

프로그램 합성

이우석

한양대학교 ERICA 소프트웨어학부

2021. 03. 19 @ 충남대학교

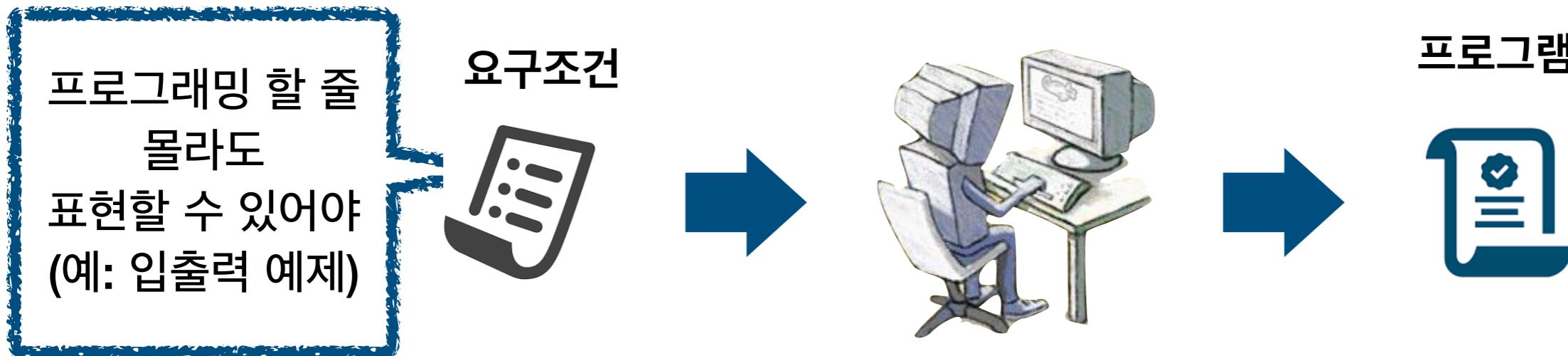


내용

- 프로그램 합성 개괄
- 귀납 합성 전략
- 연역 합성 전략

프로그램 합성

- 사용자가 원하는 프로그램을 자동으로 생성하는 기술



- 1945년부터 언급(앨런 튜링), 실현될 경우 개발의 무거운 짐을 덜어줄 것으로 기대되는 컴퓨터공학의 성배(holy grail)
- 1960년대부터 느리게 연구되어오다가 2008년 CEGIS라는 기술의 등장을 기점으로 최근까지 두드러지게 발전

프로그램 합성과 연관기술 비교 (합성 vs. 컴파일러)

```
let rec insert x xs =
  match xs with
  | [] -> x :: []
  | h :: t ->
    if x <= h then x :: xs
    else h :: (insert x t)
```

OCaml



```
append:
push ebp
mov ebp, esp
push eax
push ebx
push len
call malloc
mov ebx, [ebp + 12]
mov [eax + info], ebx
mov dword [eax + next], 0
mov ebx, [ebp + 8]
cmp dword [ebx], 0
je null_pointer
mov ebx, [ebx]

next_element:
cmp dword [ebx + next], 0
je found_last
mov ebx, [ebx + next]
jmp next_element

found_last:
push eax
push addMes
call puts
add esp, 4
pop eax
mov [ebx + next], eax

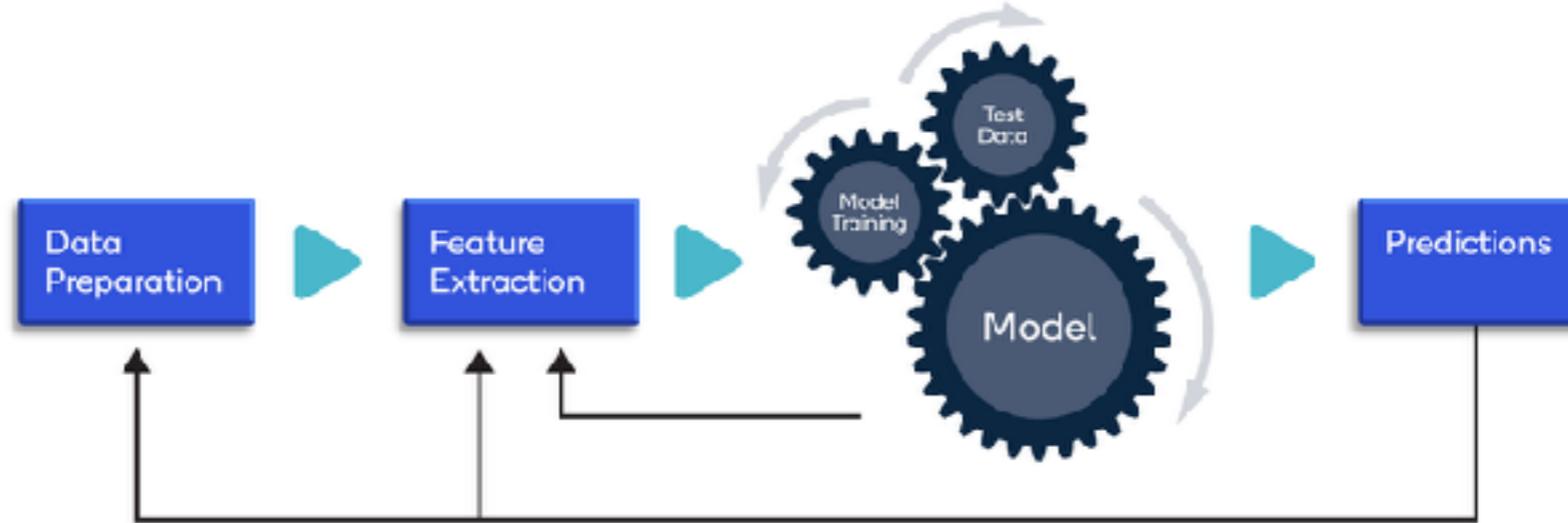
go_out:
pop ebx
pop eax
mov esp, ebp
pop ebp
ret 8

null_pointer:
push eax
push nullMes
call puts
add esp, 4
pop eax
mov [ebx], eax
jmp go_out
```

어셈블리

- **공통점:** 고 수준 (high level, 사람이 이해할 수 있는) 언어로 쓰인 개념을 실행 가능한 프로그램으로 변환
- **차이점:**
 - **컴파일러:** 탐색 과정 없이 단순 변환
 - **합성:** “어떻게” 주어진 목표를 달성할지 컴퓨터가 탐색

프로그램 합성과 연관기술 비교 (합성 vs. 기계학습)



- **공통점:** 주어진 입출력 데이터에 맞게 작동하는 함수 찾기
- **차이점:**
 - 기계학습: 투링완전(turing complete)하지 않은 대상(결정트리, 인공신경망 구조 등) 학습, (대개)결과물 해독 불가, 주어진 데이터에 꼭 맞지 않아도 됨 (노이즈 허용).
 - 합성: 일반 프로그램 대상(투링 완전), 해독 가능, 주어진 데이터에 꼭 맞아야.

프로그램 합성과 연관기술 비교 (합성 vs. 선언형 프로그래밍)

프롤로그 (Prolog) 의 예:

```
grandparent(X, Z) :- parent(X, Y), parent(Y, Z).  
parent(X, Y) :- father(X, Y).  
mother(X, Y) :- father(X, Y).
```

```
mother(Mary, Stan). Father(Stan, Alice).
```

- **공통점:** 사용자는 계산과정에 대한 얘기없이 결과물에 대한 조건만 기술
- **차이점:**
 - 선언형 프로그램: 요구조건에 모호성이 없어야 함 (예: 입출력 예제는 안됨), 결과물의 실행 성능 느림.
 - 합성: 요구조건이 모호해도 되고 결과물은 일반 프로그램

응용사례 – 문자열 변환 (FlashFill)

The screenshot shows a Microsoft Excel spreadsheet titled "dr-2 - Microsoft Excel". The ribbon at the top has the "Table Tools" tab selected. In the "Design" section of the ribbon, there is a dropdown menu with options like "Quick Fill", "Auto Fill", "Quick Layout", "Apply", "HiLight", "CurrencyWidget", "Undo", and "Commit AddressWidget". The main area of the spreadsheet contains a table named "Table116" with 26 rows and 6 columns. The first column is labeled "Column1". The second column contains names, and the third column contains addresses and phone numbers. The fourth column contains city names, the fifth column contains state abbreviations, and the sixth column contains zip codes. The "Table Tools" ribbon tab is highlighted in yellow.

Column1	A	B	C	D	E	F
Ana Trujillo	357 21th Place SE,Redmond,WA,(757) 555-1634,140-37-6064,27171	Redmond	WA	(757) 555-1634	140-37-6064	27171
Antonio Moreno	515 93th Lane ,Renton,WA,(411) 555-2786,562-87-3127,28581					
Thomas Hardy	742 17th Street NE,Seattle,WA,(412) 555-5719,921-29-4931,24607					
Christina Berglund	475 22th Lane ,Redmond,WA,(443) 555-6774,844-35-6764,30146					
Hanna Moos	785 45th Street NE,Puyallup,WA,(376) 555-2462,515-68-1285,29284					
Frédérique Citeaux	308 66th Place ,Redmond,WA,(689) 555-2770,552-23-2508,21415					
Martin Sommer	887 86th Place ,Kent,WA,(715) 555-5450,870-91-9824,21536					
Laurence Lebihan	944 13th Street NE,Redmond,WA,(620) 555-2361,649-25-5312,25252					
Elizabeth Lincoln	452 73th Lane NE,Renton,WA,(851) 555-4561,425-97-6344,22279					
Victoria Ashworth	463 16th Street ,Renton,WA,(696) 555-6044,690-29-7926,22832					
Patricia Simpson	630 20th Street ,Redmond,WA,(179) 555-3265,389-78-3236,24525					
Francisco Chang	683 49th Lane ,Seattle,WA,(272) 555-7434,665-18-6435,29453					
Yang Wang	944 28th Lane ,Redmond,WA,(151) 555-2272,846-78-8452,24388					
Pedro Afonso	411 70th Place ,Kent,WA,(170) 555-2964,774-35-2298,29485					
Elizabeth Brown	971 20th Lane ,Puyallup,WA,(373) 555-4134,476-53-7164,26417					
Sven Ottlieb	676 17th Lane NE,Redmond,WA,(828) 555-1593,548-73-8633,27440					
Janine Labrone	267 95th Place SE,Seattle,WA,(949) 555-1316,350-27-8300,28074					
Ann Devon	694 53th Place ,Kent,WA,(194) 555-8124,559-74-4016,22367					
Roland Mendel	581 12th Street NW,Kent,WA,(103) 555-2146,303-79-1328,20518					
Aria Cruz	594 85th Lane ,Renton,WA,(431) 555-1376,329-93-9992,21498					
Diego Roel	550 22th Lane ,Renton,WA,(639) 555-6238,918-34-5172,25931					
Martine Rancé	688 93th Place NW,Kent,WA,(573) 555-3571,695-94-3479,22424					
24						
25						
26						

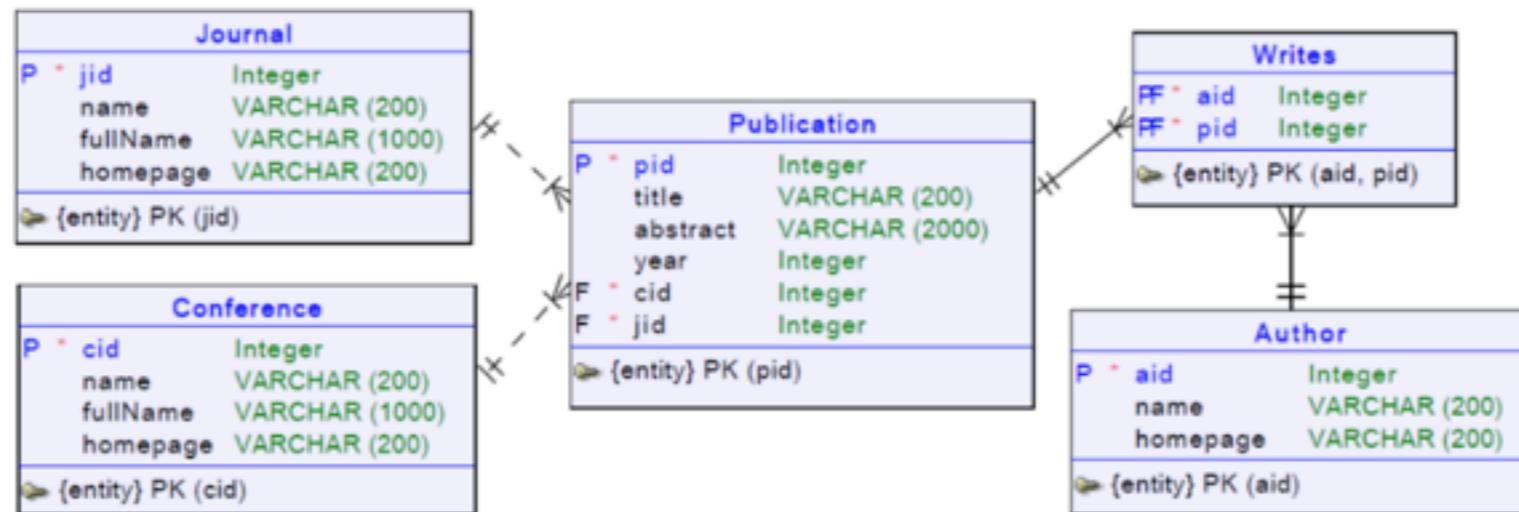
응용사례 – 문자열 변환 (FlashFill)

The screenshot shows a Microsoft Excel spreadsheet titled "dr-2 - Microsoft Excel". The ribbon at the top has tabs for File, Home, Insert, Page Layout, Formulas, Data, Review, View, Quick Code, Load Test, Team, and Table Tools. The Table Tools tab is selected. The main area displays a table named "Table116" with 23 rows of data. The first column is labeled "Column1" and contains names. The second column contains full addresses. The third column contains city names, the fourth contains state abbreviations, and the fifth contains phone numbers. The sixth column contains zip codes. The "Data" tab in the ribbon is also highlighted, indicating the use of the FlashFill feature.

1	Column1	A	B	C	D	E	F
2	Ana Trujillo	357 21th Place SE,Redmond,WA,(757) 555-1634,140-37-6064,27171	Redmond	WA	(757) 555-1634	140-37-6064	27171
3	Antonio Moreno	515 93th Lane ,Renton,WA,(411) 555-2786,562-87-3127,28581	Renton	WA	(411) 555-2786	562-87-3127	28581
4	Thomas Hardy	742 17th Street NE,Seattle,WA,(412) 555-5719,921-29-4931,24607	Seattle	WA	(412) 555-5719	921-29-4931	24607
5	Christina Berglund	475 22th Lane ,Redmond,WA,(443) 555-6774,844-35-6764,30146	Redmond	WA	(443) 555-6774	844-35-6764	30146
6	Hanna Moos	785 45th Street NE,Puyallup,WA,(376) 555-2462,515-68-1285,29284	Puyallup	WA	(376) 555-2462	515-68-1285	29284
7	Frédérique Citeaux	308 66th Place ,Redmond,WA,(689) 555-2770,552-23-2508,21415	Redmond	WA	(689) 555-2770	552-23-2508	21415
8	Martin Sommer	887 86th Place ,Kent,WA,(715) 555-5450,870-91-9824,21536	Kent	WA	(715) 555-5450	870-91-9824	21536
9	Laurence Lebihan	944 13th Street NE,Redmond,WA,(620) 555-2361,649-25-5312,2525	Redmond	WA	(620) 555-2361	649-25-5312	25252
10	Elizabeth Lincoln	452 73th Lane NE,Renton,WA,(851) 555-4561,425-97-6344,22279	Renton	WA	(851) 555-4561	425-97-6344	22279
11	Victoria Ashworth	463 16th Street ,Renton,WA,(696) 555-6044,690-29-7926,22832	Renton	WA	(696) 555-6044	690-29-7926	22832
12	Patricia Simpson	630 20th Street ,Redmond,WA,(179) 555-3265,389-78-3236,24525	Redmond	WA	(179) 555-3265	389-78-3236	24525
13	Francisco Chang	683 49th Lane ,Seattle,WA,(272) 555-7434,665-18-6435,29453	Seattle	WA	(272) 555-7434	665-18-6435	29453
14	Yang Wang	944 28th Lane ,Redmond,WA,(151) 555-2272,846-78-8452,24388	Redmond	WA	(151) 555-2272	846-78-8452	24388
15	Pedro Afonso	411 70th Place ,Kent,WA,(170) 555-2964,774-35-2298,29485	Kent	WA	(170) 555-2964	774-35-2298	29485
16	Elizabeth Brown	971 20th Lane ,Puyallup,WA,(373) 555-4134,476-53-7164,26417	Puyallup	WA	(373) 555-4134	476-53-7164	26417
17	Sven Ottlieb	676 17th Lane NE,Redmond,WA,(828) 555-1593,548-73-8633,27440	Redmond	WA	(828) 555-1593	548-73-8633	27440
18	Janine Labrunie	267 95th Place SE,Seattle,WA,(949) 555-1316,350-27-8300,28074	Seattle	WA	(949) 555-1316	350-27-8300	28074
19	Ann Devon	694 53th Place ,Kent,WA,(194) 555-8124,559-74-4016,22367	Kent	WA	(194) 555-8124	559-74-4016	22367
20	Roland Mendel	581 12th Street NW,Kent,WA,(103) 555-2146,303-79-1328,20518	Kent	WA	(103) 555-2146	303-79-1328	20518
21	Aria Cruz	594 85th Lane ,Renton,WA,(431) 555-1376,329-93-9992,21498	Renton	WA	(431) 555-1376	329-93-9992	21498
22	Diego Roel	550 22th Lane ,Renton,WA,(639) 555-6238,918-34-5172,25931	Renton	WA	(639) 555-6238	918-34-5172	25931
23	Martine Rancé	688 93th Place NW,Kent,WA,(573) 555-3571,695-94-3479,22424	Kent	WA	(573) 555-3571	695-94-3479	22424
24							
25							
26							

응용사례 – SQL 합성 (SQLizer)

Problem: “Find the number of papers in OOPSLA 2010”

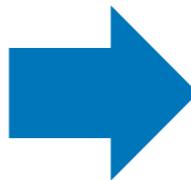


Output:

```
SELECT count(Publication.pid)
FROM Publication JOIN Conference ON Publication.cid = Conference.cid
WHERE Conference.name = "OOPSLA" AND Publication.year = 2010
```

응용사례 – 수퍼최적화 (STOKE)

```
1 # gcc -O3
2
3 movq rsi, r9
4 movl ecx, ecx
5 shrq 32, rsi
6 andl 0xffffffff, r9d
7 movq rcx, rax
8 movl edx, edx
9 imulq r9, rax
10 imulq rdx, r9
11 imulq rsi, rdx
12 imulq rsi, rcx
13 addq rdx, rax
14 jae .L0
15 movabsq 0x100000000, rdx
16 addq rdx, rcx
17 .L0:
18 movq rax, rsi
19 movq rax, rdx
20 shrq 32, rsi
21 salq 32, rdx
22 addq rsi, rcx
23 addq r9, rdx
24 adcq 0, rcx
25 addq r8, rdx
26 adcq 0, rcx
27 addq rdi, rdx
28 adcq 0, rcx
29 movq rcx, r8
30 movq rdx, rdi
```

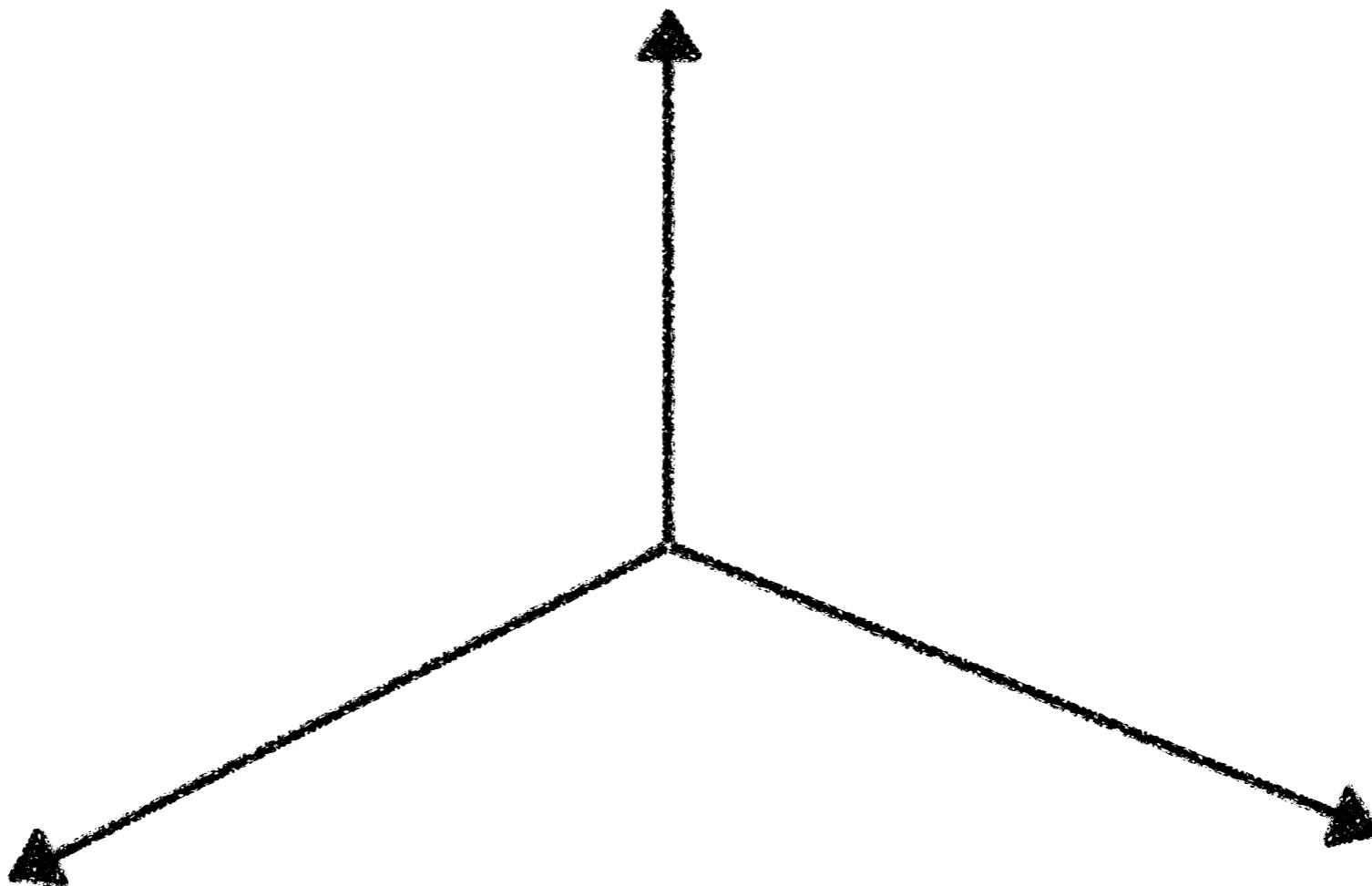


```
1 # STOKE
2
3 sh1q 32, rcx
4 movl edx, edx
5 xorq rdx, rcx
6 movq rcx, rax
7 mulq rsi
8 addq r8, rdi
9 adcq 0, rdx
10 addq rdi, rax
11 adcq 0, rdx
12 movq rdx, r8
13 movq rax, rdi
```

Montgomery multiplication kernel from the OpenSSL RSA library. Compilations shown for gcc -O3 (left) and a stochastic optimizer (right).

