



CoqPilot

a plugin for LLM-based generation of proofs

Andrei Kozyrev

Gleb Solovev

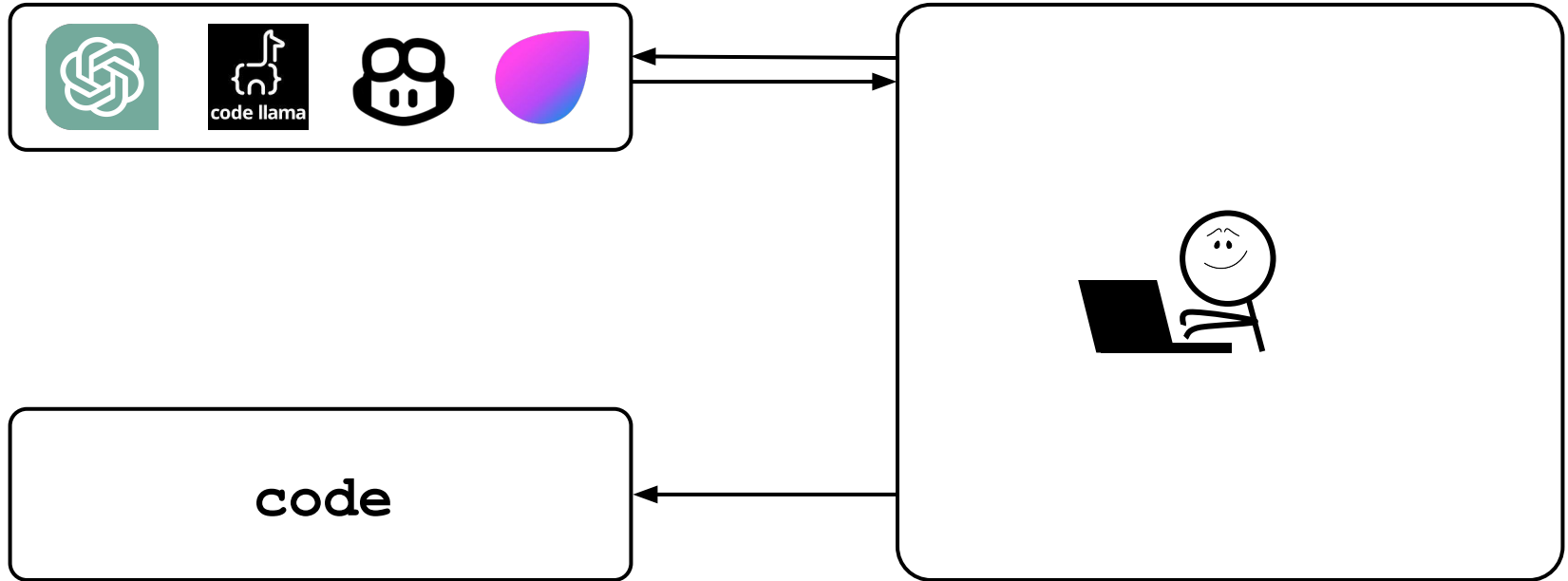
Nikita Khramov

Anton Podkopaev

Programming Languages and Program Analysis Lab (PLAN), JetBrains Research

September, 2024

LLMs are used more and more for code generation



Do Users Write More Insecure Code with AI Assistants?

Neil Perry *
Stanford University

Megha Srivastava *
Stanford University

Deepak Kumar
Stanford University

Dan Boneh
Stanford University

Do Users Write More Insecure Code with AI Assistants?

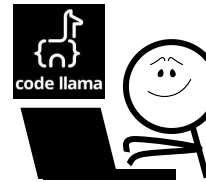
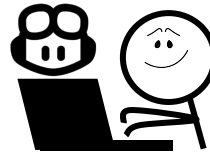
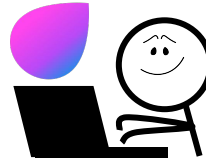
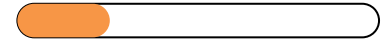
Neil Perry *
Stanford University

Megha Srivastava *
Stanford University

Deepak Kumar
Stanford University

Dan Boneh
Stanford University

introduce security *vulnerabilities*



Do Users Write More Insecure Code with AI Assistants?

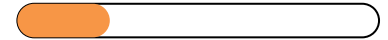
Neil Perry *
Stanford University

Megha Srivastava *
Stanford University

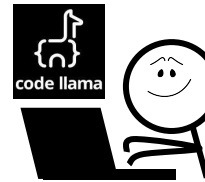
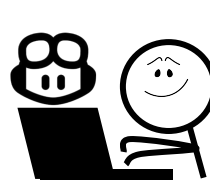
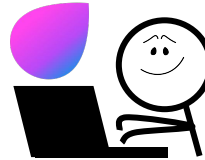
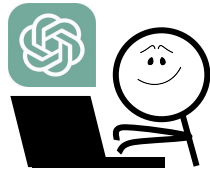
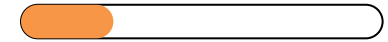
Deepak Kumar
Stanford University

