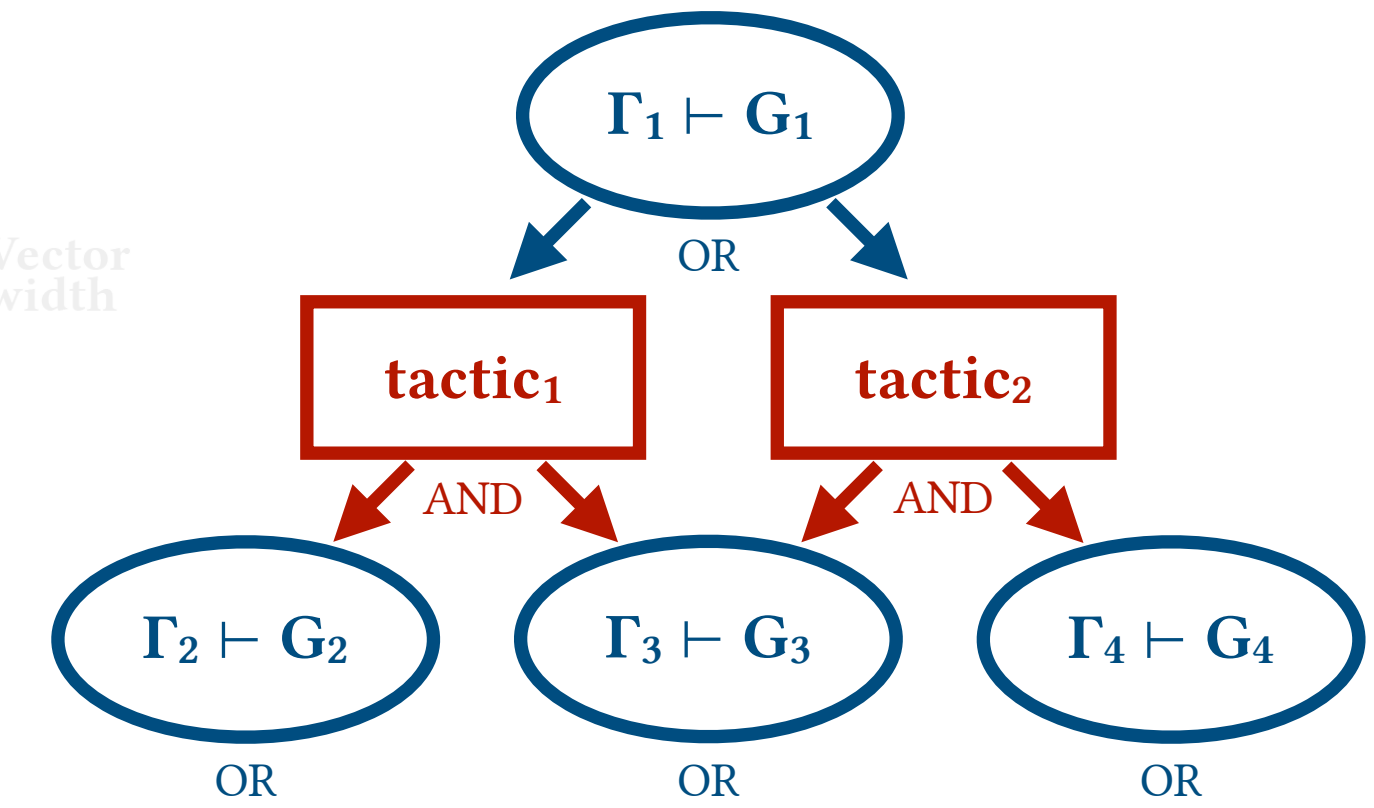
Programming by Navigation Ingredients

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor [?], Memory layout [?], Vector width [?], Dataflow ordering [?], Memory layout [?]



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

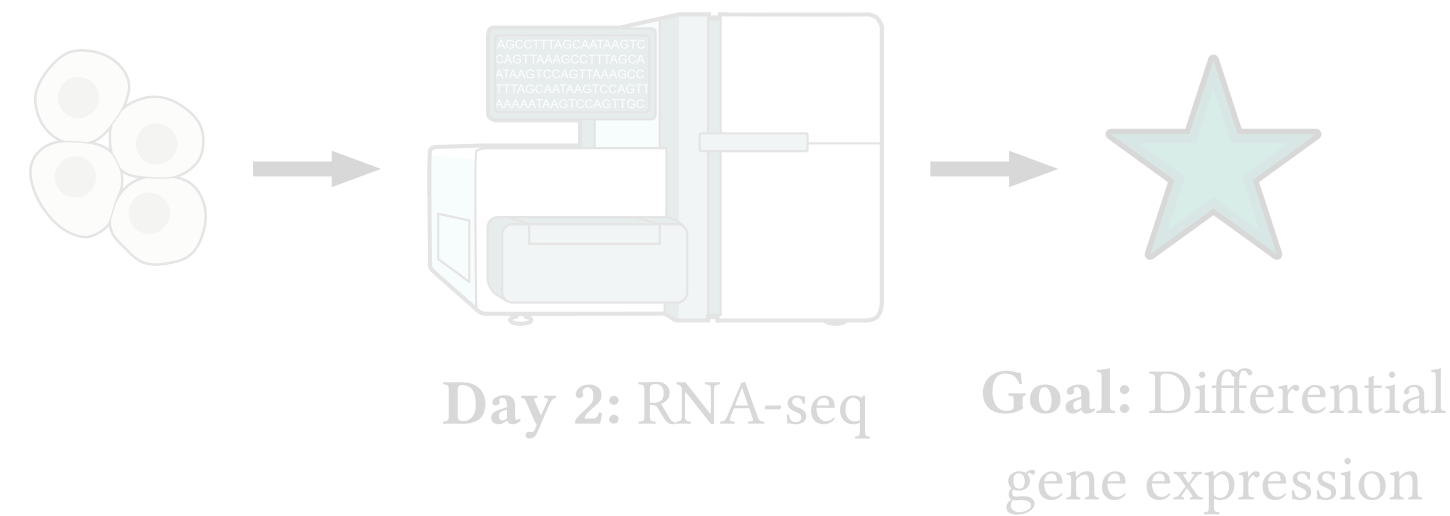
Step provider

Top-down classical-constructive synthesis

Big-step synthesis

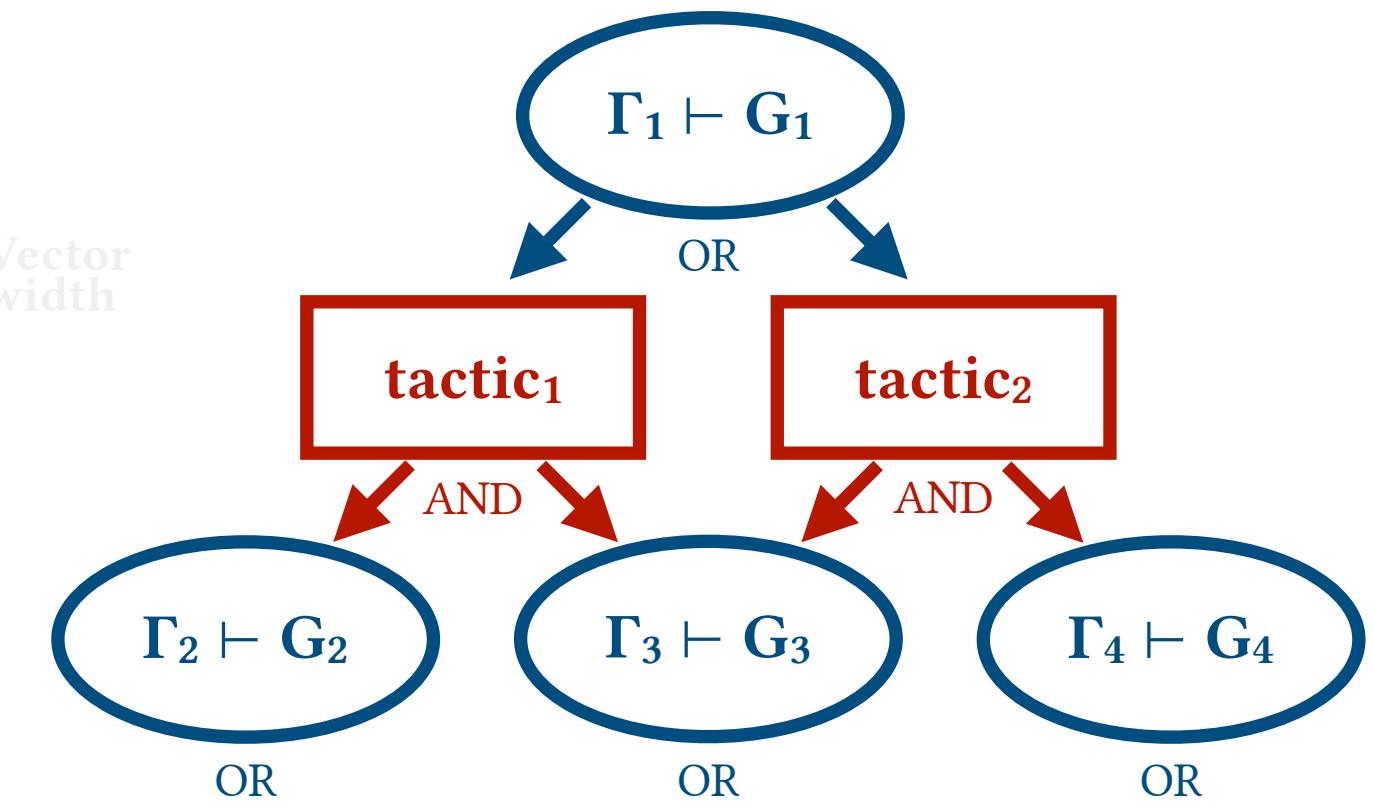
Programming by Navigation Ingredients

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor [?], Memory layout [?], Vector width [?], Dataflow ordering [?], Memory layout [?]



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

Modification proves final goal node

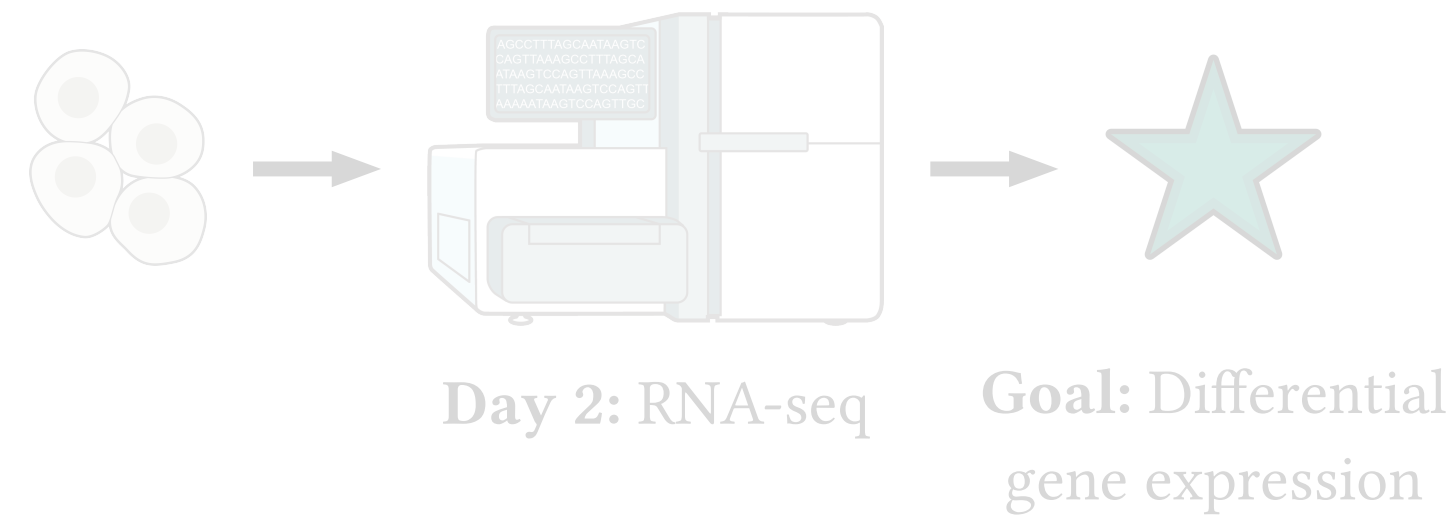
Step provider

Top-down classical-constructive synthesis

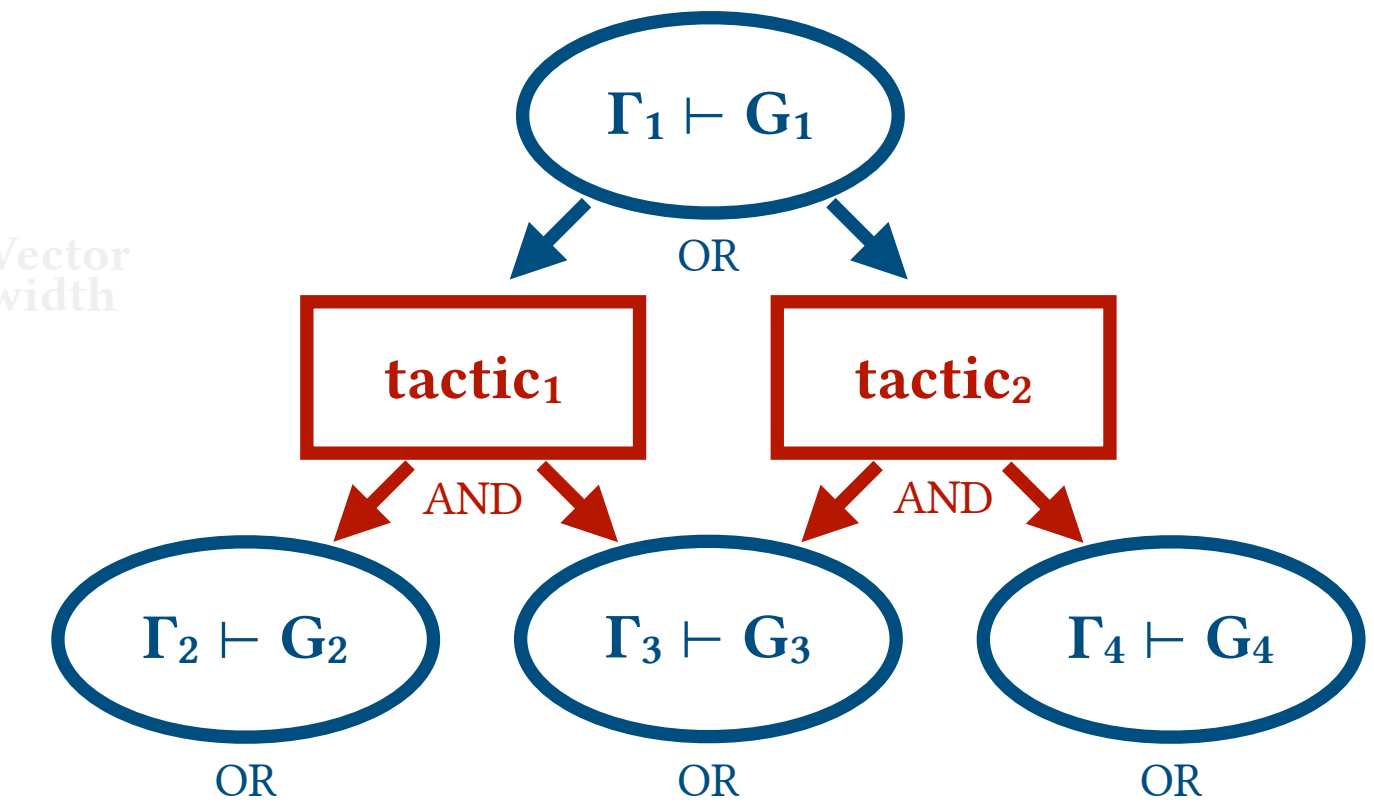
Big-step synthesis

Programming by Navigation Ingredients

Lacouture *et al.* FuseFlow: A Fusion-Centric Compilation Framework for Sparse Deep Learning on Streaming Dataflow (ASPLOS 2026).



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \quad \text{Parallelization factor} \quad \text{Memory layout} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \quad \text{Dataflow ordering} \quad \text{Memory layout} \end{array} \right.$$



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

Modification proves final goal node

Step provider

Top-down classical-constructive synthesis

Big-step synthesis

?

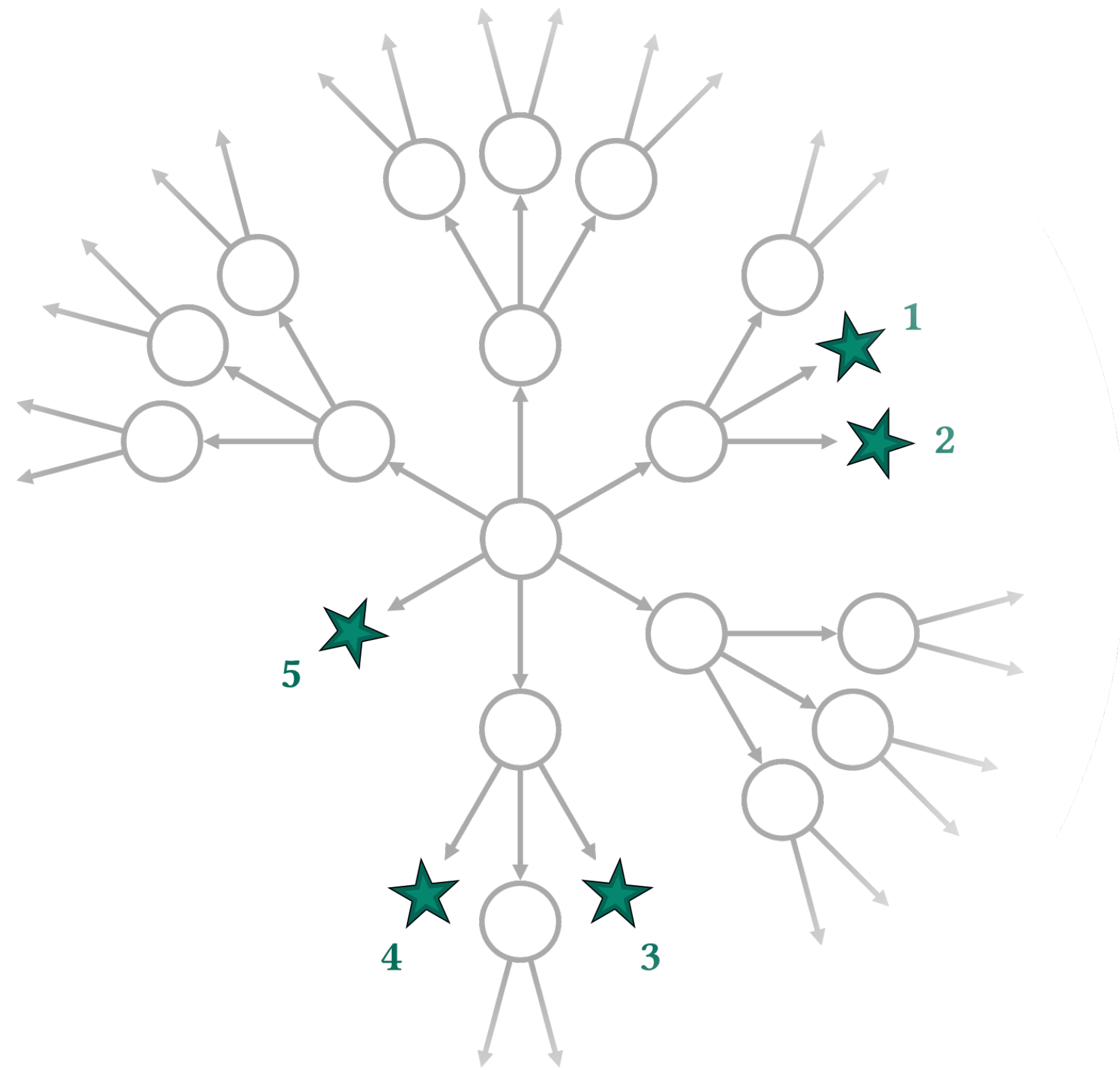
Helpful building block: Nonempty-completion oracle

Helpful building block: Nonempty-completion oracle

$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$

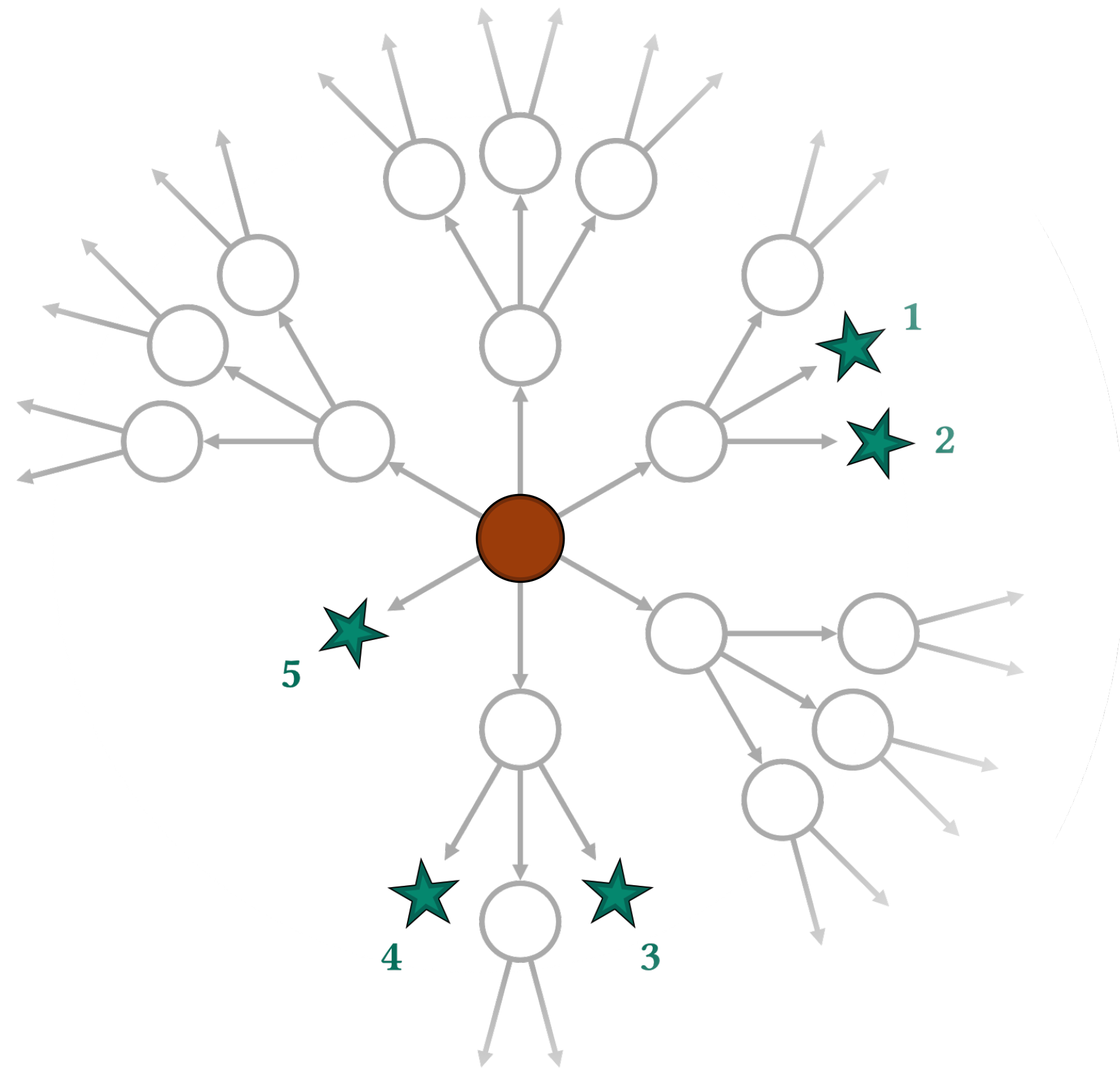
Helpful building block: Nonempty-completion oracle

$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



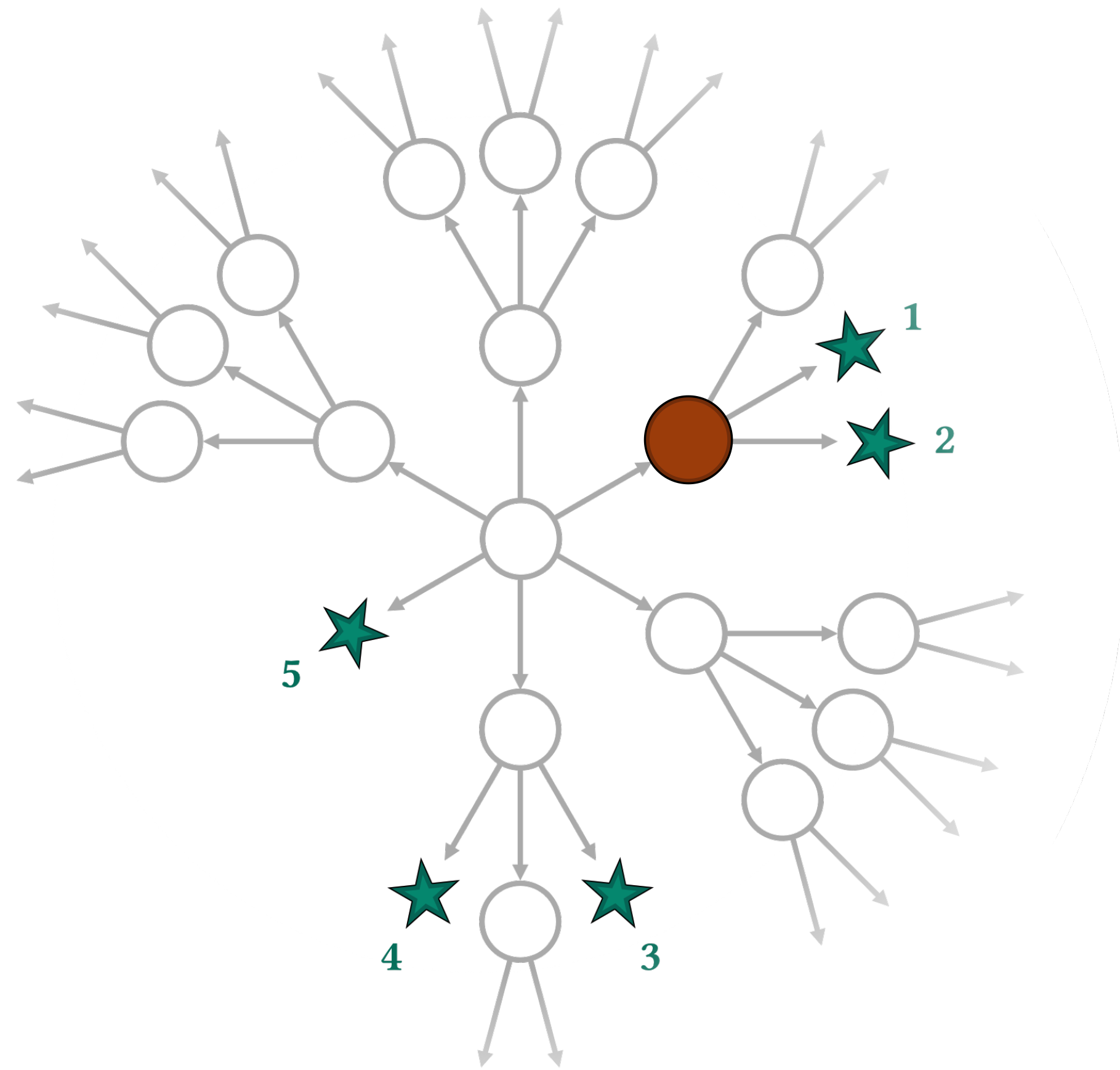
Helpful building block: Nonempty-completion oracle

$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



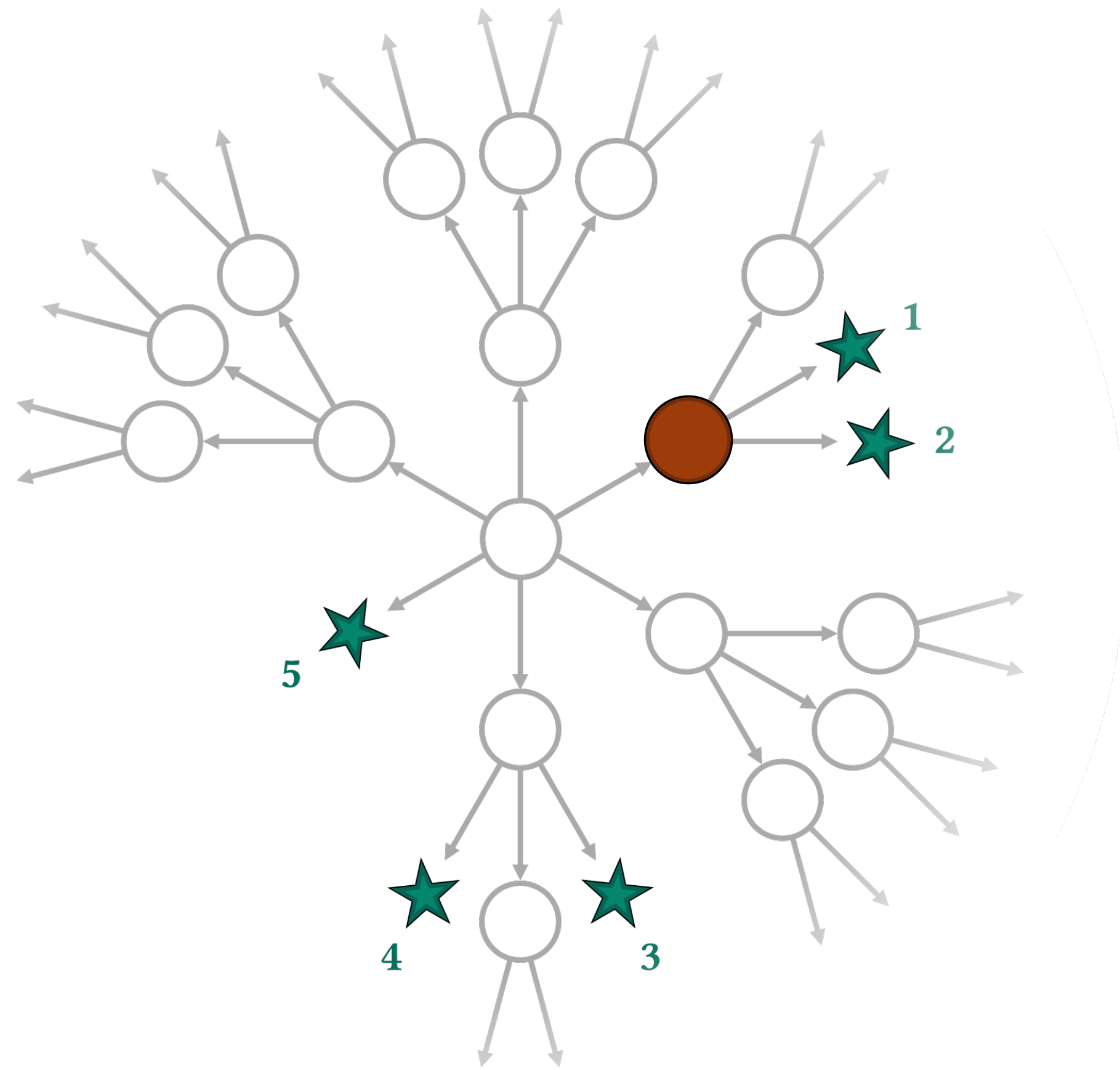
Helpful building block: Nonempty-completion oracle