프로그램 합성의 차원

제약 조건: 프로그램이 달성할 목표를 어떻게 기술할지

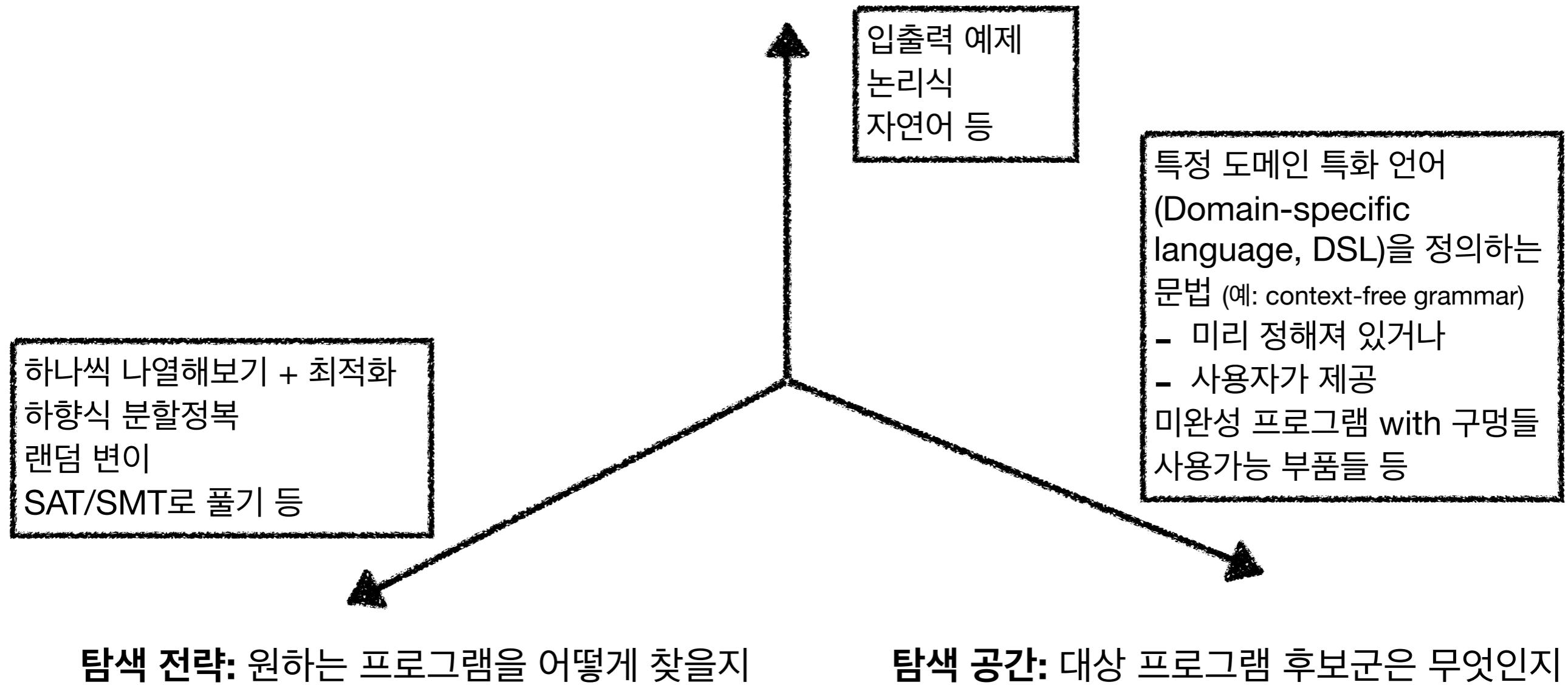


탐색 전략: 원하는 프로그램을 어떻게 찾을지

탐색 공간: 대상 프로그램 후보군은 무엇인지

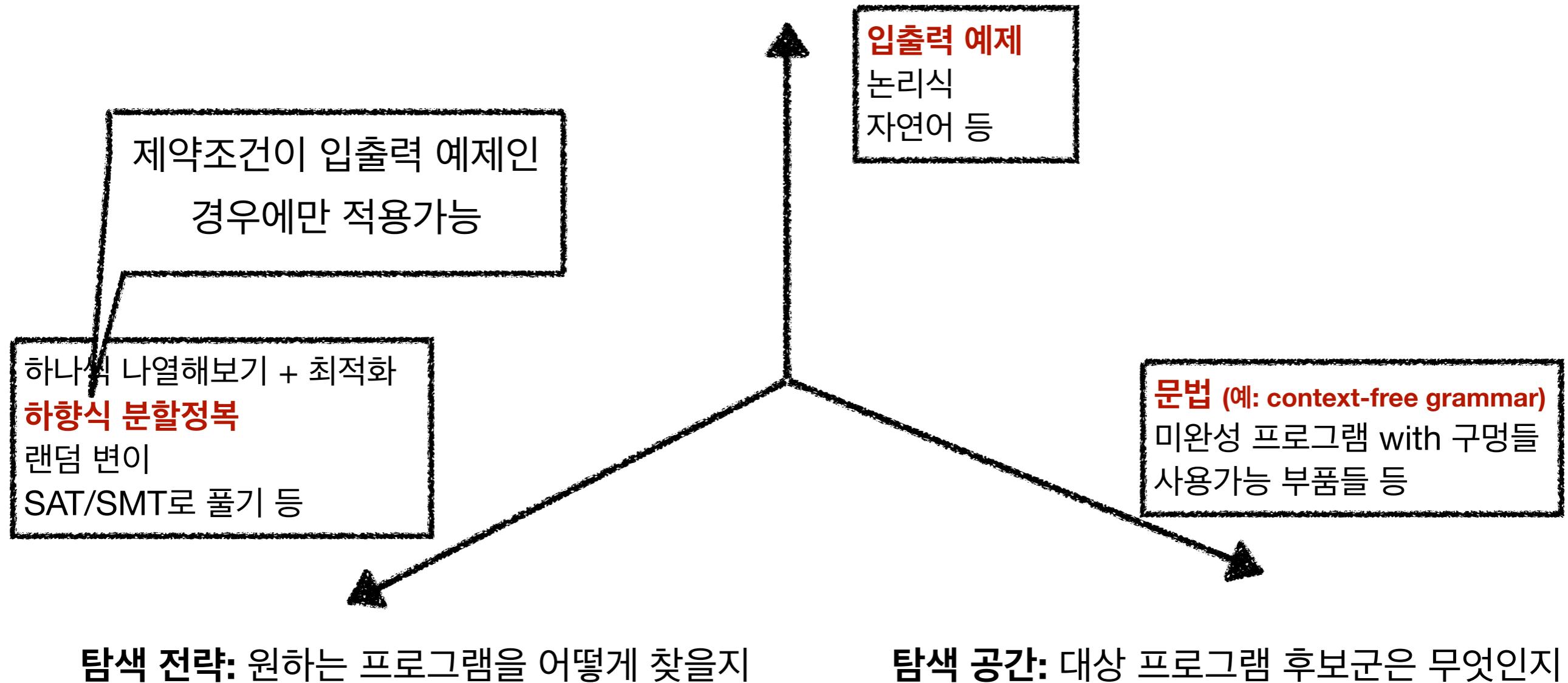
프로그램 합성의 차원

제약 조건: 프로그램이 달성할 목표를 어떻게 기술할지



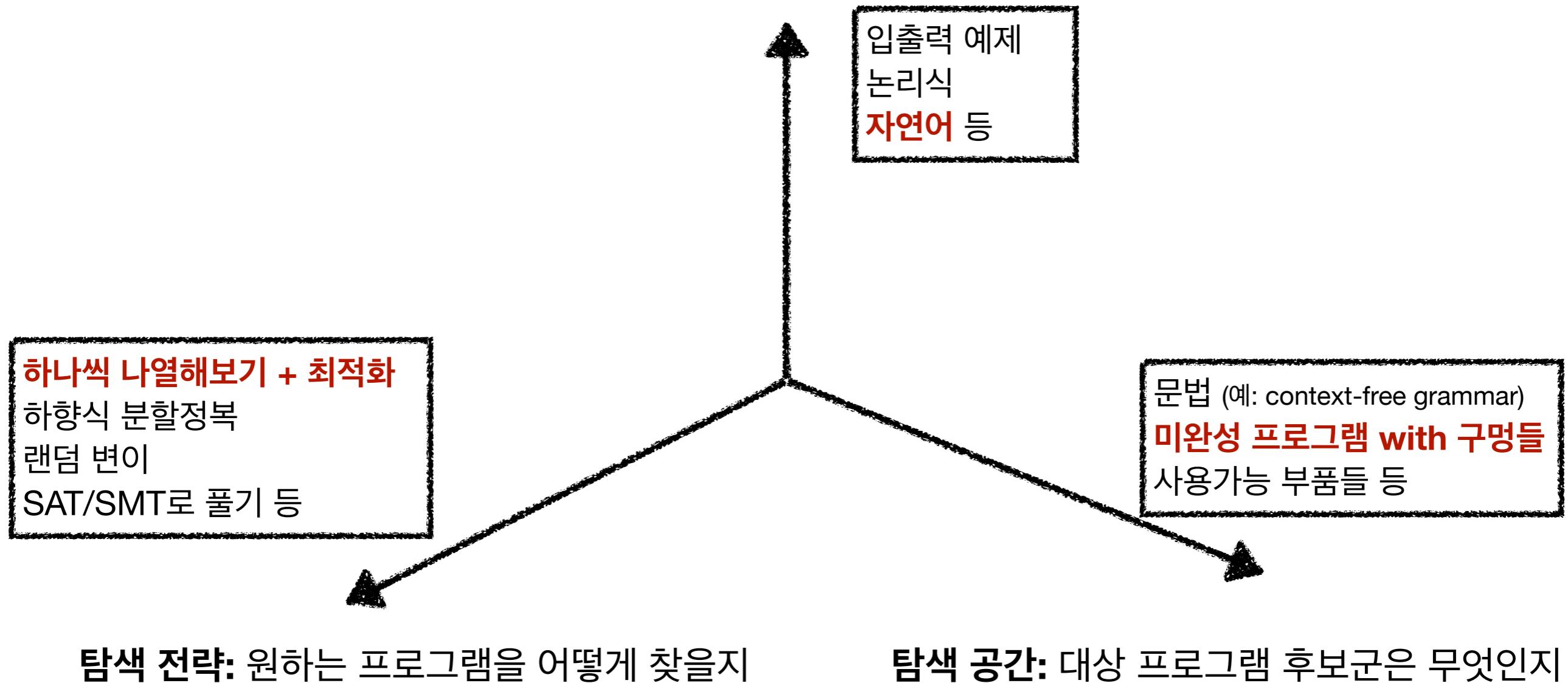
FlashFill (Excel 문자열 변환)

제약 조건: 프로그램이 달성할 목표를 어떻게 기술할지



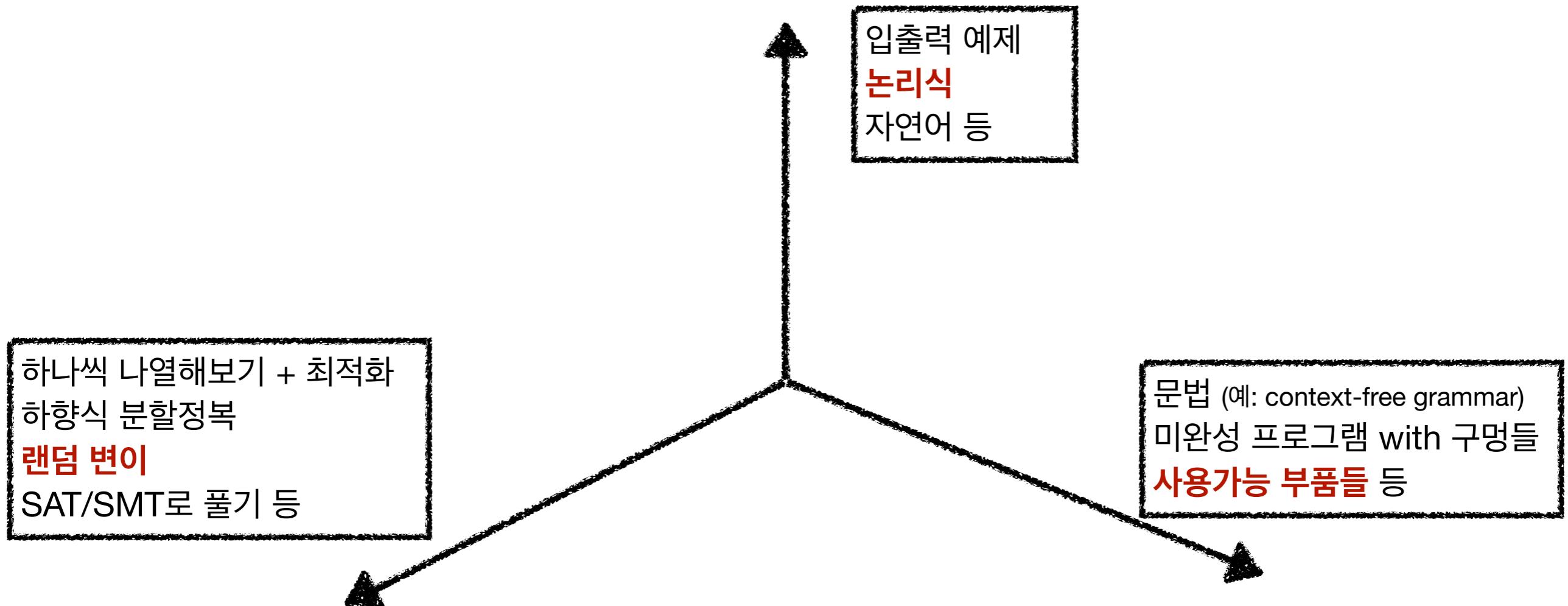
SQLizer (DB질의문 합성)

제약 조건: 프로그램이 달성할 목표를 어떻게 기술할지



STOKE (프로그램 최적화)

제약 조건: 프로그램이 달성할 목표를 어떻게 기술할지



탐색 전략: 원하는 프로그램을 어떻게 찾을지

탐색 공간: 대상 프로그램 후보군은 무엇인지

너무 다양한 합성 문제들

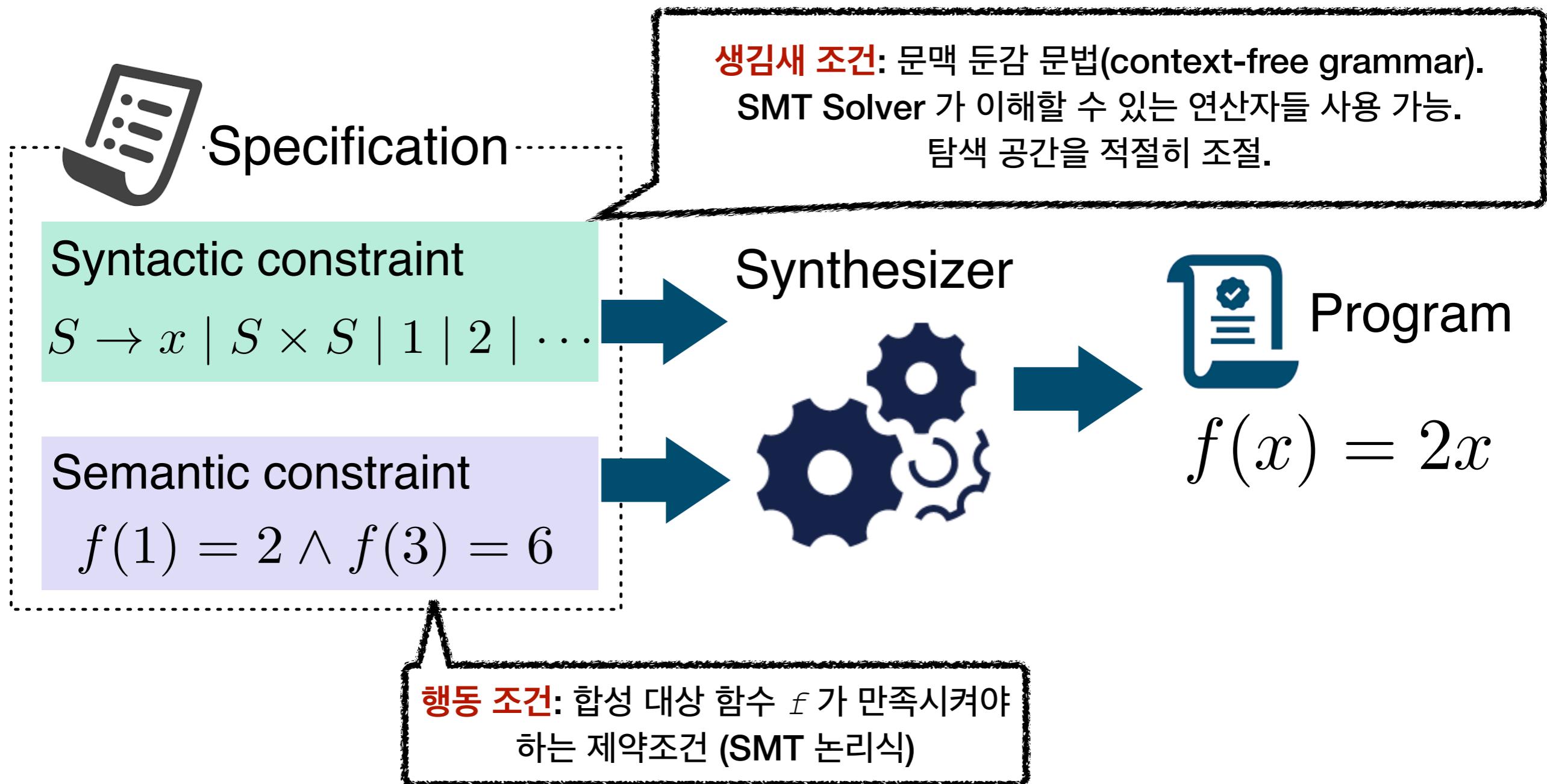
- 각 응용 도메인에 대해서 따로따로 개발되는 합성 기술들의 비교가 어려움.
- 관찰:
 - 대부분의 합성 문제가 문맥 둔감 문법(context-free grammar)로 기술되는 프로그램 공간에서 하나의 해를 찾는 문제임.
 - 또한, 원하는 프로그램에 대한 제약조건이 대개의 경우 SMT 논리식으로 작성될 수 있음.

너무 다양한 합성 문제들

- 결과:
 - 합성 문제를 기술하는 표준화된 포맷 (Syntax-guided Synthesis, SyGuS) 고안.
 - 이 포맷으로 기술된 합성 문제를 푸는데 사용될 수 있는 다목적(general-purpose) 합성 전략들 등장
 - 공통 포맷으로 인해 서로다른 합성기들 성능비교 가능해짐. 매년 합성기 경진대회 (SyGuS competition) 개최.

Syntax-Guided Program Synthesis (SyGuS)[†]

다양한 종류의 합성 문제들을 기술하기 위한 표준화된 포맷



[†]<http://www.sygus.org>

합성전략의 두 갈래

- 제약조건의 기술 방법에 따라 전략이 두 갈래로 갈림.
- 귀납 합성 (inductive synthesis) — 제약조건이 입출력 예제 (a.k.a. Programming by example)
 - 합성된 프로그램의 올바름 체크: 입출력 매치
- 연역 합성 (deductive synthesis) — 제약조건이 논리식
 - 합성된 프로그램의 올바름 체크: SAT/SMT solver 사용

내용

- 프로그램 합성 개괄
- 귀납 합성 전략
- 연역 합성 전략

귀납 합성 전략

- 하나씩 나열해보기 + 최적화 (일반적으로 적용 가능, 느림)
 - 하향식 (Top-Down) / 상향식 (Bottom-Up)
 - 최적화: 통계모델 사용 (PLDI'18), 프로그램 분석 등
- 하향식 분할정복 (특정 문제들에만 적용 가능, 빠름)
 - 제품에 탑재될 정도로 빠름 (Excel, PowerShell, Visual Studio)
- 상향식 나열 + 효율적 자료구조 쓰기 (일반적으로 적용 가능, 빠름)
 - 최근의 연구 성과 (POPL'21)

귀납 합성 전략

- 하나씩 나열해보기
- 하향식 (**Top-Down**)
 - 상향식 (Bottom-Up)
 - 하향식 + 통계모델 사용 (PLDI'18)
- 하향식 분할정복
- 상향식 나열 + 효율적 자료구조 쓰기 (POPL'21)

하향식(Top-Down) 나열

- 시작 비말단기호(start non-terminal)에서 시작
- 비말단 기호하나를 (보통 가장 왼쪽의 것) 선택하여 문법 규칙을 적용하여 확장
- 예제: 합성대상 $f : \mathbb{Z} \rightarrow \mathbb{Z}$

생김새 조건: $S \rightarrow x \mid 1 \mid -S \mid S + S \mid S \times S$

행동조건: $f(1) = 2$

하향식(Top-Down) 나열

Iter 0	S				
Iter 1	x	1	$-S$	$S + S$	$S \times S$
Iter 2	$-x$	-1	$-(-S)$	$x + S \dots$	
Iter 3	$-(-x)$	$-(-1)$	$x + 1$		

하향식 나열 최적화

- 탐색 중 미완성 프로그램 (예: $x + S$) 들 등장.
- 미완성 프로그램의 행동을 미리 어림잡아 (프로그램 분석)
해가 될 가능성이 없는 것들은 미리 가지치기 가능
- 예: 입출력 예제 조건이 $f(2) = 3$ 일 경우 $x \times S$ 는 후보에
서 제외하고 더 탐색 진행 안함
- 타입/SMT제약식 풀기/요약해석 등 다양하게 이용.

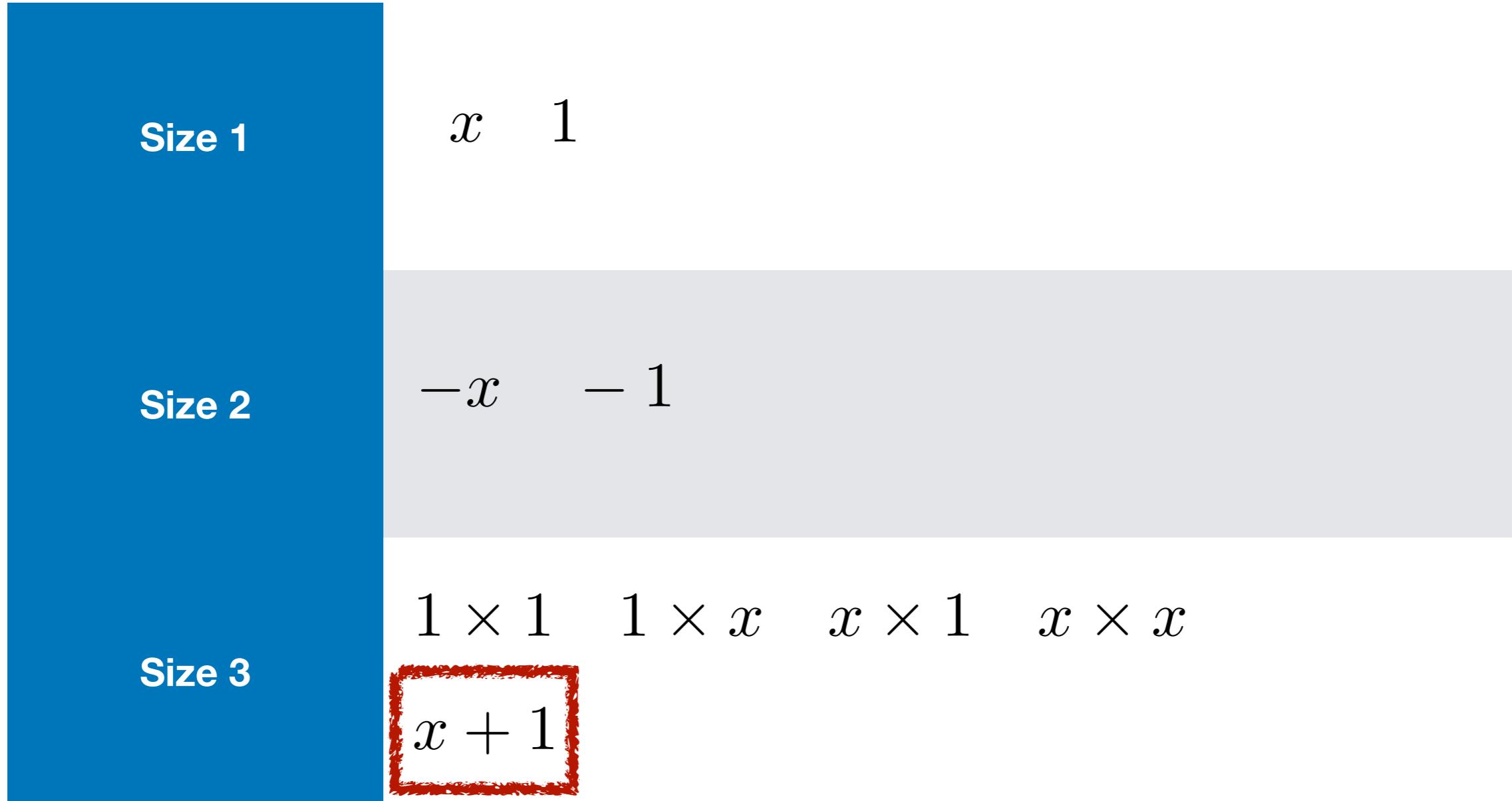
귀납 합성 전략

- 하나씩 나열해보기
 - 하향식 (Top-Down)
- 상향식 (**Bottom-Up**)
 - 하향식 + 통계모델 사용 (PLDI'18)
- 하향식 분할정복
- 상향식 나열 + 효율적 자료구조 쓰기 (POPL'21)

상향식(Bottom-Up) 나열

- 말단기호(terminal)에서 시작
- 가장 작은 프로그램들 부터 만들고, 그것들을 조합하여 더 큰 프로그램을 만듦.

상향식(Bottom-Up) 나열

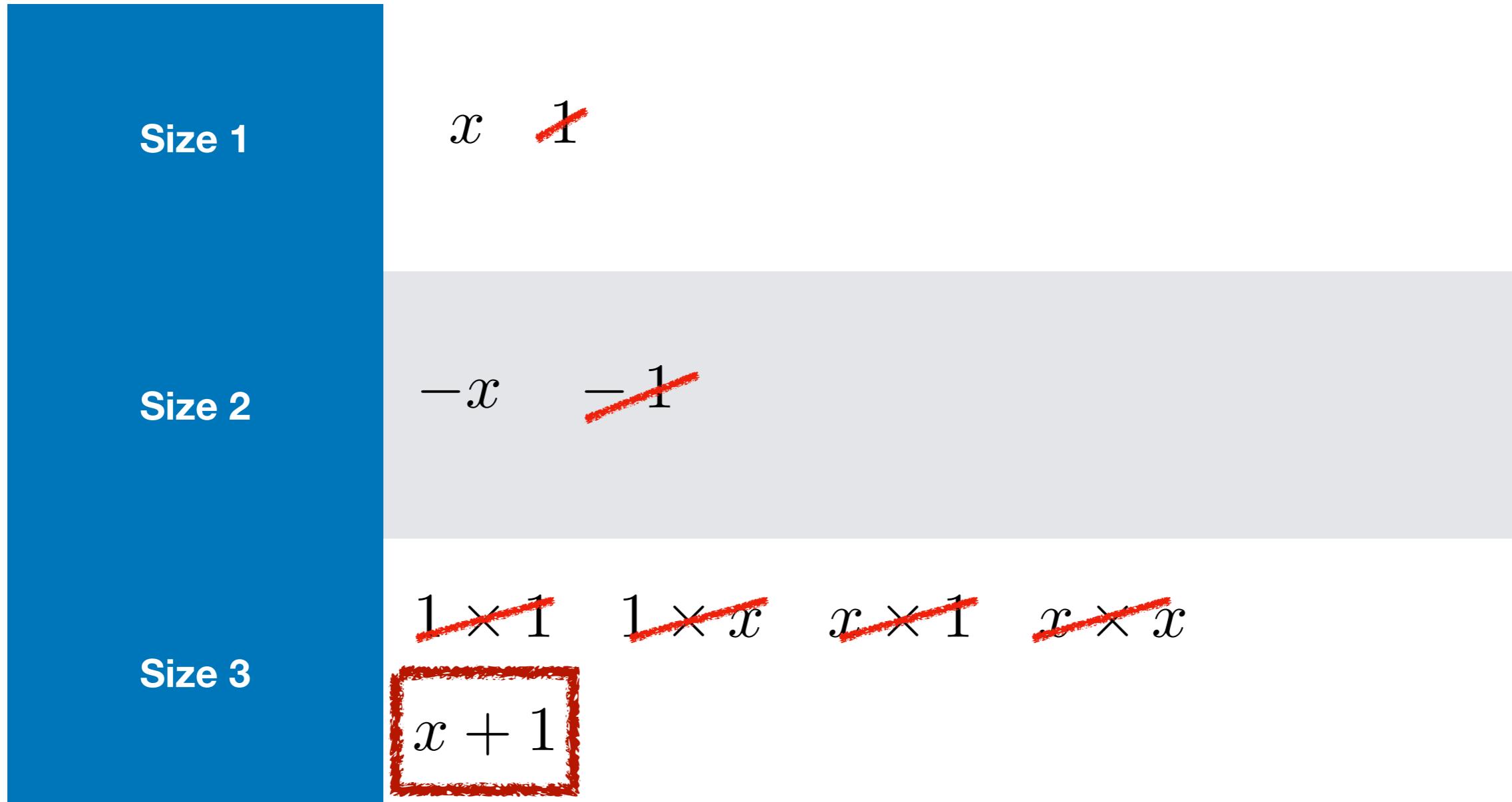


상향식 나열 최적화

- 탐색 중 완성 프로그램 (예: $I + I$) 들 등장. 완성 프로그램들은 실행이 가능
- 실행 결과가 같은 부품들 중 하나만 남겨놓고 더 큰 프로그램 생성에 사용.
- 이를 같은 것들 (observational equivalence) 뭉뚱그리기라고 함.

상향식(Bottom-Up) 나열

$$x = 1 \because f(1) = 2$$

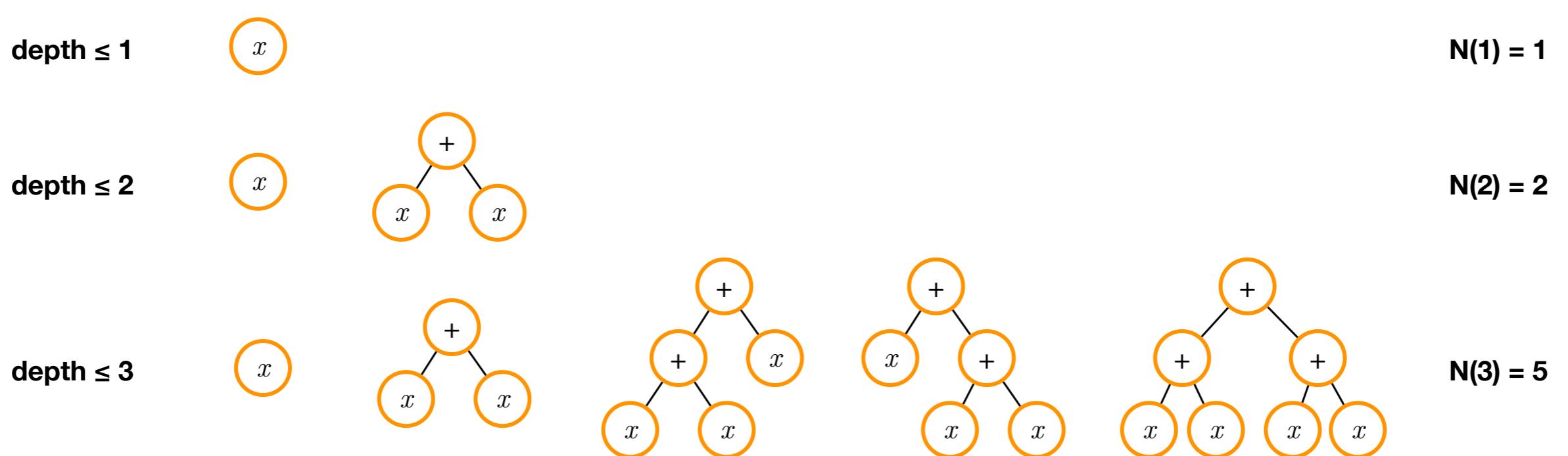


하나씩 나열하기

- 장점: 귀납 합성 뿐 아니라 일반적으로 거의 모든 종류의 프로그램 합성 문제에 사용 가능
- 단점: 탐색 공간이 너무나 크므로, 여러가지 최적화(예: 같은 것 뭉뚱그리기)에도 불구하고 큰 프로그램 찾기 불가능

탐색공간 얼마나 큰가?

$$S \rightarrow x \mid S + S$$



$$N(d) = 1 + N(d - 1)^2$$

*Examples from Nadia Polikarpova's slides

탐색공간 얼마나 큰가?

$$S \rightarrow x \mid S + S$$

$$N(d) = 1 + N(d - 1)^2$$

$N(1) = 1$
 $N(2) = 2$
 $N(3) = 5$
 $N(4) = 26$
 $N(5) = 677$
 $N(6) = 458330$
 $N(7) = 210066388901$
 $N(8) = 44127887745906175987802$
 $N(9) = 1947270476915296449559703445493848930452791205$
 $N(10) = 3791862310265926082868235028027893277370233152247388584761734150717768254410341175325352026$

!!

*Examples from Nadia Polikarpova's slides

귀납 합성 전략

- 하나씩 나열해보기
 - 하향식 (Top-Down)
 - 상향식 (Bottom-Up)
- 하향식 + 통계모델 사용 (**PLDI'18**)
- 하향식 분할정복
- 상향식 나열 + 효율적 자료구조 쓰기 (**POPL'21**)

통계 모델을 이용한 탐색 가이드

- 하나씩 나열하기의 성능 문제는 탐색 공간에 존재하는 모든 프로그램을 모두 “공평”하게 대하기 때문.
- 모든 프로그램이 “그럴싸”한 것은 아님. 프로그램에는 반복적이고 예측 가능한 패턴이 존재.

```
for (i = 0; i < 100; ??)
```

- 프로그램의 통계적 모델: 프로그램에 대한 확률 분포 결정

$$Pr(\text{??} \rightarrow i++ \mid \text{for}(i=0; i<100; \text{??})) = 0.85$$

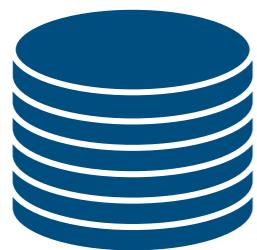
$$Pr(\text{??} \rightarrow i-- \mid \text{for}(i=0; i<100; \text{??})) = 0.01$$

통계 모델을 이용한 탐색 가이드

- 통계 프로그램 모델은 역난독화, IDE의 코드 완성, 프로그램 자동 수정 등에 이미 사용되어옴
 - N-gram, Probabilistic context-free grammar (PCFG) 등
- 아이디어: 프로그램 크기 오름차순이 아닌, 그럴싸함 확률 내림차순으로 나열

통계적 프로그램 모델

- 사람이 쓴, 혹은 합성된 코드로부터 학습 가능



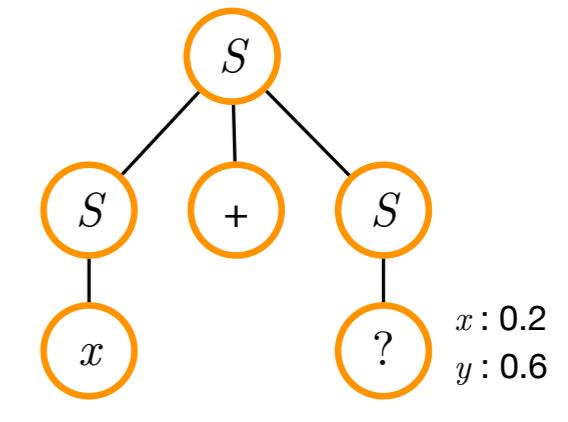
$$Pr(S \rightarrow S + S) = 0.3$$

$$Pr(S \rightarrow x | S + S) = 0.8$$

$$Pr(S \rightarrow x | \dots x + S) = 0.2$$

$$Pr(S \rightarrow y | x + S) = 0.6$$

...



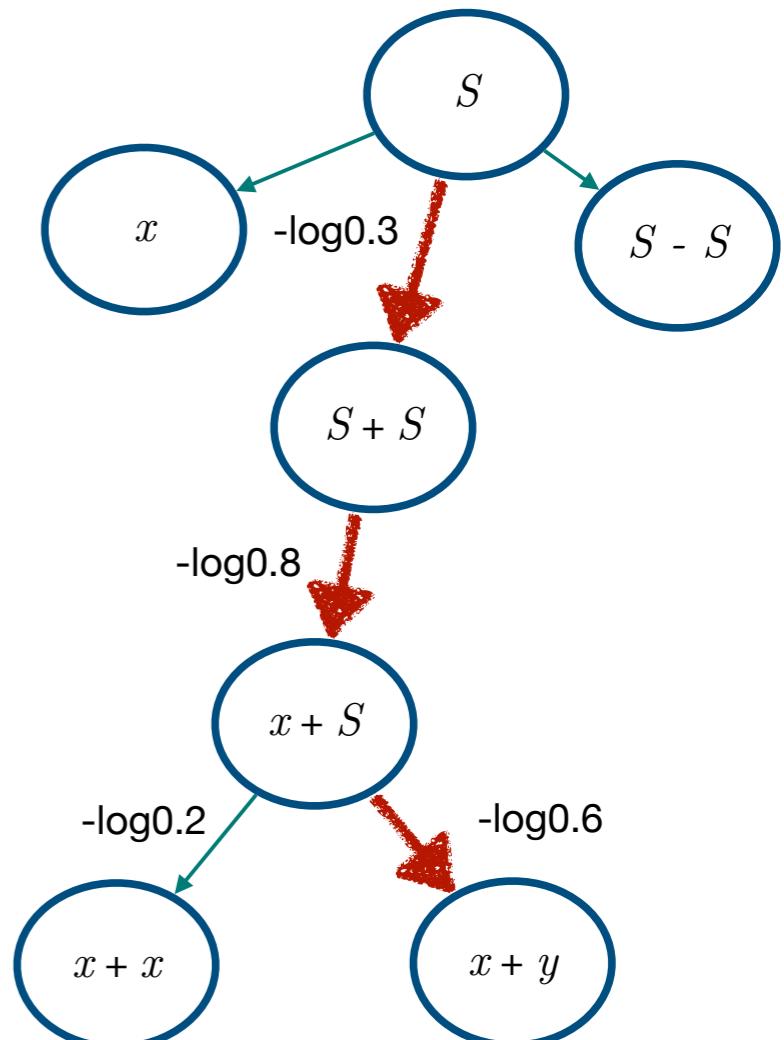
Program Corpus

Learned Probabilistic Model

Probability of Programs

통계 모델을 이용한 탐색 가이드

- 통계 모델로부터 그래프 구축
 - 노드: 시작 비밀단 기호로부터 도출 가능한 (부분)프로그램들
 - 간선: 문법 규칙 1번 적용에 대응
 - 가중치: 문법 규칙이 적용될 확률에 $-\log$ 취한 값
 - 시작 비밀단 기호부터 가장 거리가 짧은 순으로 프로그램 나열
 - 이 때 최단거리 계산은 A* 알고리즘 사용



$Pr(S \rightarrow S + S) = 0.3$
$Pr(S \rightarrow x S + S) = 0.8$
...
$Pr(S \rightarrow x x + S) = 0.2$
$Pr(S \rightarrow y x + S) = 0.6$
...

실험 벤치마크

	A	B	C	D
1	Number	Phone		
2	02082012225	020-8201-2225		
3	02072221236	020-7222-1236		
4	0208123654	020-8123-654		
5	0207236523	020-7236-523		
6	02082012222	020-8201-2222		
7				
8				

STRING: End-user Programming
205 problems

complement

$\sim 01010001110101110000000000001111$
 $10101110001010001111111111110000$

bitwise and

$01010001110101110000000000001111$
 $\& 00110001011011100011000101101110$
 $00010001010001100000000000001110$

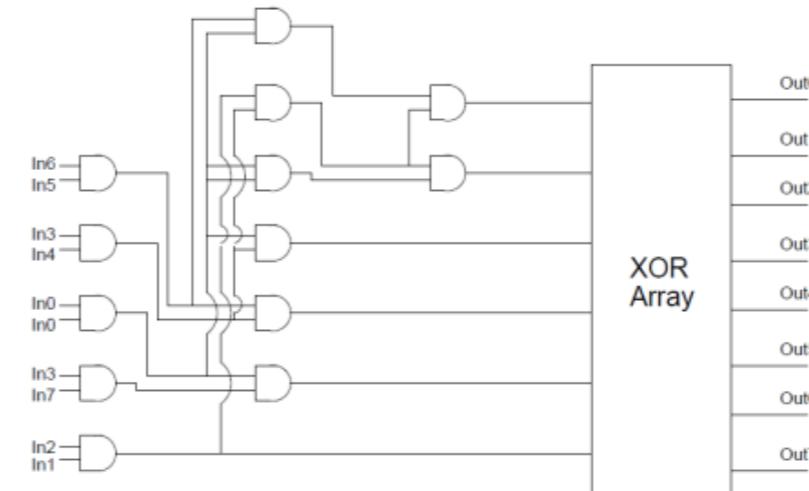
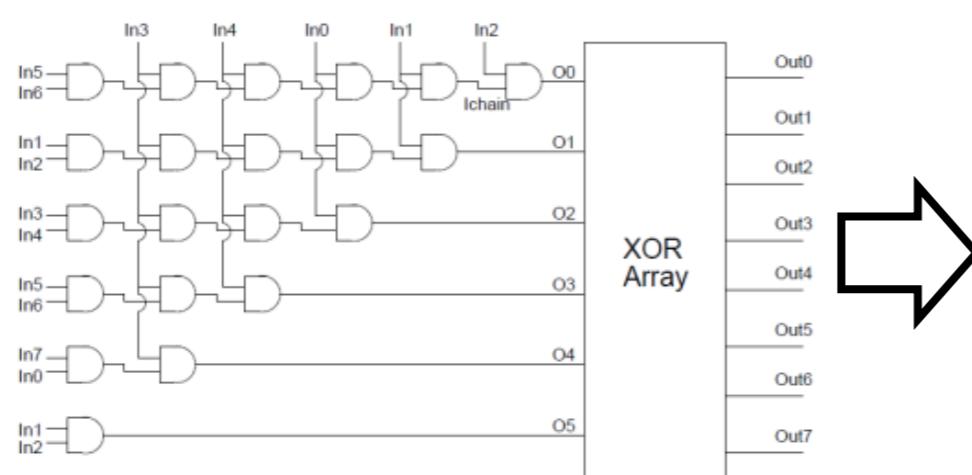
bitwise or

$01010001110101110000000000001111$
 $| 00110001011011100011000101101110$
 $0111000111111110011000101101111$

bitwise xor

$01010001110101110000000000001111$
 $\wedge 00110001011011100011000101101110$
 $01100000101110010011000101100001$

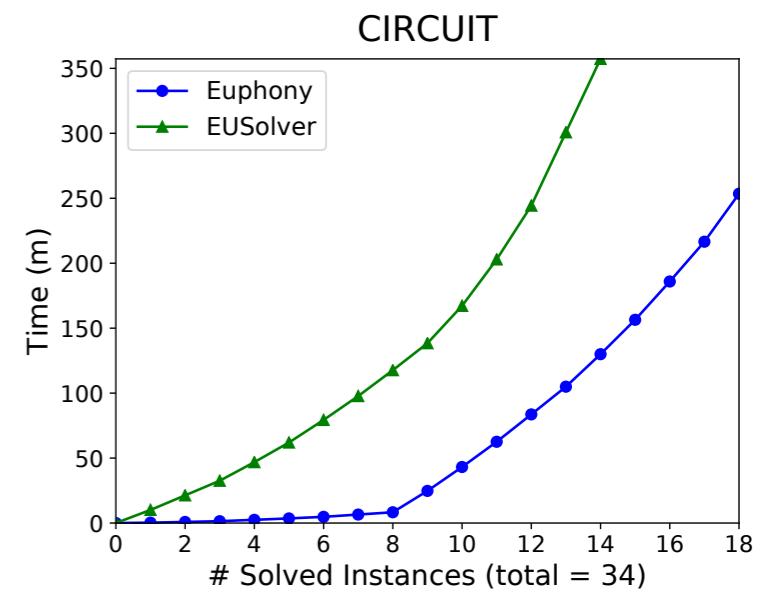
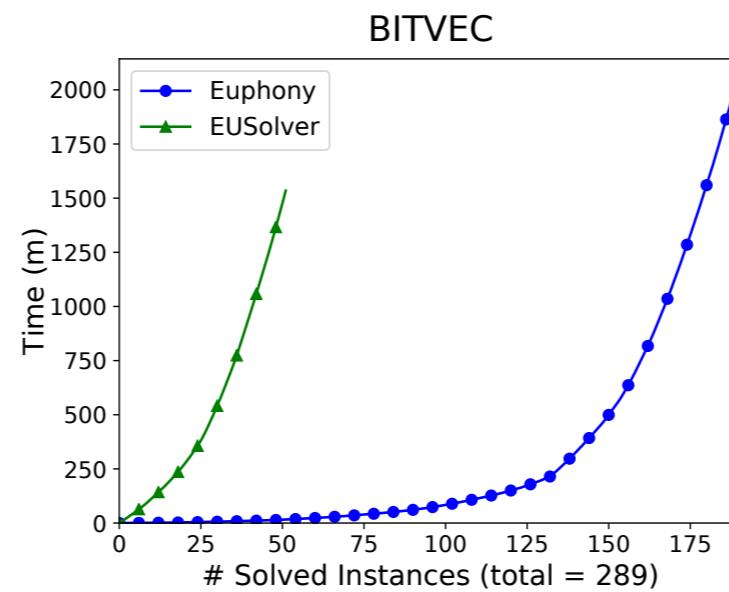
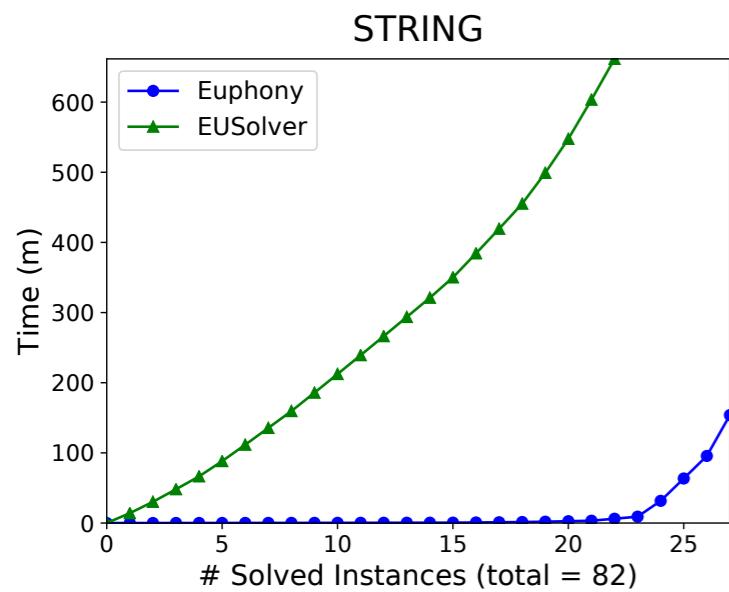
BITVEC: Efficient low-level algorithm
750 problems



CIRCUIT: Attack-resistant crypto circuits
212 problems

실험 결과

- EU Solver(확률모델 가이드 없는 상향식 나열 전략)와 비교
 - 학습: EU Solver 가 10분안에 푼 문제들의 솔루션
 - 테스트셋: 405 (timeout: 1시간)



귀납 합성 전략

- 하나씩 나열해보기
 - 하향식 (Top-Down)
 - 상향식 (Bottom-Up)
 - 하향식 + 통계모델 사용 (PLDI'18)
- 하향식 분할정복
- 상향식 나열 + 효율적 자료구조 쓰기 (POPL'21)