$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



Helpful building block: Nonempty-completion oracle

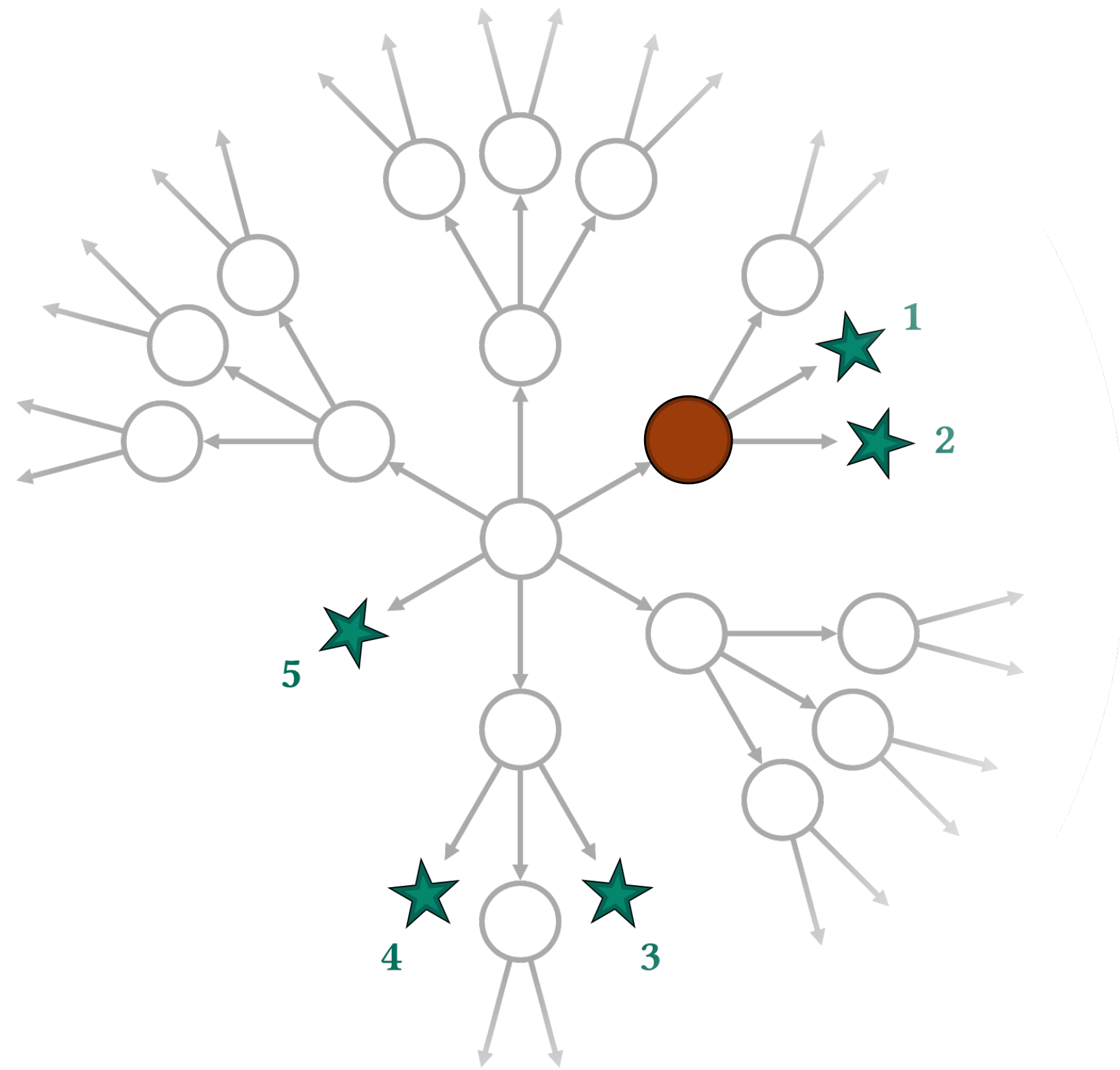
$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



Define via compilation φ to SAT:

Helpful building block: Nonempty-completion oracle

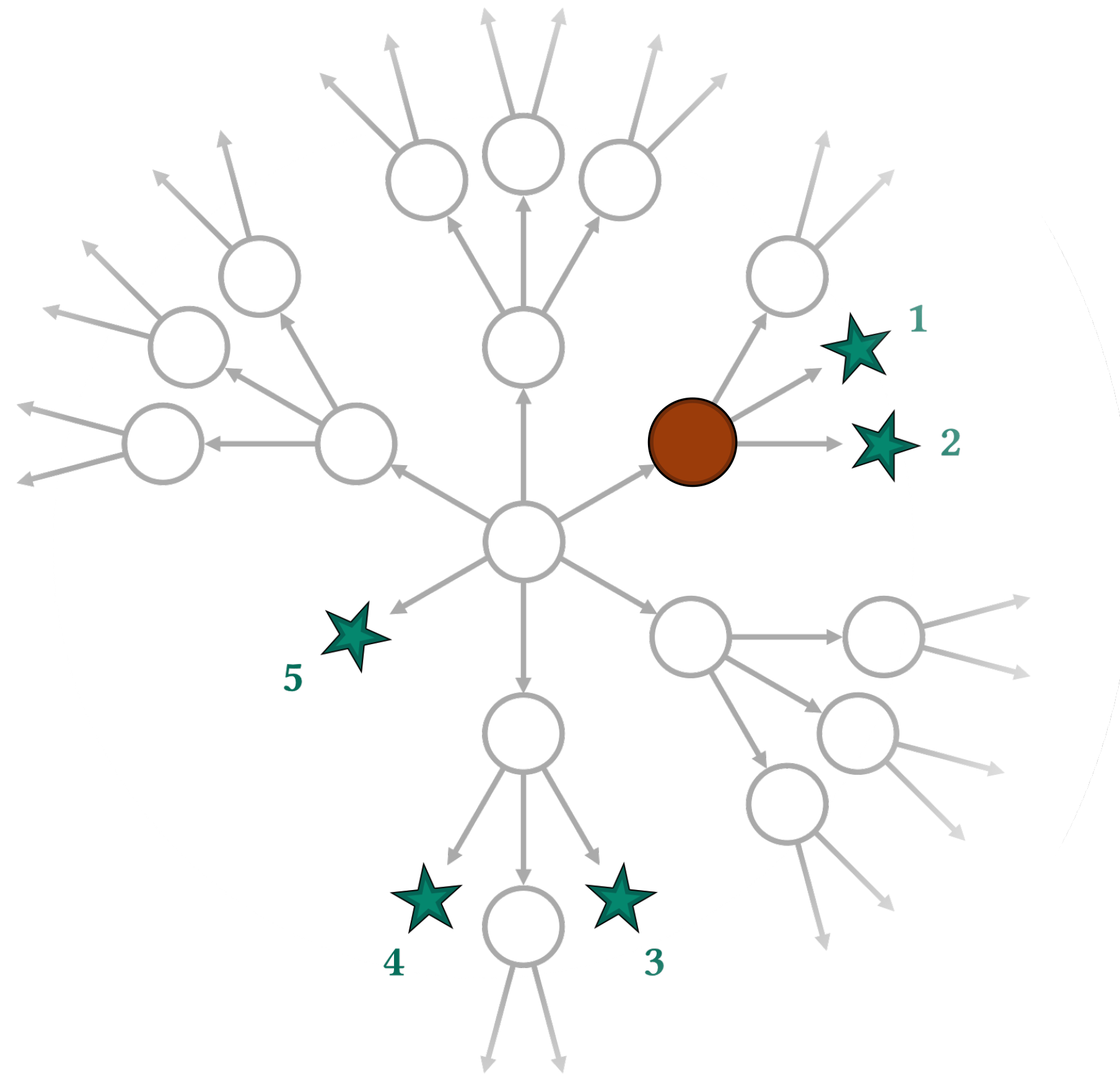
$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



Define via compilation φ to SAT: $\text{SAT}(\varphi(e))$ iff $C(e) \neq \emptyset$

Helpful building block: Nonempty-completion oracle

$$C(e) = \{e' \mid e \preceq e' \wedge e' \text{ valid}\}$$



266:28

Justin Lubin, Marlena Preigh, Max Willsey, and Sarah E. Chasins

B SAT-Based Oracle to Check for Nonempty Completions

Semantic Formulas for Propositions. Propositions can either be assumed, in which case they are automatically true, or else they are subject to the standard interpretation of Horn clauses in a proof system: a proposition is true iff there exists a rule that fires with that proposition as a conclusion. This is captured by the following two formulas for each $Q \in \mathcal{A}$:

$$\text{assume}_Q \Rightarrow \text{true}_Q \quad \text{and} \quad \neg \text{assume}_Q \Rightarrow \left(\text{true}_Q \iff \bigvee_{R \text{ with conclusion } Q} \text{true}_R \right)$$

Activity Formulas for Propositions. We start with two necessary conditions for a proposition $Q \in \mathcal{A}$ to be active. First, Q can be active only if it is actually true ($\text{active}_Q \Rightarrow \text{true}_Q$). Moreover, if Q is not the goal node, then it must have an active node that uses it:

$$\text{active}_Q \Rightarrow \bigvee_{R \text{ with premise } Q} \text{active}_R \quad (\text{only for non-goal nodes})$$

Only a single rule that provides a proposition can be active at a given time.

$$\sum_{R \text{ with conclusion } Q} \text{active}_R \leq 1$$

This constraint uses a *cardinality constraint* to limit the number of literals set to true in a given set, which is often directly supported by off-the-shelf SAT solvers; alternatively, since the upper bound of the constraint is 1, it can easily be desugared into standard boolean SAT.

Finally, if Q is to be assumed, then none of the rules that provide Q as a conclusion can be active. Otherwise, Q will be active iff there exists such a rule that is active:

$$\begin{aligned} \text{assume}_Q &\Rightarrow \bigwedge_{R \text{ with conclusion } Q} \neg \text{active}_R \\ \neg \text{assume}_Q &\Rightarrow \left(\text{active}_Q \iff \bigvee_{R \text{ with conclusion } Q} \text{active}_R \right) \end{aligned}$$

Consistency Formulas for Propositions. For each proposition $Q \in \mathcal{A}$, we add a simple consistency formula that enforces the desired meanings of the label partitions:

\top	if $e(Q) = \perp$,
$\neg \text{assume}_Q$	if $e(Q) = ?$,
$\neg \text{true}_Q \wedge \neg \text{assume}_Q \wedge \neg \text{active}_Q$	if $e(Q) = \text{F}$,
true_Q	if $e(Q) = \text{T/A}$,
$\text{true}_Q \wedge \text{active}_Q$	if $e(Q) = \text{T/A!}$,
$\text{true}_Q \wedge \neg \text{assume}_Q$	if $e(Q) = \text{T}$,
$\text{true}_Q \wedge \neg \text{assume}_Q \wedge \text{active}_Q$	if $e(Q) = \text{T!}$,
$\text{true}_Q \wedge \text{assume}_Q$	if $e(Q) = \text{A}$,
$\text{true}_Q \wedge \text{assume}_Q \wedge \text{active}_Q$	if $e(Q) = \text{A!}$.

Semantic and Activity Formulas for Rules. Because rules do not have partition labels or assume-semantic, the formulas for them are much simpler:

$$\text{true}_R \iff \bigwedge_{Q \text{ is premise of } R} \text{true}_Q \quad \text{and} \quad \text{active}_R \Rightarrow \text{true}_R \quad \text{and} \quad \text{active}_R \Rightarrow \bigwedge_{Q \text{ is premise of } R} \text{active}_Q$$

The first formula is standard AND semantics, the second captures that a rule can only be active if it is true, and the third captures that a rule can only be active if all its premises are active.

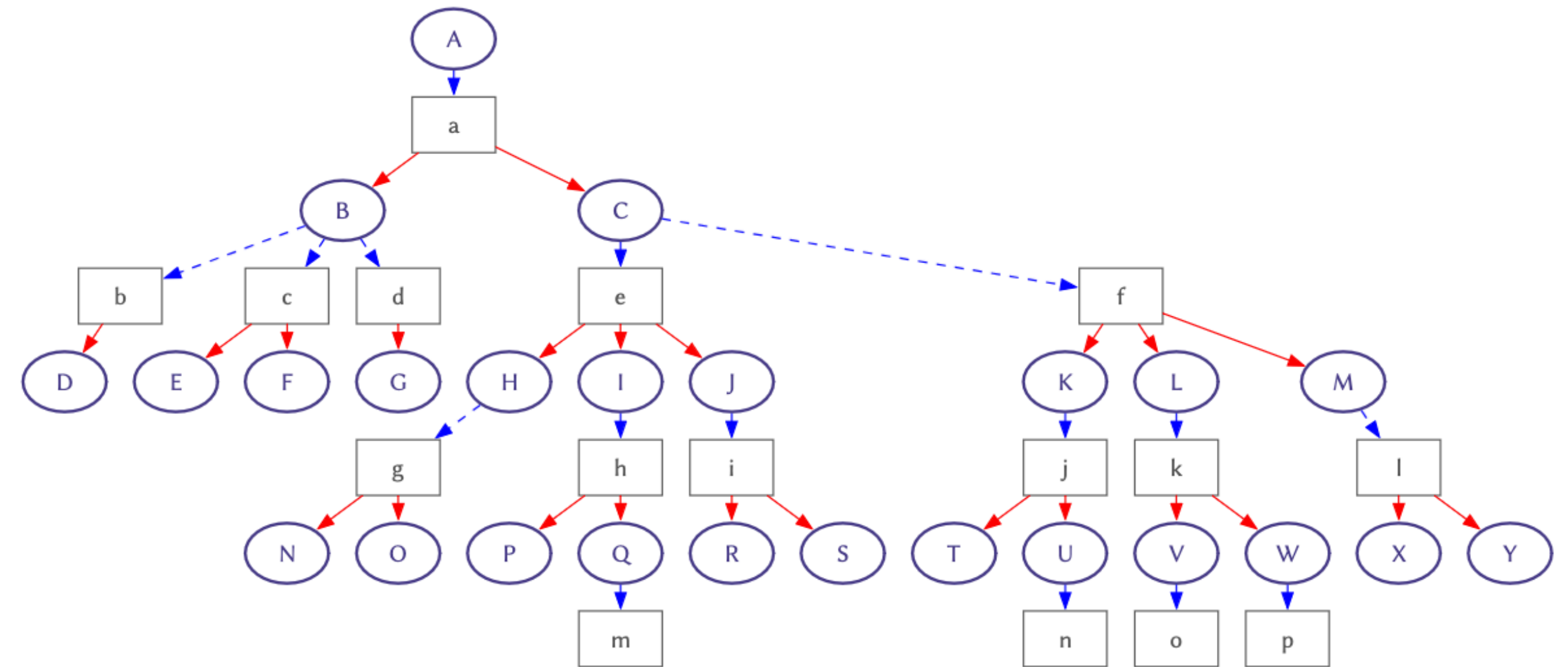
Define via compilation φ to SAT: $\text{SAT}(\varphi(e))$ iff $C(e) \neq \emptyset$

Warm-up: TOP-DOWN step provider

```

1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 

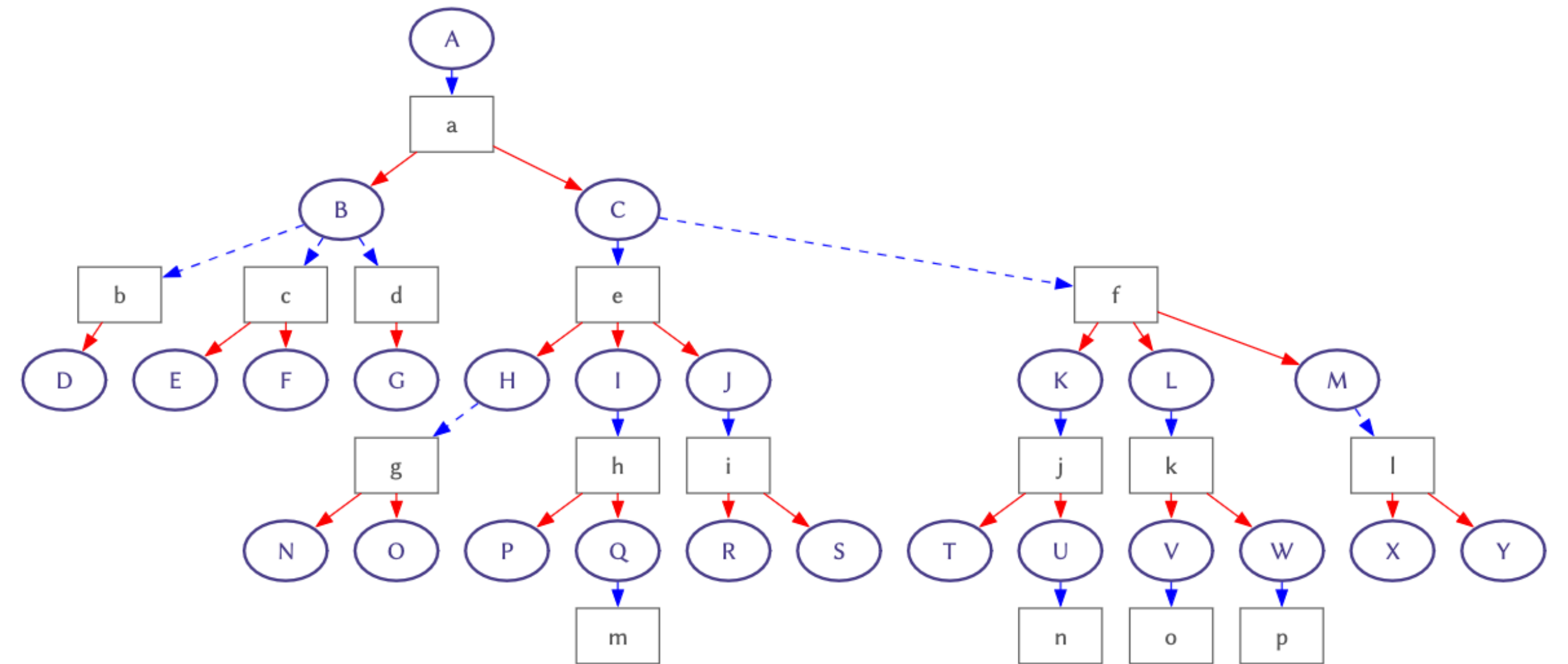
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

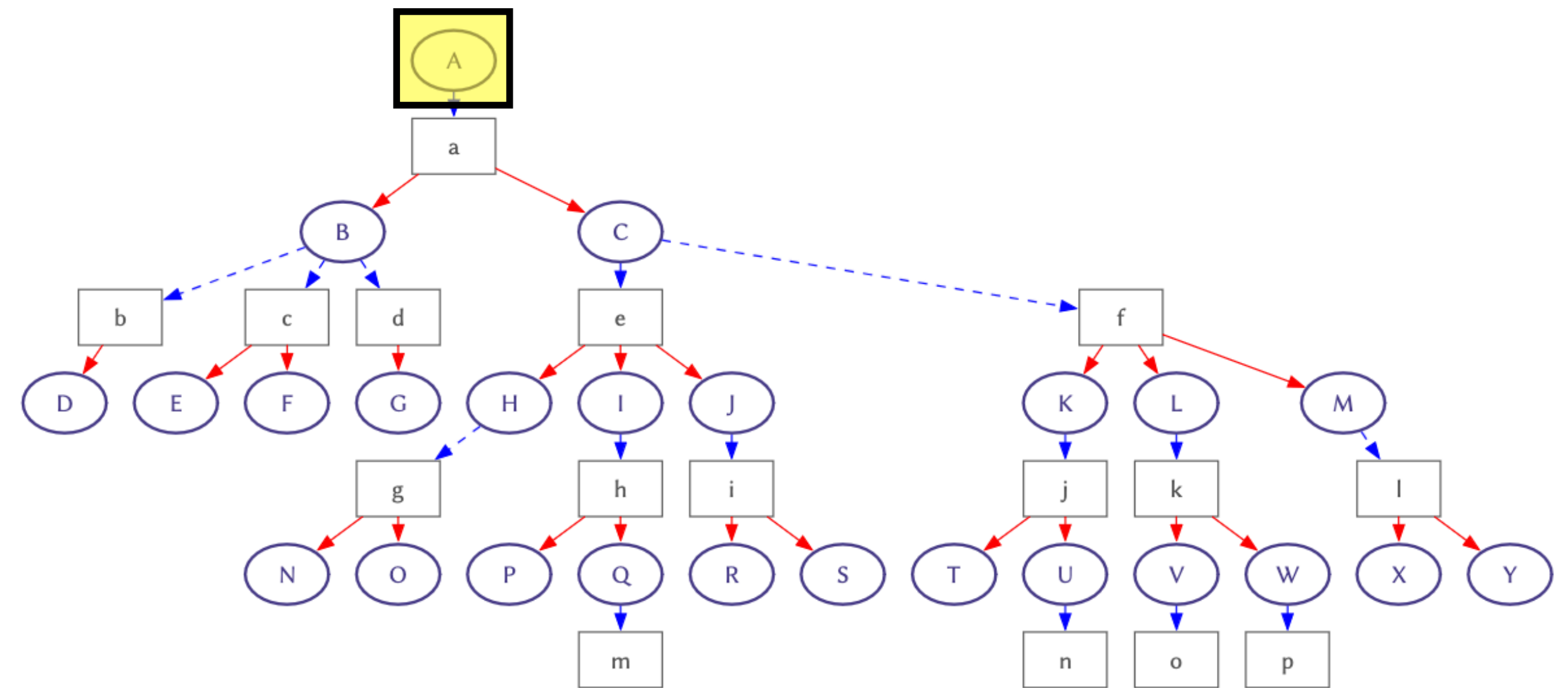
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

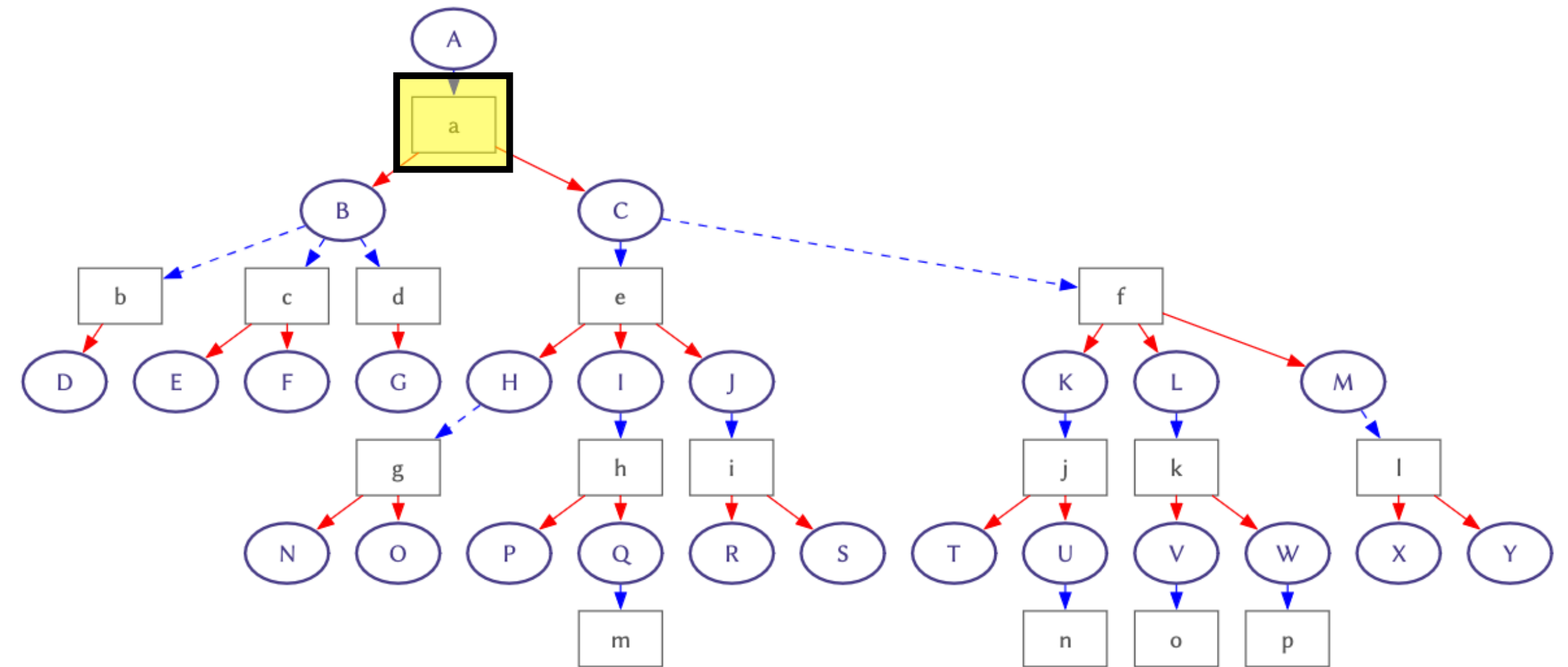
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

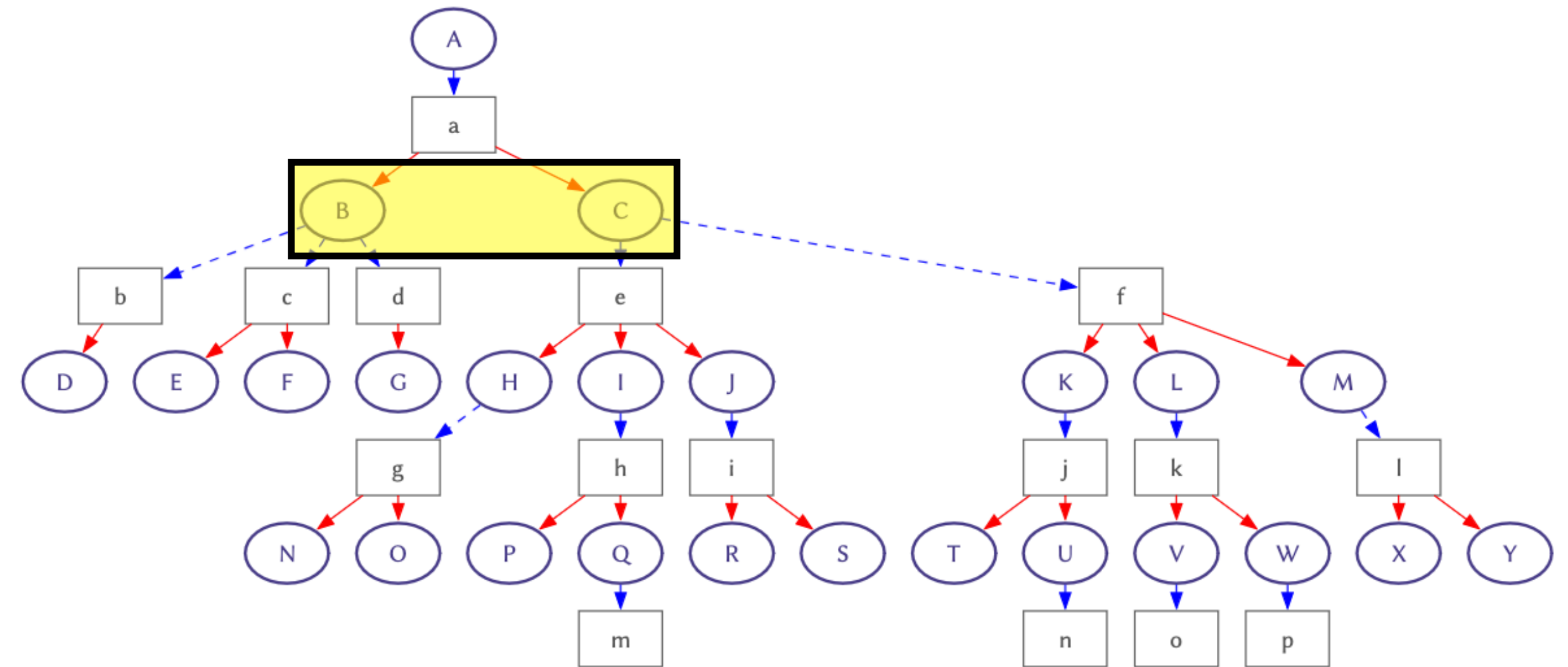
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

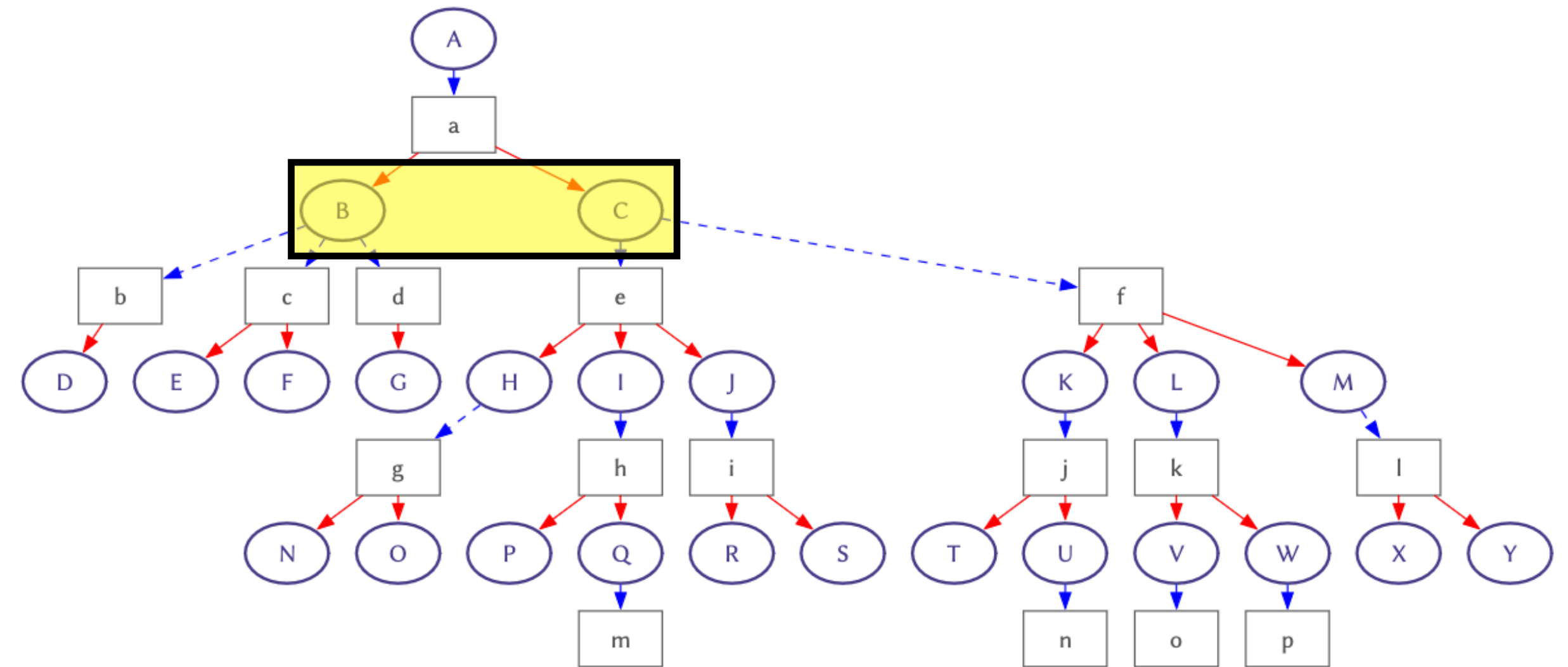
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

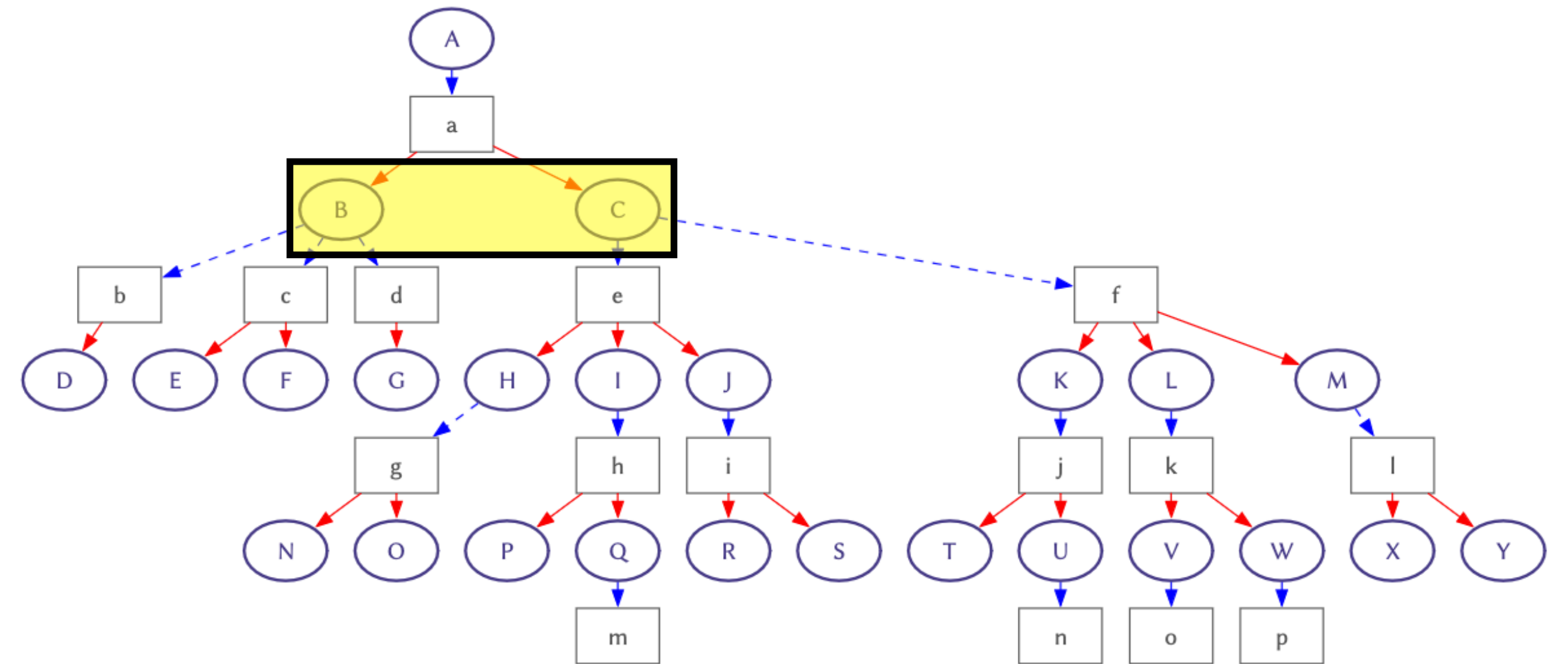


(User selection)

Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

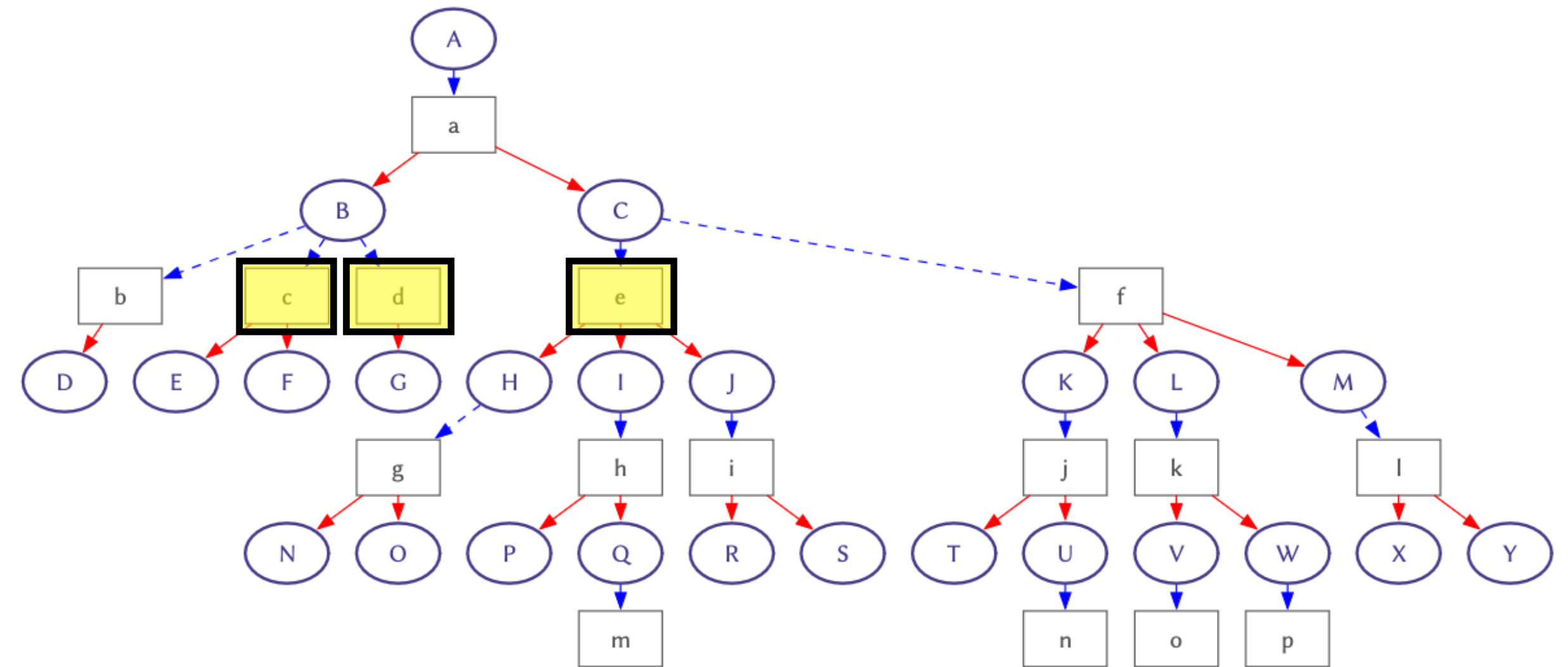
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

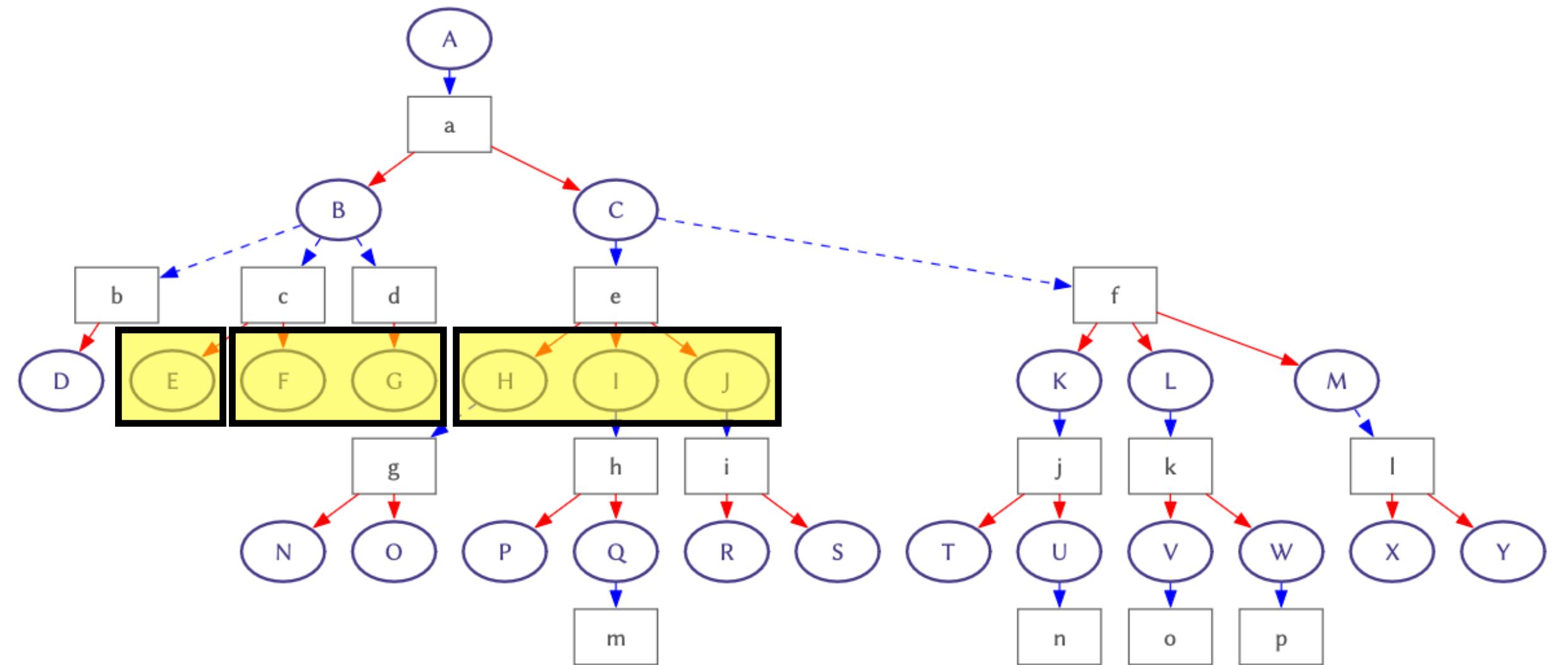
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

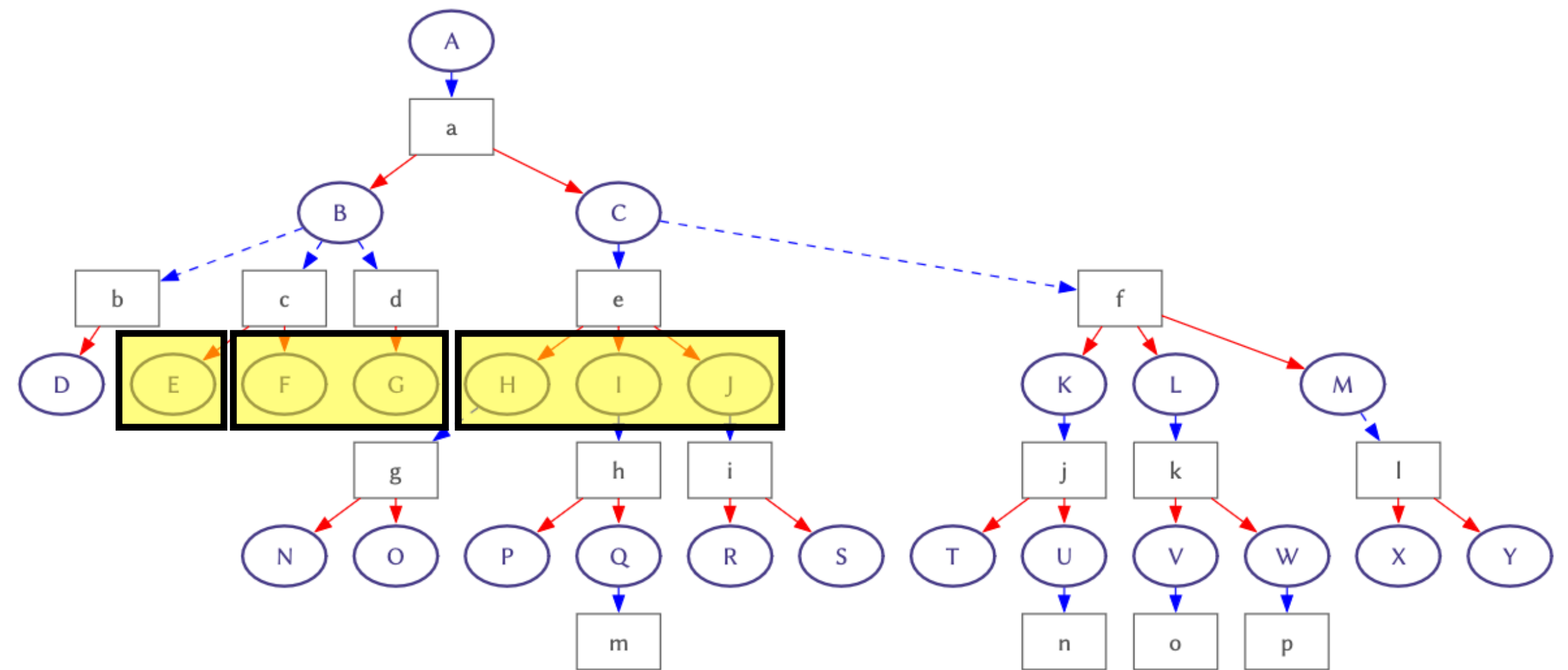
```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Warm-up: TOP-DOWN step provider

(i.e., we don't lose anything with Programming by Navigation)

```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

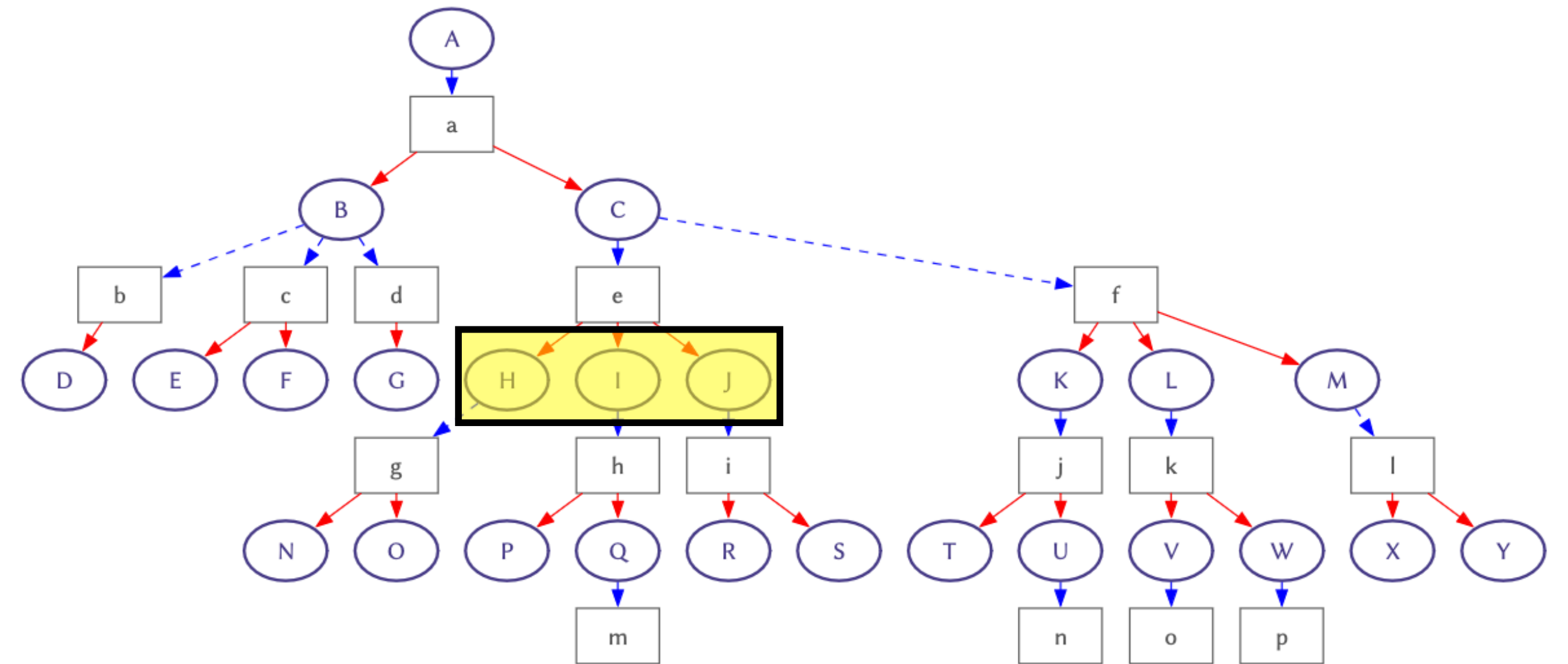


(User selection)

Warm-up: TOP-DOWN step provider

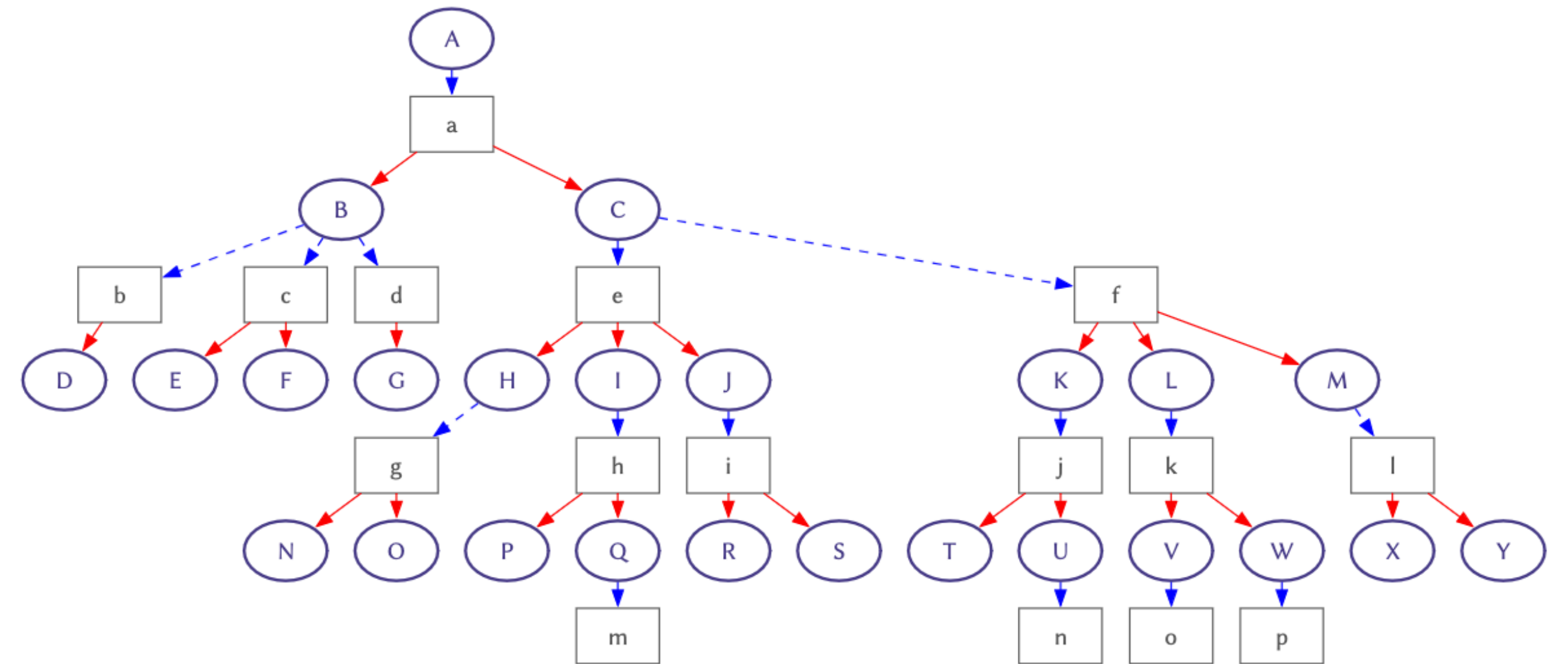
(i.e., we don't lose anything with Programming by Navigation)

```
1: procedure TOPDOWNINVERSION( $e$ )
2:   yield all from COMMIT( $e$ )
3:   for  $Q \in e^{-1}[\{T/A!, T!\}]$  do
4:     for  $R \in \mathcal{R}$  with conclusion  $Q$  do
5:        $P \leftarrow$  premises of  $R$  with label  $\perp$ 
6:       if  $P$  empty then continue
7:        $\sigma \leftarrow P_1 := T/A!; \dots; P_{|P|} := T/A!$ 
8:       if  $e(Q) = T/A!$  then  $\sigma \leftarrow \sigma; Q := T!$ 
9:       if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



Non-topological: RANDOM step provider

- 1: **procedure** RANDOM(e)
- 2: $Q \leftarrow$ random element of $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$
- 3: **for** $\ell > e(Q)$ **do**
- 4: $\sigma \leftarrow Q := \ell$
- 5: **if** SAT($\varphi(\sigma e)$) **then yield** σ



Non-topological: RANDOM step provider

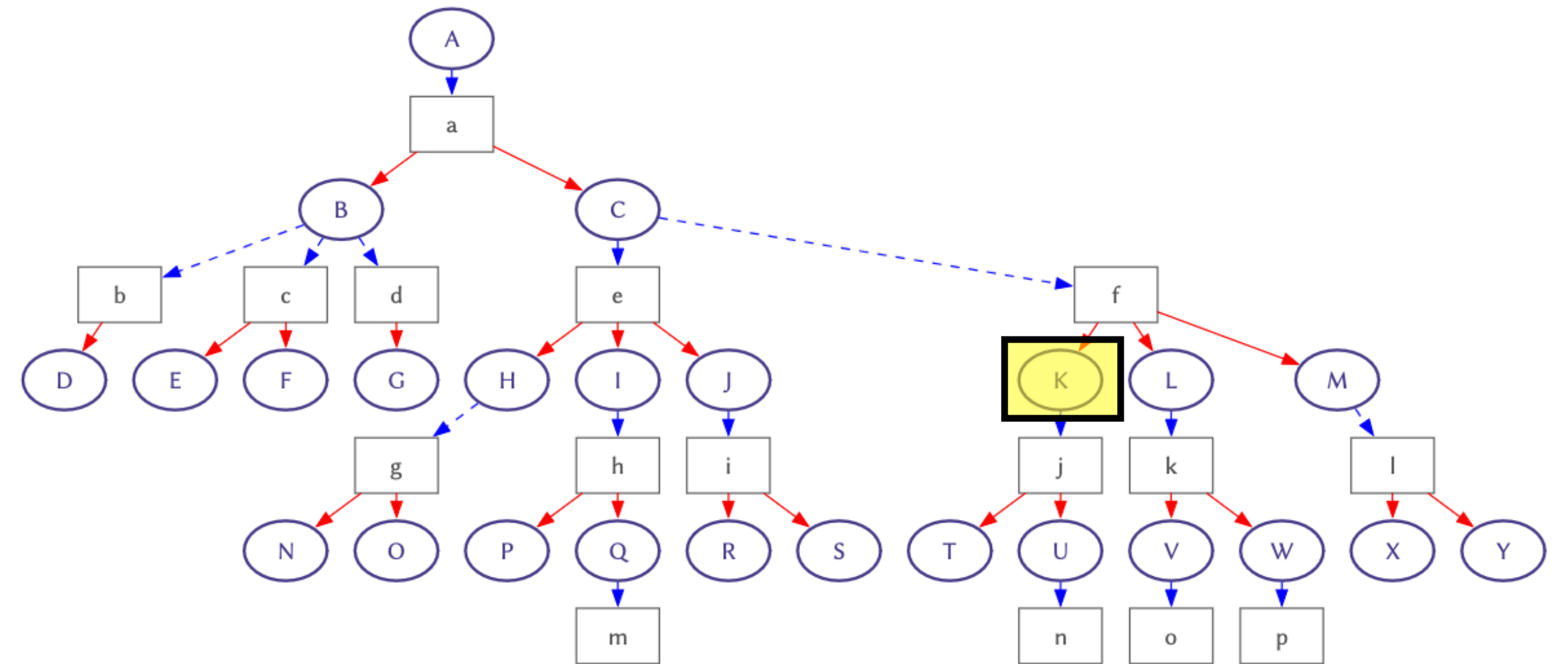
1: **procedure** RANDOM(e)

2: $Q \leftarrow$ random element of $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$

3: **for** $\ell > e(Q)$ **do**

4: $\sigma \leftarrow Q := \ell$

5: **if** SAT($\varphi(\sigma e)$) **then yield** σ



Non-topological: RANDOM step provider

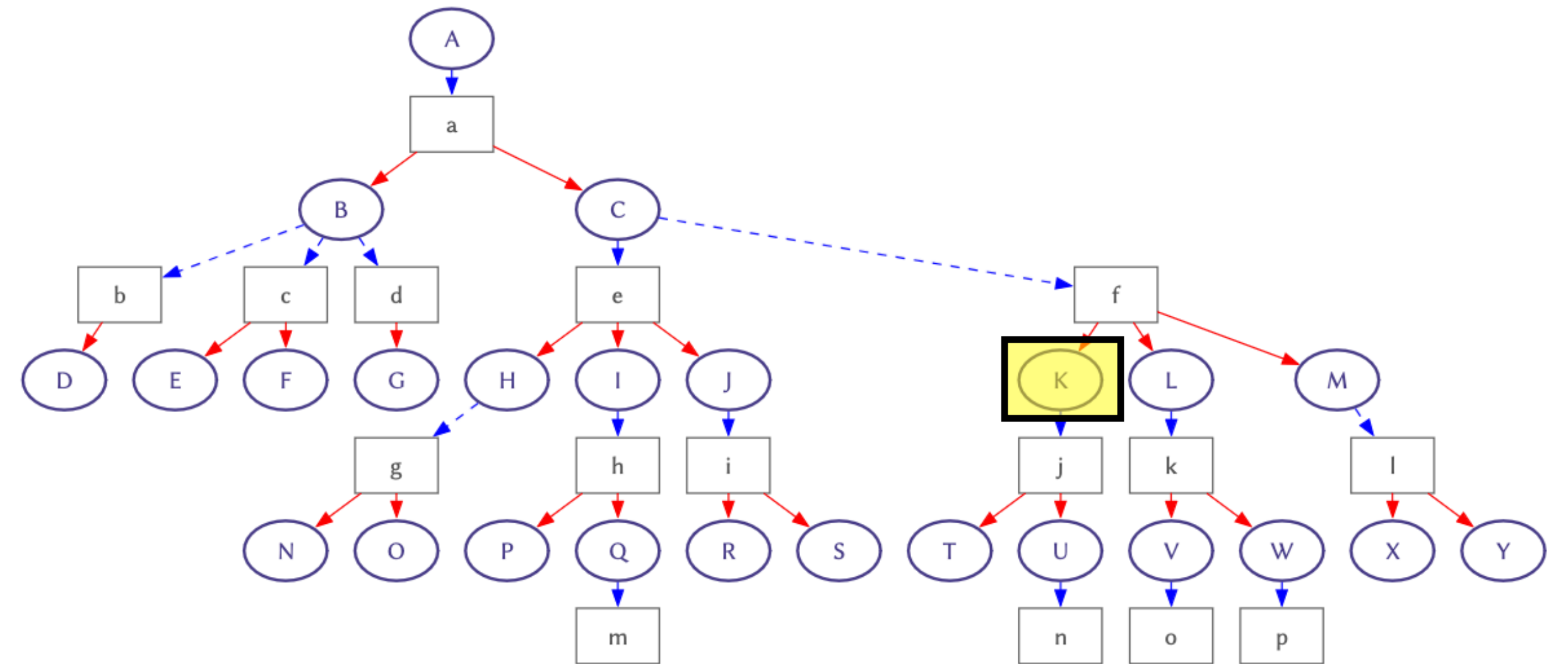
1: **procedure** RANDOM(e)

2: $Q \leftarrow$ random element of $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$

3: **for** $\ell > e(Q)$ **do**

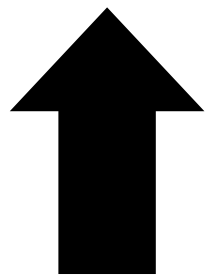
4: $\sigma \leftarrow Q := \ell$

5: **if** SAT($\varphi(\sigma e)$) **then yield** σ

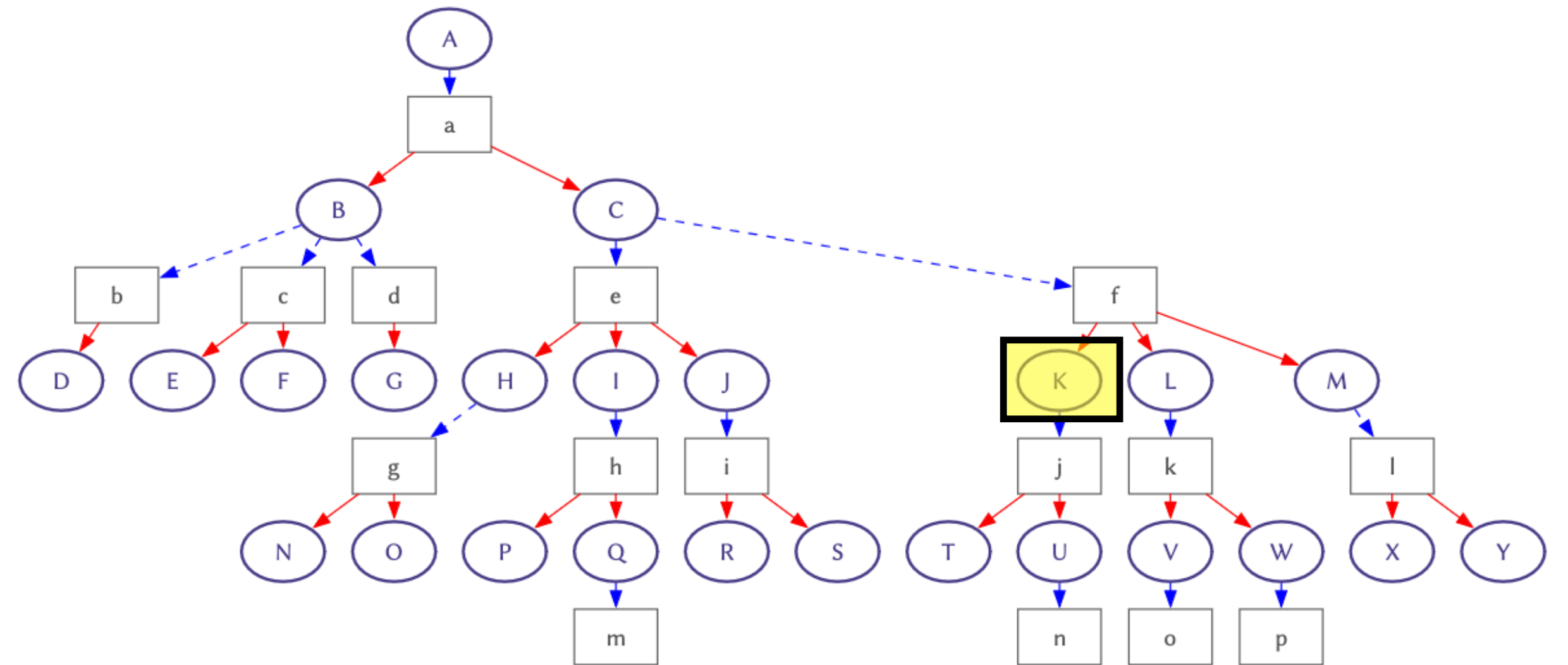


Non-topological: RANDOM step provider

```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

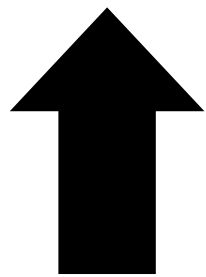


Even though the **propositions** are randomly selected,
the **labels** are still filtered to satisfy **Strong Soundness**

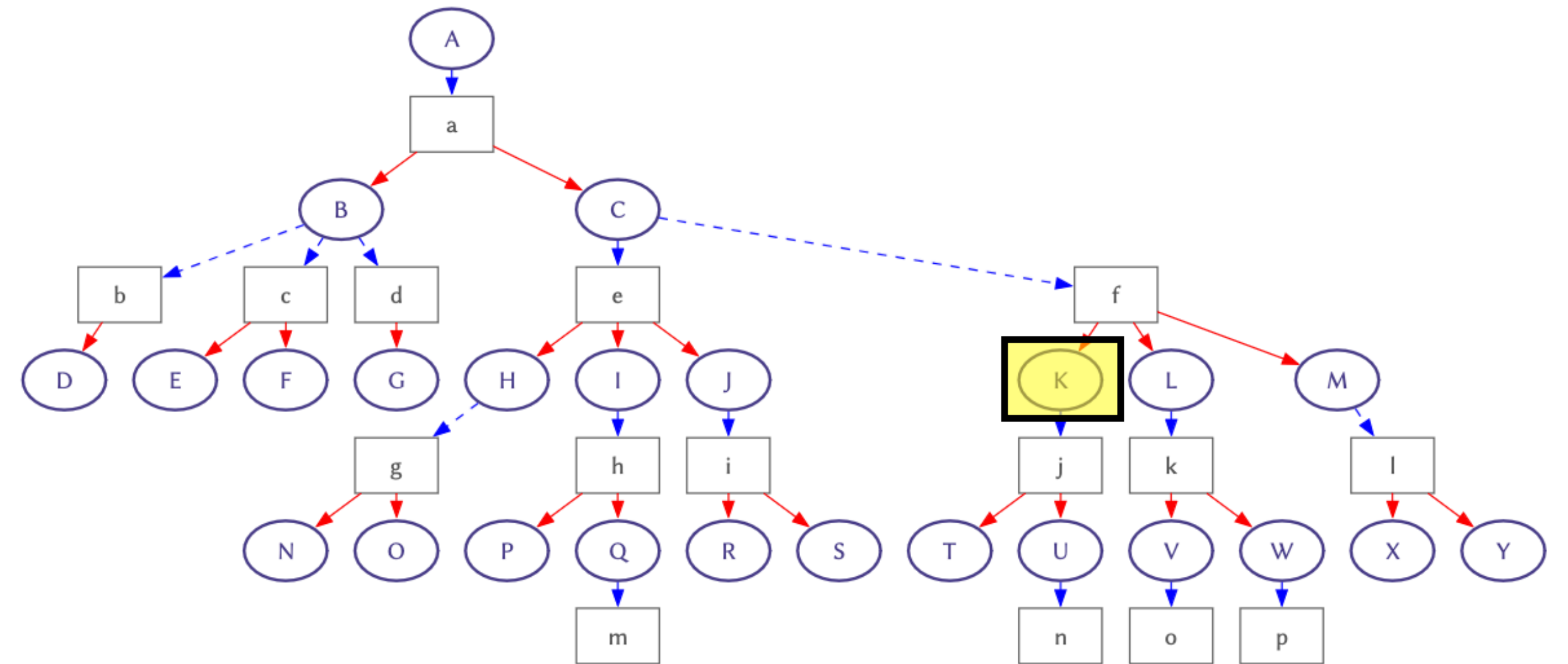


Non-topological: RANDOM step provider