예제 합성 문제

- Goal: function f converting a given phone number x into another format

-  Specification

Syntactic: $S \rightarrow x \mid "+" \mid "-" \mid ":" \mid \text{ConCat}(S, S)$
 $\mid \text{SubStr}(S, I, I)$
 $\mid \text{IntToStr}(I)$
 $\mid \text{Replace}(S, S, S)$

$I \rightarrow 1 \mid \dots \mid 9 \mid I - I \mid \text{Length}(S)$

String concatenation

Semantic:

$$f("1_2-3.") = "1_2-3." \wedge f("4_56.78") = "4_56-78."$$

예제 합성 문제

- Goal: function f converting a given phone number x into another format

- S $\text{SubStr}(str, start_pos, len)$

e.g., $\text{SubStr}("abc", 0, 2) = "ab"$

Symbole:

“-” | “.” | $\text{ConCat}(S, S)$

| $\text{SubStr}(S, I, I)$

| $\text{IntToStr}(I)$

| $\text{Replace}(S, S, S)$

$I \rightarrow 1 | \dots | 9 | I - I | \text{Length}(S)$

Semantic:

$$f(" _ + 1 _ 2.3") = "1 _ 2-3." \wedge f(" _ + 4 _ 56.78") = "4 _ 56-78."$$

예제 합성 문제

- Goal: function f converting a given phone number x into another format

-  Specification

Syntactic: $S \rightarrow$ Replace($str, match, replacement$) $S, S)$

| e.g., Replace("aba", "a", "b") = "bba"

| introduction

| Replace(S, S, S)

$I \rightarrow 1 | \dots | 9 | I - I | \text{Length}(S)$

Semantic:

$f(" _ + 1_2.3") = "1_2-3." \wedge f(" _ + 4_56.78") = "4_56-78."$

예제 합성 문제

-  Solution: Size : 12 AST nodes

SubStr(Replace(ConCat(x , “.”), “.”, “-”), 2, Length(x) – 1)

하향식 분할정복 (Top-Down Propagation)

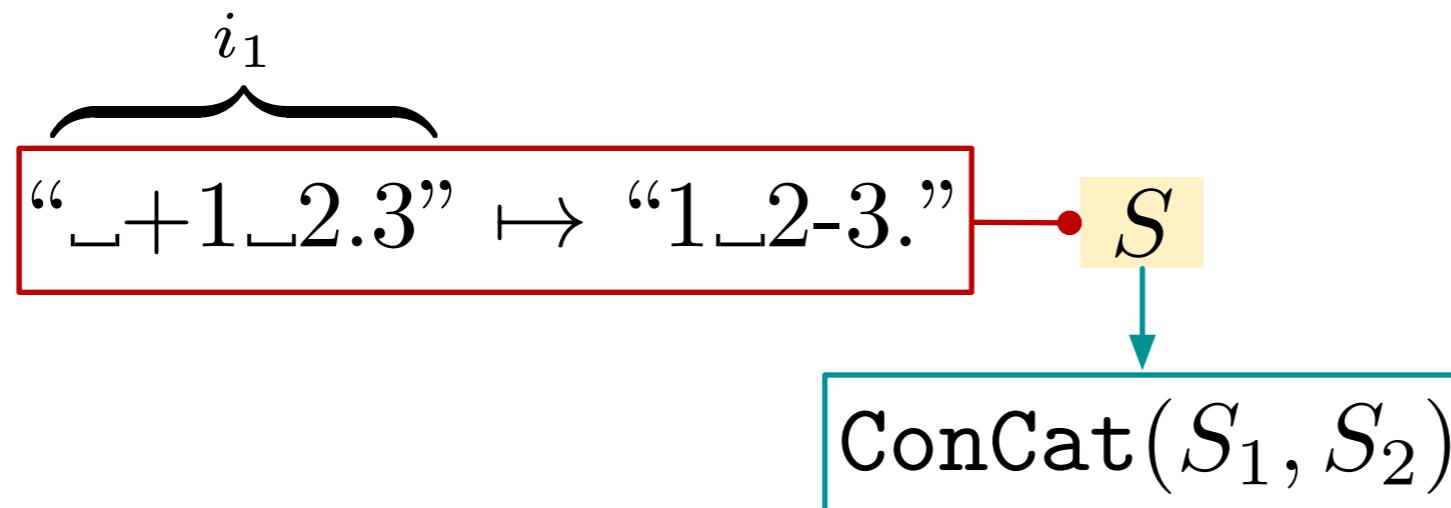
- **분할 정복:** 만약 $F(e_1, \dots, e_k)$ 꼴인 프로그램이 어떤 입력에 대해서 o 를 출력한다면, e_1, \dots, e_k 는 같은 입력에 대해 어떤 값을 출력해야 하는가?
- **역함수:** 인자들이 만족시킬 입출력 조건 유추

$$F^{-1}(o) = \{(a_1, \dots, a_k) \mid F(a_1, \dots, a_k) = o\}$$

- 가능 입력들 (*inverse-set or pre-image*):

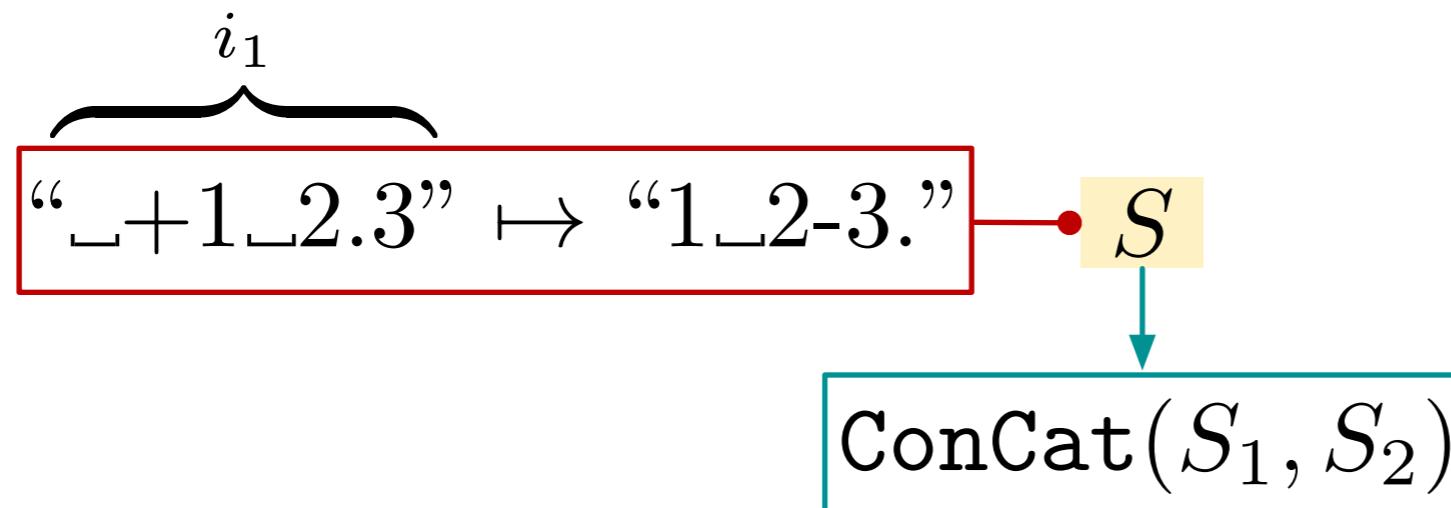
e.g., $\text{ConCat}^{-1}(\text{"USA"}) = \{(\text{"U"}, \text{"SA"}), (\text{"US"}, \text{"A"})\}$

하향식 분할정복 (Top-Down Propagation)



- Nonterminal Symbol
- Constraint on a nonterminal symbol
- One step derivation
- Solution directly derived

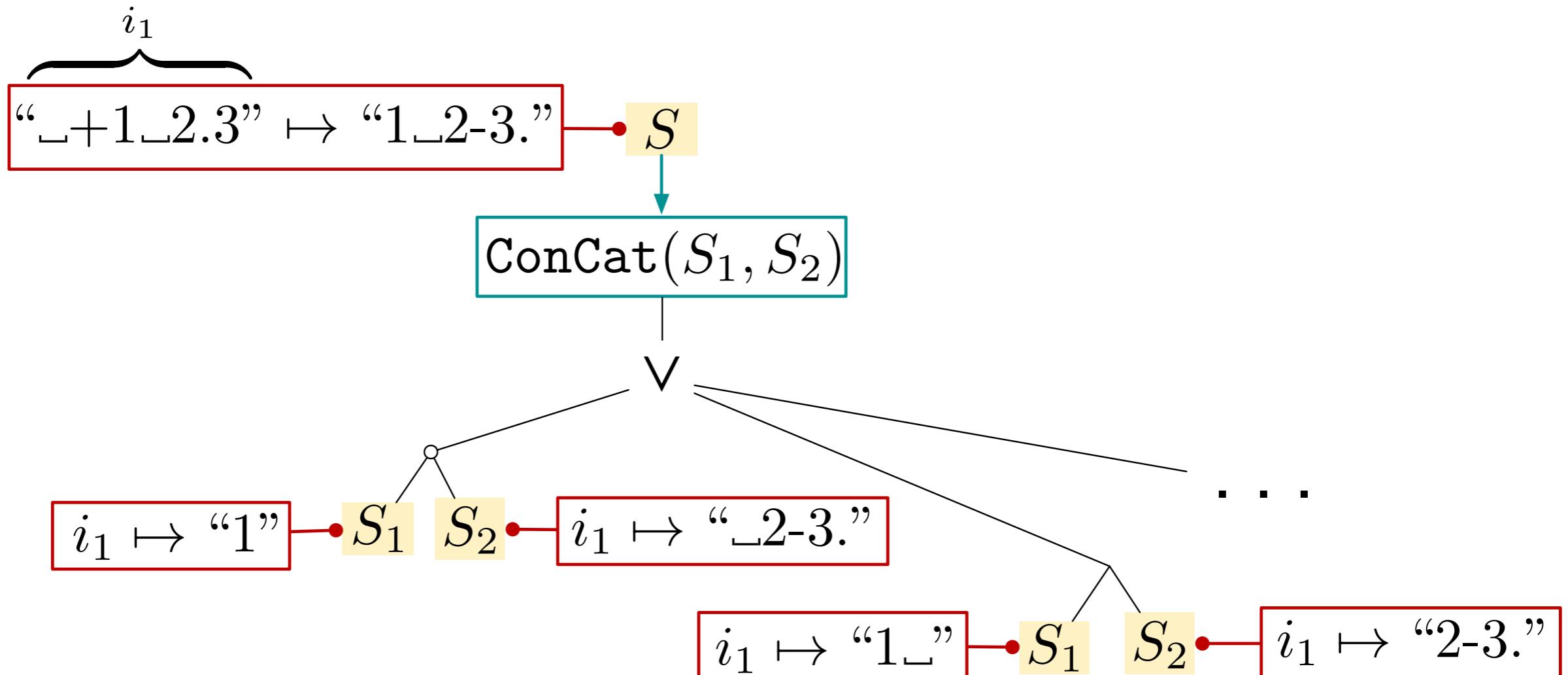
하향식 분할정복 (Top-Down Propagation)



$$\text{ConCat}^{-1}(\text{“1_2-3.”}) = \{(\text{“1”}, \text{“_2-3.”}), (\text{“1_”}, \text{“2-3.”}), \dots\}$$

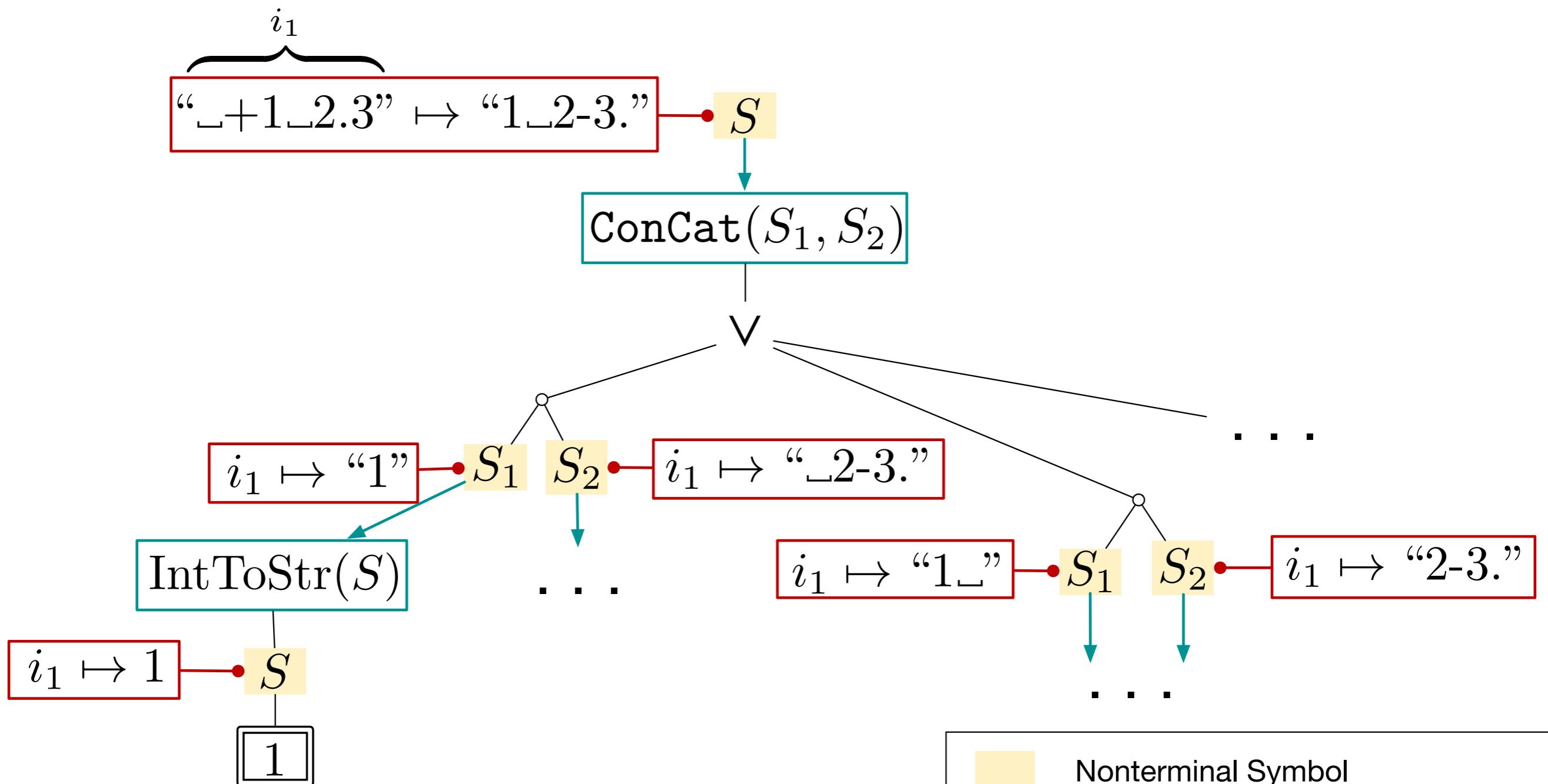
	Nonterminal Symbol
	Constraint on a nonterminal symbol
	One step derivation
	Solution directly derived

하향식 분할정복 (Top-Down Propagation)



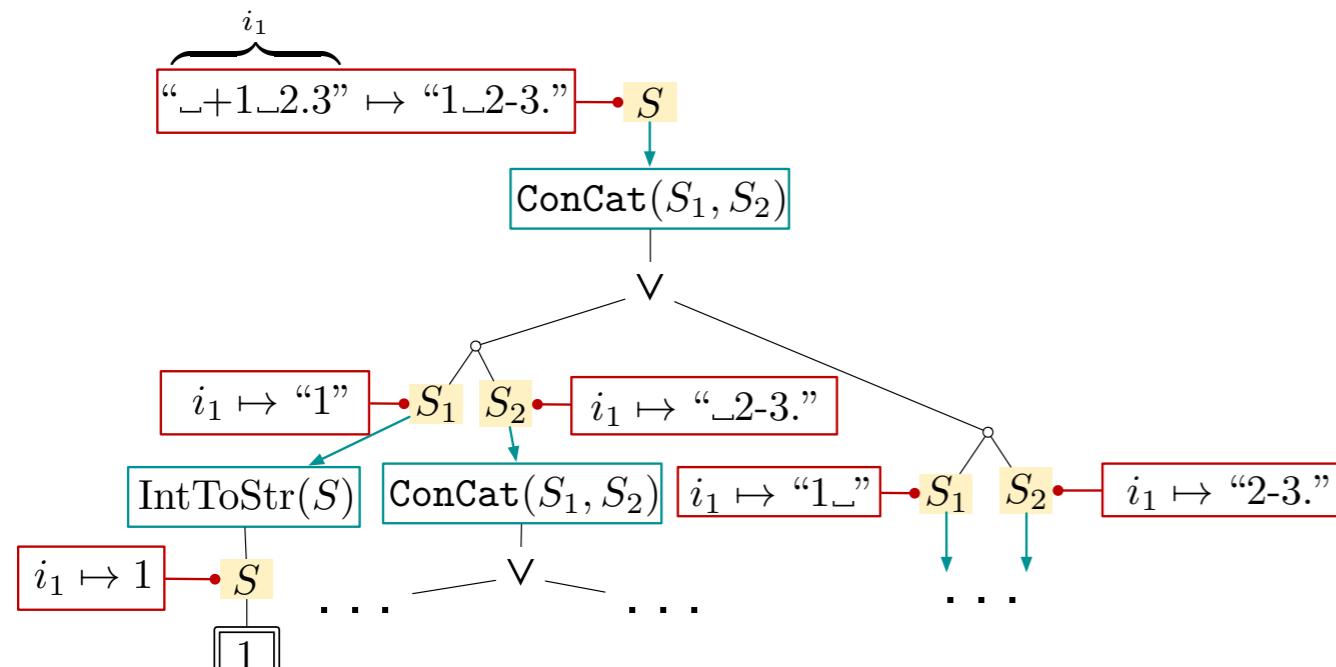
- Nonterminal Symbol
- Constraint on a nonterminal symbol
- One step derivation
- Solution directly derived

하향식 분할정복 (Top-Down Propagation)

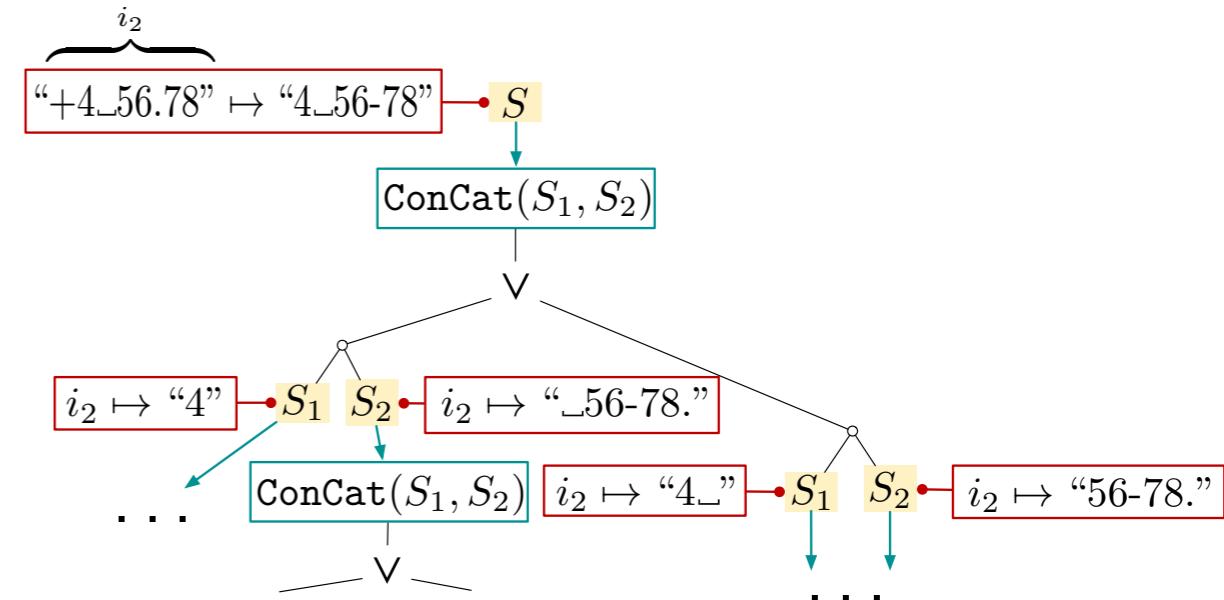


- Nonterminal Symbol
- Constraint on a nonterminal symbol
- One step derivation
- Solution directly derived

하향식 분할정복 (Top-Down Propagation)



입출력 예제1에 맞는
모든 프로그램들



입출력 예제2에 맞는
모든 프로그램들

$\{ \text{ConCat}(\text{ConCat}(\text{IntToStr}(1), _), _), \dots,$
 $\text{ConCat}(\text{ConCat}(\text{IntToStr}(1), \text{ConCat}(_, _)), \dots),$
 $\dots \}$



$\{ \text{ConCat}(\text{ConCat}(\text{IntToStr}(4), _), _), \dots,$
 $\text{ConCat}(\text{ConCat}(\text{IntToStr}(4), \text{ConCat}(_, _)), \dots),$
 $\dots \}$

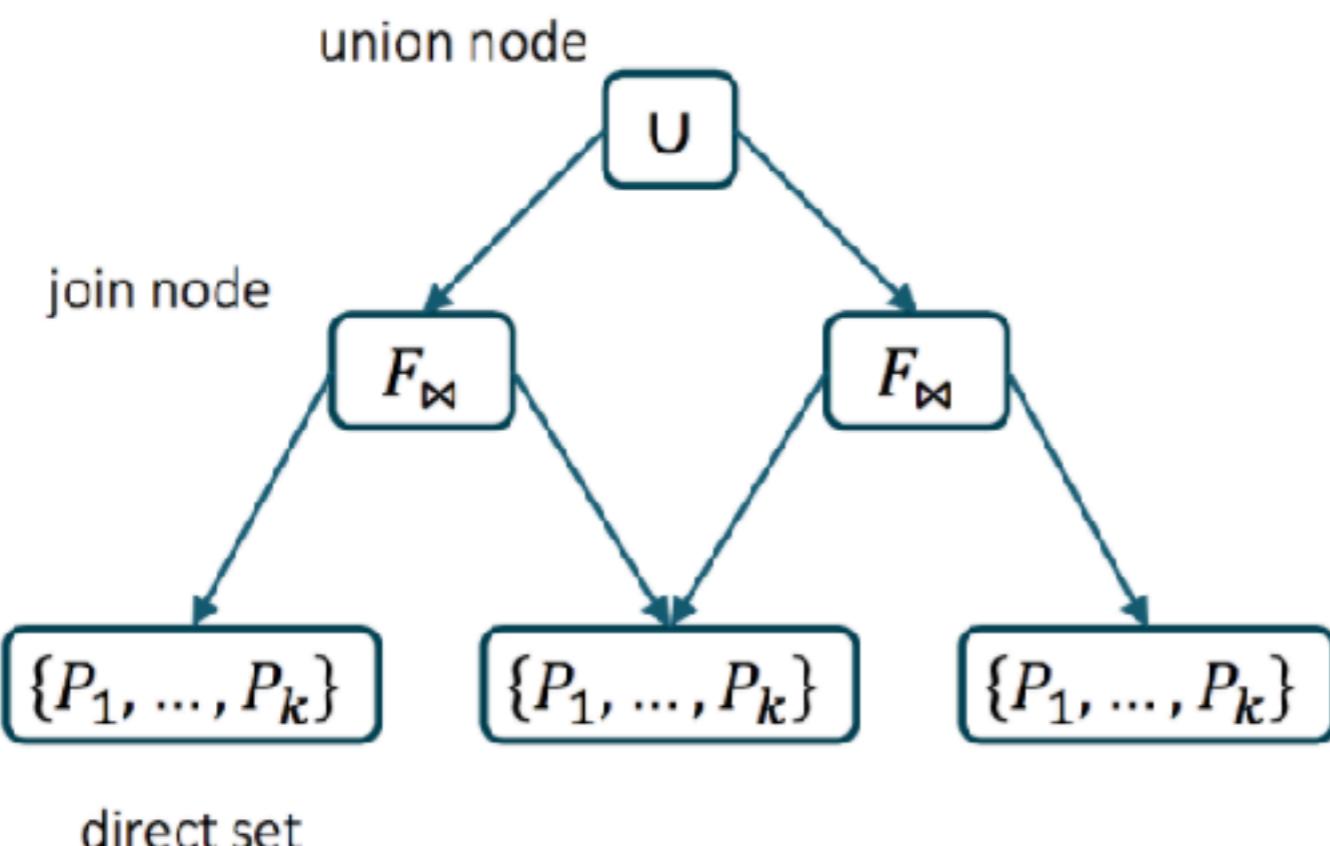
공간 효율적 자료구조에 저장됨.