```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



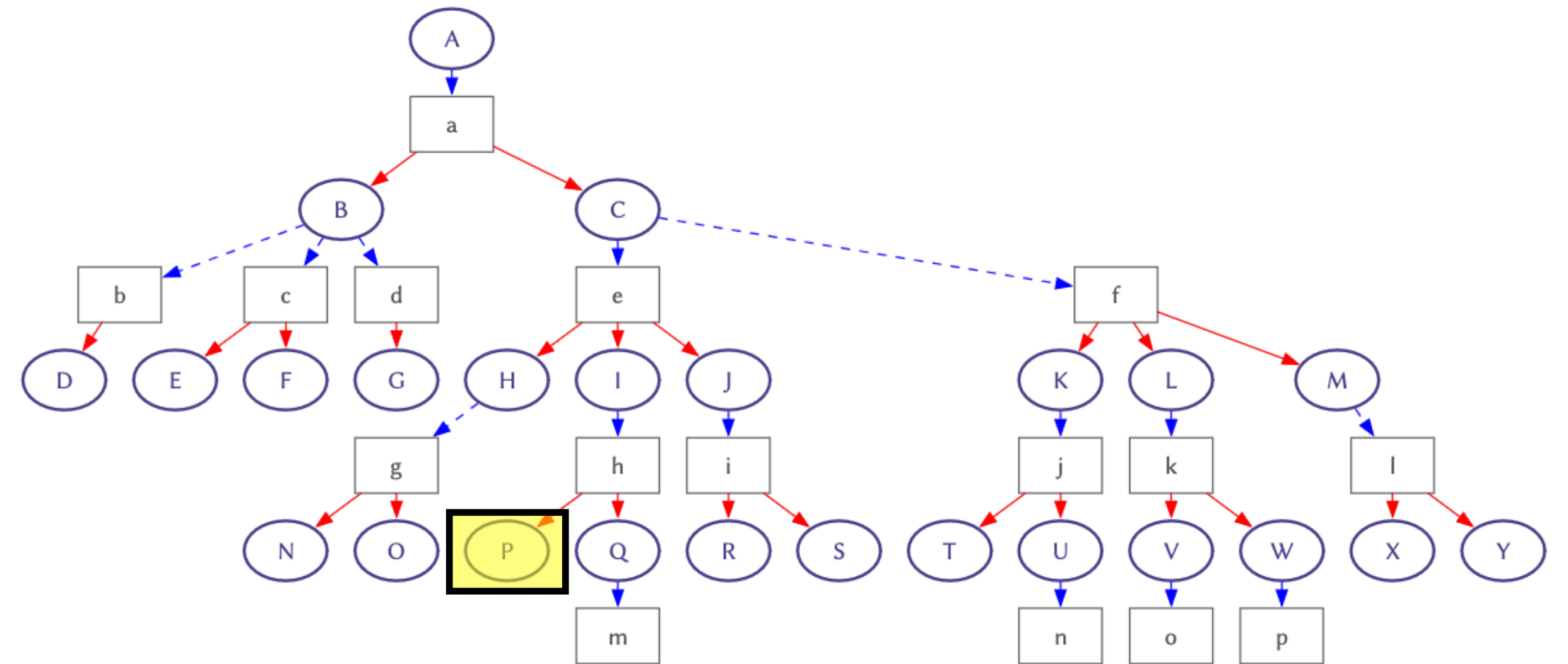
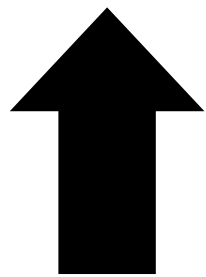
Even though the **propositions** are randomly selected,
the **labels** are still filtered to satisfy **Strong Soundness**



(User selection)

Non-topological: RANDOM step provider

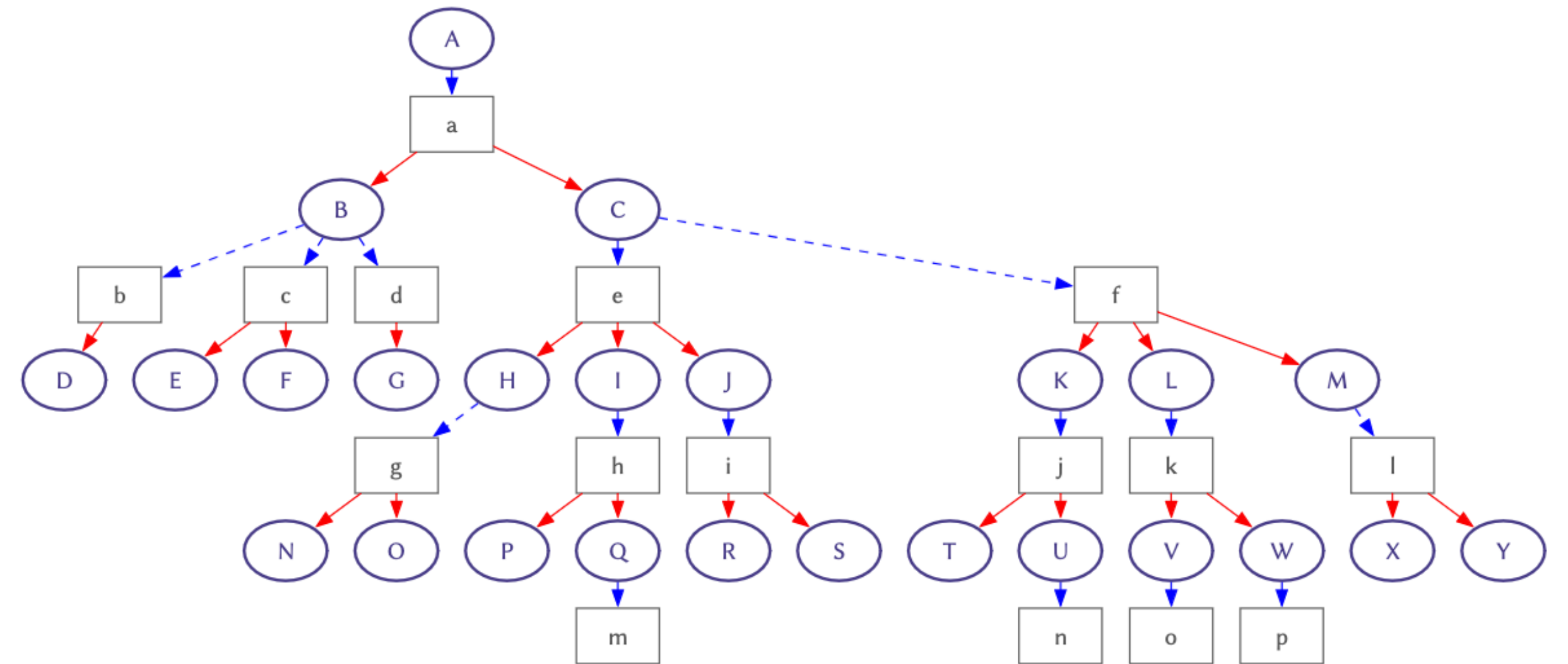
```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```



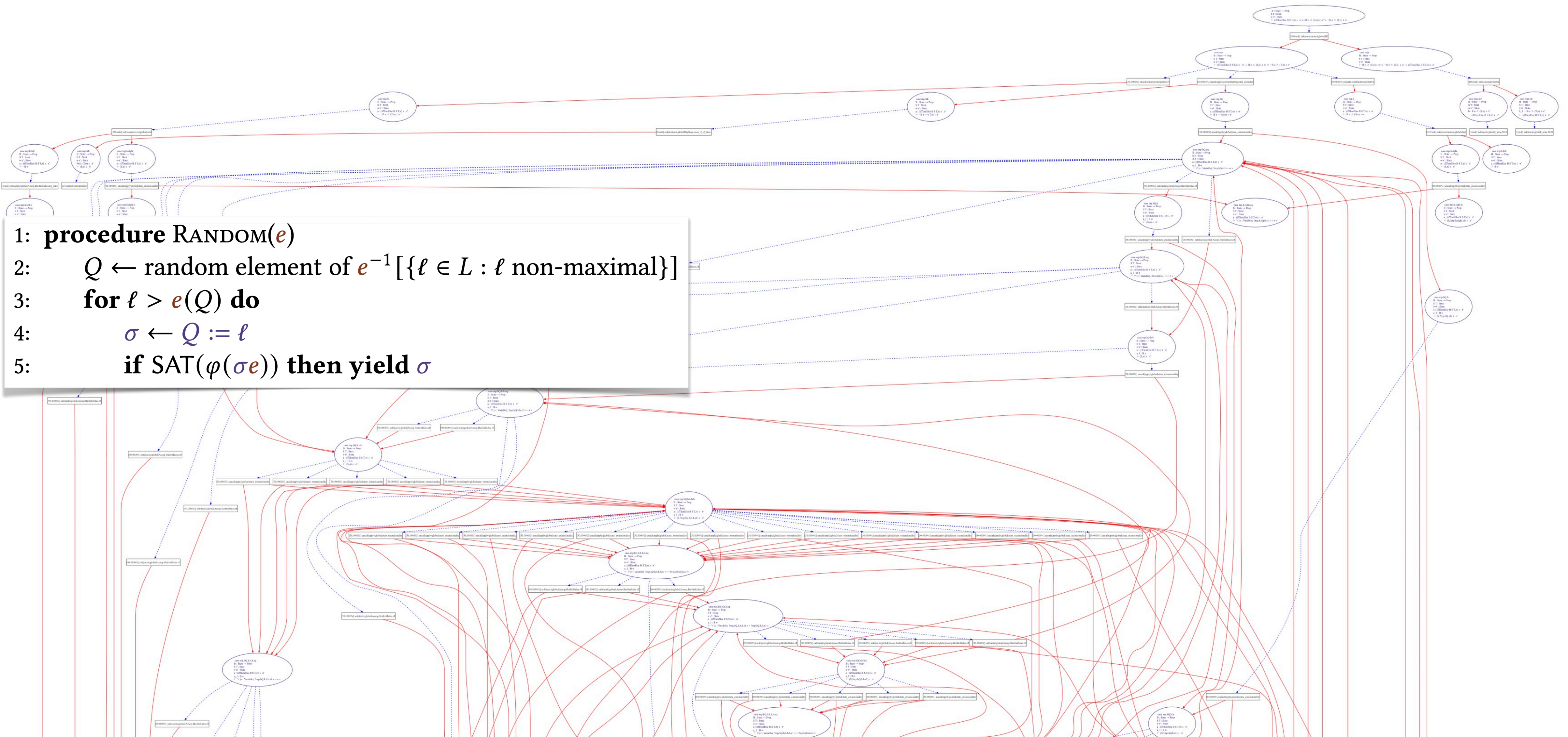
Even though the **propositions** are randomly selected, the **labels** are still filtered to satisfy **Strong Soundness**

Non-topological: CUT step provider

- 1: **procedure** RANDOM(e)
- 2: $Q \leftarrow$ random element of $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$
- 3: **for** $\ell > e(Q)$ **do**
- 4: $\sigma \leftarrow Q := \ell$
- 5: **if** SAT($\varphi(\sigma e)$) **then yield** σ

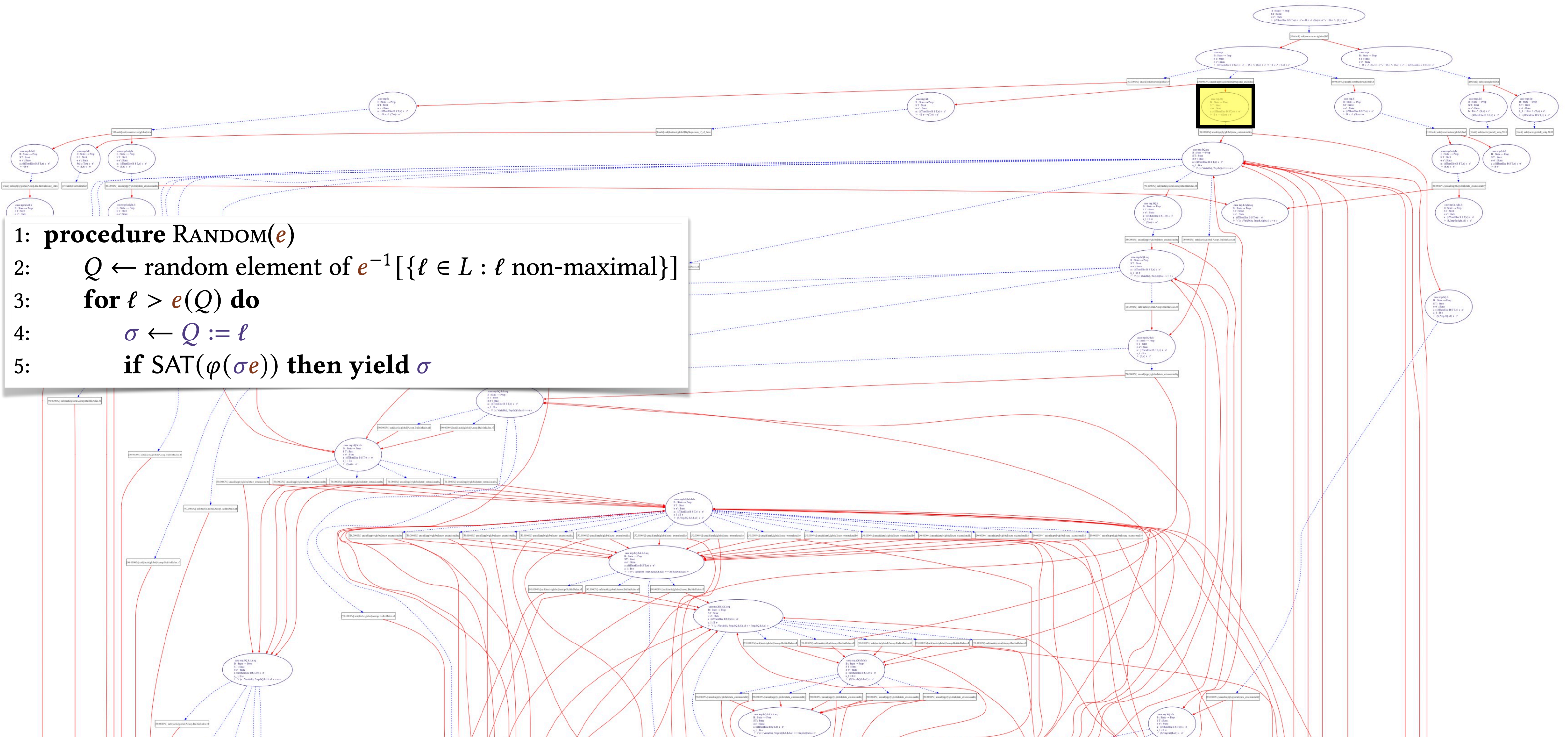


Non-topological: CUT step provider



```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

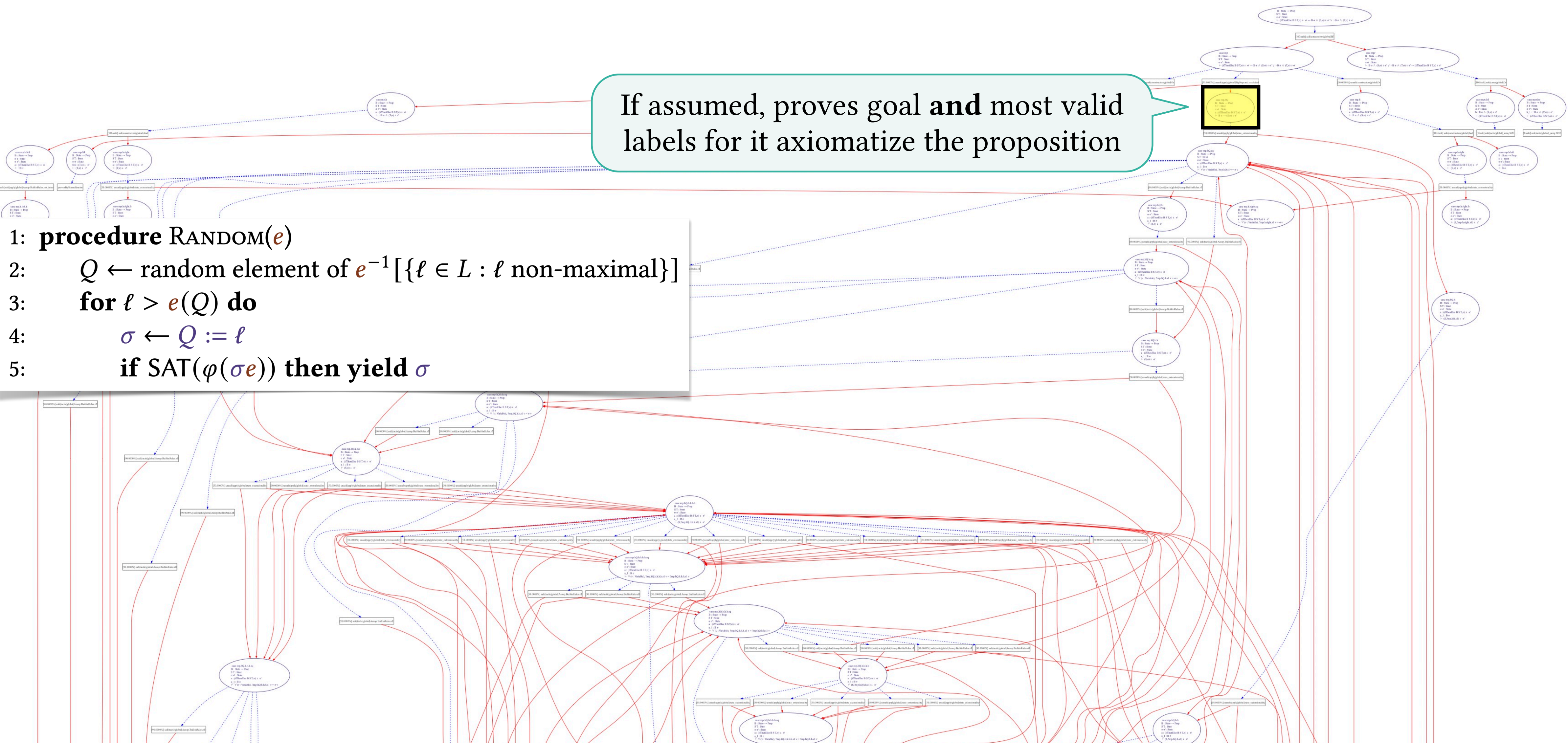
Non-topological: CUT step provider



```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

Non-topological: CUT step provider

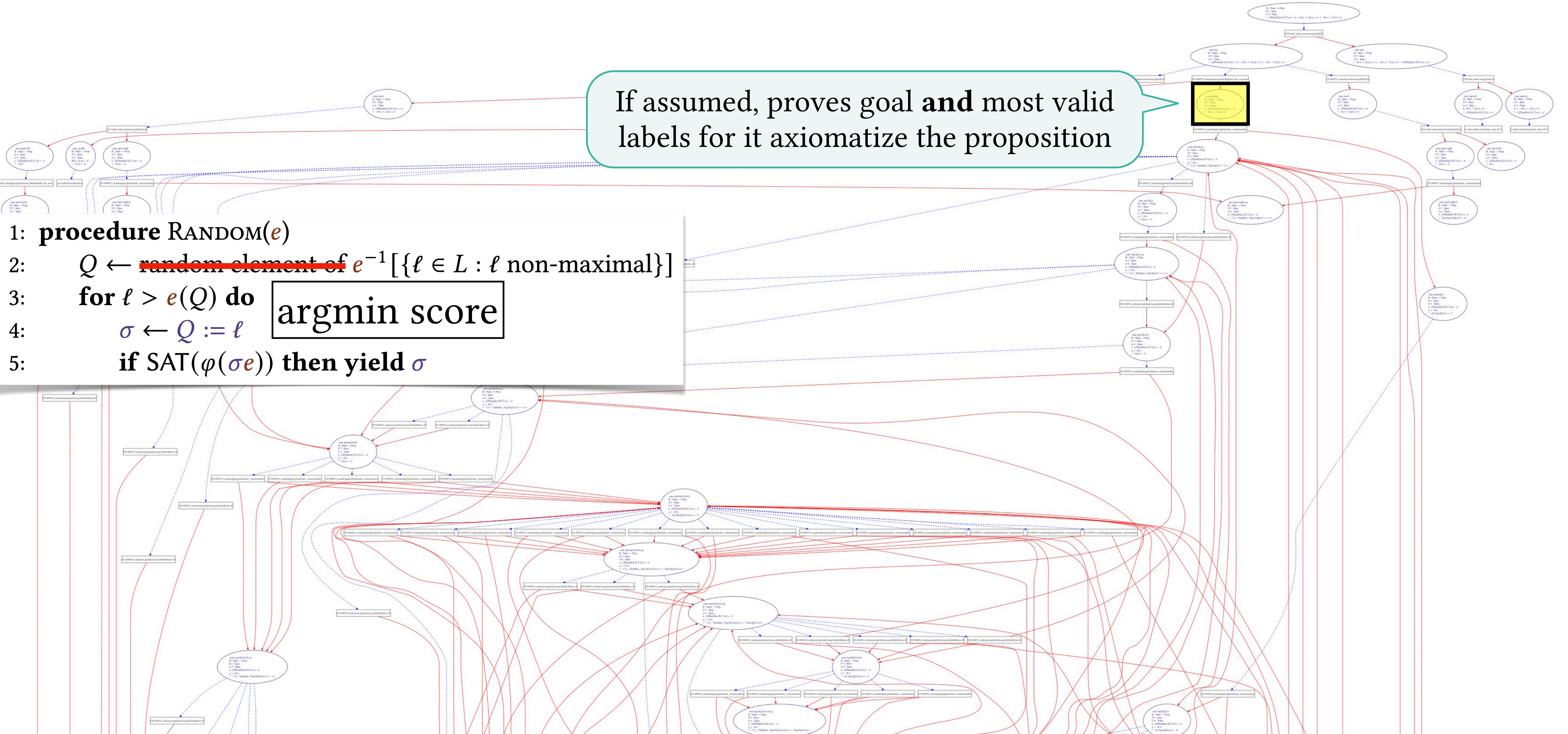
If assumed, proves goal **and** most valid labels for it axiomatize the proposition



```
1: procedure RANDOM( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

Non-topological: CUT step provider

If assumed, proves goal **and** most valid labels for it axiomatize the proposition



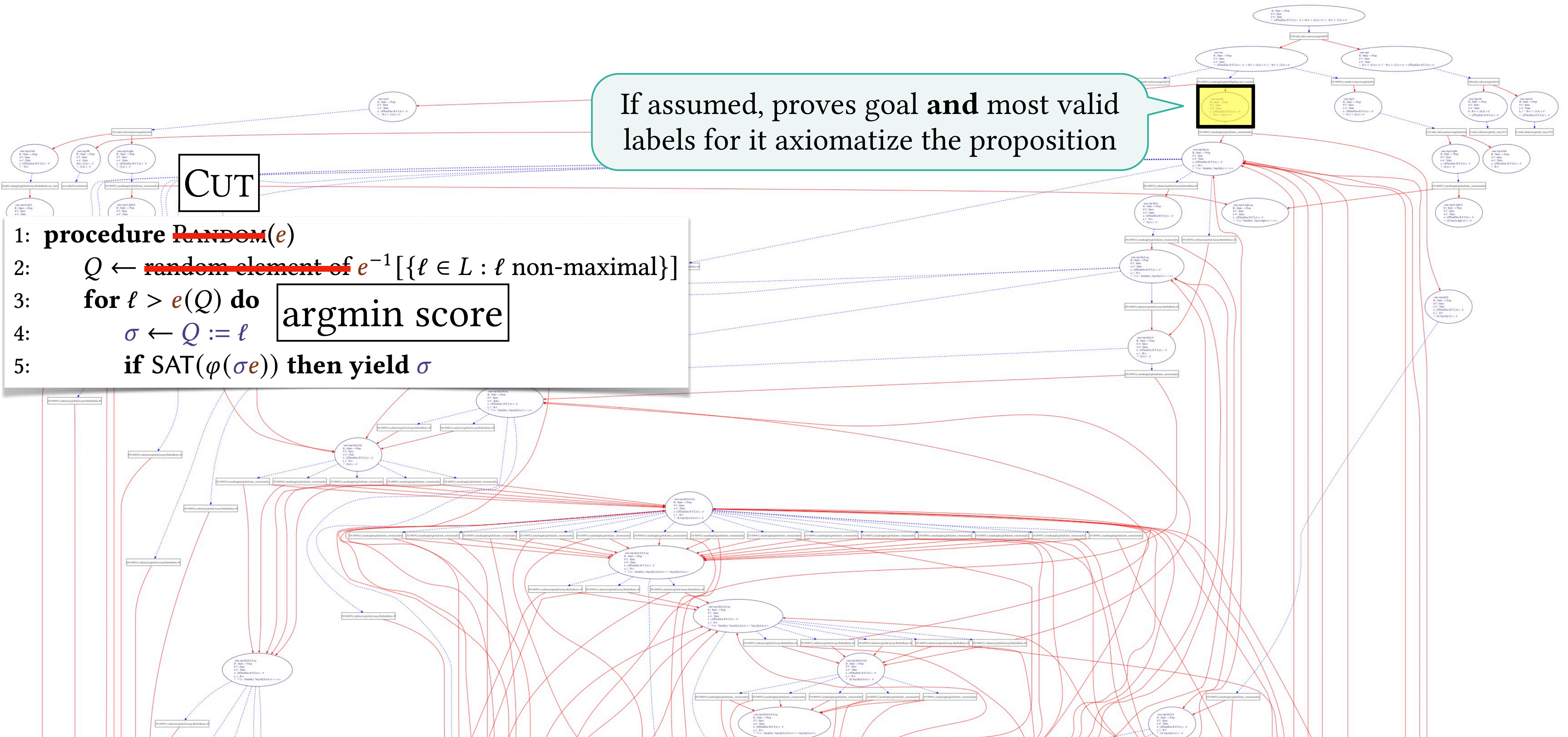
```
1: procedure RANDOM(e)
2:   Q ← random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do argmin score
4:      $\sigma \leftarrow Q := \ell$ 
5:   if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

Non-topological: CUT step provider

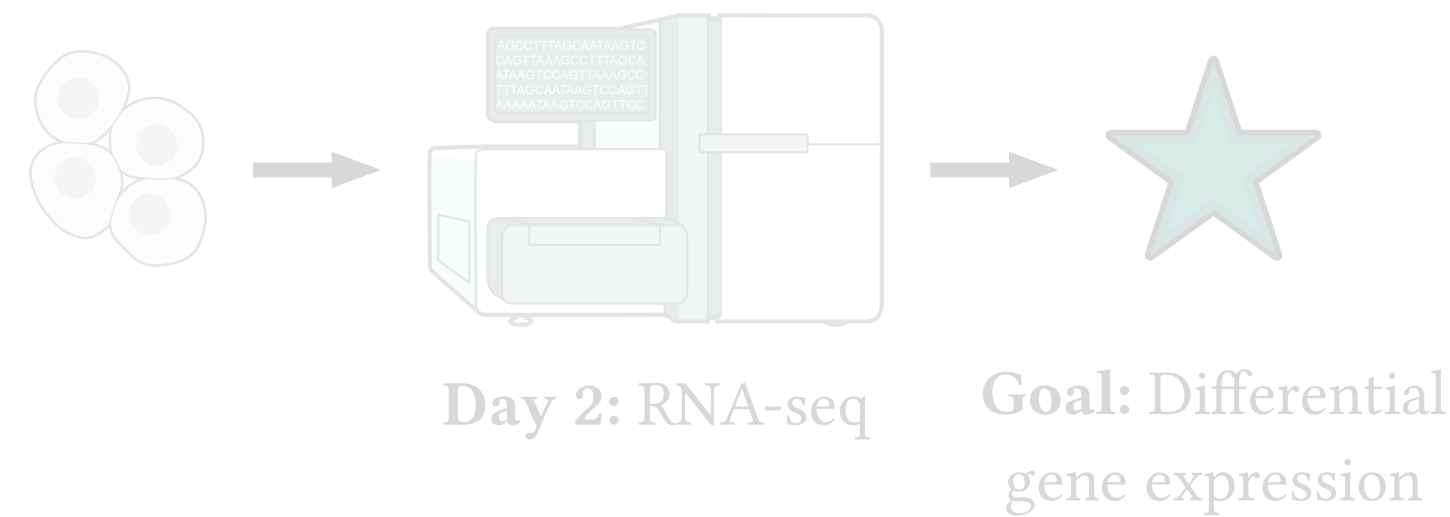
If assumed, proves goal **and** most valid labels for it axiomatize the proposition

CUT

```
1: procedure RANDOM(e)
2:   Q ← random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do argmin score
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

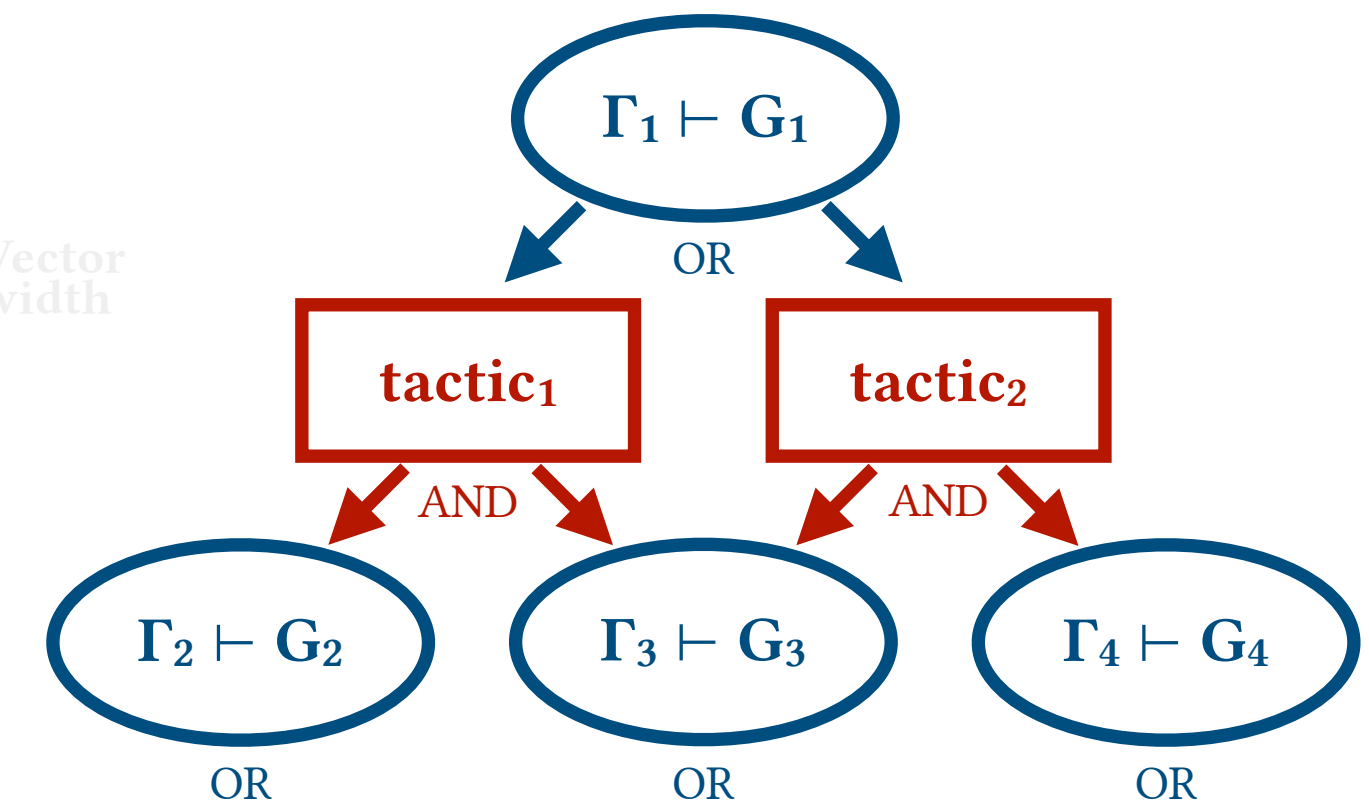


Programming by Navigation Ingredients



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor, Memory layout, Vector width, Dataflow ordering, Memory layout



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

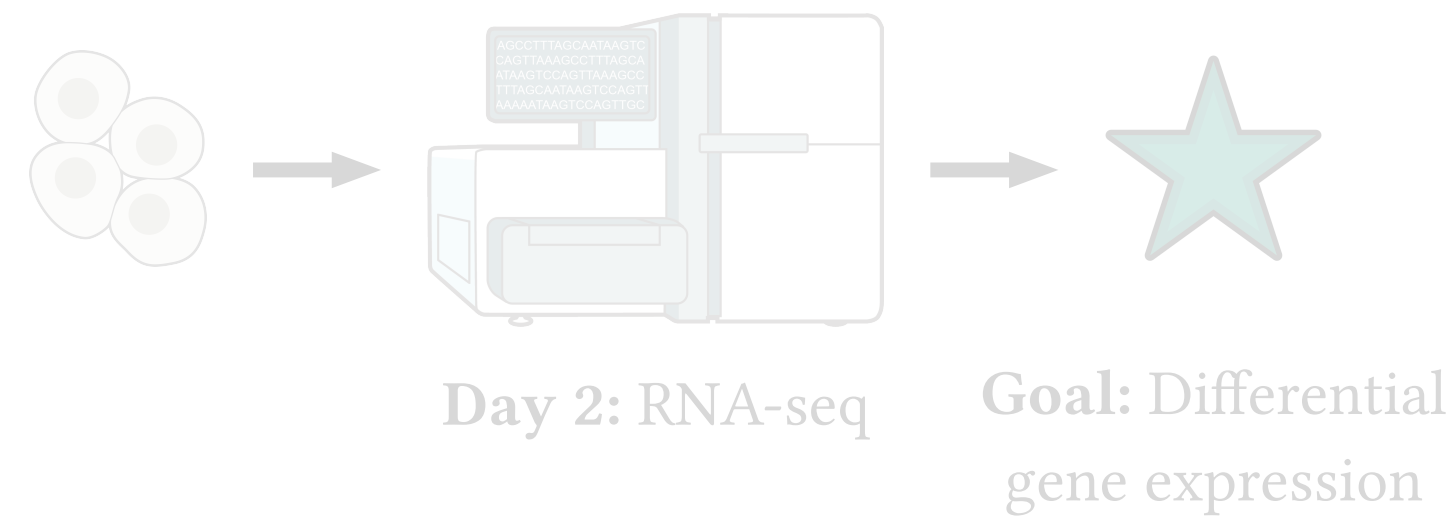
Modification proves final goal node

Step provider

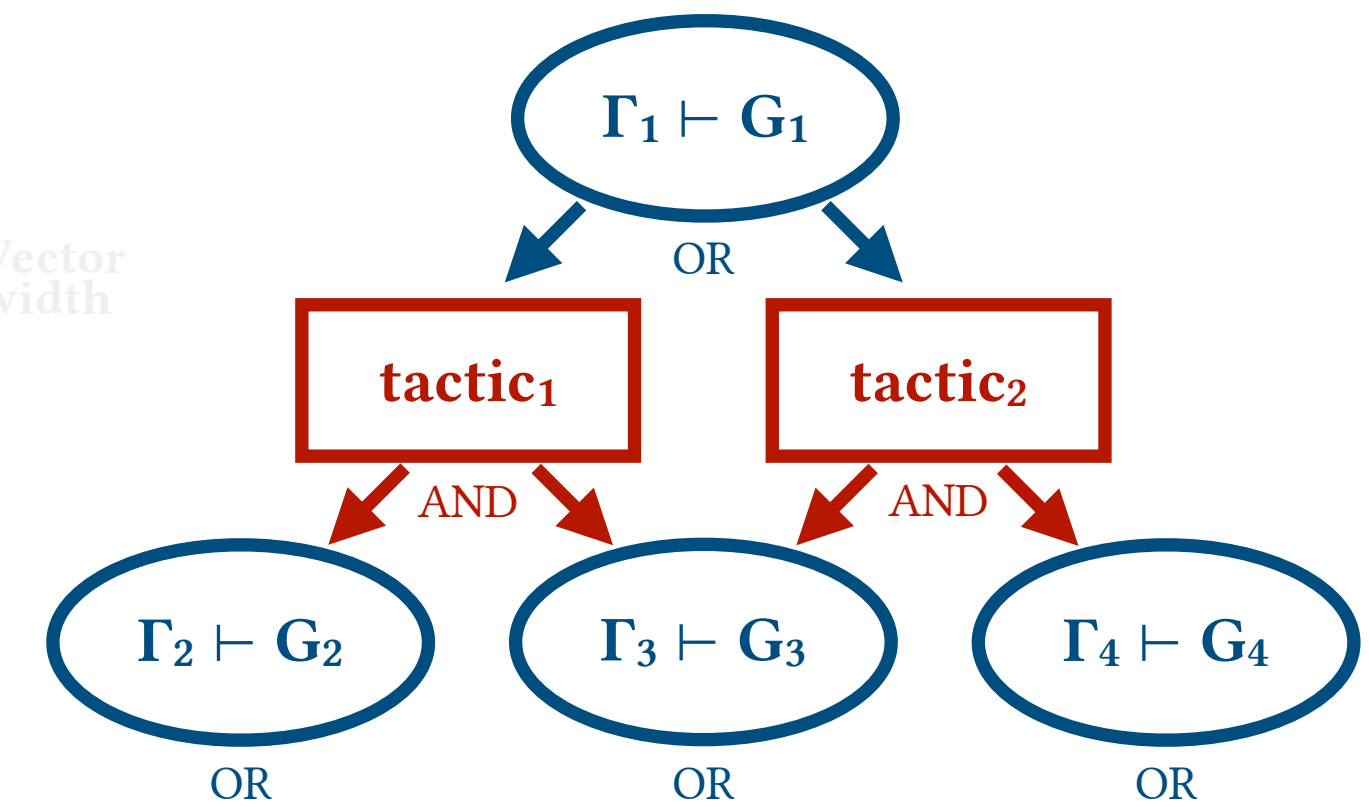
Top-down classical-constructive synthesis

Big-step synthesis

Programming by Navigation Ingredients



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \quad \text{Parallelization factor } \boxed{?} \quad \boxed{?} \quad \text{Memory layout} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \quad \text{Dataflow ordering} \quad \boxed{?} \quad \text{Memory layout} \end{array} \right\}$$



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

Modification proves final goal node

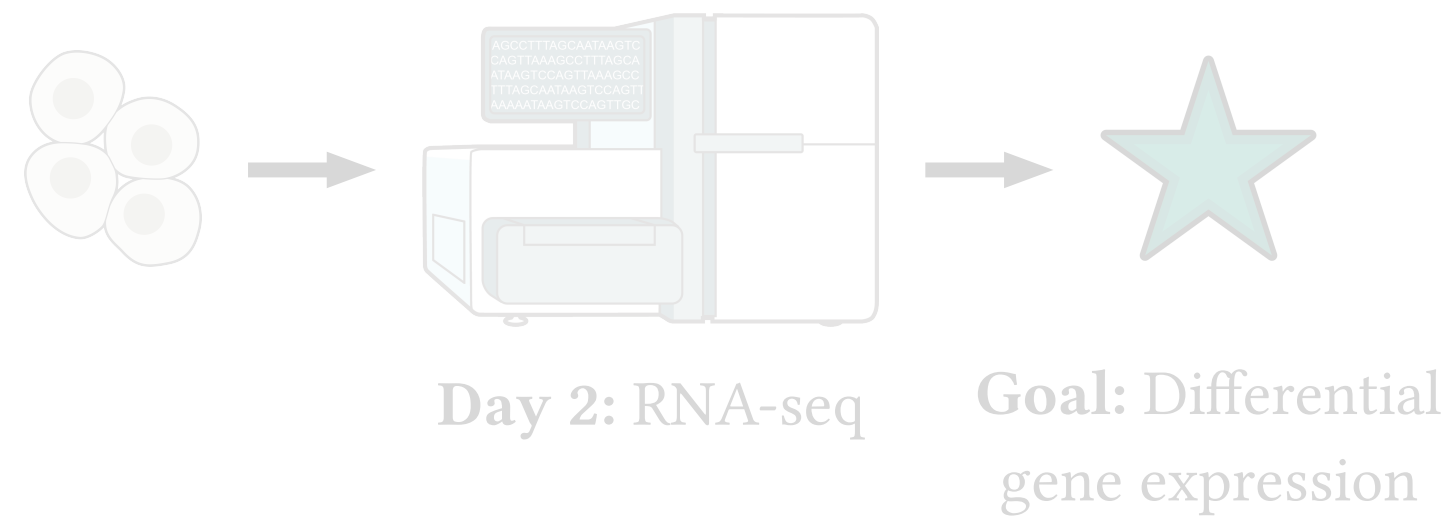
Step provider

Top-down classical-constructive synthesis

Big-step synthesis

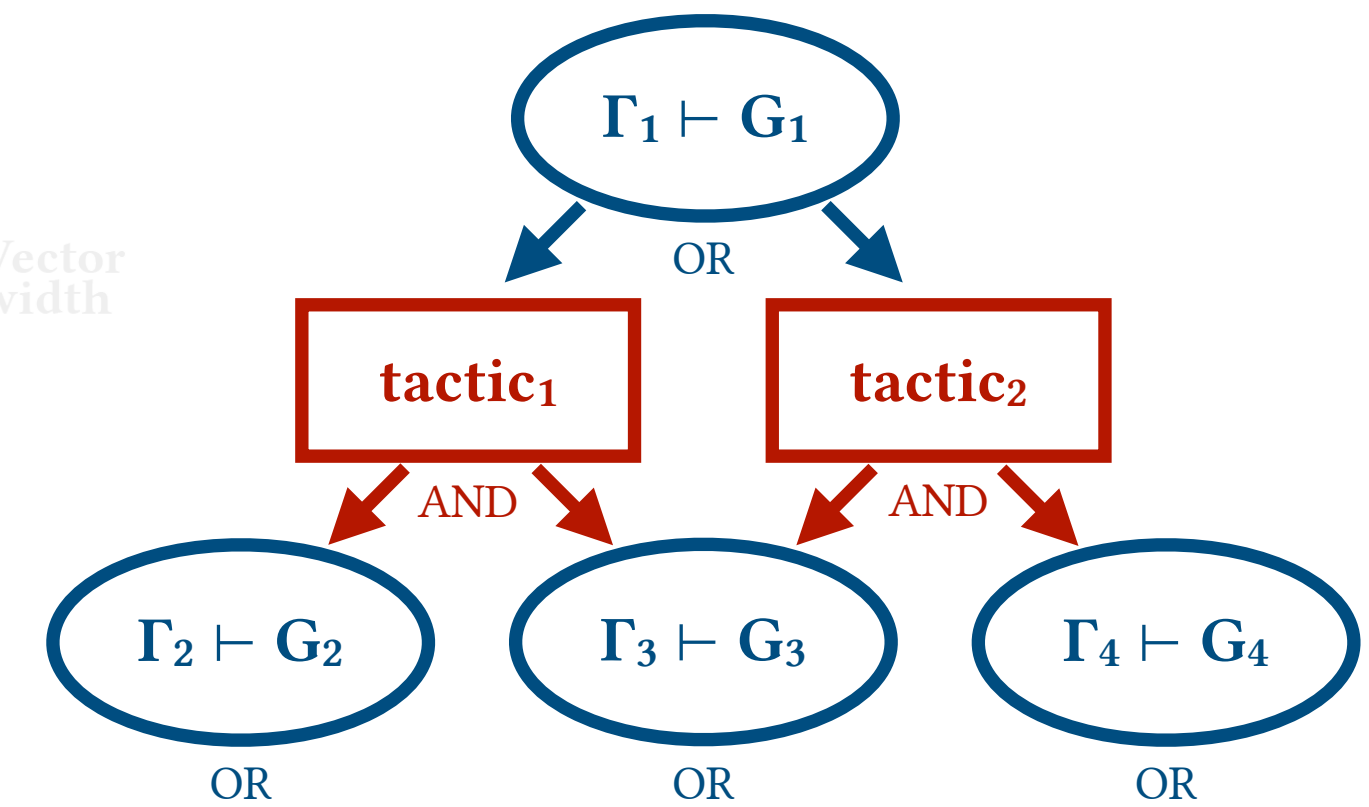
Maximum Information Gain, Cut, Random, Top-Down, Bottom-Up, Minimal Leaf Heuristic, ...

Programming by Navigation Ingredients



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{ij} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor, Memory layout, Vector width, Dataflow ordering, Memory layout.



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

Modification proves final goal node

Step provider

Top-down classical-constructive synthesis

Big-step synthesis

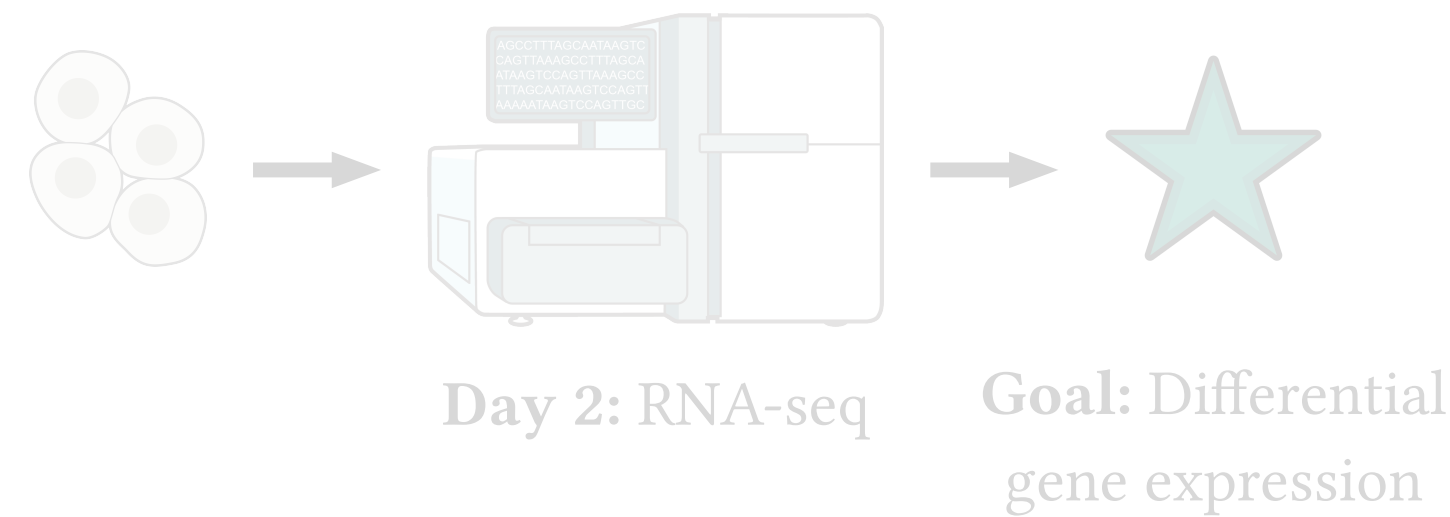
Maximum Information Gain, Cut, Random, Top-Down, Bottom-Up, Minimal Leaf Heuristic, ...

Nonempty-completion oracle

Datalog

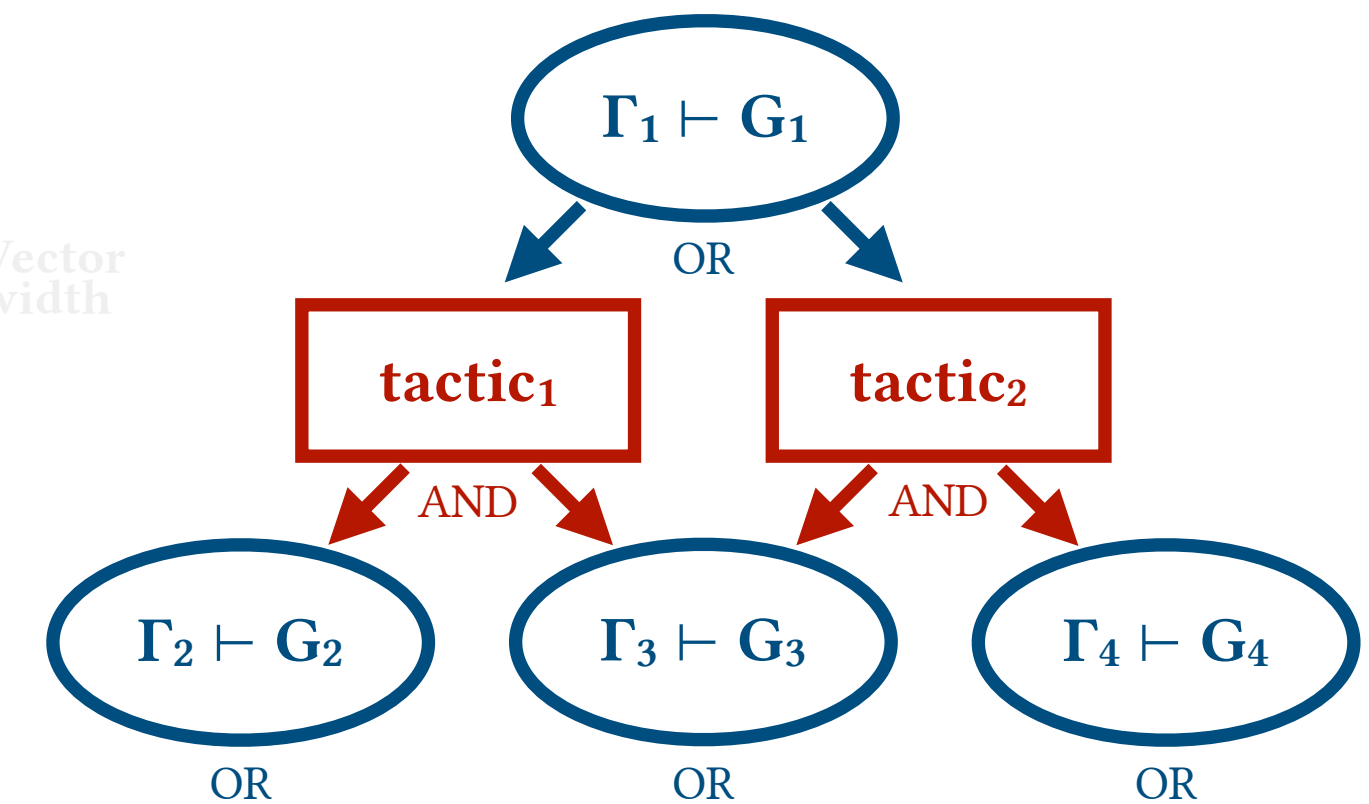
FuseFlow compiler + Datalog

Programming by Navigation Ingredients



$$\text{Fuse} \left\{ \begin{array}{l} \forall_{ikj} A_{ij} = B_{ik} C_{kj} \\ \forall_{?} E_{ij} = A_{ik} D_{kj} \end{array} \right\}$$

Annotations: Parallelization factor, Memory layout, Vector width, Dataflow ordering, Memory layout.



Expressions

Bioinformatics pipelines

FuseFlow compiler schedules

Proof system modifications

Validity

Meets user-defined validity conditions

Meets FuseFlow-defined validity conditions

Modification proves final goal node

Step provider

Top-down classical-constructive synthesis

Big-step synthesis

Maximum Information Gain, Cut, Random, Top-Down, Bottom-Up, Minimal Leaf Heuristic, ...

Nonempty-completion oracle

Datalog

FuseFlow compiler + Datalog

SAT

Navigating AND-OR Graph Modifications to Debug Failing Proof Search



Marlena Preigh



Max Willsey



Sarah E. Chasins

justinlubin@berkeley.edu

Expressions

Proof system modifications

Validity

Modification proves final goal node

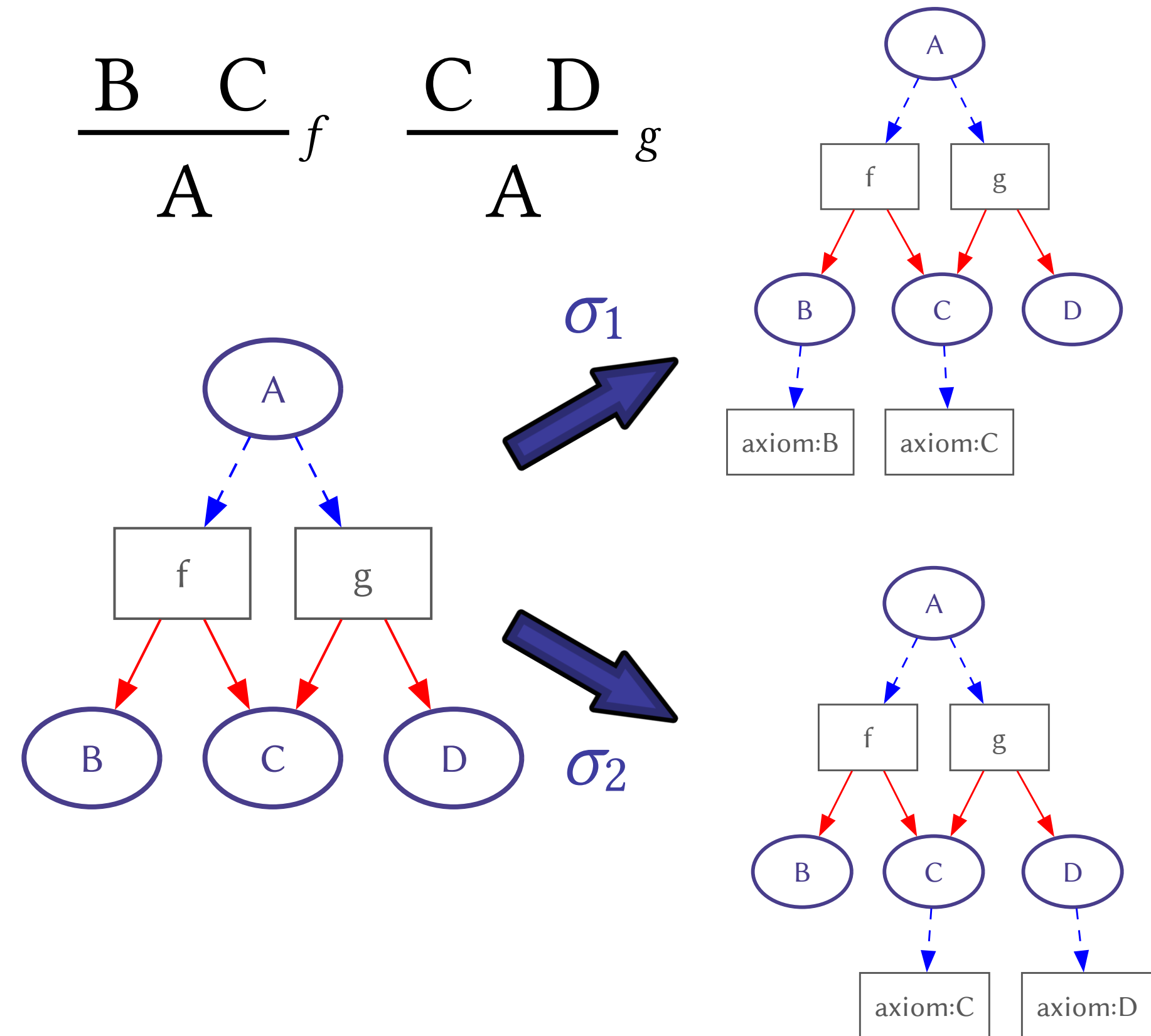
Step provider

Maximum Information Gain, Cut, Random, Top-Down, Bottom-Up, Minimal Leaf Heuristic, ...

Nonempty-completion oracle

SAT

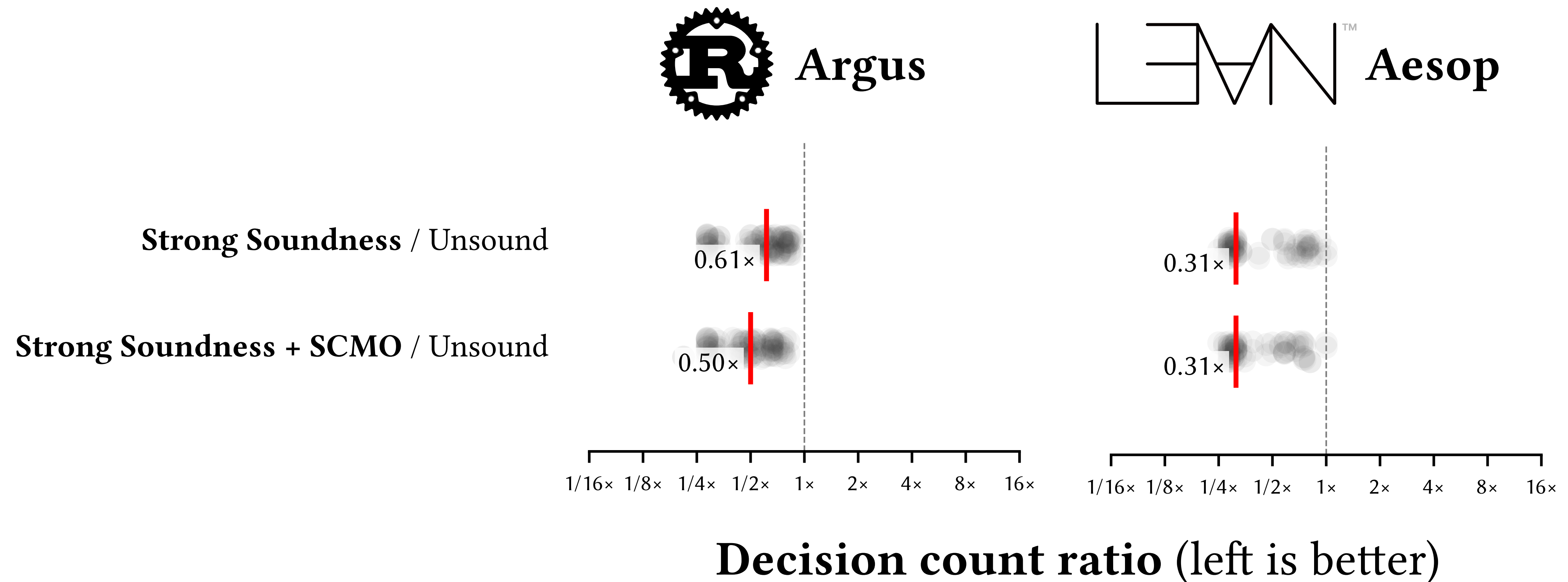
$$\frac{B \quad C}{A} f \quad \frac{C \quad D}{A} g$$





Extra Slides

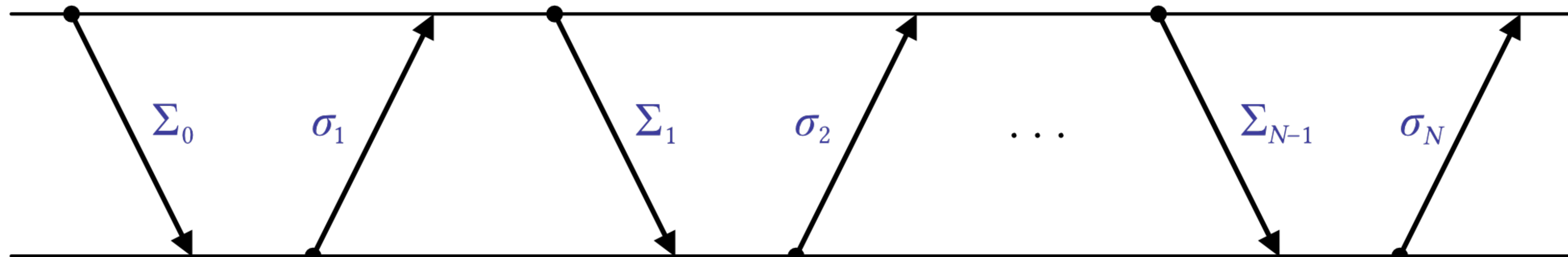
Strong Completeness Modulo Observability (SCMO) reduces decision counts



Gray, Crichton, and Krishnamurthi. An Interactive Debugger for Rust Trait Errors (*PLDI 2025*).

Limperg and From. Aesop: White-Box Best-First Proof Search for Lean (*CPP 2023*).

Step Provider



Step Decider

CUT

```
1: procedure Random( $e$ )
2:    $Q \leftarrow$  random element of  $e^{-1}[\{\ell \in L : \ell \text{ non-maximal}\}]$ 
3:   for  $\ell > e(Q)$  do argmin score
4:      $\sigma \leftarrow Q := \ell$ 
5:     if SAT( $\varphi(\sigma e)$ ) then yield  $\sigma$ 
```

Syntax, semantics, and steps for assume-partitions

- Expression e : Partition of propositions
 - *Notation*: $e(Q)$ is label of proposition Q in e
 - Label of goal is **T!** (must be true and used, but not assumed)

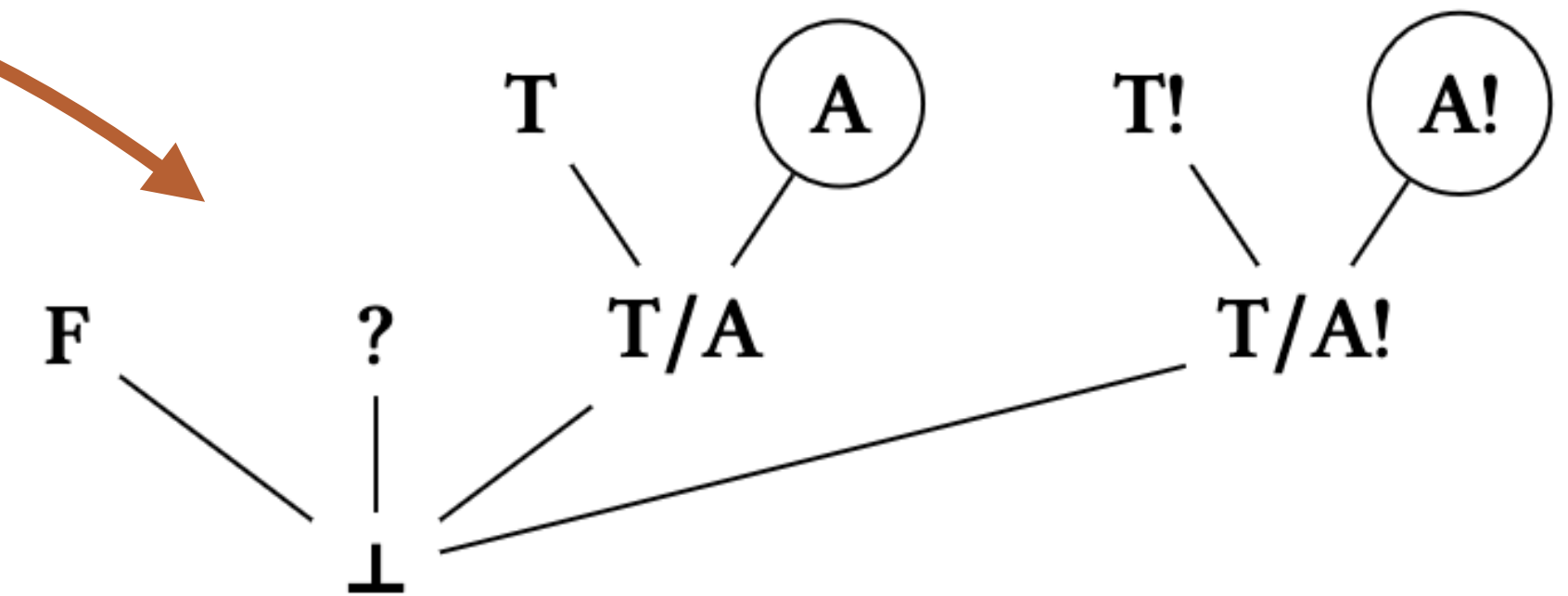
- Semantics $\llbracket e \rrbracket$: Axiomatizes nodes labeled with **A** or **A!**

- e **valid** iff:

1. $e(Q)$ maximal for all Q
2. $\llbracket e \rrbracket \mathbf{P} \vdash Q$ for $e(Q) = \mathbf{T}$ or $e(Q) = \mathbf{T!}$
3. $\llbracket e \rrbracket \mathbf{P} \not\vdash Q$ for $e(Q) = \mathbf{F}$
4. There exists a derivation of goal containing all propositions Q with $e(Q) = \mathbf{T!}$ or $e(Q) = \mathbf{A!}$

- Step $\sigma ::= Q := l \mid \sigma_1 ; \sigma_2$

with these labels



$$\frac{e(Q) < \ell}{e \xrightarrow{Q := \ell} [Q := \ell]e} \text{ STEP/SET}$$

$$\frac{e_1 \xrightarrow{\sigma_1} e_2 \quad e_2 \xrightarrow{\sigma_2} e_3}{e_1 \xrightarrow{\sigma_1 ; \sigma_2} e_3} \text{ STEP/SEQ}$$

with these semantics

Model proof search debugging as interactive program synthesis (specification refinement)

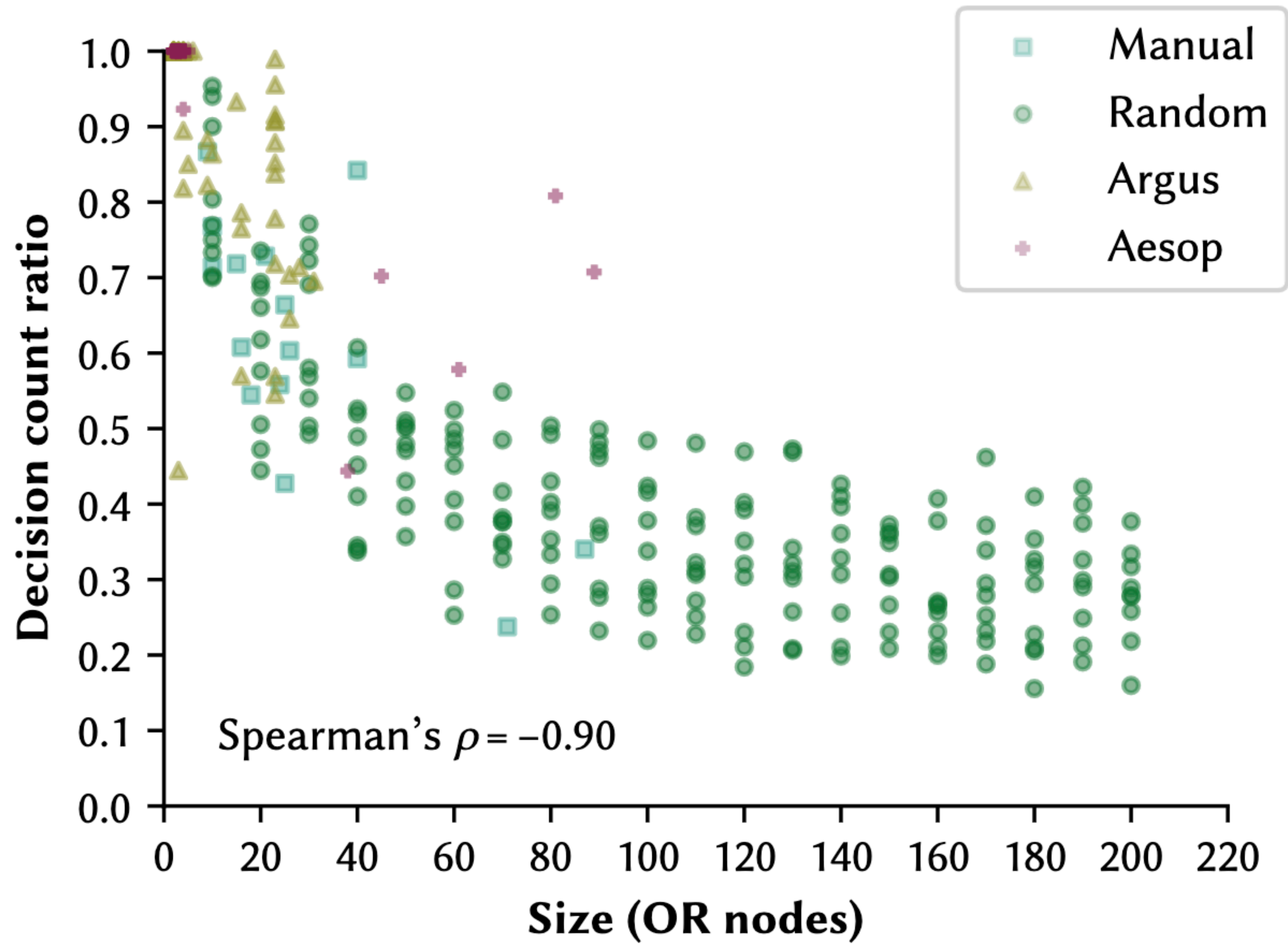
Goal: Synthesize modification to proof system that proves the goal... but not just any modification!