최종 솔루션 집합

하향식 분할정복 + 효율적 자료구조

효율적 자료구조: 버전 공간 표현 (Version Space Algebra)

기계학습에서: 입출력 예제
에 맞는 모델 집합



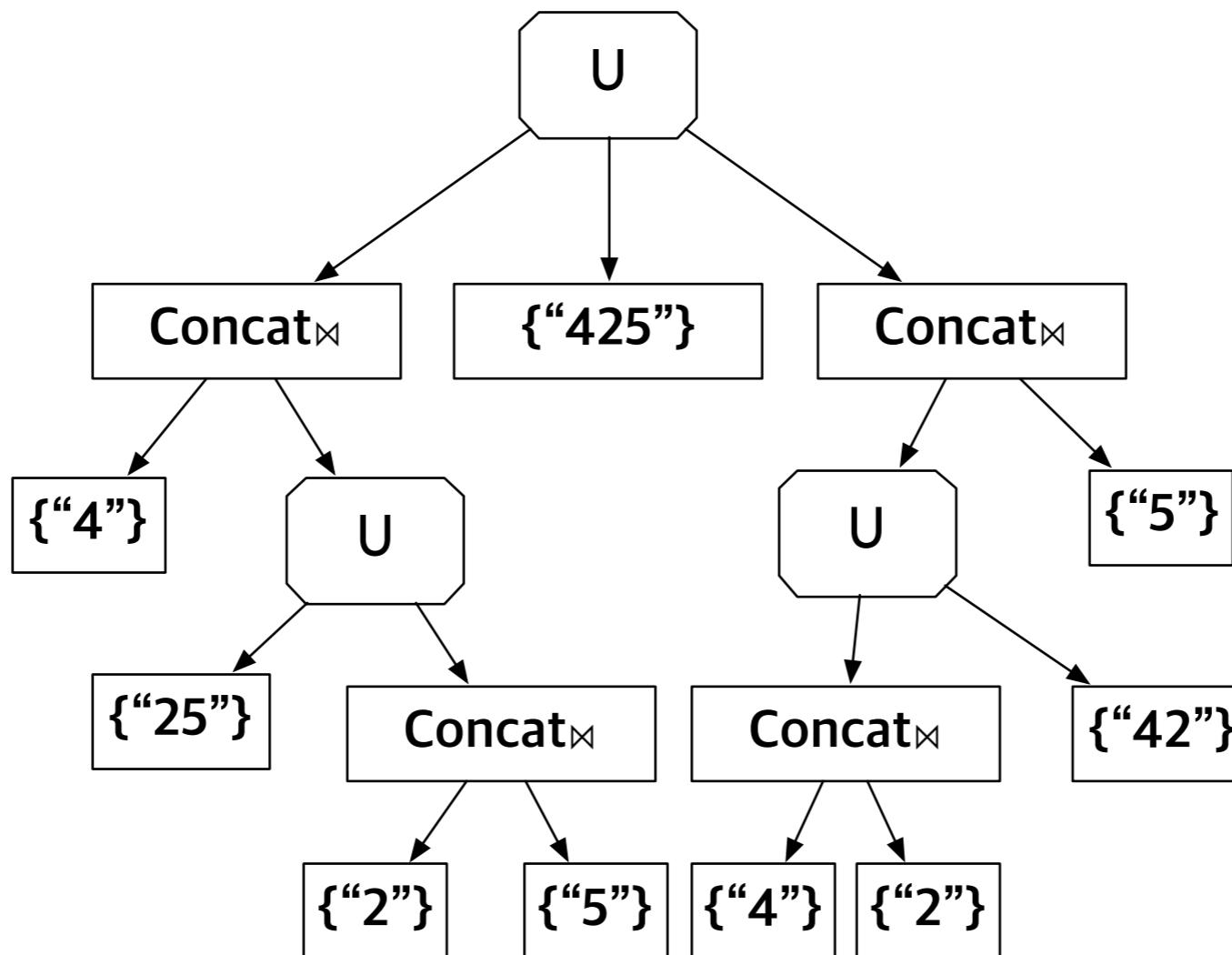
Volume of a VSA
(the number of nodes) $V(VSA)$

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

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

버전 공간 표현 (Version Space Algebra)

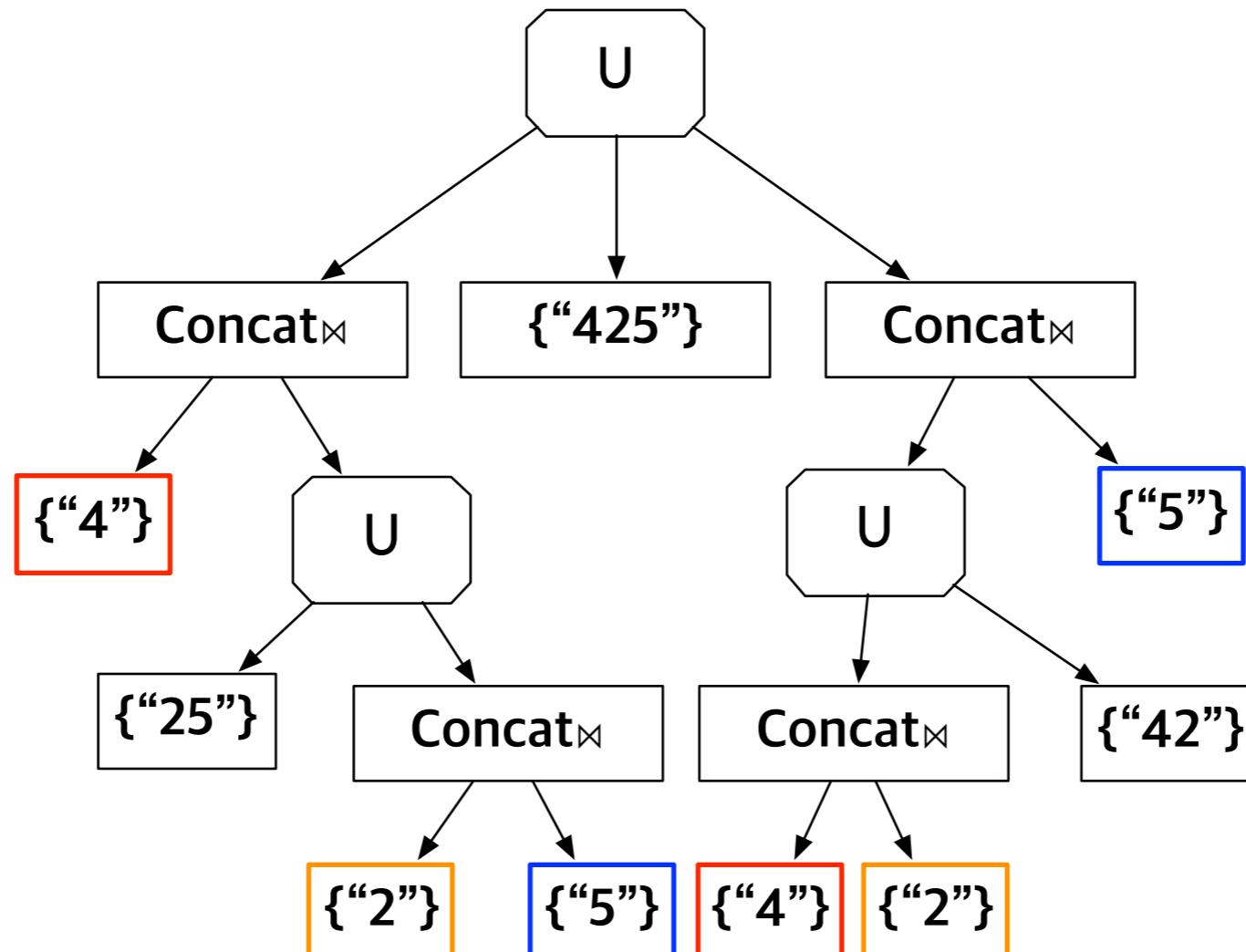
“425”를 출력하는 프로그램 집합을 표현하는 VSA



$$\text{Concat}_{\bowtie}(\mathbf{VS}_1, \mathbf{VS}_2) = \{\text{Concat}(P_1, P_2) \mid P_1 \in \mathbf{VS}_1, P_2 \in \mathbf{VS}_2\}$$

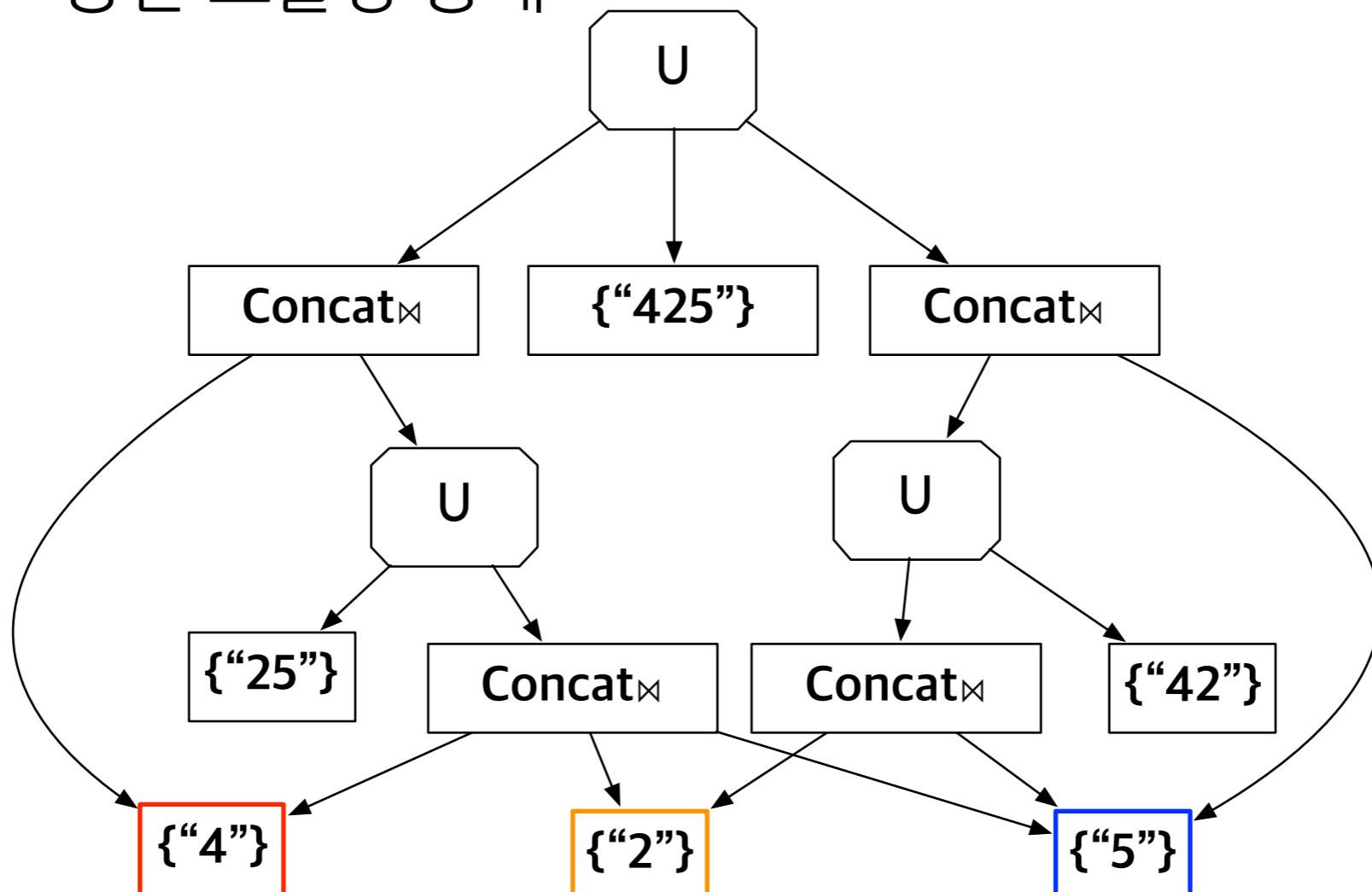
버전 공간 표현 (Version Space Algebra)

중복된 노드들 다수 존재

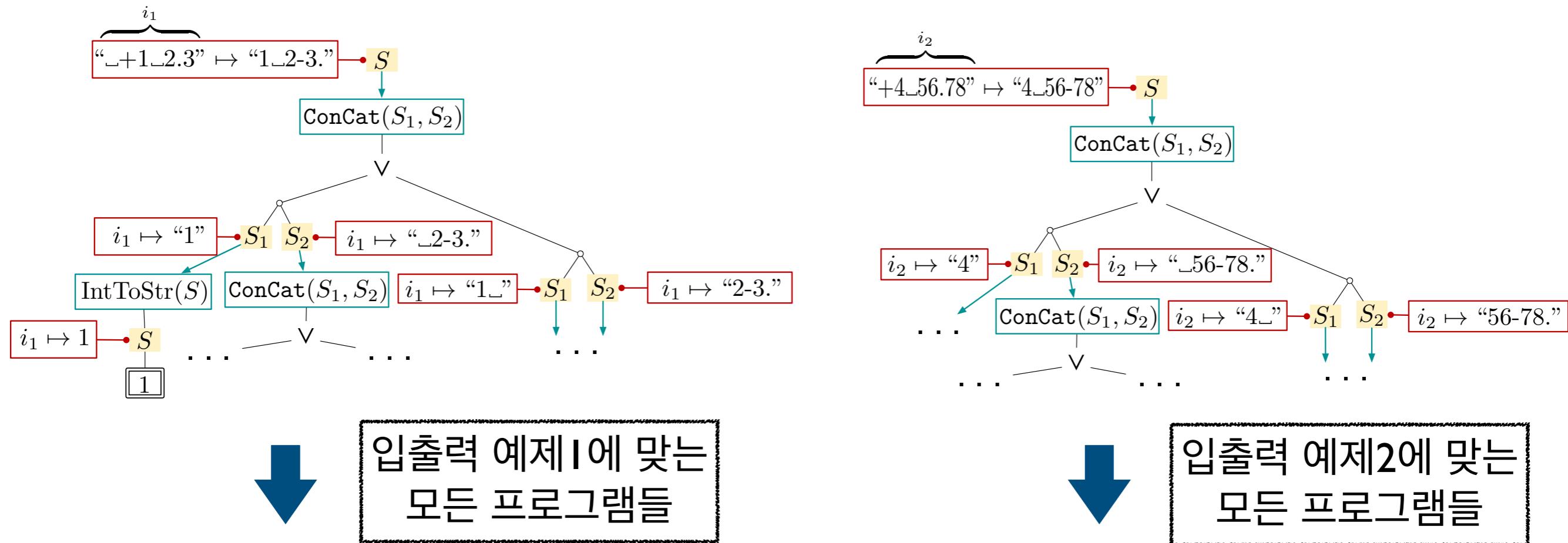


버전 공간 표현 (Version Space Algebra)

중복 제거 → 공간 효율성 증대



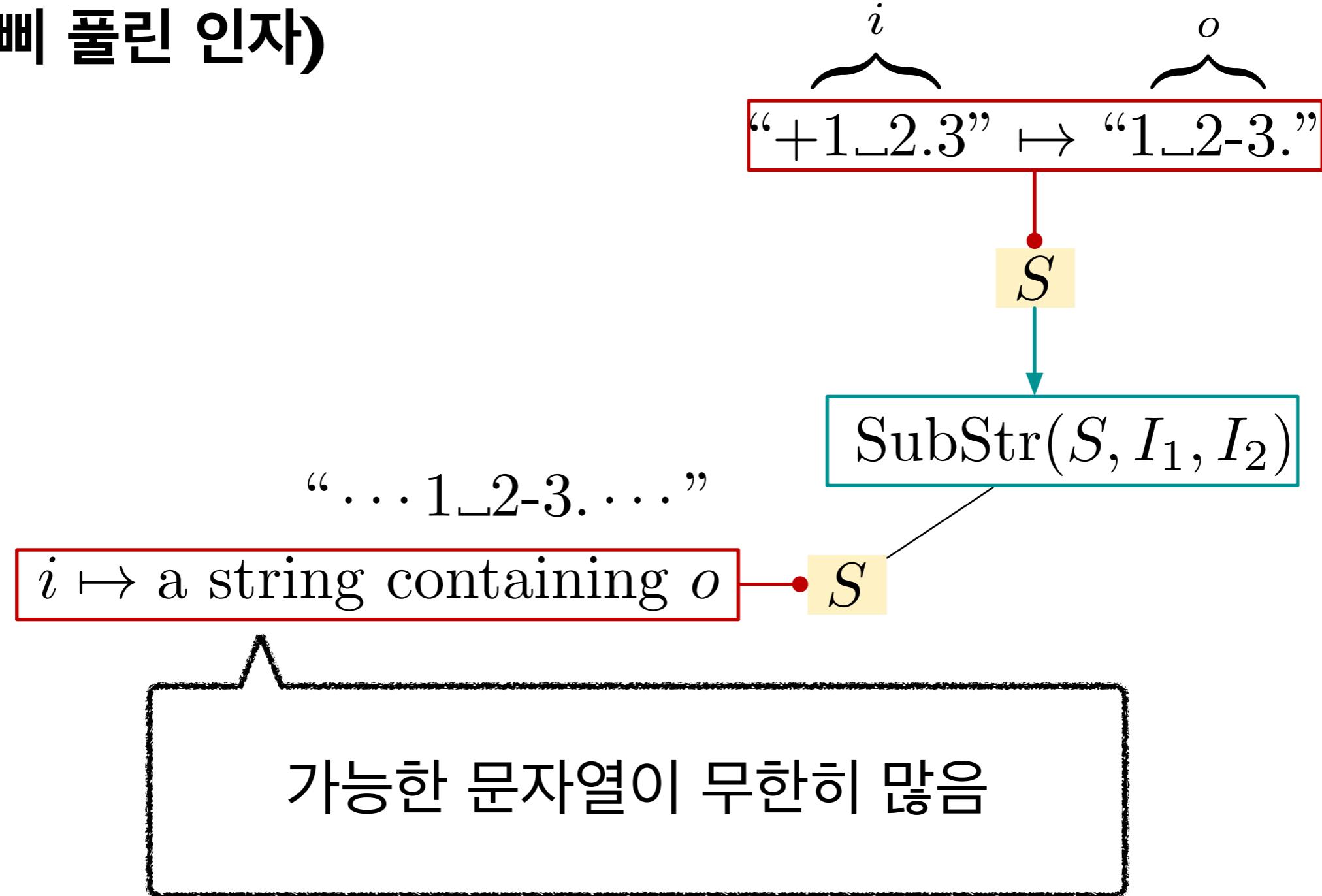
하향식 분할정복 (Top-Down Propagation)



- + 매우 효율적 (분할정복, 최종목표 달성에 필요한 문제들만 풀),
- 입출력 예제 갯수↑ → 성능 ↓
- 항상 적용 가능하지는 않음

하향식 분할정복의 문제점

가능한 입력이 무한한, 자유도가 지나친 인자들
(고삐 풀린 인자)

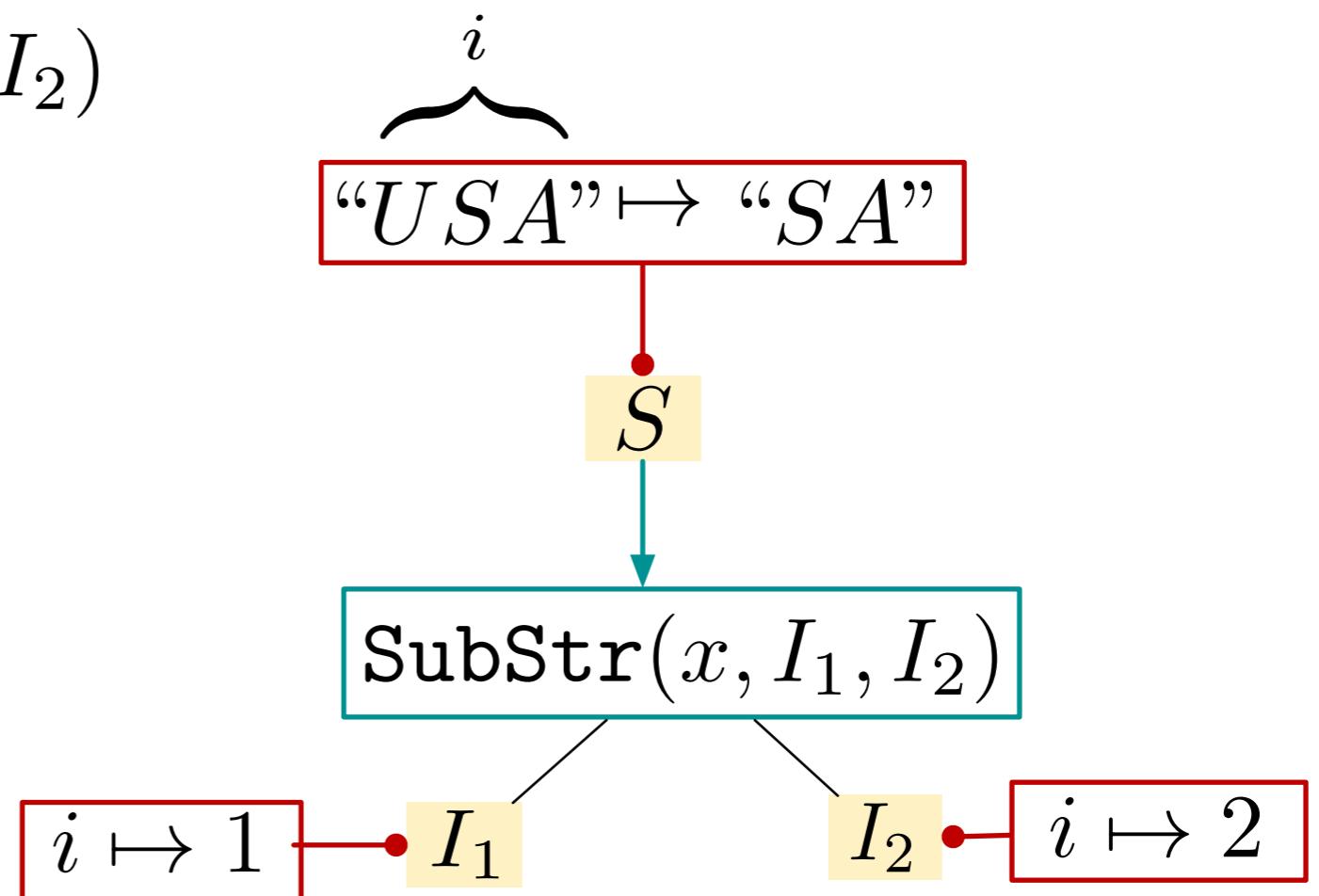


하향식 분할정복의 문제점

기존의 미봉책: 문법 제한하기

$$S \rightarrow \dots \mid \text{SubStr}(x, I_1, I_2)$$

Input example



DSL 정의에 “장인정신”이 필요 :(

“DSL design = Art + *Lots* of iterations”

“The DSL should be *expressive enough* to represent various tasks... and *restricted enough* to allow efficient search.”