Ingredients

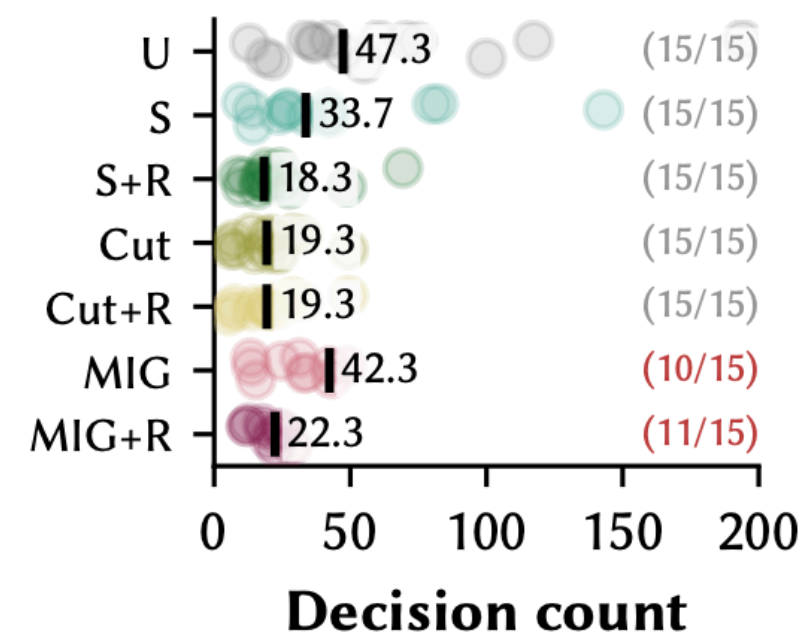
- Expressions e
- Semantics $\llbracket e \rrbracket$, maps proof systems to proof systems
 - *E.g.*, axiomatize a proposition, introduce new rules, remove rules, *etc.*
- For a proof system \mathbf{P} and goal proposition G , a notion of validity (written e valid) such that e valid implies $\llbracket e \rrbracket \mathbf{P} \vdash G$.
 - Specifies “successful” modifications

Goal restated: Synthesize e such that e valid

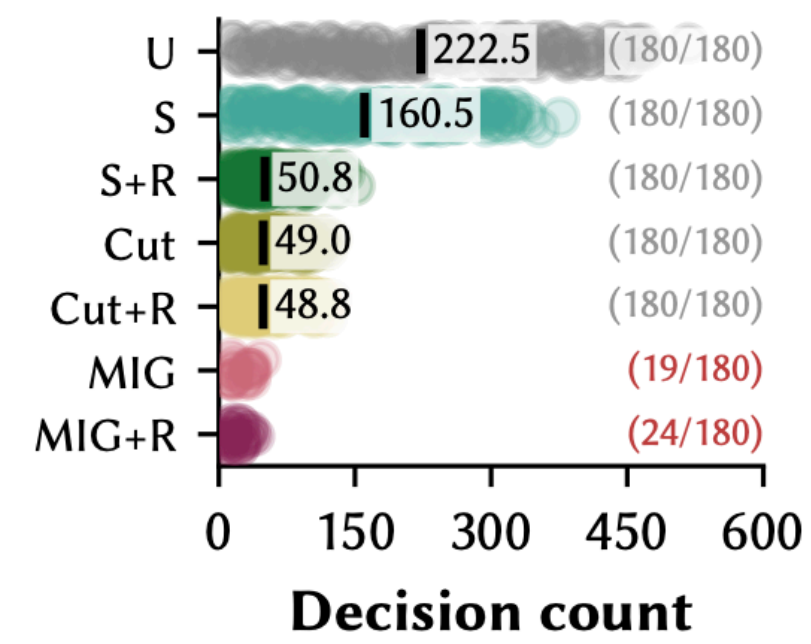
... but this is an **underspecification**; need to interactively **refine** specification with user



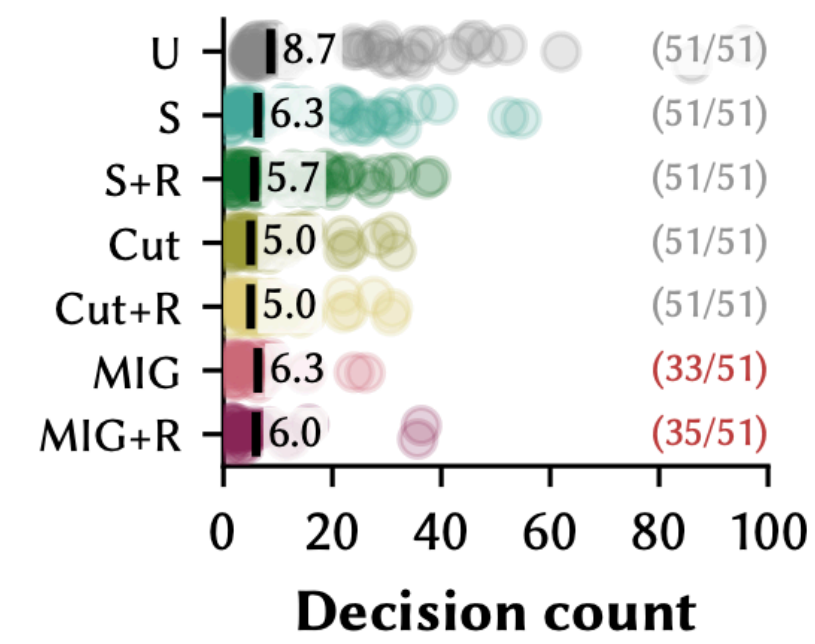
MANUAL



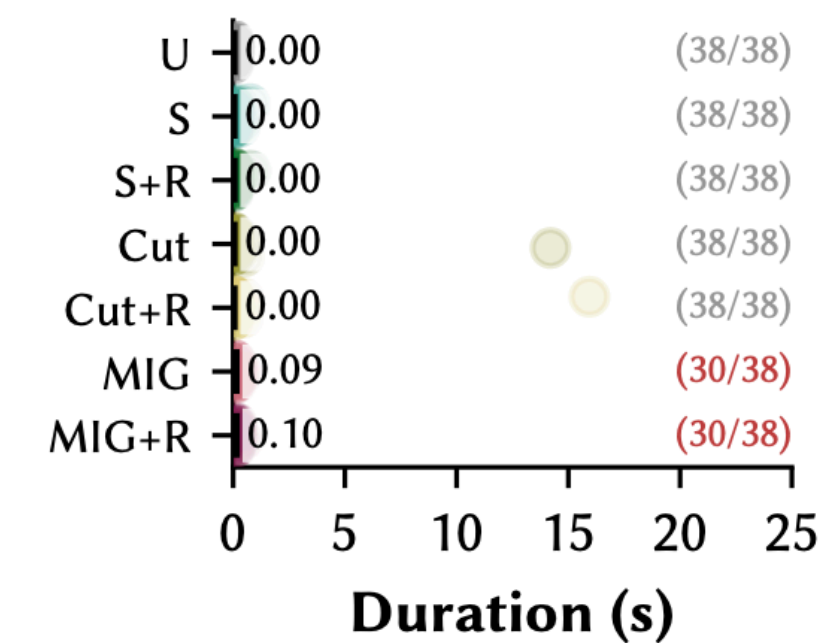
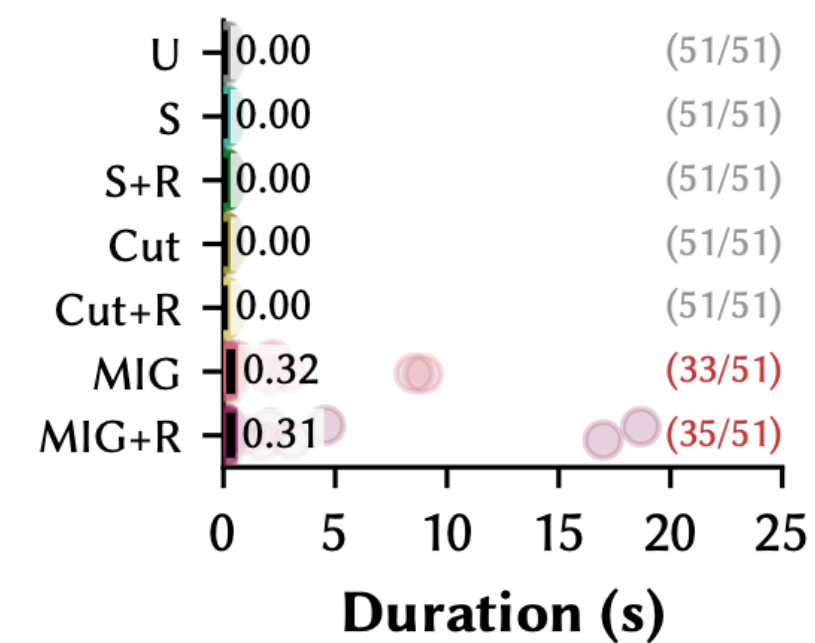
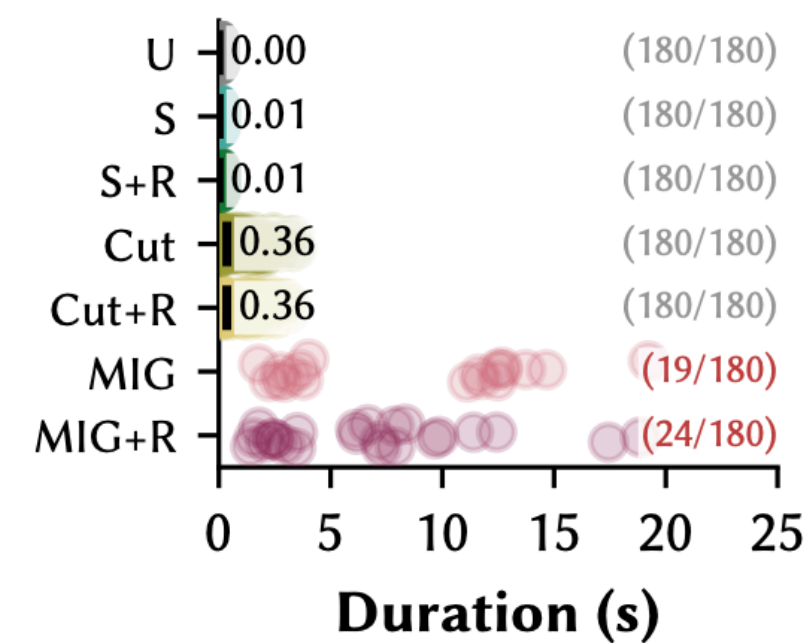
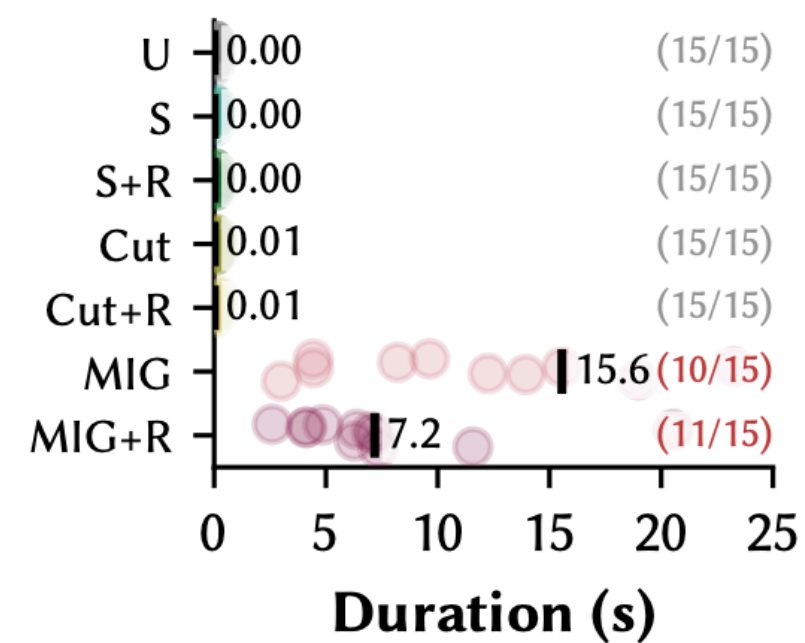
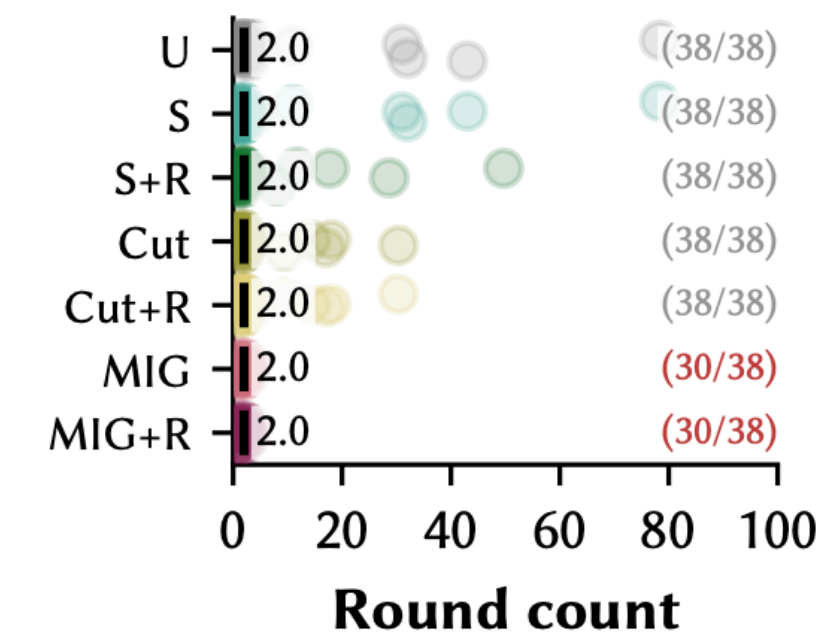
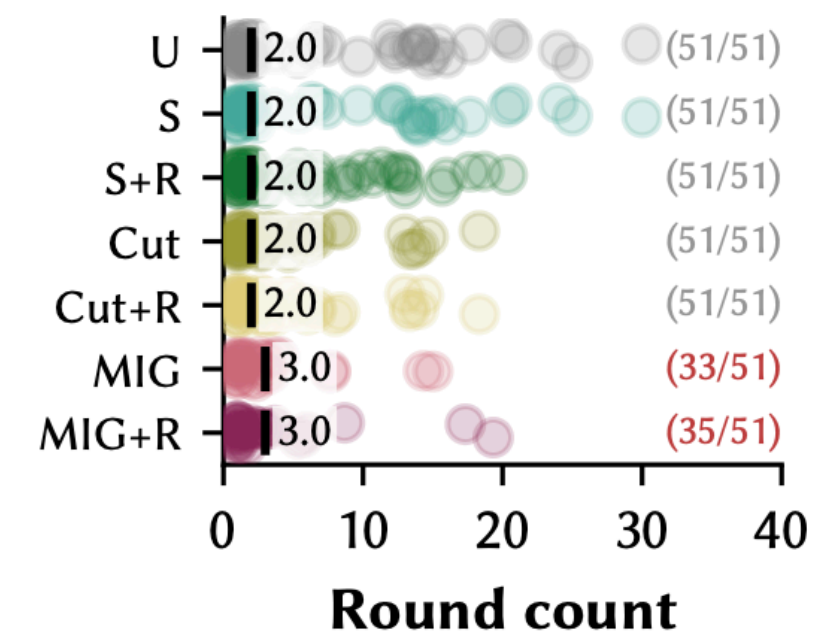
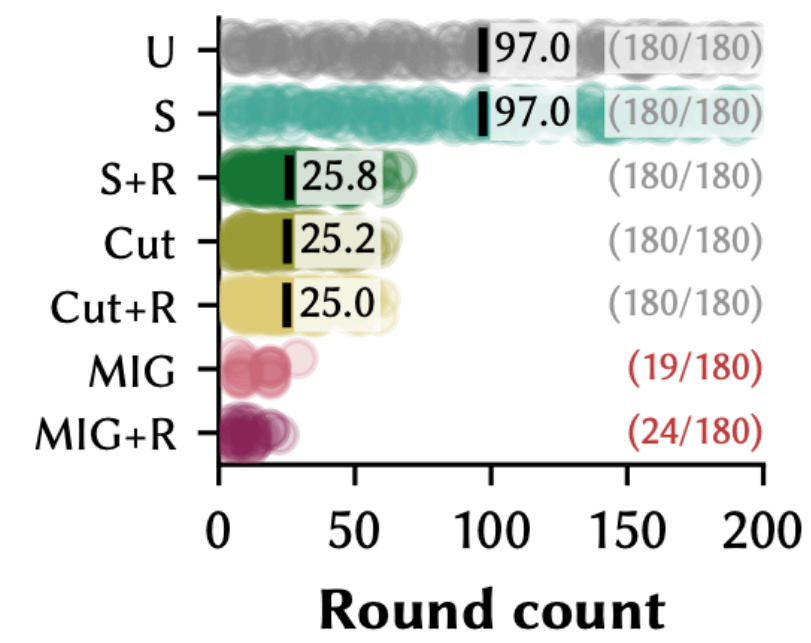
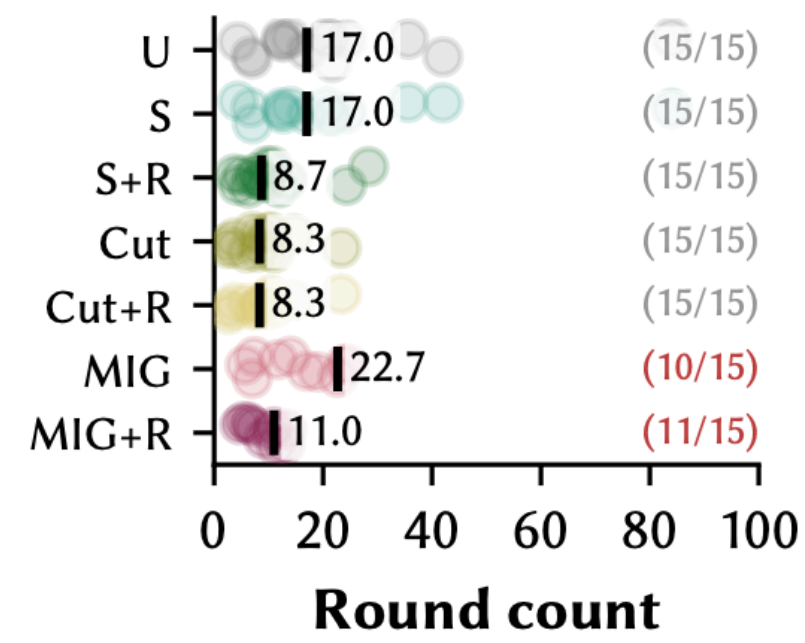
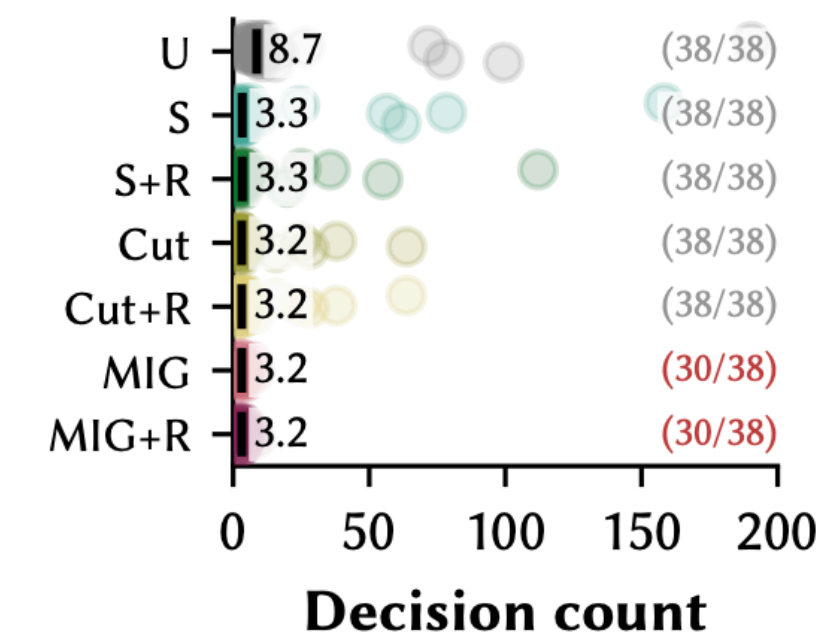
RANDOM