Gulwani et al., Program synthesis

Polozov et al., FlashMeta: a framework for inductive program synthesis

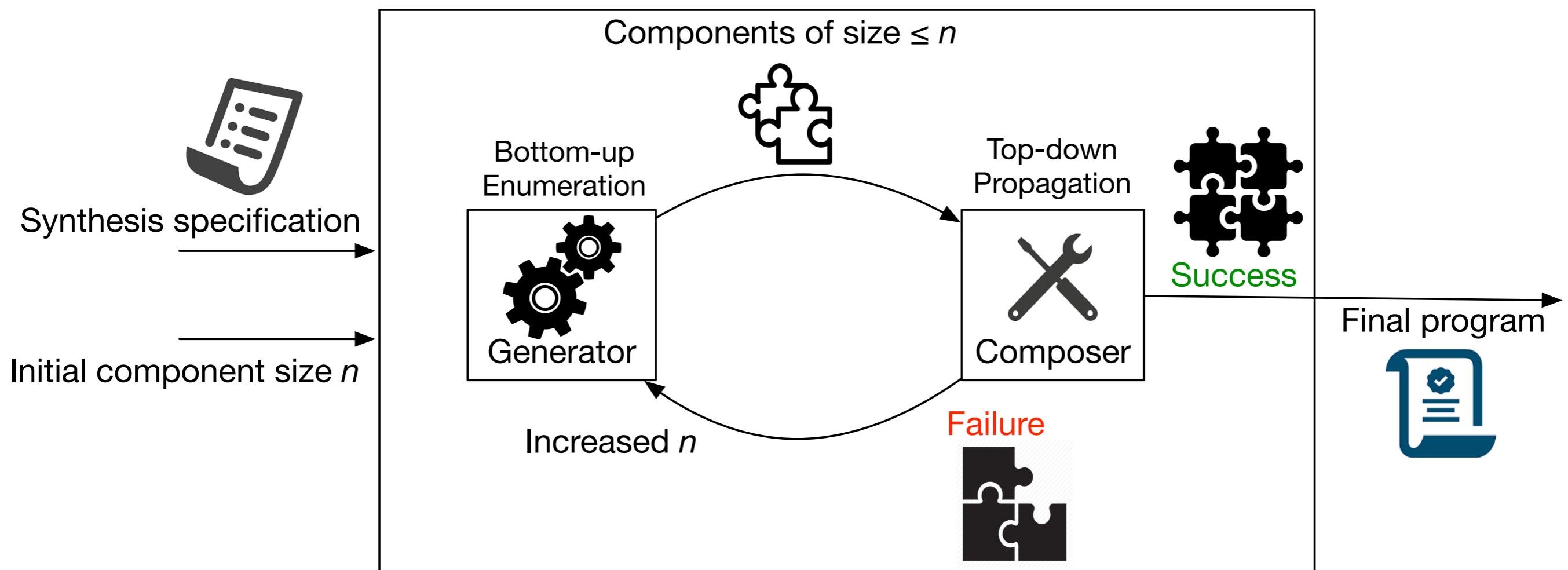
Gulwani et al., Programming by example (and its application to data wrangling)

귀납 합성 전략

- 하나씩 나열해보기
 - 하향식 (Top-Down)
 - 상향식 (Bottom-Up)
 - 하향식 + 통계모델 사용 (PLDI'18)
- 하향식 분할정복
- 상향식 나열 + 효율적 자료구조 쓰기 (**POPL'21**)

DUET: 양방향식 프로그램 합성⁺

하향식 분할정복(일반성 부족), 상향식 나열(느린 성능)의 한계 모두 해결.



+ Domain-unaware inductive synthesis combining Enumeration and Top-down Propagation

아이디어

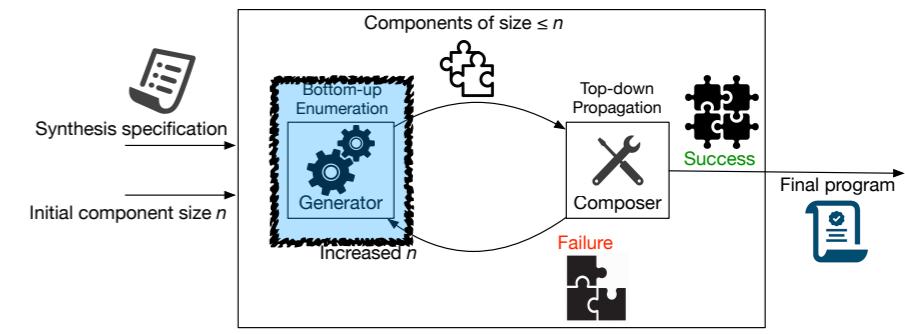
- 상향식 나열은 느리지만 “핵심”부분 표현식들은 빨리 찾을 수 있음.

SubStr(Replace(ConCat(x , “.”), “.”, “-”), 2, Length(x) – 1)

고빼풀린 인자들 자리에 들어갈 표현식

- 우리의 하향식 분할정복: 고빼 풀린 인자들 자리에 상향식 나열로 만든 프로그램들 끼워넣기.

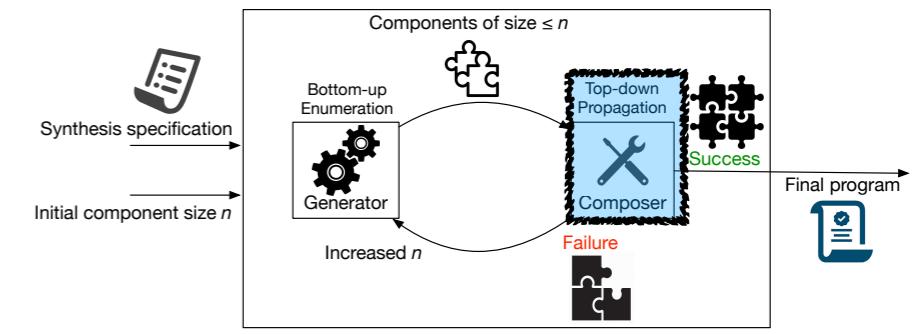
부품 생성



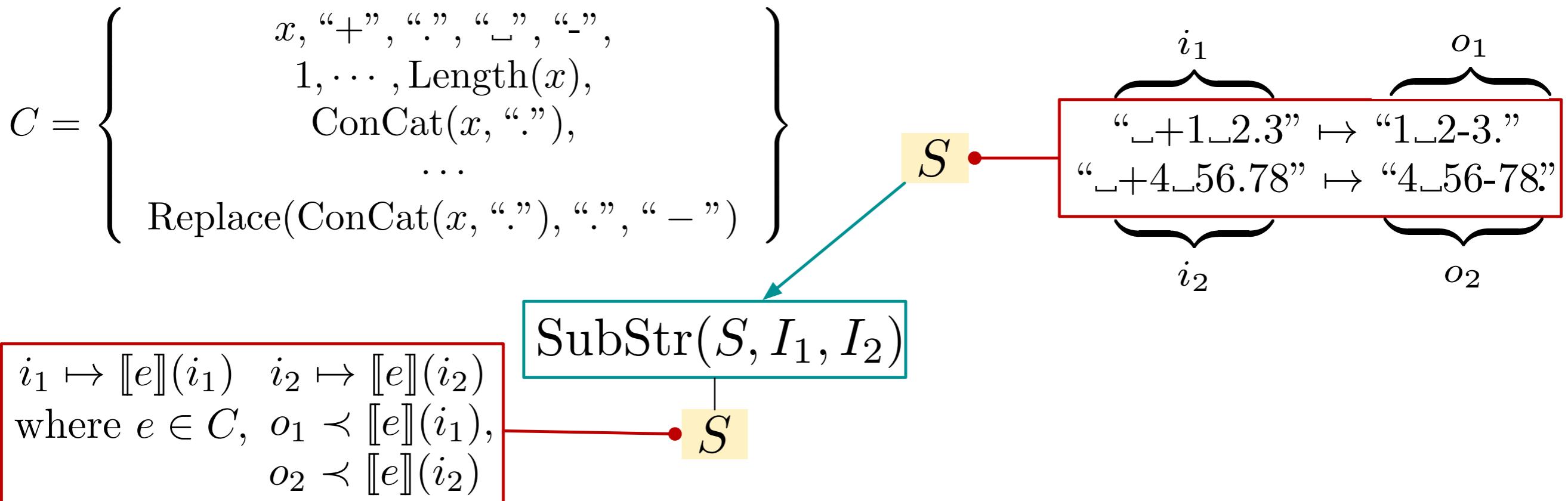
처음 부품 크기가 6으로 주어졌다고 가정
상향식 나열로 다음 부품들 생성

$$C = \left\{ \begin{array}{l} x, "+", "\cdot", "-", \\ 1, \dots, \text{Length}(x), \\ \text{ConCat}(x, ".") , \\ \dots \\ \text{Replace}(\text{ConCat}(x, "."), "\cdot", "-") \end{array} \right\}$$

가능입력들 ⊂ 부품의 출력들

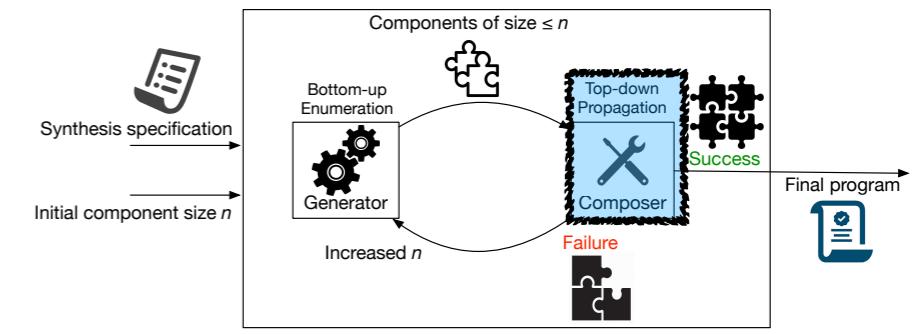


Component Pool



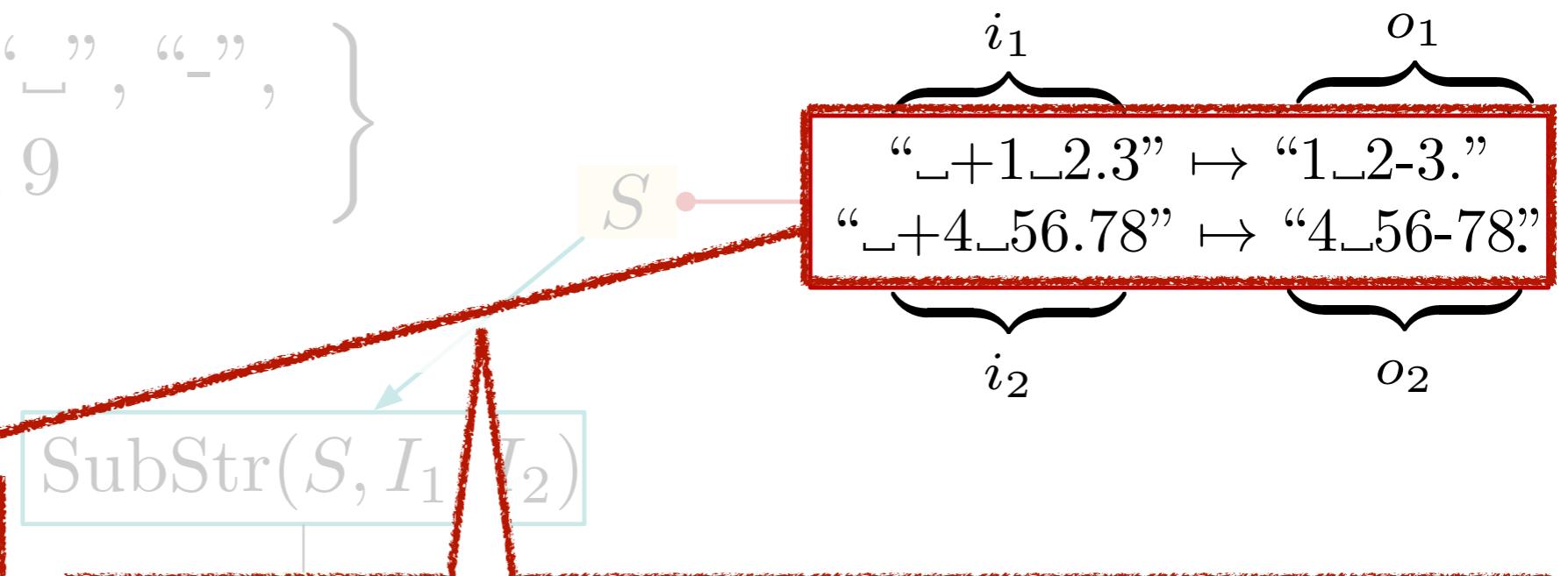
출력 문자열들을 ‘포함’하는 문자열들을 생성하는 부품 C 에서 찾아서 있으면 쓰기

가능입력들 ⊂ 부품의 출력들



Component Pool

$$C = \left\{ x, "+", "\cdot", "\lfloor", "\rfloor", "-", 1, \dots, 9 \right\}$$

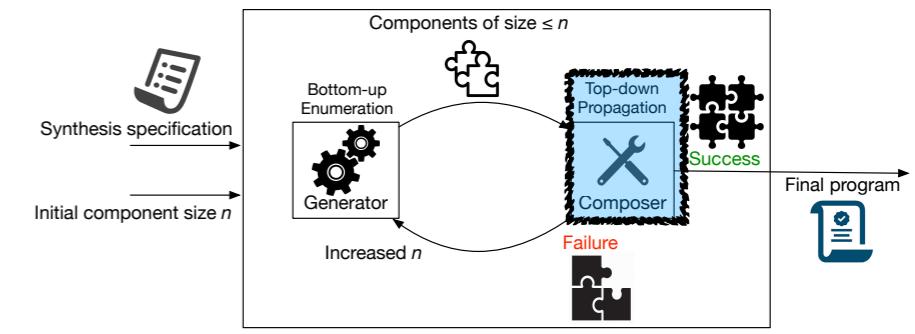


$i_1 \mapsto [e](i_1) \quad i_2 \mapsto [e](i_2)$
where $e \in C, o_1 \prec [e](i_1), o_2 \prec [e](i_2)$

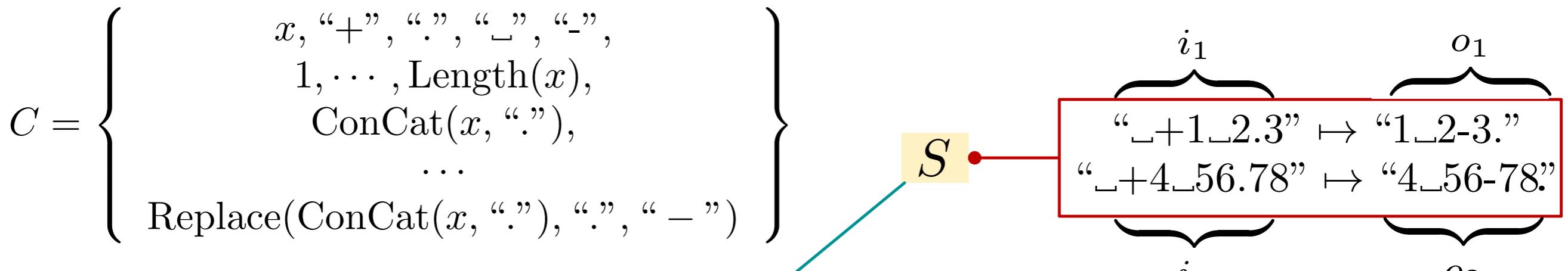
부분 문제 생성시 모든 입출력 예제 한꺼번에 고려:

- 최종 교집합 연산 필요 없음
- 입출력 예제 갯수가 늘어나도 성능 부하 적음.

가능입력들 ⊂ 부품의 출력들



Component Pool

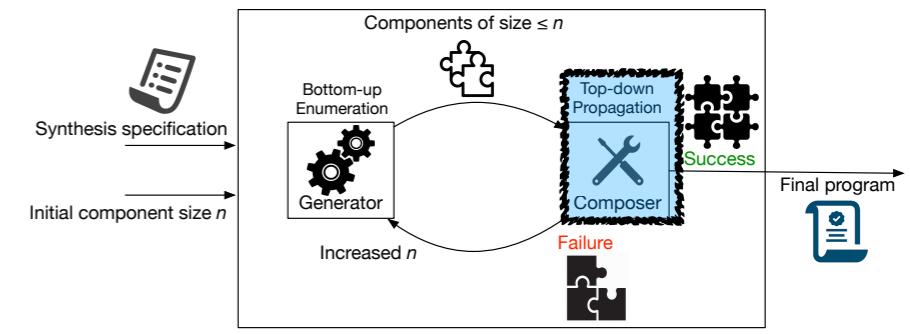


$i_1 \mapsto [e](i_1) \quad i_2 \mapsto [e](i_2)$
where $e \in C, o_1 \prec [e](i_1), o_2 \prec [e](i_2)$

$\text{SubStr}(S, I_1, I_2)$

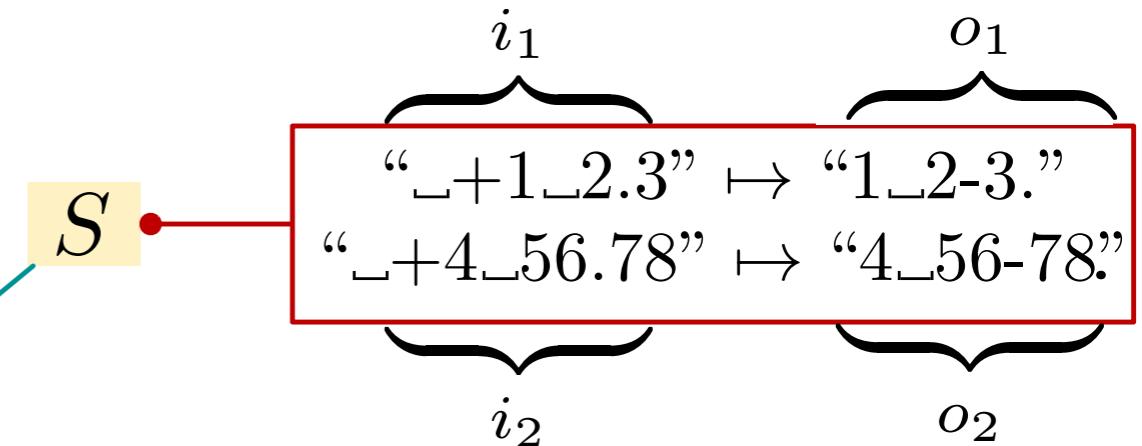
$\text{Replace}(\text{ConCat}(x, "."), ".", "-")$
가 바로 그런 부품! (입력 1,2에 대해 각각
“ +1 2-3 . ”, “ +4 56-78 . ” 출력)

가능입력들 ⊂ 부품의 출력들



Component Pool

$$C = \left\{ \begin{array}{l} x, "+", ".", "-", \\ 1, \dots, \text{Length}(x), \\ \text{ConCat}(x, ".") , \\ \dots \\ \text{Replace}(\text{ConCat}(x, "."), ".", "-") \end{array} \right\}$$



$i_1 \mapsto [e](i_1)$ $i_2 \mapsto [e](i_2)$
where $e \in C$, $o_1 \leftarrow [e](i_1)$,
 $o_2 \leftarrow [e](i_2)$

SubStr(S, I_1, I_2)

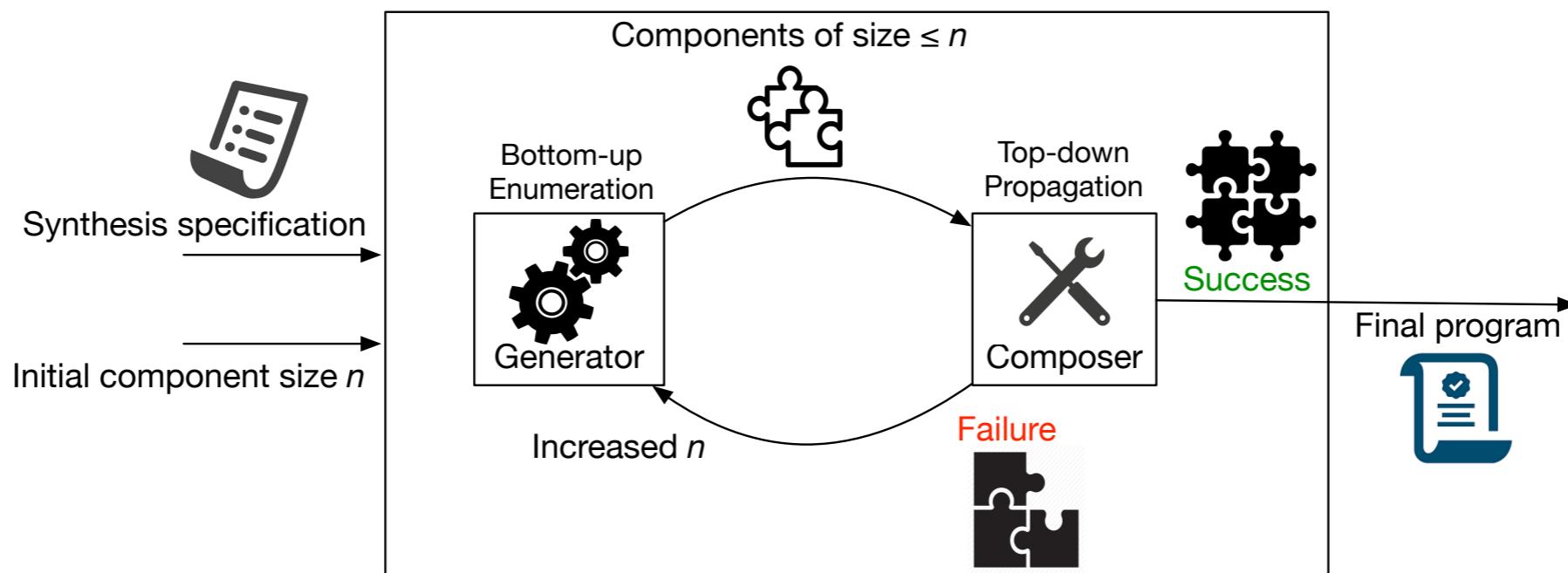
S

첫번째 인자가 정해진 후, 나머지 인자에 적절한 부품
들 또한 C에서 찾음.

(I_1 은 2, I_2 은 $\text{Length}(x) - 1$)

탐색 완전성

- 자유로운 인자 자리에 부품값만 들어갈 수 있음에도 불구하고
- 해가 탐색 공간에 존재할 경우, 이를 항상 찾아냄



- 부품 크기를 계속 늘리다 보면 솔루션이 포함될 것이기 때문

Evaluation Setup

- Benchmarks: **1,536** SyGuS problems
 - **1,167** from the SyGuS annual competitions + **369** from optimization tasks for homomorphic evaluation [Lee et al. PLDI'20]
- Comparison to three baselines (Timeout 1 hour):
 - EUSolver: winner of 2016 SyGuS competition
 - CVC4: winner of 2017 - 2019 SyGuS competition
 - Euphony [Lee et al. PLDI'18]: statistical model-guided synthesizer

실험 벤치마크

- **1,536** SyGuS 문제들
 - 경진대회 **1,167** 문제 + 프로그램 합성 기반 동형암호 프로그램 최적화기에서 **369** 문제 [Lee et al. PLDI'20]
- 경쟁 대상들
 - EUSolver: 2016 SyGuS 경진대회 우승
 - CVC4: 2017년 아래 계속 SyGuS 경진대회 우승
 - Euphony [Lee et al. PLDI'18]: 통계모델 기반 합성기

실험 벤치마크

A	B
1 Email	Column 2
2 Nancy.FreeHafer@fourthcoffee.com	nancy freehafer
3 Andrew.Cencici@northwindtraders.com	andrew cencici
4 Jan.Kotas@litwareinc.com	jan kotas
5 Mariya.Sergienko@gradicdesigninstitute.com	mariya sergienko
6 Steven.Thorpe@northwindtraders.com	steven thorpe
7 Michael.Neipper@northwindtraders.com	michael neipper
8 Robert.Zare@northwindtraders.com	robert zare
9 Laura.Giussani@adventure-works.com	laura giussani
10 Anne.HL@northwindtraders.com	anne hl
11 Alexander.David@contoso.com	alexander david
12 Kim.Shane@northwindtraders.com	kim shane
13 Manish.Chopra@northwindtraders.com	manish chopra
14 Gerwald.Oberleitner@northwindtraders.com	gerwald oberleitner
15 Amr.Zaki@northwindtraders.com	amr zaki
16 Yvonne.McKay@northwindtraders.com	yvonne mckay
17 Amanda.Pinto@northwindtraders.com	amanda pinto

complement

$\sim \underline{01010001110101110000000000001111}$
 $\underline{10101110001010001111111111110000}$

bitwise and

$01010001110101110000000000001111$
 $\& \underline{0011000101101100011000101101110}$
 $\underline{00010001010001100000000000001110}$

bitwise or

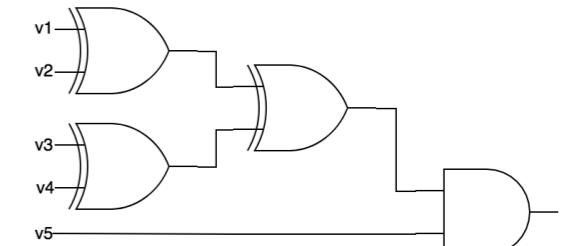
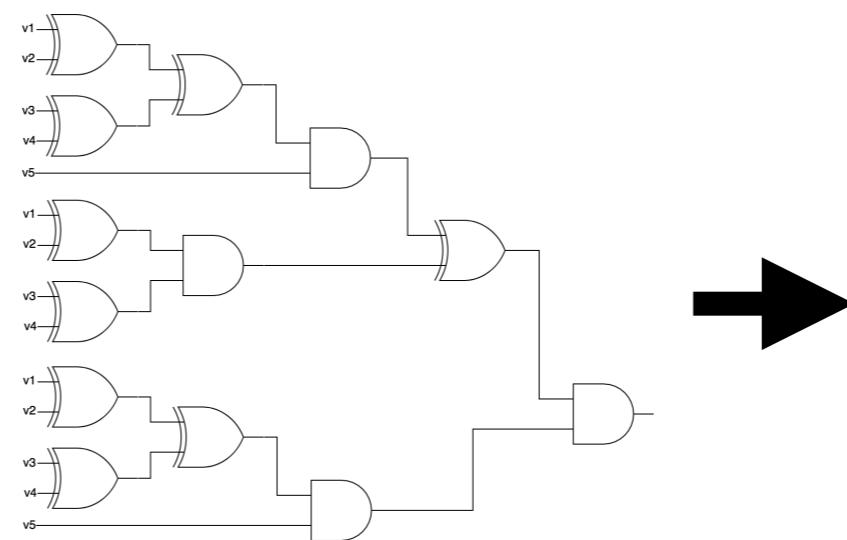
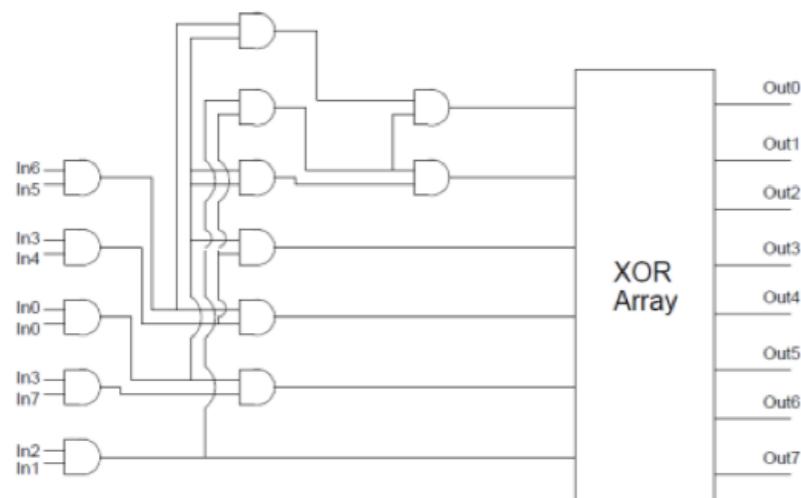
$01010001110101110000000000001111$
 $| \underline{0011000101101100011000101101110}$
 $\underline{0111000111111110011000101101111}$

bitwise xor

$01010001110101110000000000001111$
 $\wedge \underline{0011000101101100011000101101110}$
 $\underline{0110000101110010011000101100001}$

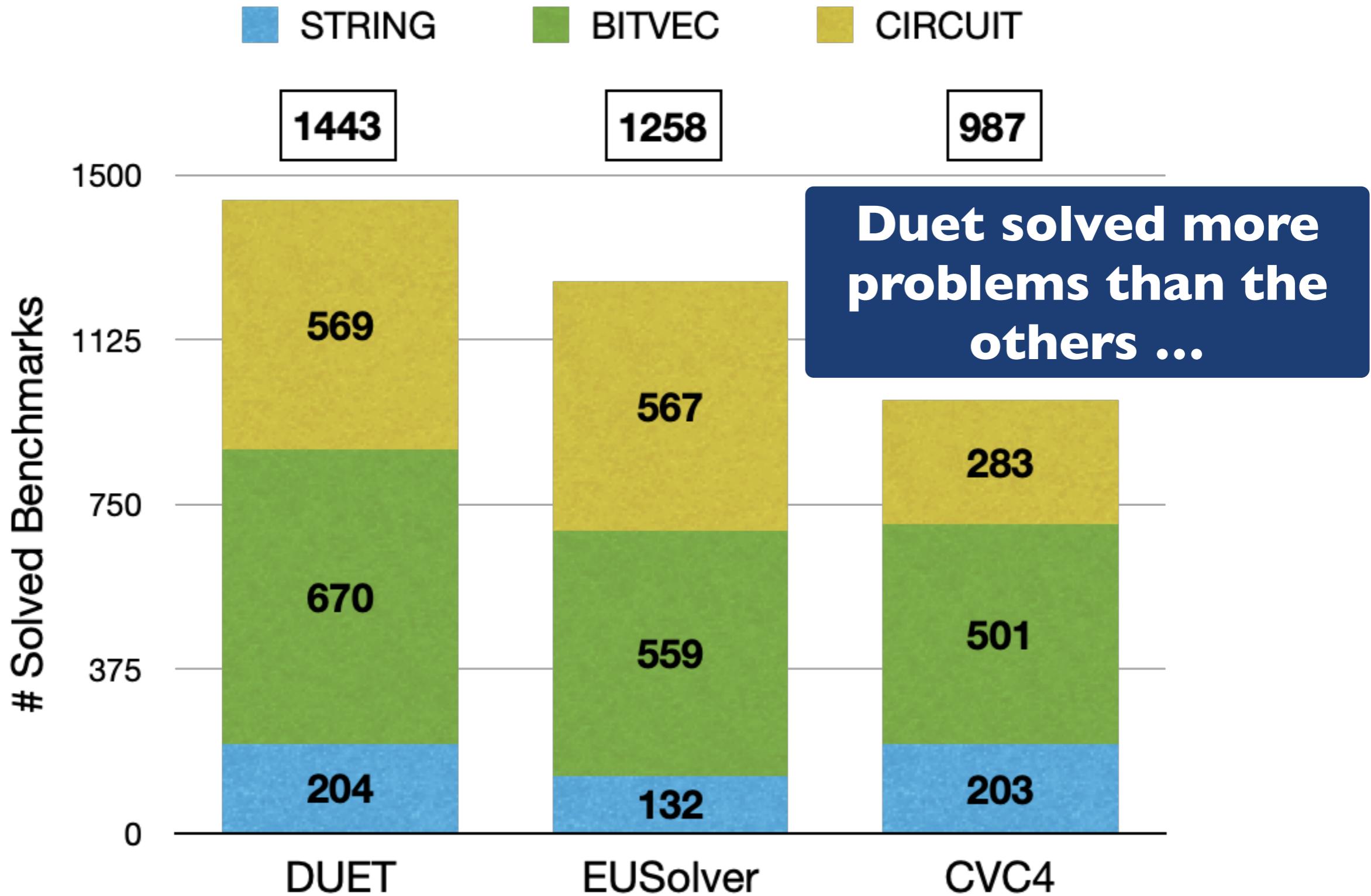
STRING: End-user programming
205 problems

BITVEC: Efficient low-level algorithm
750 problems

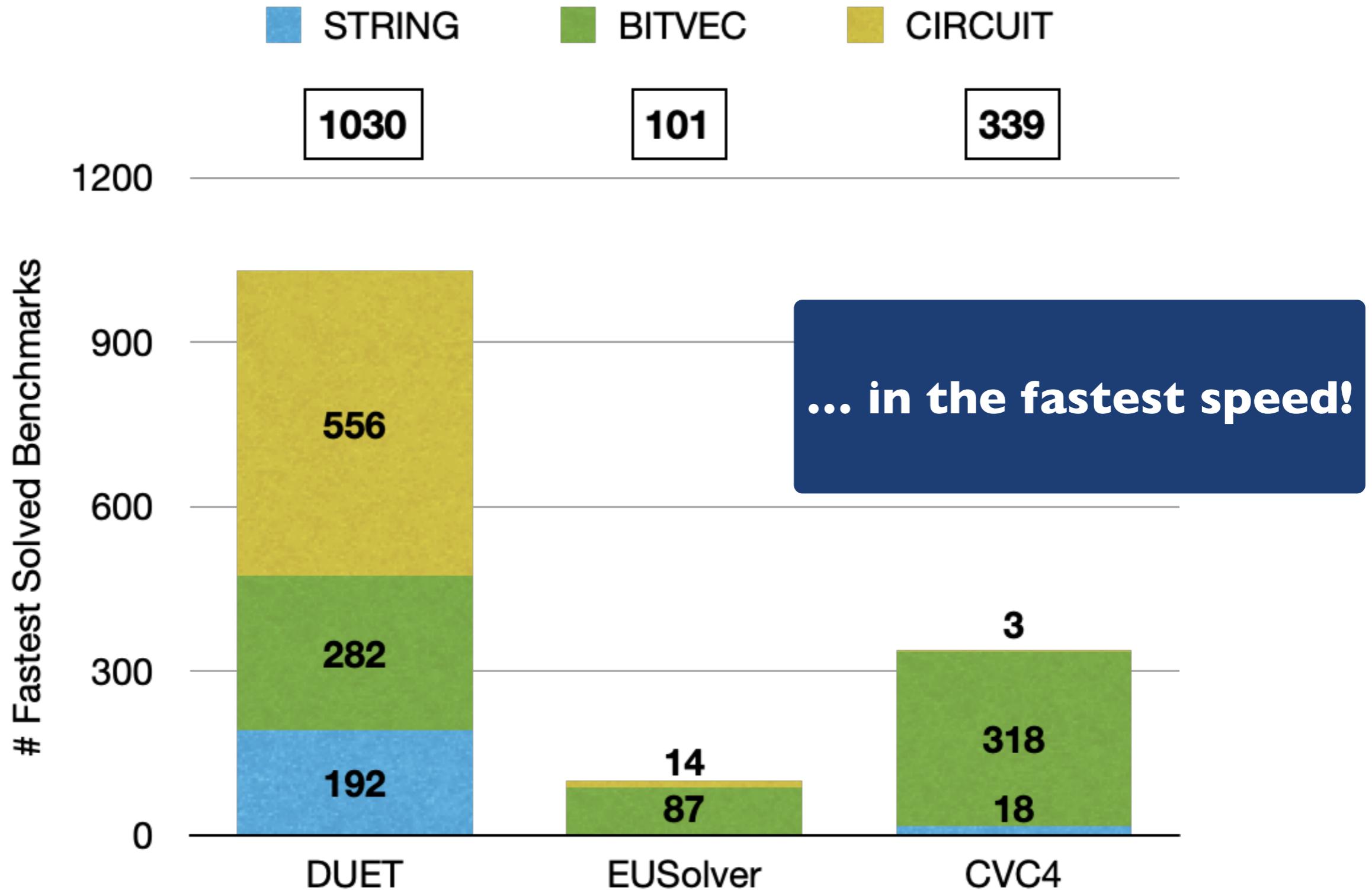


CIRCUIT: Attack-resilient crypto circuits + Optimized homomorphic evaluation circuits
581 problems

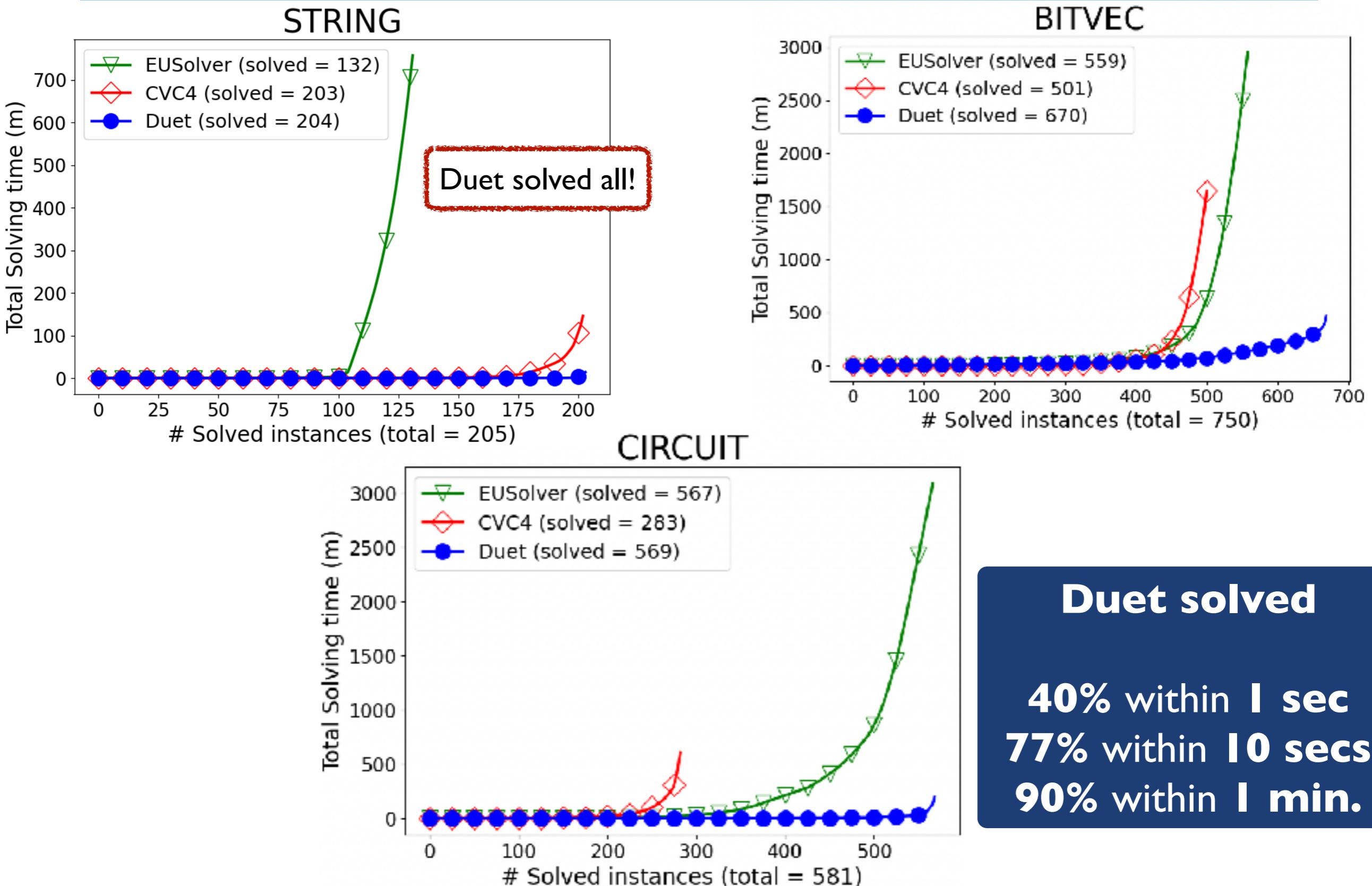
실험 결과



실험 결과

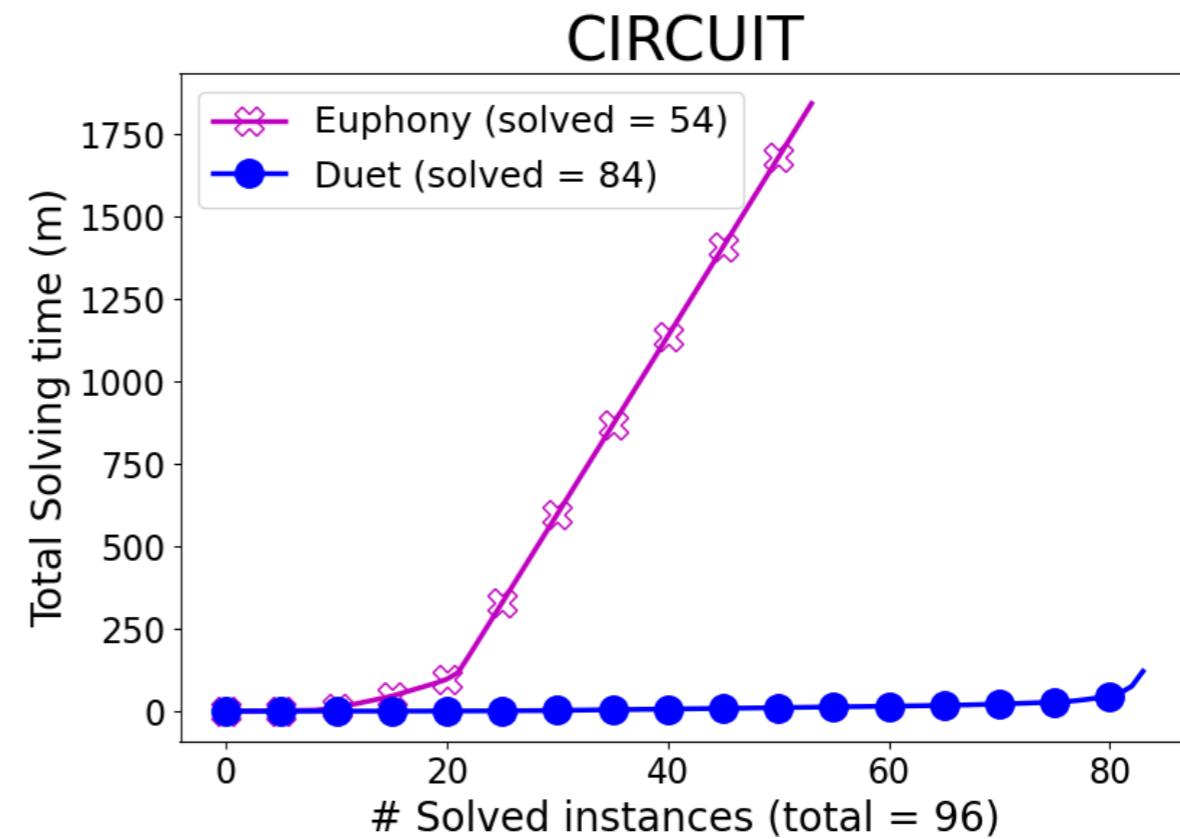
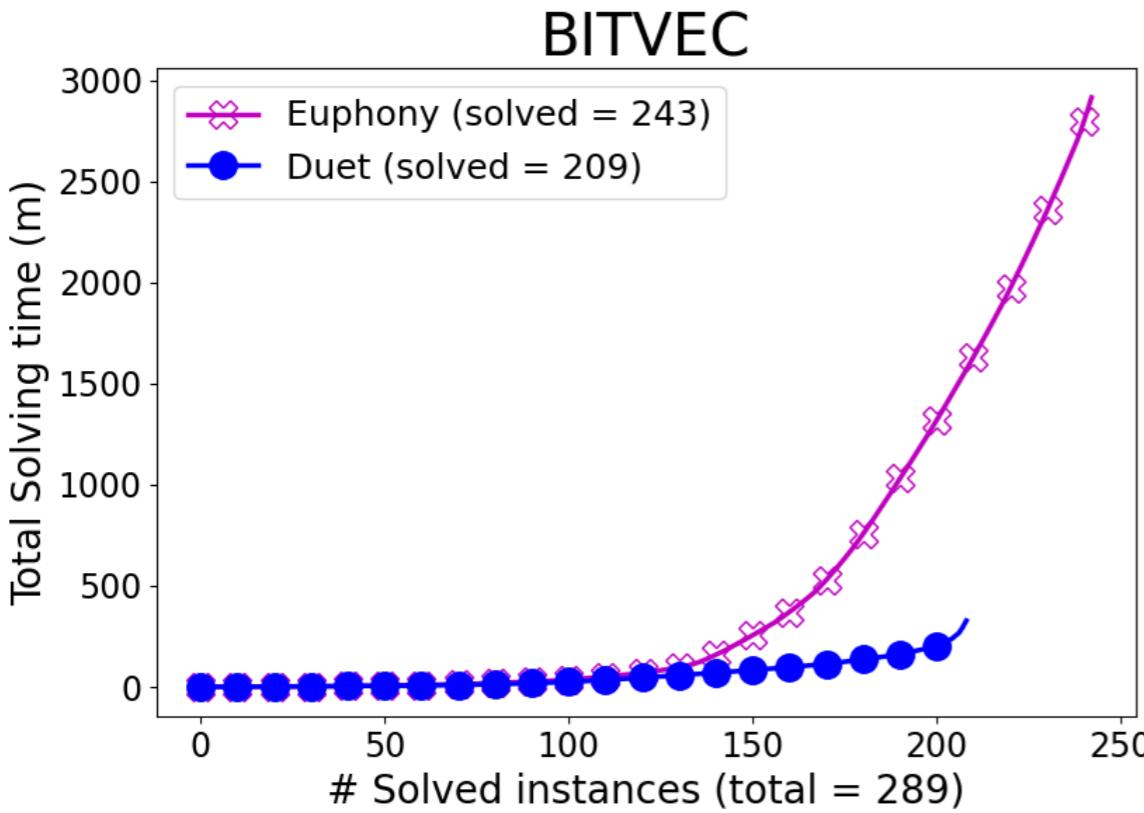
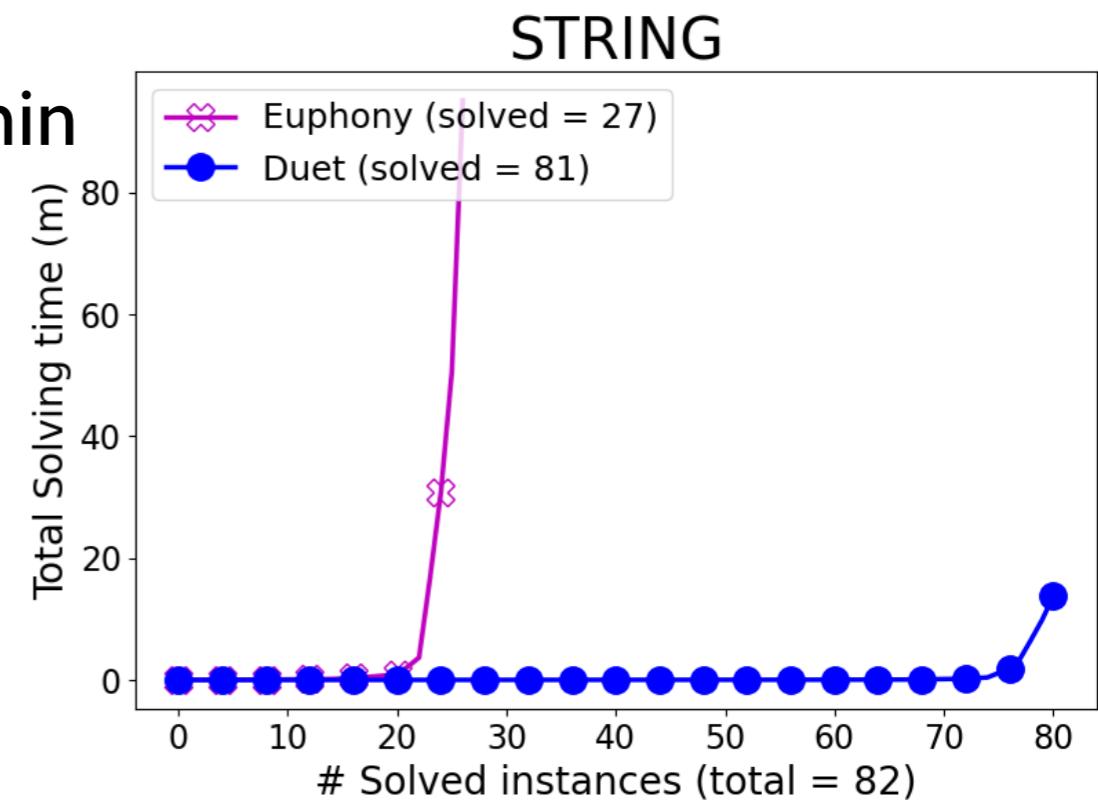