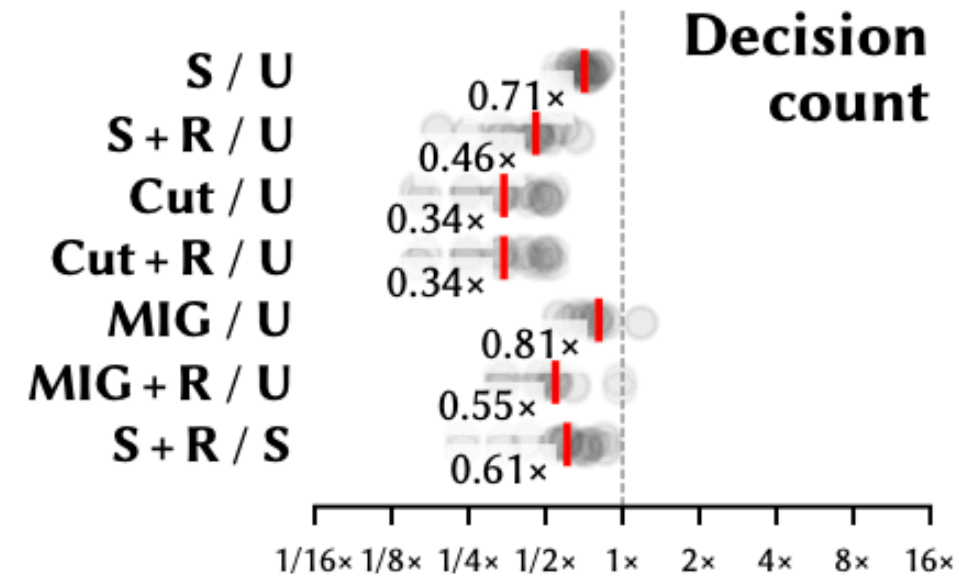
ARGUS



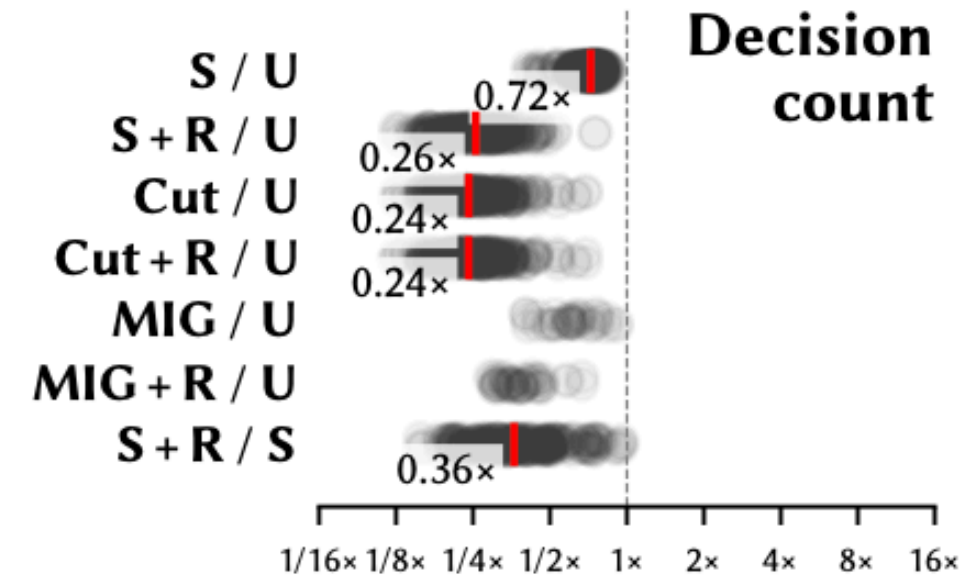
AESOP



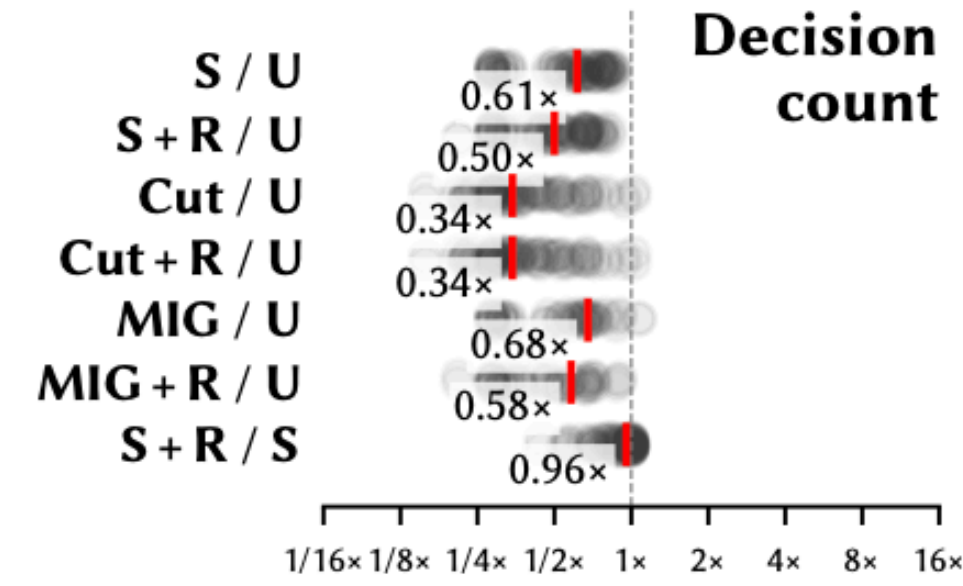
MANUAL



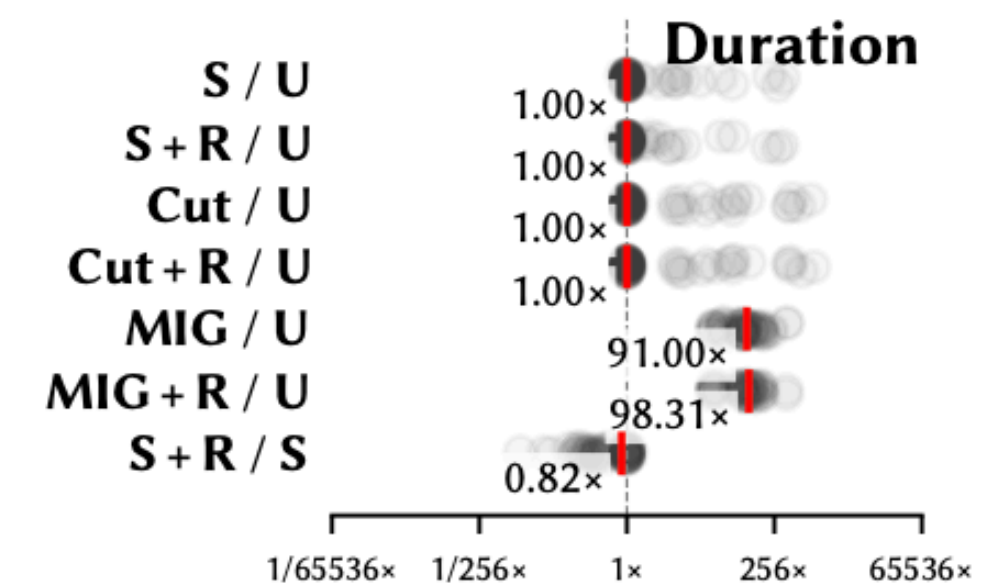
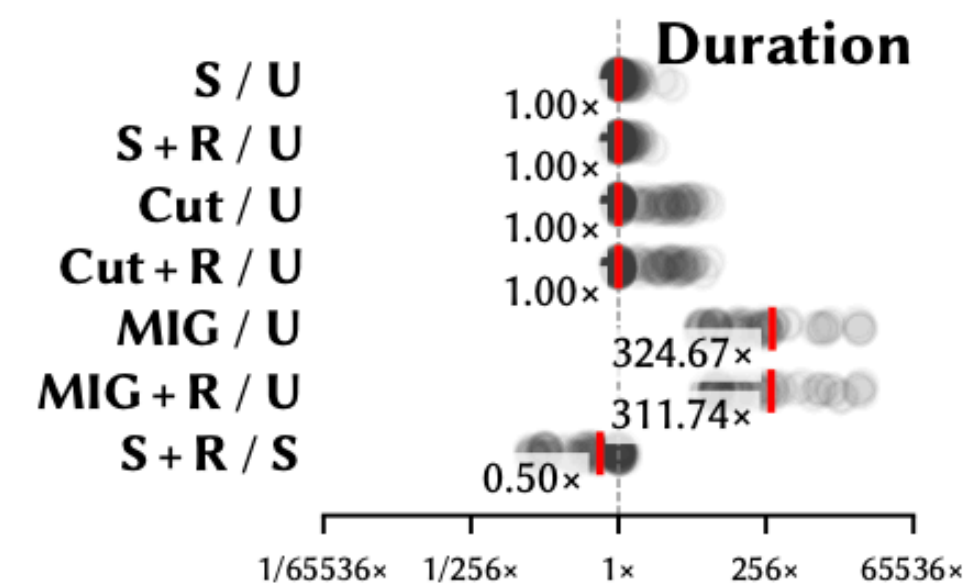
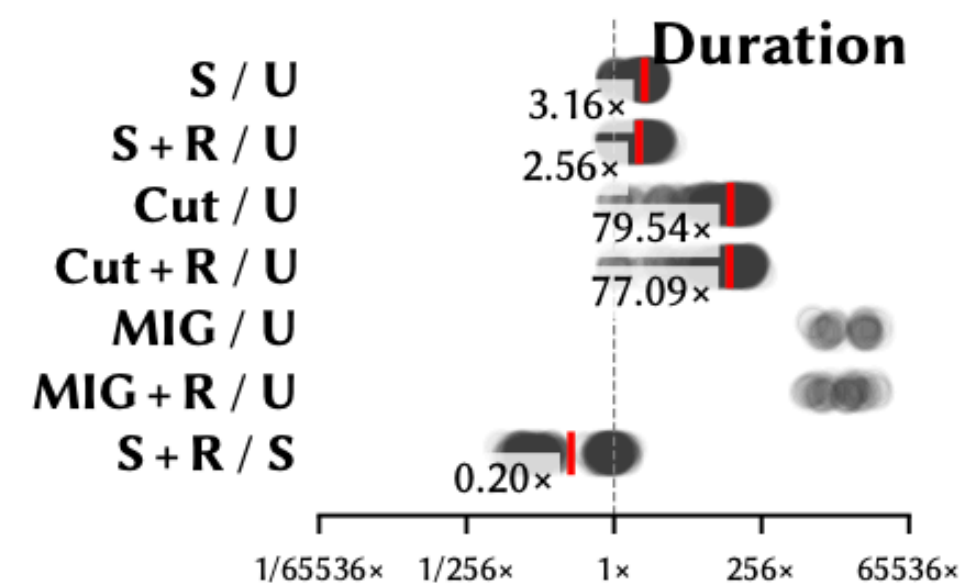
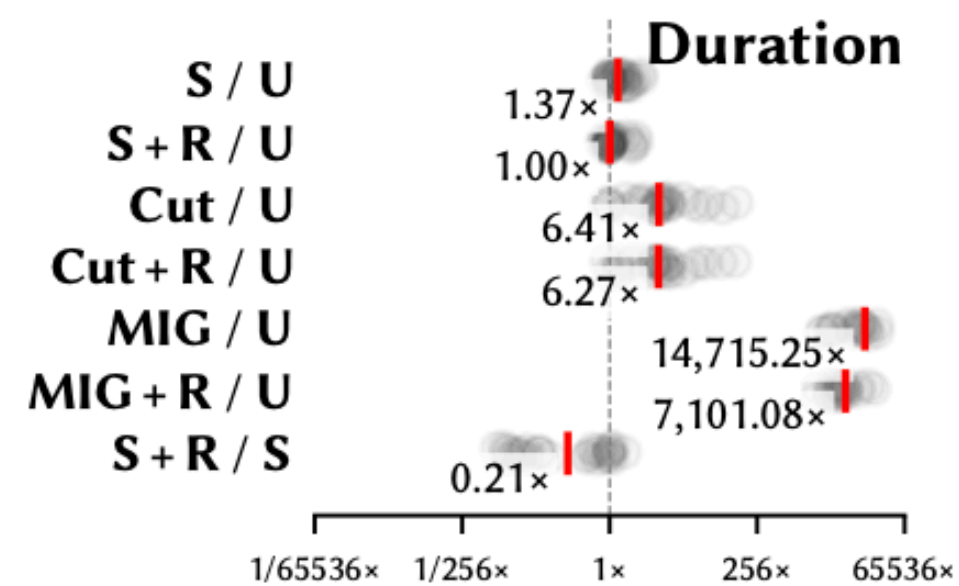
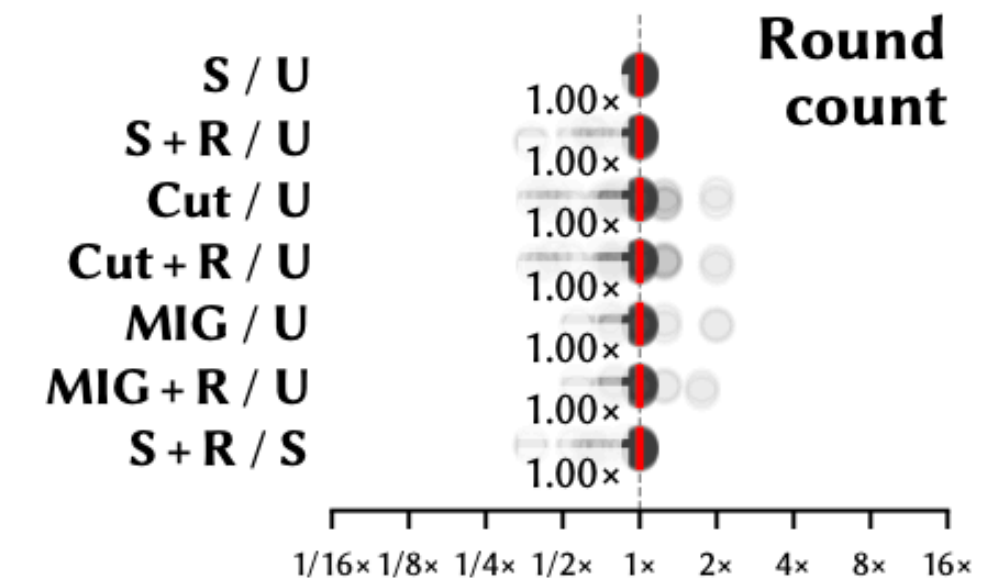
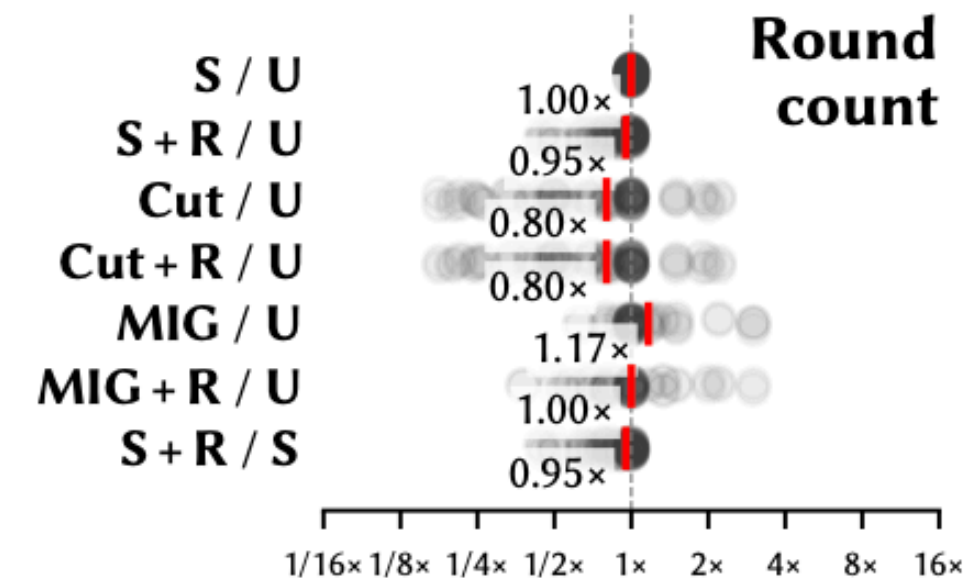
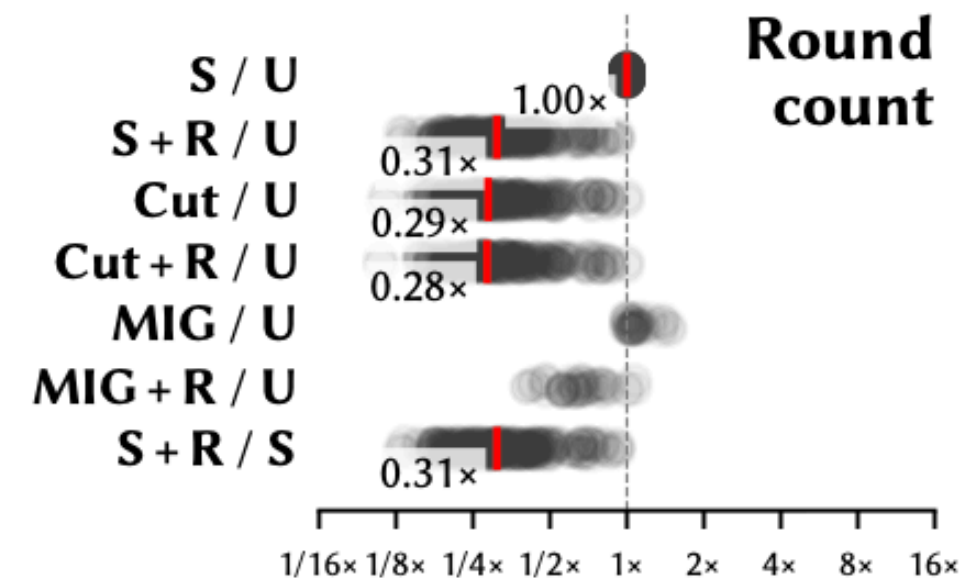
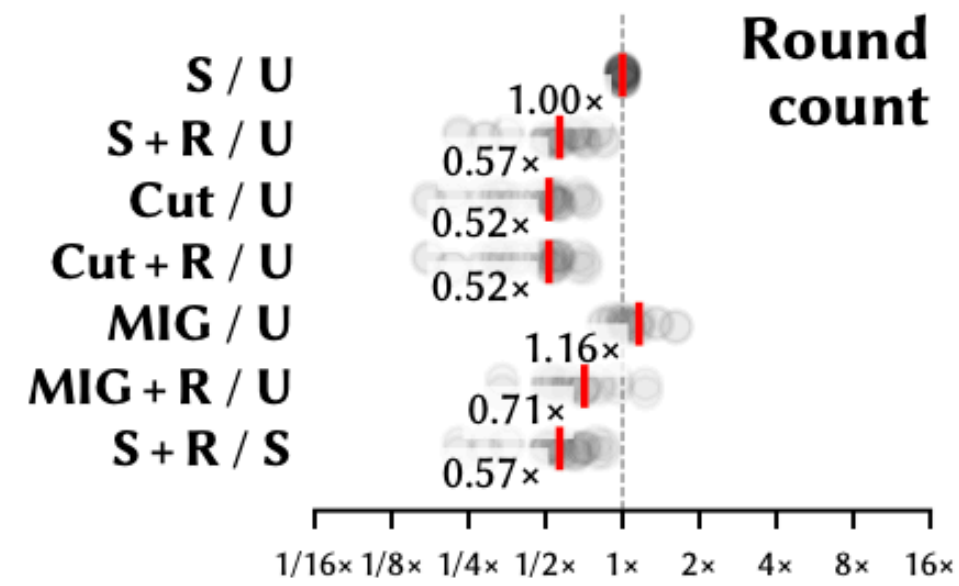
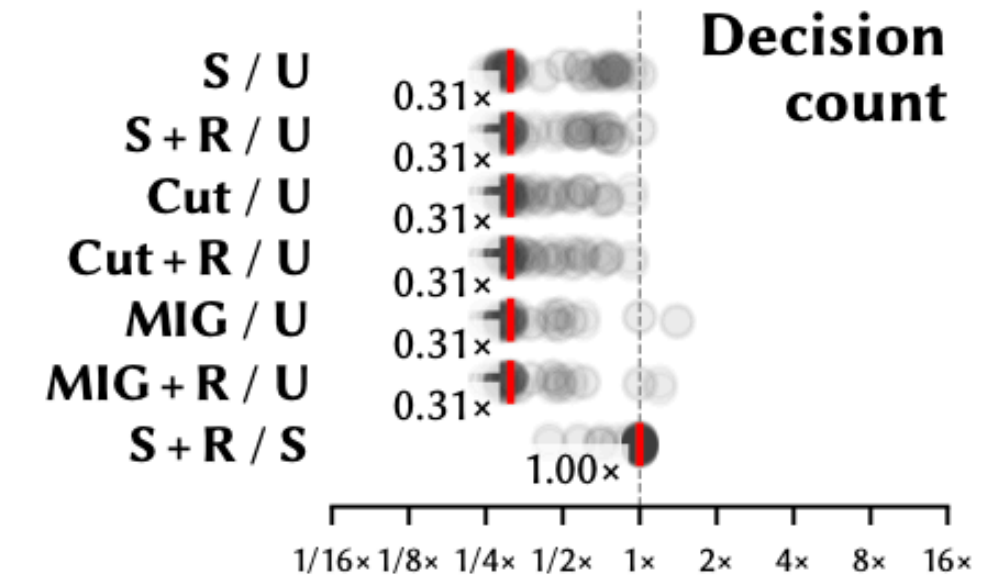
RANDOM



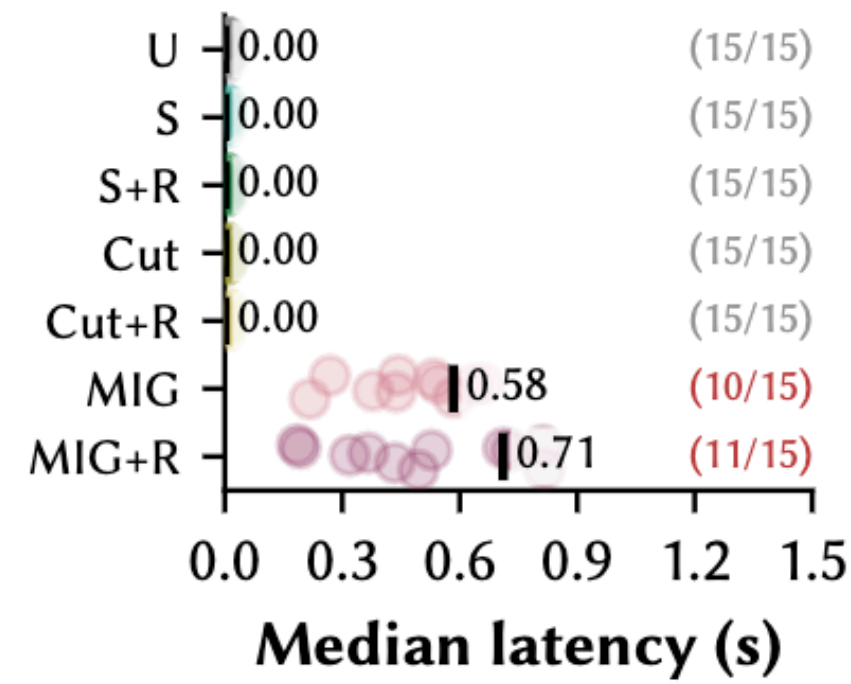
ARGUS



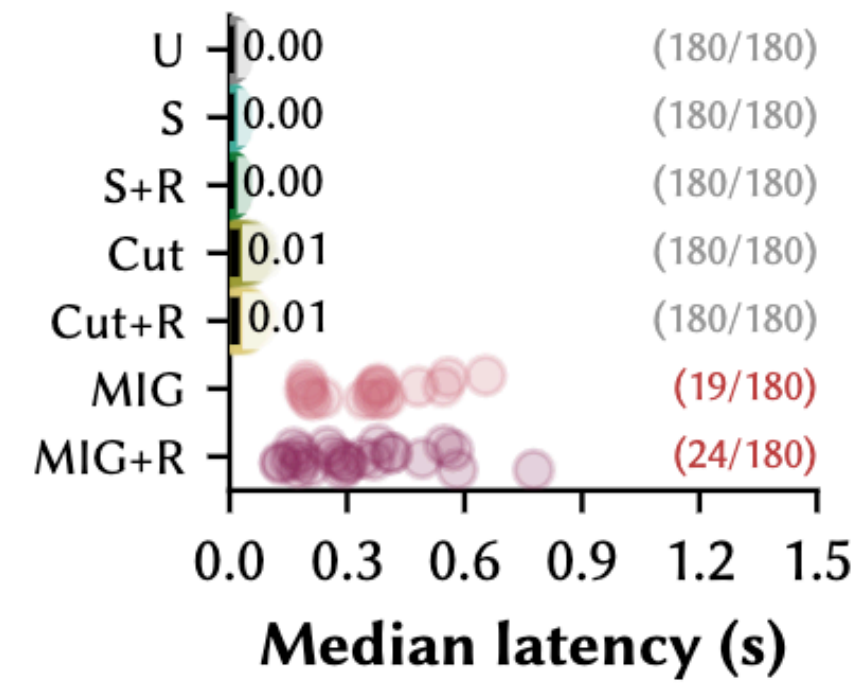
AESOP



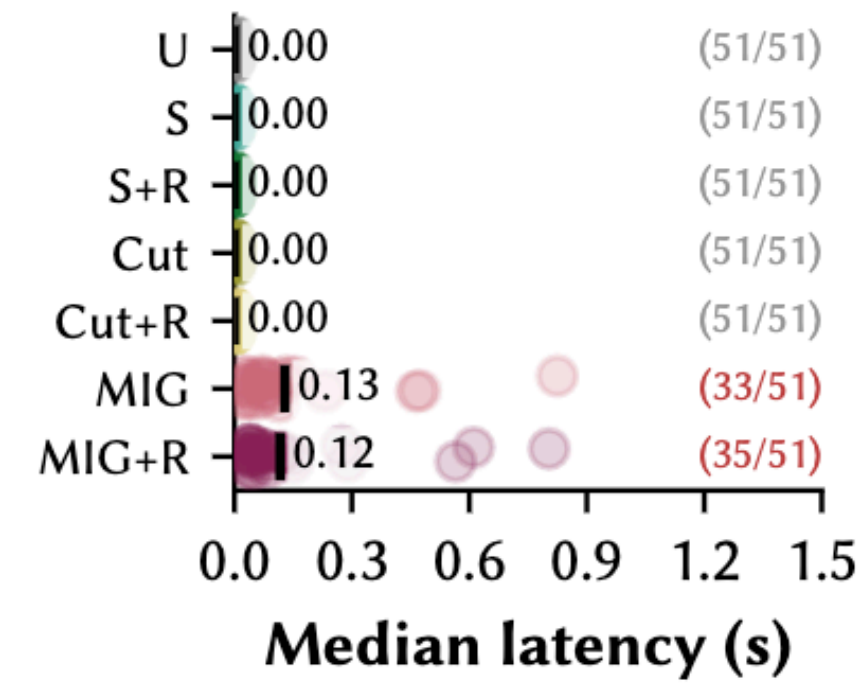
MANUAL



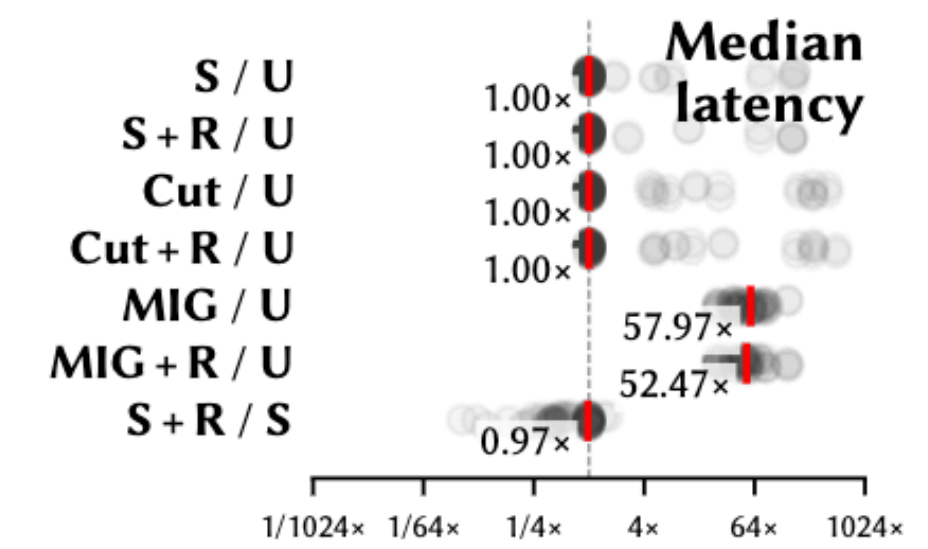
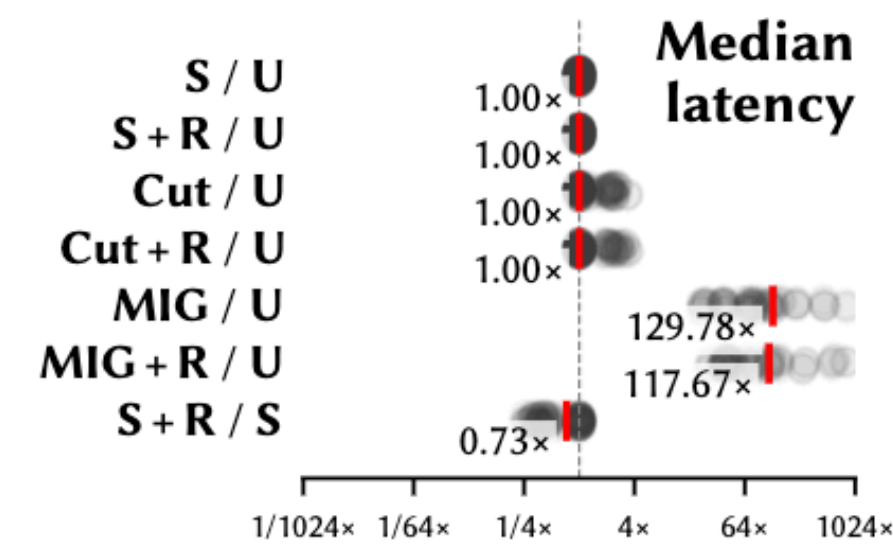
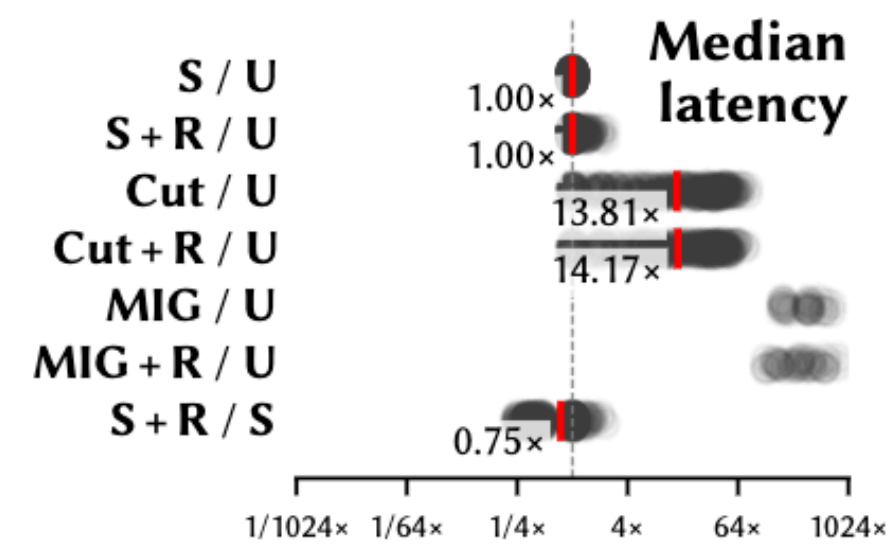
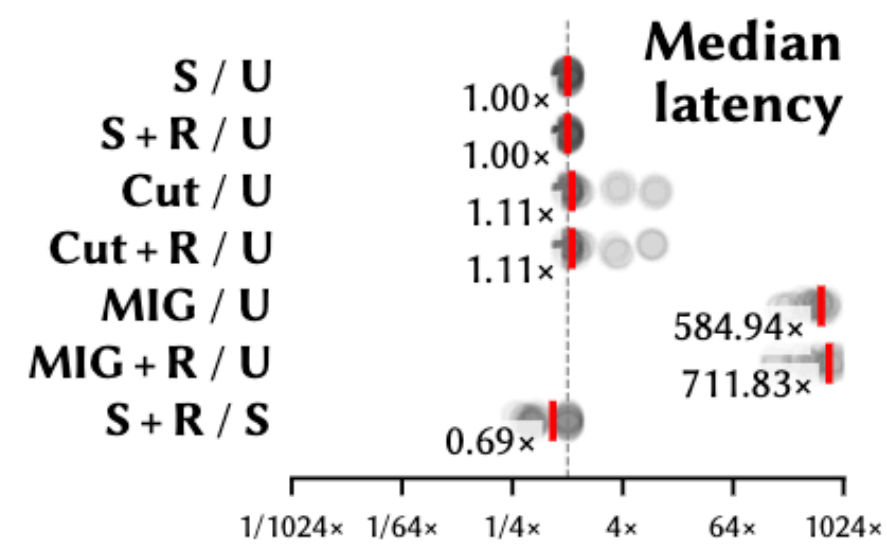
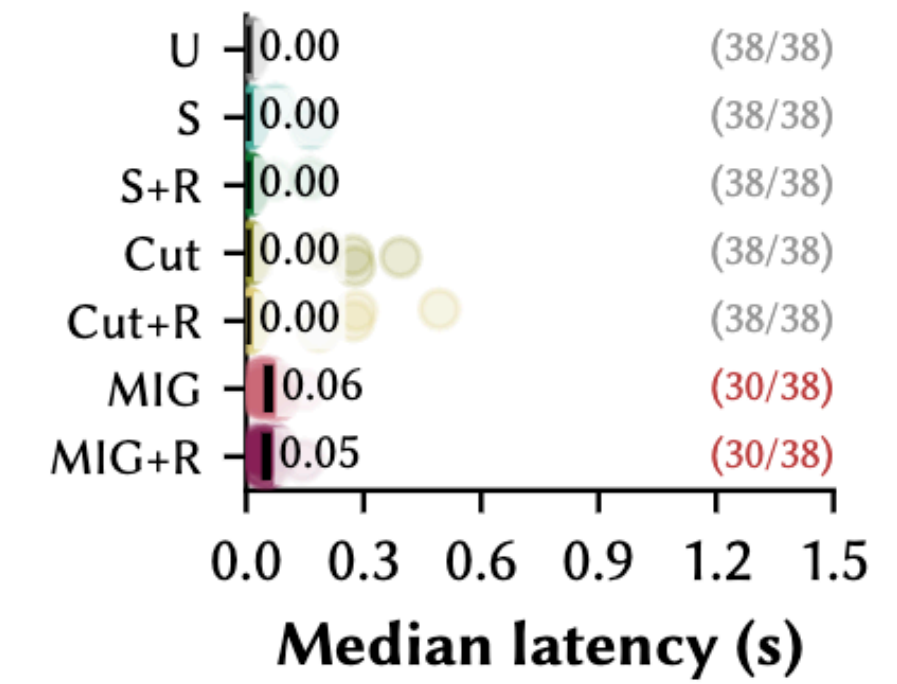
RANDOM



ARGUS



AESOP



Incremental SAT Solving

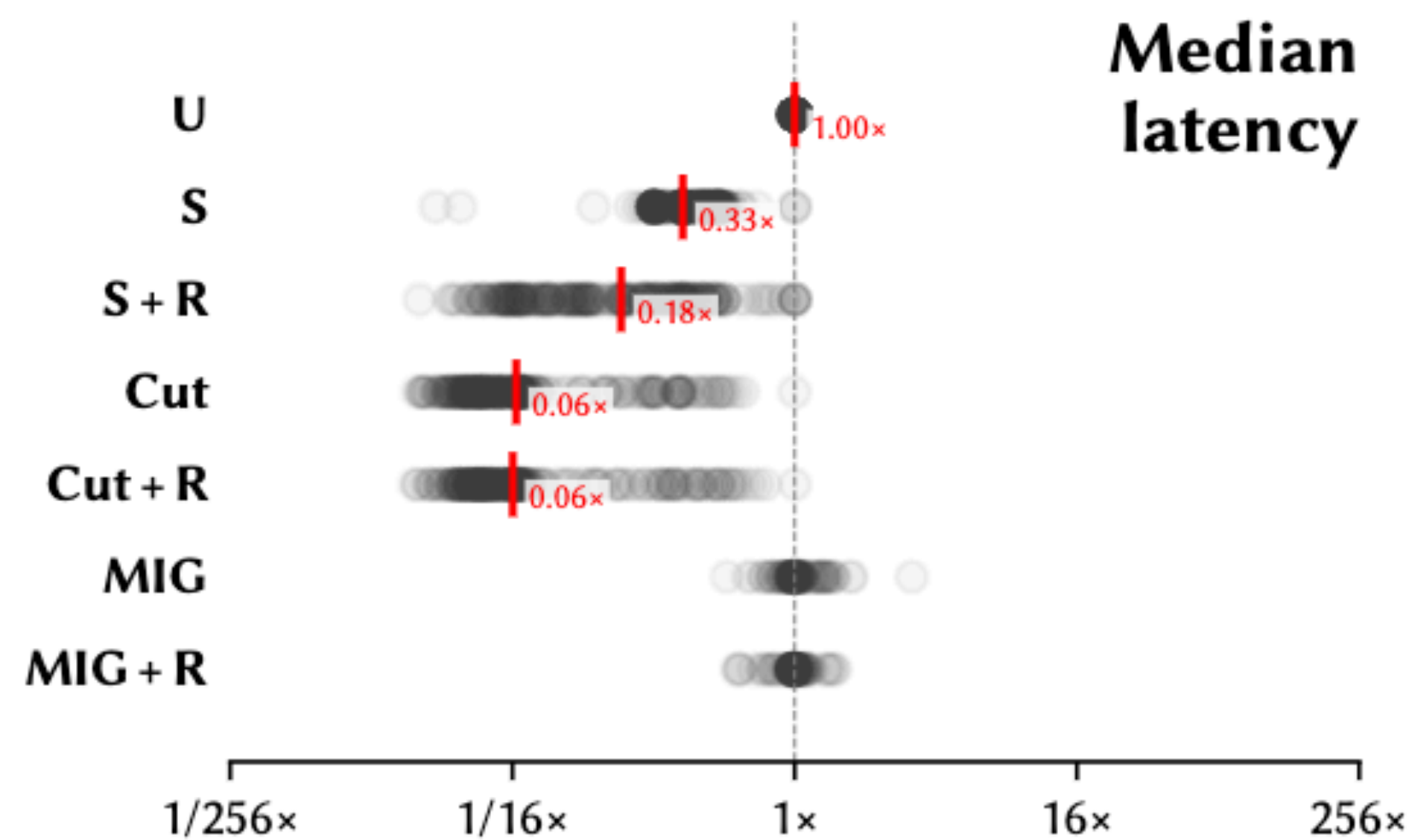
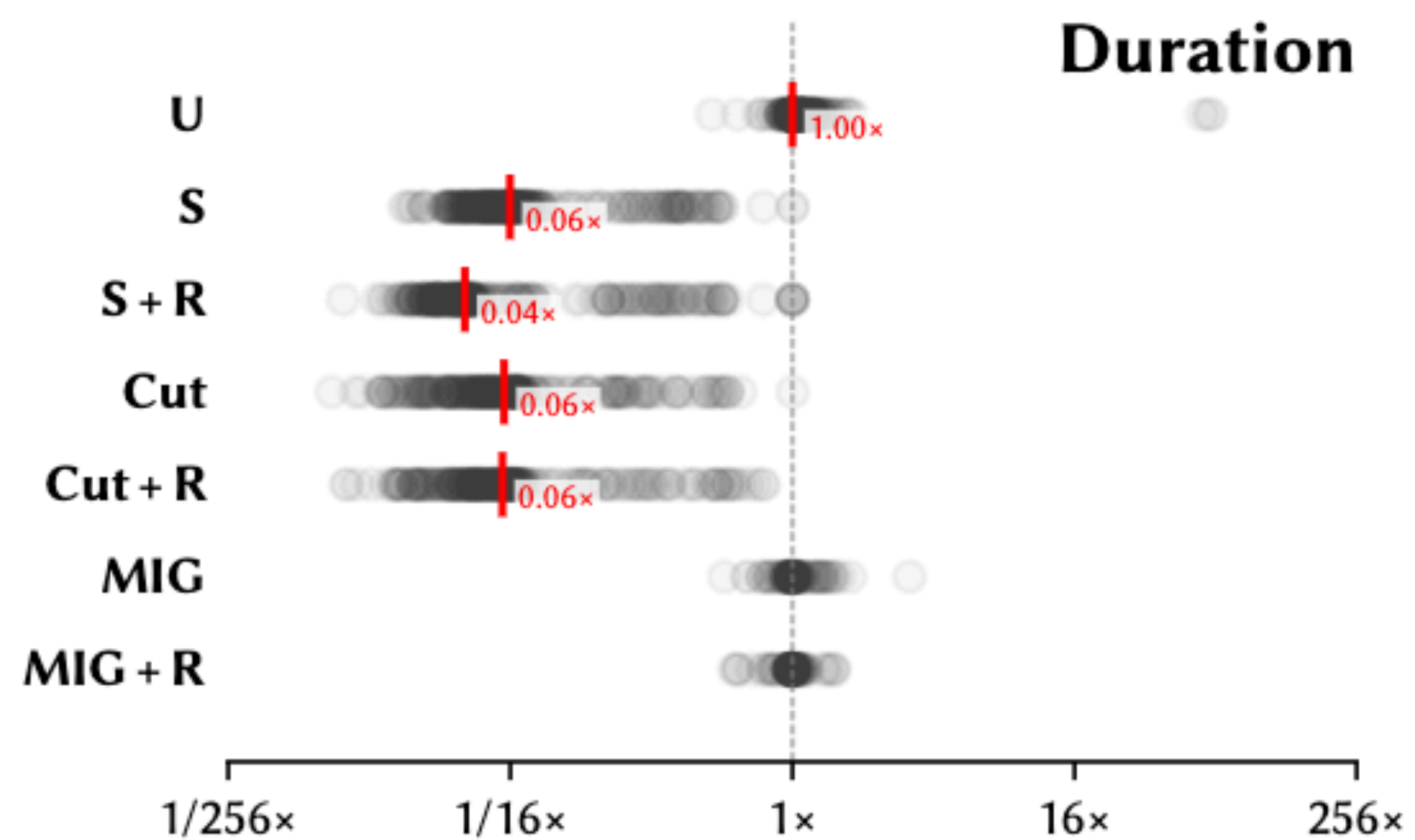