실험 결과



실험 결과 (vs. Euphony)

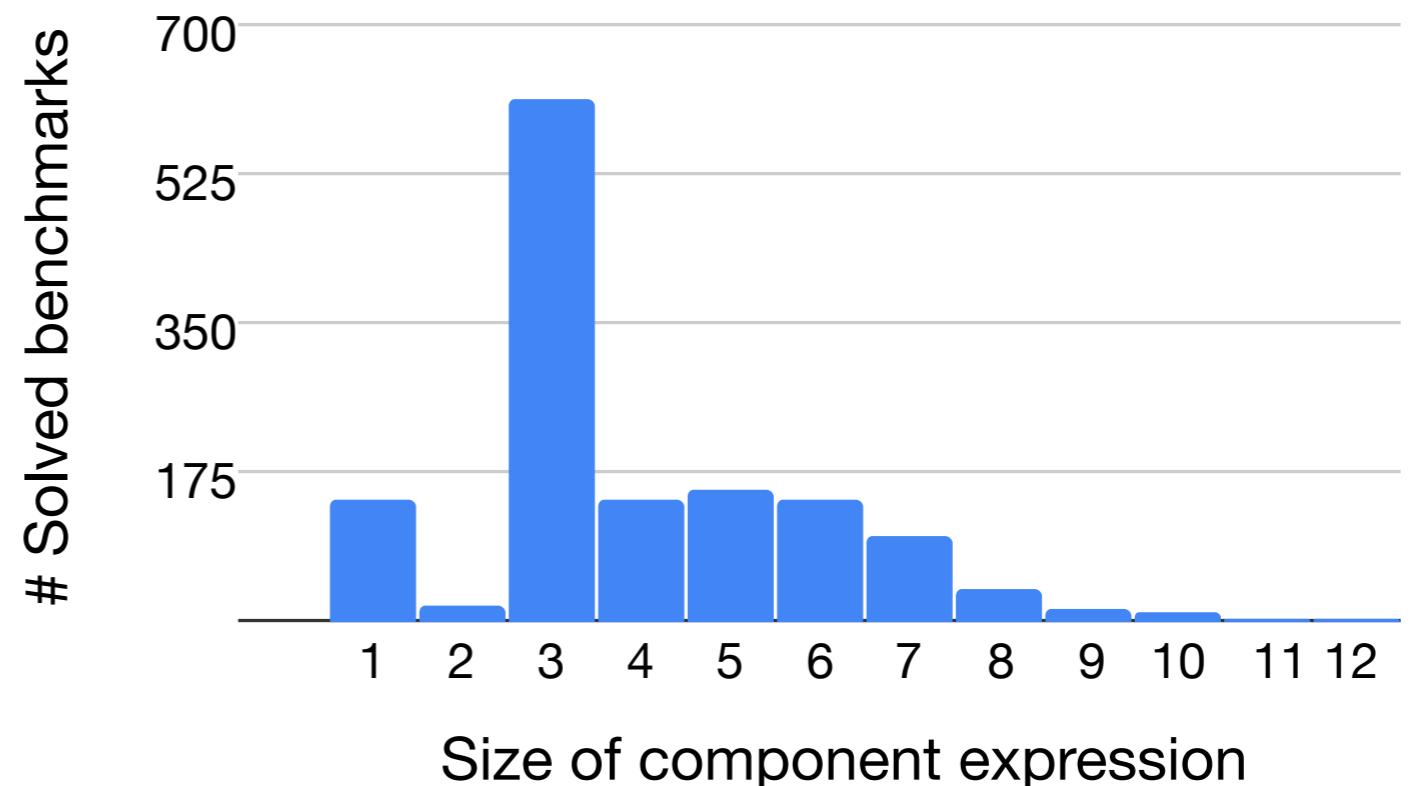
- Training: 1,069 solved EU Solver in 10 min
- Testing: 467
- # solved: Duet: 374, Euphony 324



부품 크기 분석

- 해를 찾는데 필요한 부품 크기

- STRING: 1 ~ 6 (avg: 1.8)
- BITVEC: 3 ~ 5 (avg: 3.5)
- CIRCUIT: 3 ~ 12 (avg: 5.9)



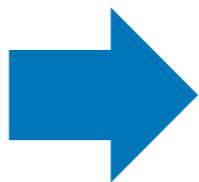
- 92% 의 문제가 크기 7 이하 부품으로 해결 가능.

내용

- 프로그램 합성 개괄
- 귀납 합성 전략
- 연역 합성 전략

입출력 예제는 종종 (제약조건으로) 너무 느슨

```
1 # gcc -O3
2
3 movq rsi, r9
4 movl ecx, ecx
5 shrq 32, rsi
6 andl 0xffffffff, r9d
7 movq rcx, rax
8 movl edx, edx
9 imulq r9, rax
10 imulq rdx, r9
11 imulq rsi, rdx
12 imulq rsi, rcx
13 addq rdx, rax
14 jae .L0
15 movabsq 0x100000000, rdx
16 addq rdx, rcx
17 .L0:
18 movq rax, rsi
19 movq rax, rdx
20 shrq 32, rsi
21 salq 32, rdx
22 addq rsi, rcx
23 addq r9, rdx
24 adcq 0, rcx
25 addq r8, rdx
26 adcq 0, rcx
27 addq rdi, rdx
28 adcq 0, rcx
29 movq rcx, r8
30 movq rdx, rdi
```



```
1 # STOKE
2
3 sh1q 32, rcx
4 movl edx, edx
5 xorq rdx, rcx
6 movq rcx, rax
7 mulq rsi
8 addq r8, rdi
9 adcq 0, rdx
10 addq rdi, rax
11 adcq 0, rdx
12 movq rdx, r8
13 movq rax, rdi
```

Montgomery multiplication kernel from the OpenSSL RSA library. Compilations shown for
gcc -O3 (left) and a stochastic optimizer (right).

입출력 예제는 종종 (제약조건으로) 너무 느슨

- 원본 프로그램 f 과 의미가 같으면서 최적화된 프로그램 f' 을 원할 경우,
 - 귀납합성으로 간단히 생각 할 수 있는 방법: 랜덤 입력들을 생성하여 원본 프로그램에 대해서 돌린 후 출력값 얻은 후, 입출력 예제들을 사용
 - 문제: 올바름 조건 100% 보장 안됨
- 제대로 된 올바름 조건: $\forall x. f(x) = f'(x)$

입출력 예제는 종종 (제약조건으로) 너무 느슨

- 올바름 조건이 논리식으로 사전/사후(pre/post) 조건으로 기술된 경우에
도 기존 귀납 합성 전략 먹통.

```
int W = 32;

void swap(ref bit[W] x, ref bit[W] y) {
    if(??){ x = x ^ y; }else{ y = x ^ y; }
    if(??){ x = x ^ y; }else{ y = x ^ y; }
    if(??){ x = x ^ y; }else{ y = x ^ y; }
}

harness void main(bit[W] x, bit[W] y) {
    bit[W] tx = x; bit[W] ty = y;
    swap(x, y);
    assert x==ty && y == tx;
}
```

- Program sketch: <https://people.csail.mit.edu/asolar/sketch-1.7.6.tar.gz>

프로그램 합성 문제의 일반적 기술

$$\exists f. \forall x. P(f, x)$$



There exists a function f such that for all x , property P holds

- 모든 가능한 입력에 대해 ($\forall x$) 어떤 성질을 만족하는 함수 찾기
- 기존 SMT solvers로 풀 수 있는가? $\Rightarrow \text{NO!}$ 2차 변수 f 와 \forall (universal quantifier) 이 붙은 1차 변수가 있는 식 잘 못 다룸.

프로그램 합성 문제의 일반적 기술

$$\exists f. \forall x. P(f, x)$$



There exists a function f such that for all x , property P holds

- 그냥 가능한 f 를 나열하면서 찾아보기? \Rightarrow 2차 변수가 사라져서 SMT solver가 다룰 수 있으나, SMT solver 호출(계산비용 비쌈) 횟수가 너무 많아짐.

Counterexample guided Inductive Synthesis (CEGIS)

- 2008년 Armando Solar-Lezama 가 제안한 알고리즘.
- 기존의 귀납 합성 전략을 일반적인 제약조건을 다루는 연역 합성에 사용할 수 있게 해 줌.
- 합성 대상 함수의 후보를 제안하는 탐색파트(generator)와, 제안된 후보가 정말 맞는지 체크하는 검증파트(verifier)로 나뉨.

Example

- Goal: a function f computes the maximum of integers x and y

-  Specification

if C then S else S

Syntactic specification:

$$S \rightarrow x \mid y \mid n \mid \text{ite}(C, S, S)$$

$$C \rightarrow S \leq S \mid \neg C \mid C \& C$$

Semantic specification:

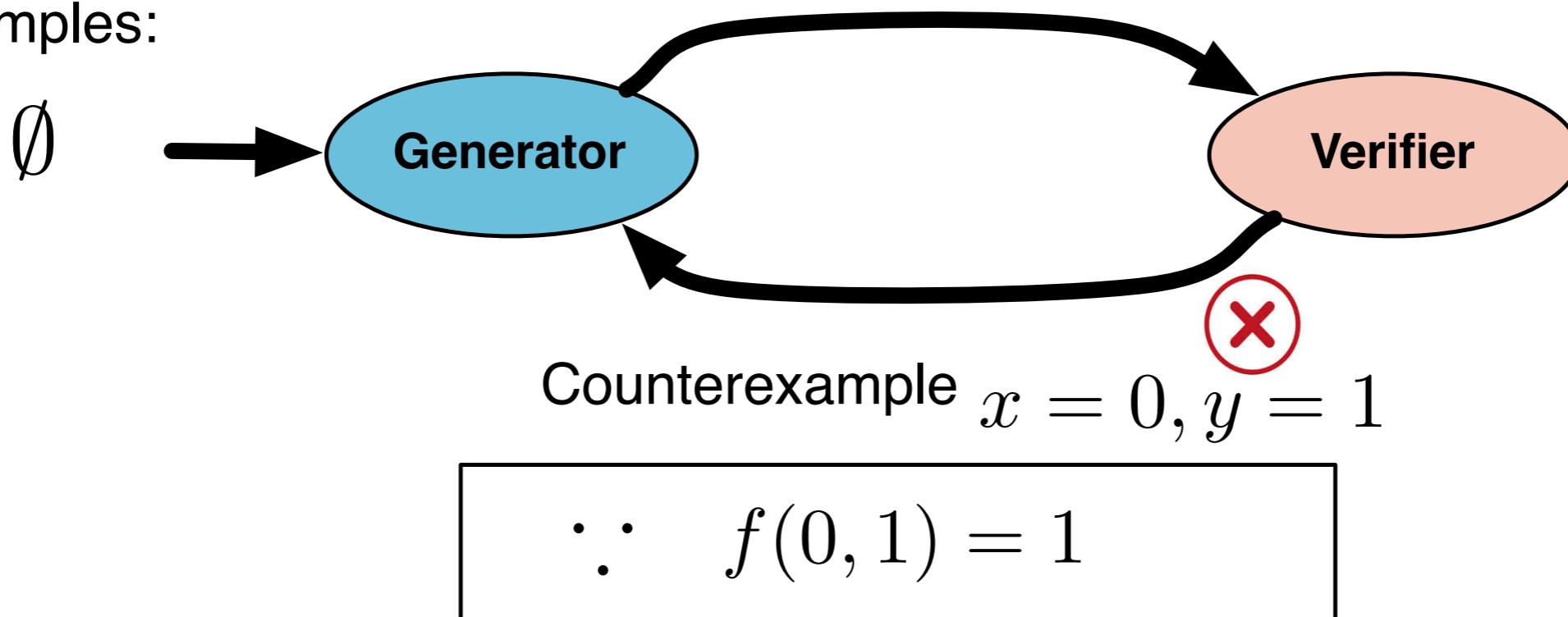
$$\forall x, y. f(x, y) \geq x \wedge f(x, y) \geq y \wedge (f(x, y) = x \vee f(x, y) = y)$$

-  Solution $f(x, y) = \text{ite}(x > y, x, y)$

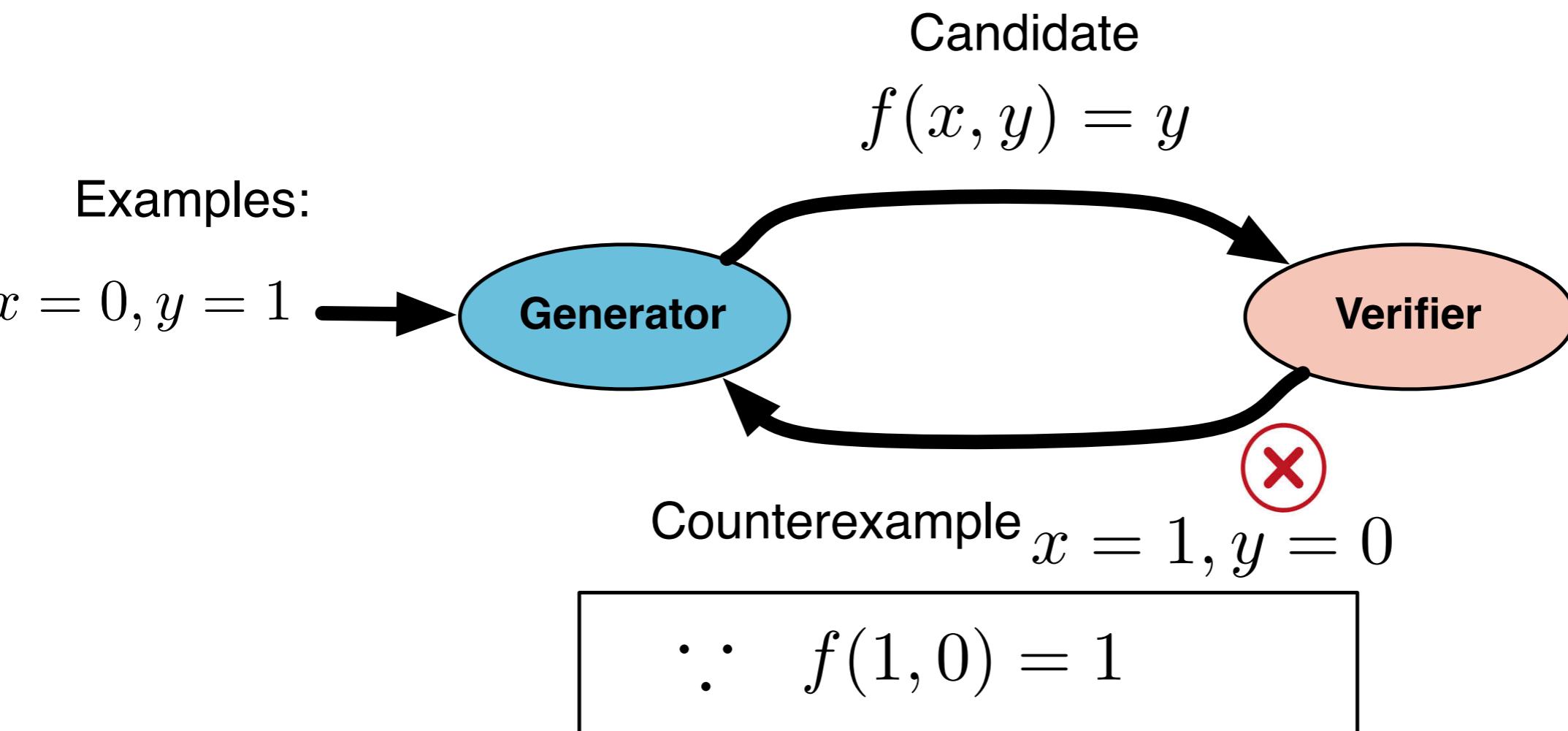
CEGIS + 상향식 나열

Candidate
 $f(x, y) = x$

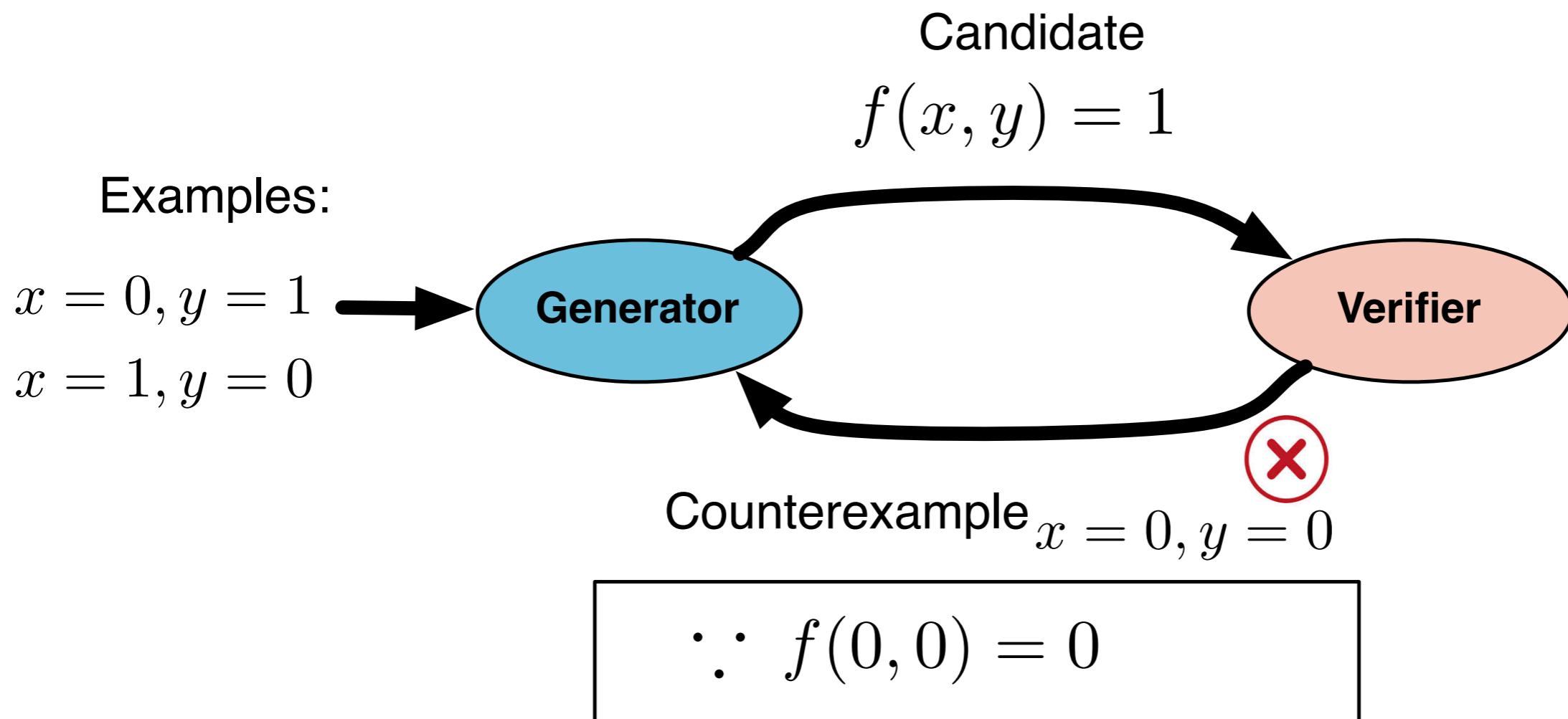
Examples:



CEGIS + 상향식 나열



CEGIS + 상향식 나열



CEGIS + 상향식 나열

Candidate

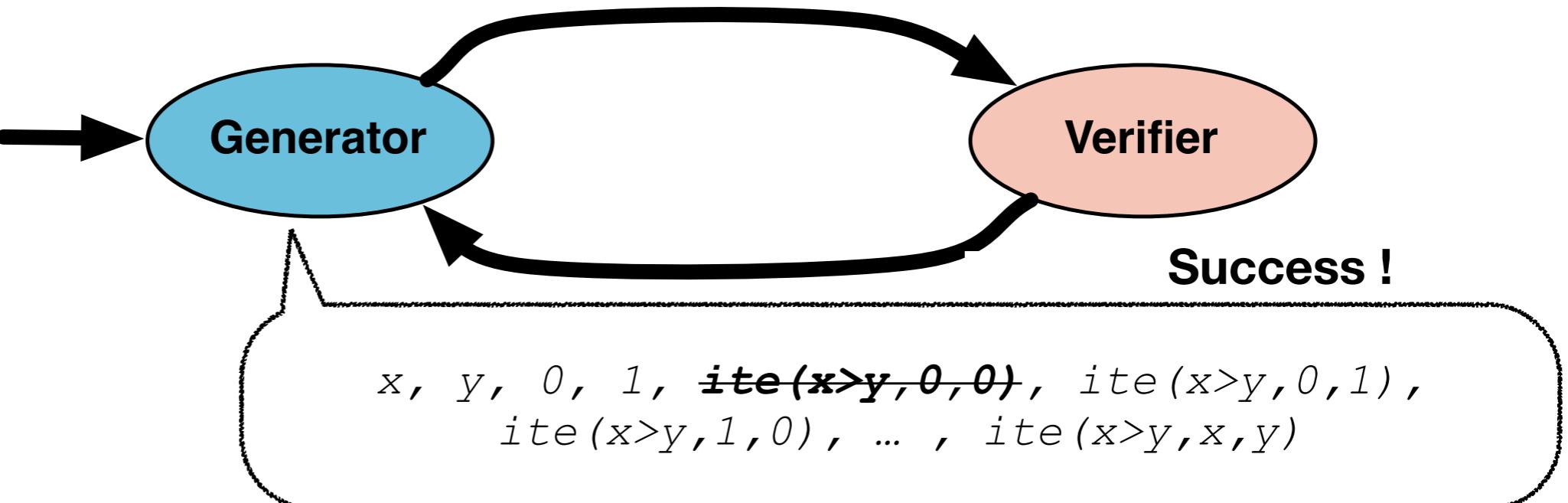
$$f(x, y) = \text{ite}(x > y, x, y)$$

Examples:

$$x = 0, y = 1$$

$$x = 1, y = 0$$

$$x = 0, y = 0$$



최적화: $\text{ite}(x > y, 0, 0)$ 는 무시됨.

$\because \text{ite}(x > y, 0, 0) \approx 0$ wrt examples

CEGIS가 가져다주는 이점

- 검증기와 탐색기는 각기 따로따로 설계되어도 됨 (탐색기는 임의의 귀납 합성전략, 검증기는 임의의 검증 방식 ...)
- 검증기는 2차변수 없이 1차변수만 다루어도 됨.
- 탐색기는 2차변수를 다루지만 universal quantifier 는 다룰 필요가 없음.

정리

- 프로그램 합성 문제 = 제약조건 + 탐색공간
- 가능 탐색 전략은 제약조건, 탐색공간에 따라 결정됨.
- 제약조건이 입출력예제인 경우 귀납합성전략 사용
- 제약조건이 논리식인 경우 귀납합성전략 + CEGIS 사용