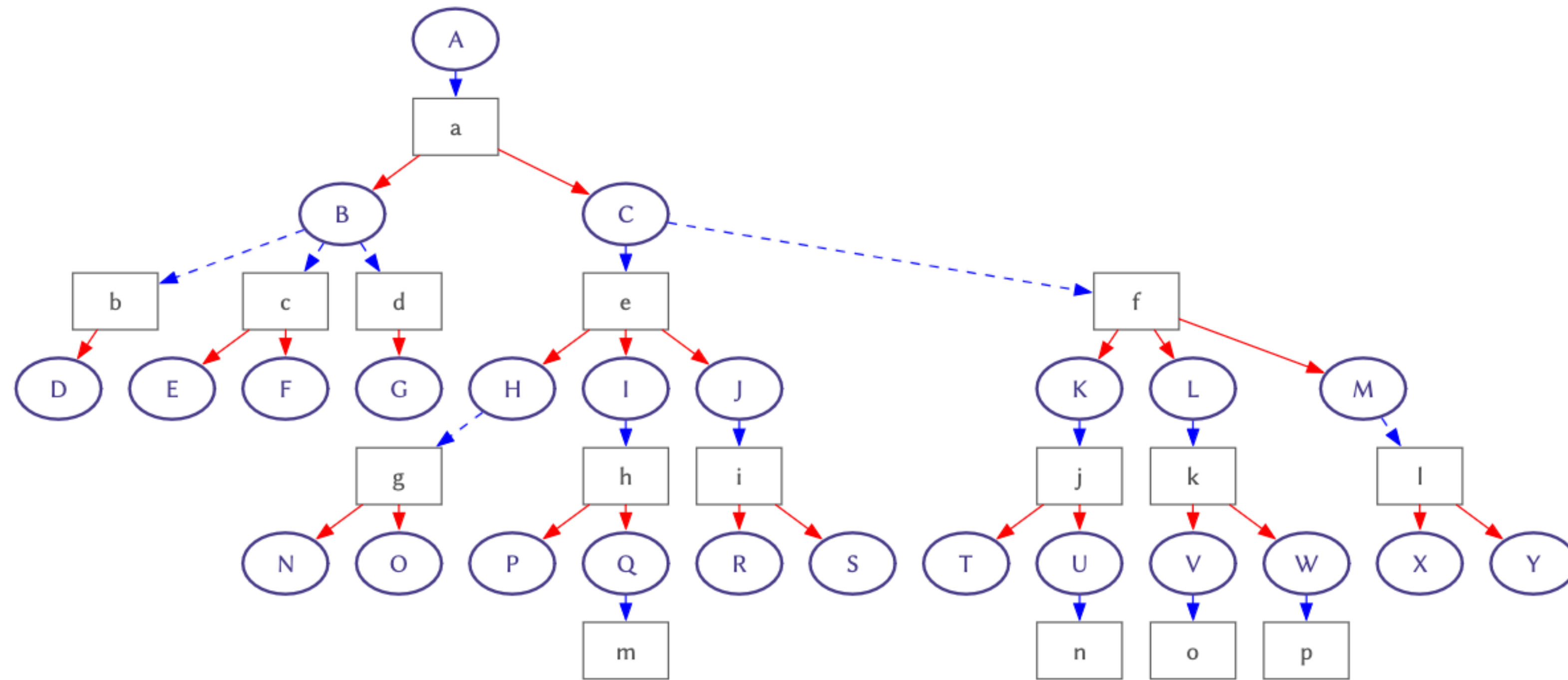
Table 1. **Benchmark summary statistics.** Each entry describes the median and IQR. **Depth** is the maximum graph distance from the goal to any other node in the same connected component. **ORs** and **ANDs** are the OR and AND node count. **Consumers** is the median number of parents for OR nodes, **Providers** is the median number of children for OR nodes, and **Premises** is the median number of children for AND nodes.

SUITE	DEPTH	ORs	ANDs	CONSUMERS	PROVIDERS	PREMISES
MANUAL (15 entries)	8 (6–12)	24 (16–40)	12 (9–32)	1 (1–1)	0 (0–1)	2 (1–2)
RANDOM (180 entries)	8 (6–12)	105 (60–150)	30 (16–48)	1 (1–1)	0 (0–0)	3 (2–4)
ARGUS (51 entries)	4 (2–4)	5 (3–23)	3 (2–18)	1 (0.5–1)	0 (0–0.5)	1 (1–4)
AESOP (38 entries)	2 (2–4)	3 (3–4)	1 (1–5)	1 (1–1)	0.5 (0–1)	1 (1–2)

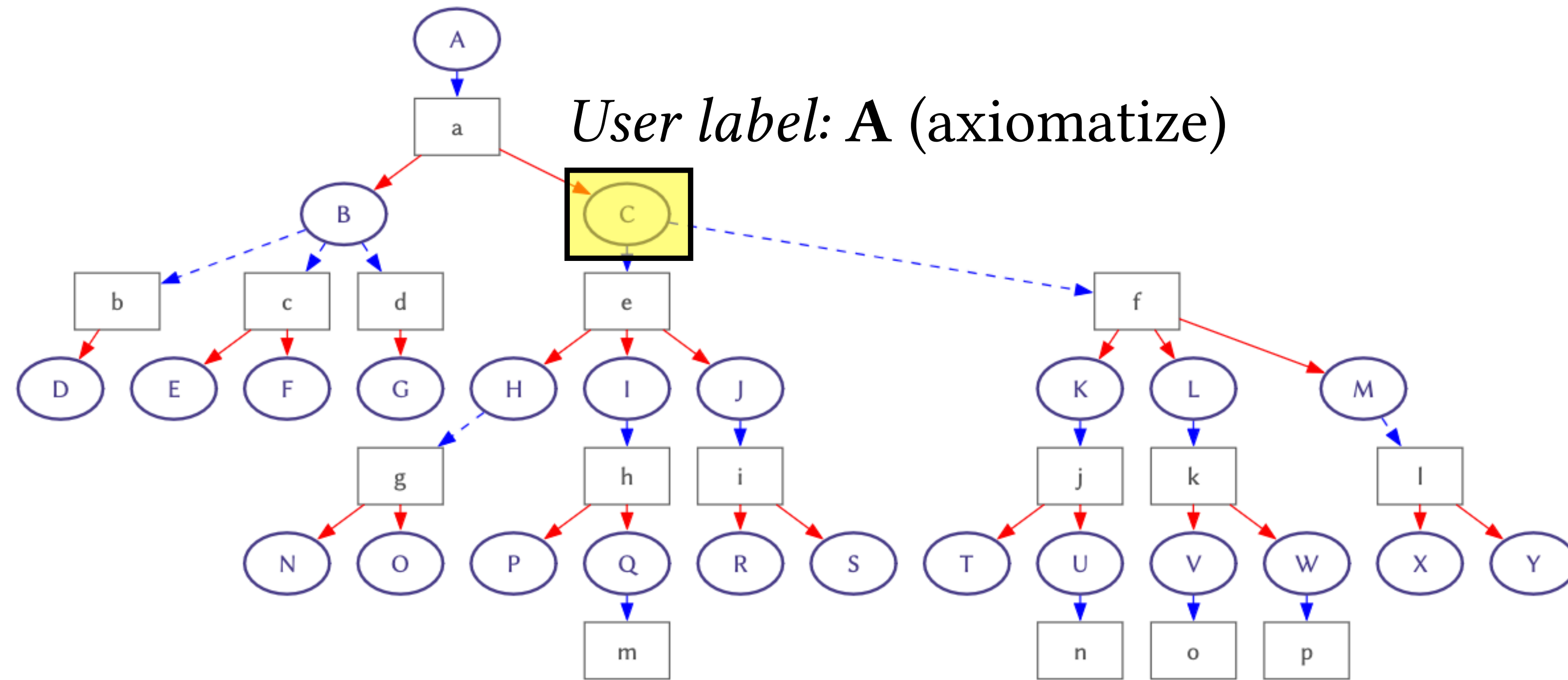
Table 2. **Maximum latencies.** Each entry is the maximum latency of the step provider on the given suite.

	MANUAL	RANDOM	ARGUS	AESOP
U	<1 ms	<1 ms	<1 ms	<1 ms
S	<1 ms	1 ms	4 ms	224 ms
S+R	<1 ms	6 ms	<1 ms	201 ms
Cut	16 ms	92 ms	4 ms	907 ms
Cut+R	15 ms	83 ms	4 ms	1084 ms

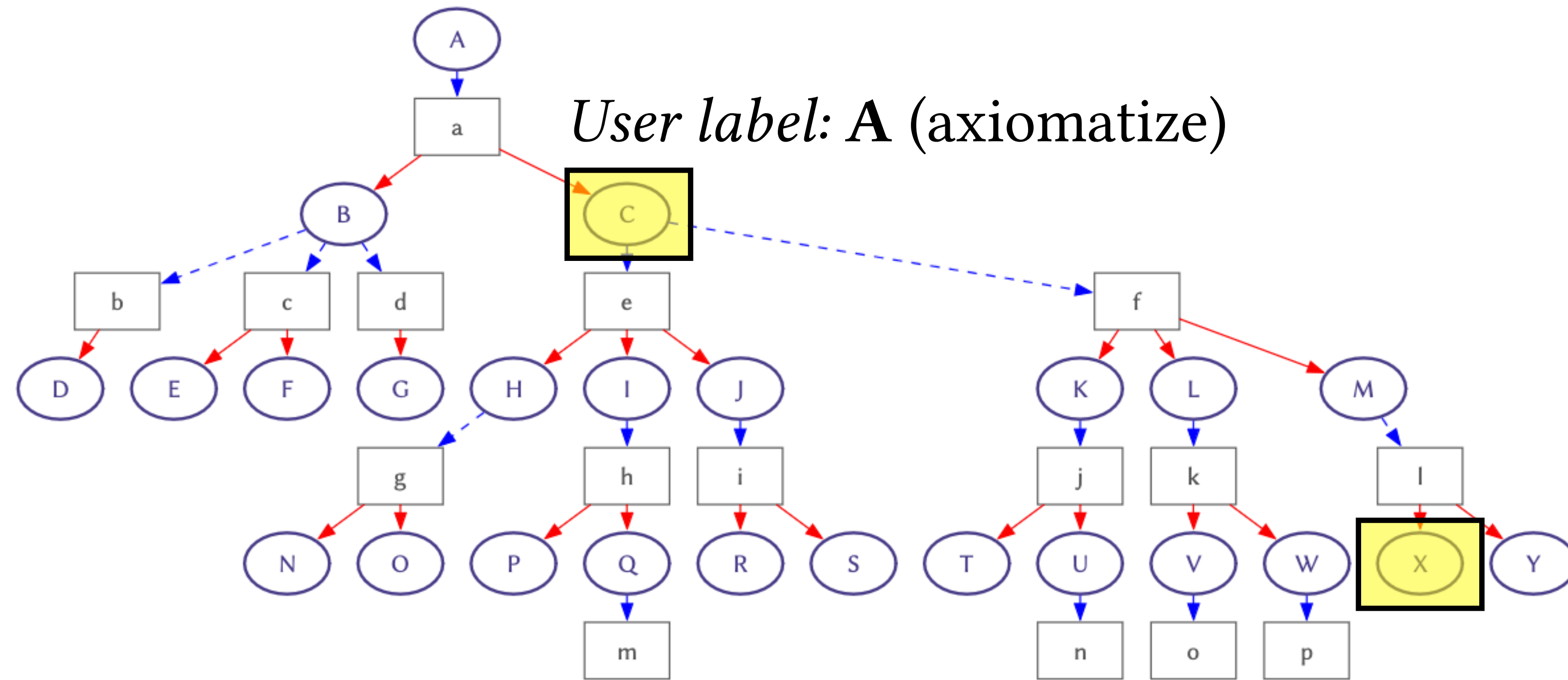
One last technique: We don't need to reach every partition



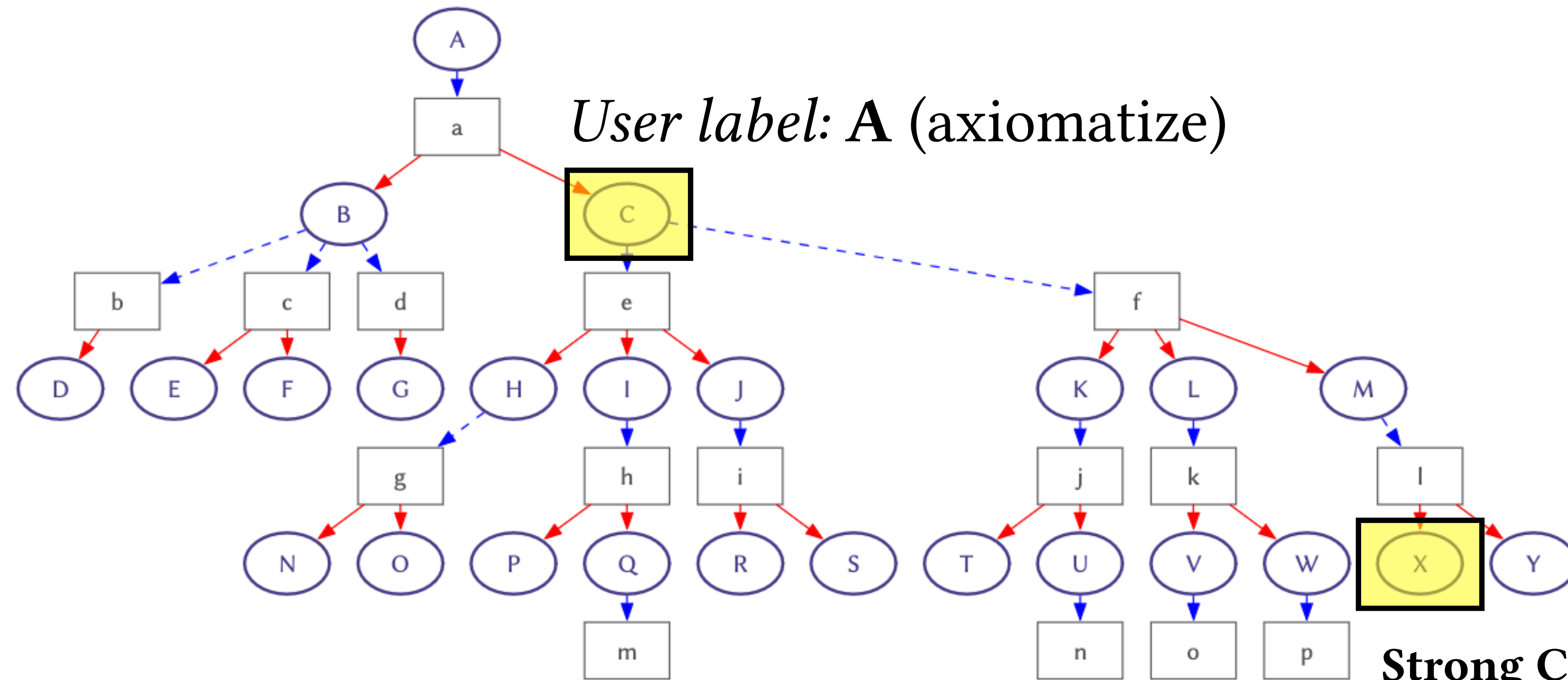
One last technique: We don't need to reach every partition



One last technique: We don't need to reach every partition



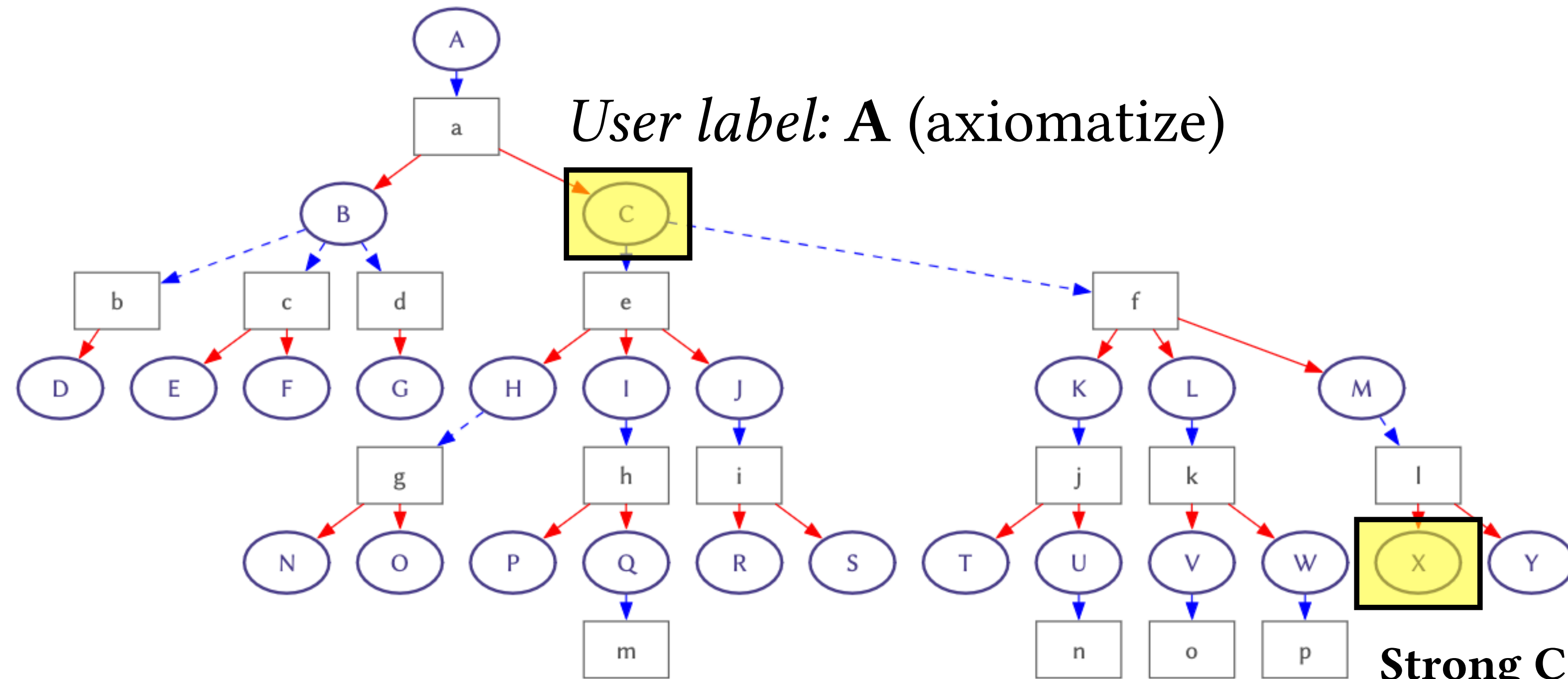
One last technique: We don't need to reach every partition



User label: A (axiomatize)

Strong Completeness
requires us to still be able
to label this as, e.g., **A**

One last technique: We don't need to reach every partition



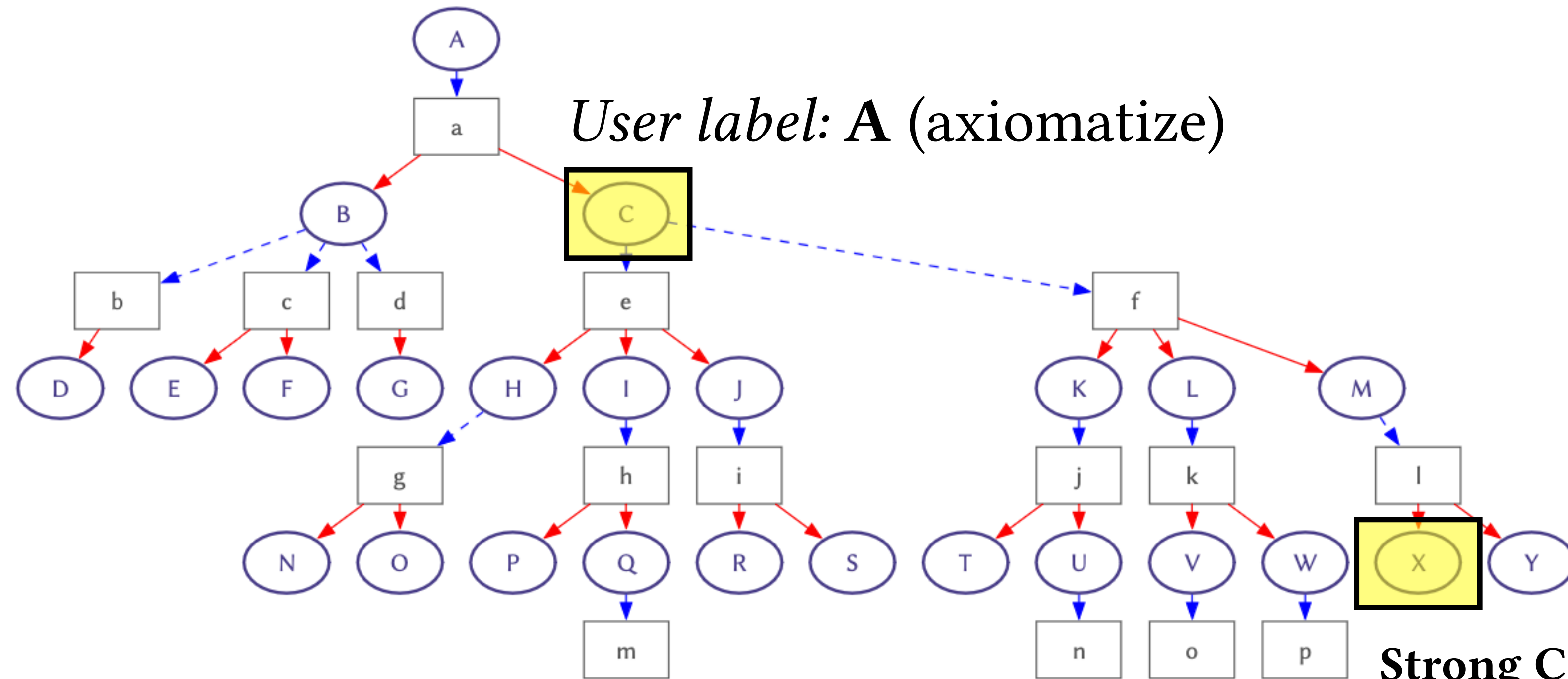
User label: **A** (axiomatize)

**Strong Completeness
Modulo Observability**

$$\bigcup_{\sigma \in \Sigma} [C(\sigma e)]_{\approx} \supseteq [C(e) \setminus \{e\}]_{\approx}$$

Strong Completeness
requires us to still be able
to label this as, e.g., **A**

One last technique: We don't need to reach every partition

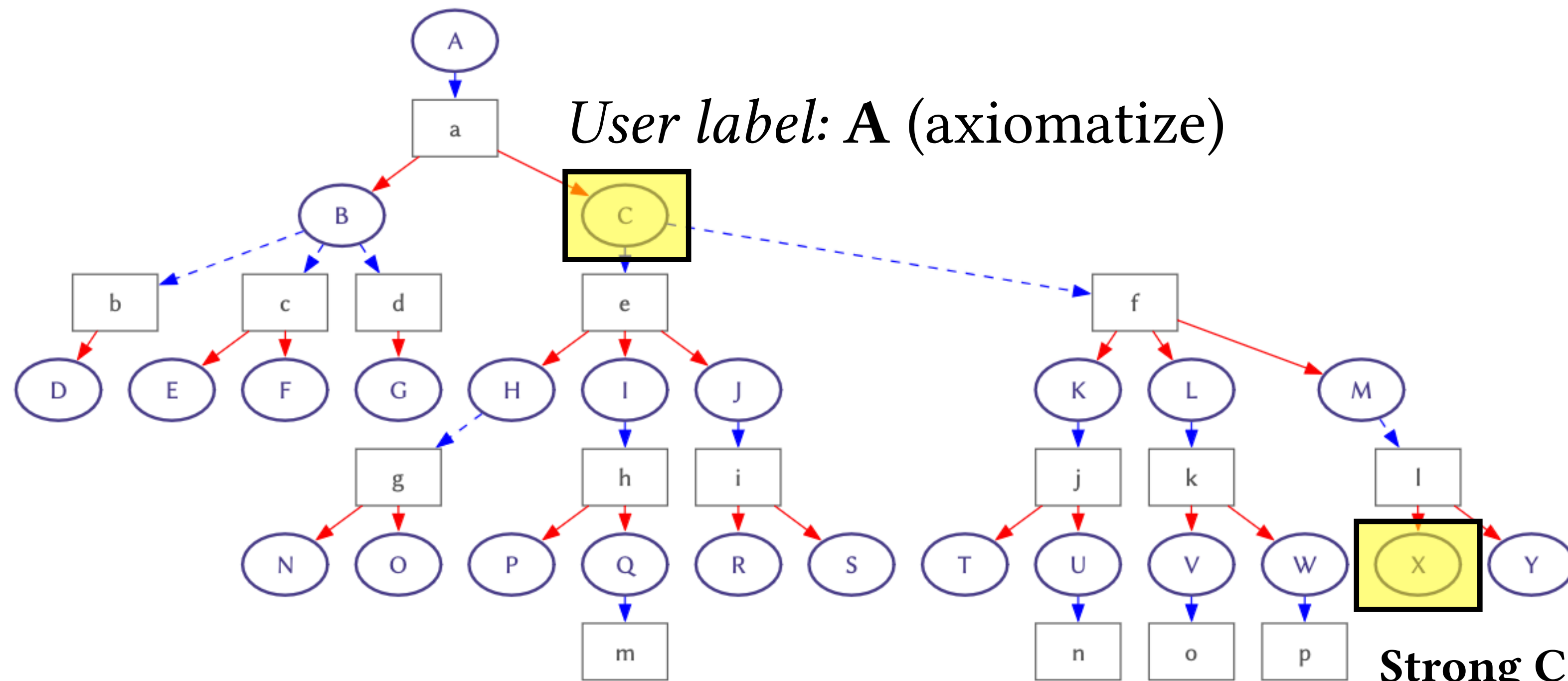


**Strong Completeness
Modulo Observability**

$$\bigcup_{\sigma \in \Sigma} [C(\sigma e)]_{\approx} \supseteq [C(e) \setminus \{e\}]_{\approx}$$

Strong Completeness
requires us to still be able
to label this as, e.g., A

One last technique: We don't need to reach every partition



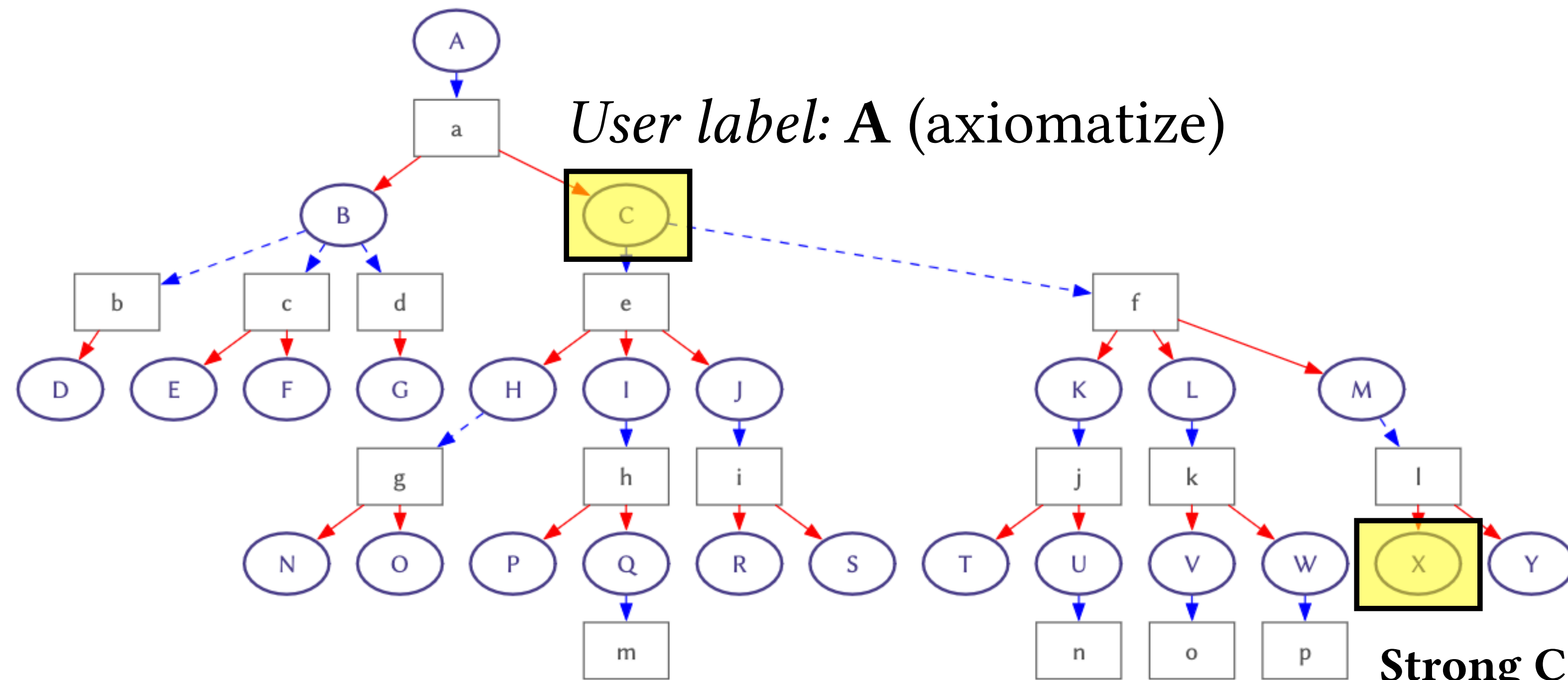
**Strong Completeness
Modulo Observability**

$$\bigcup_{\sigma \in \Sigma} [C(\sigma e)]_{\approx} \supseteq [C(e) \setminus \{e\}]_{\approx}$$

Strong Completeness
requires us to still be able
to label this as, e.g., A

Standard observational equivalence works ($e_1 \approx e_2$ iff $\llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket$), but we can go coarser:

One last technique: We don't need to reach every partition



User label: **A** (axiomatize)

Strong Completeness Modulo Observability

$$\bigcup_{\sigma \in \Sigma} [C(\sigma e)] \approx \supseteq [C(e) \setminus \{e\}] \approx$$

Strong Completeness requires us to still be able to label this as, e.g., **A**

Standard observational equivalence works ($e_1 \approx e_2$ iff $\llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket$), but we can go coarser:

Special type of “why-provenance”

where

$$e_1 \approx e_2 \text{ iff } \mathbf{Assumes}(e_1; \mathcal{P} \vdash G) = \mathbf{Assumes}(e_2; \mathcal{P} \vdash G),$$

$$\mathbf{Assumes}(e; \mathcal{P} \vdash G) = \{ \{Q \in e^{-1}[L_{\text{assume}}] : Q \text{ appears in } P\} : P \text{ a proof of } G \text{ in } \llbracket e \rrbracket \mathcal{P} \text{ respecting } e \}$$