Dan Boneh
Stanford University

introduce security *vulnerabilities*



rate their *insecure* code as *secure*





Generate a list sorting function in Python
without using any library functions



```
def custom_sort(arr):  
    n = len(arr)  
    for i in range(n):  
        # Last i elements are already sorted  
        for j in range(0, n - i - 1):  
            # Traverse the array from 0 to n-i-1  
            # Swap if the element found is greater  
            # than the next element  
            if arr[j] > arr[j + 1]:  
                arr[j], arr[j + 1] = arr[j + 1], arr[j]  
    return arr
```



Generate a list sorting function in Python
without using any library functions



```
def custom_sort(arr):
```

Not only familiar code is needed

```
    # Last i elements are already sorted
    for j in range(0, n - i - 1):
        # Traverse the array from 0 to n-i-1
        # Swap if the element found is greater
        # than the next element
        if arr[j] > arr[j + 1]:
            arr[j], arr[j + 1] = arr[j + 1], arr[j]
    return arr
```



Generate a list sorting function in Python
without using any library functions



```
def custom_sort(arr):
```

Not only **familiar** code is needed

Testing may **not** be **sufficient** (eg, concurrency)

```
    # Last i elements are already sorted
    for j in range(0, n - i - 1):
        # Traverse the array from 0 to n-i-1
        # Swap if the element found is greater
        # than the next element
        if arr[j] > arr[j + 1]:
            arr[j], arr[j + 1] = arr[j + 1], arr[j]
    return arr
```




Generate a list sorting function in Python
without using any library functions



```
def custom_sort(arr):
```

Not only **familiar** code is needed

Testing may **not** be **sufficient** (eg, concurrency)

Plain **English** is **hard** to debug and **imprecise**

```
# Traverse the array from 0 to n-1-1
# Swap if the element found is greater
# than the next element
if arr[j] > arr[j + 1]:
    arr[j], arr[j + 1] = arr[j + 1], arr[j]
return arr
```

CoqPilot



```

src > s.v
89 Definition p1 {l} (x : {l' | Permutation l l' & is_so
90   destruct x as [l']. exact l.
91 Defined.
92
93 Definition sort l : {l' | Permutation l l' & is_sorte
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHl as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105 Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108
109
110
111
112
113
114
115
116
117

```



Proof

Main 1 Shelved 0 Given up 1

Goal 1

a : nat
 l, l' : list nat
 p : Permutation l l'
 i : is_sorted l'

(1/1) _____
 is_sorted l'

Messages

src > s.v

```

43 Lemma is_inserted_perm a l l' (INS : is_inserted a l
44 (* Hint: perm_swap *))
45 Proof.
46   generalize dependent l'.
47   generalize dependent a.
48   induction l; ins; inv INS.
49   apply IHL in INS0.
50   etransitivity.
51   { by apply perm_swap. }
52   by constructor.
53 Qed.
54
55 Lemma insert_sorted a l (SORT : is_sorted l) :
56   {l' | is_inserted a l l' & is_sorted l'}.
57 (* Hint: le_gt_dec *)
58 Proof.
59   induction l; eauto with myconstr.
60   edestruct IHL as [l'].
61   { clear -SORT. inv SORT. auto with myconstr. }
62   destruct (le_gt_dec a a0).
63   { exists (a::a0::l); auto with myconstr.
64     apply sorted_cons; auto.
65     eapply smallest_head; eauto.
66     inv SORT. auto with myconstr. }
67   exists (a0::l'); auto. constructor; auto.
68
69   clear -SORT i i0 g.
70   induction i; auto.

```



Proof

Main 1 Shelved 0 Given up 1

Goal 1

```

a : nat
l, l' : list nat
p : Permutation l l'
i : is_sorted l'

```

(1/1) is_sorted l'

Messages

```

43 Lemma is_inserted_perm a l l' (INS : is_inserted a l
44 (* Hint: perm_swap *))
45 Proof.
46   generalize dependent l'.
47   generalize dependent a.
48   induction l; ins; inv INS.
49   apply IHL in INS0.
50   etransitivity.
51   { by apply perm_swap. }
52   | by constructor.
53 Qed.
54
55 Lemma insert_sorted a l (SORT : is_sorted l) :
56   {l' | is_inserted a l l' & is_sorted l'}.
57 (* Hint: le_gt_dec *)
58 Proof.
59   induction l; eauto with myconstr.
60   edestruct IHL as [l'].
61   { clear -SORT. inv SORT. auto with myconstr. }
62   destruct (le_gt_dec a a0).
63   { exists (a::a0::l); auto with myconstr.
64     apply sorted_cons; auto.
65     eapply smallest_head; eauto.
66     inv SORT. auto with myconstr. }
67   exists (a0::l'); auto. constructor; auto.
68
69   clear -SORT i i0 g.
70   induction i; auto.

```



Proof
Main 1 Shelved 0 Given up 1

Goal 1

```

a : nat
l, l' : list nat
p : Permutation l l'
i : is_sorted l'

```

(1/1) _____
is_sorted l'

Messages

```

43 Lemma is_inserted_perm a l l' (INS : is_inserted a l
44 (* Hint: perm_swap *)
45 Proof.
46   generalize dependent l'.
47   generalize dependent a.
48   induction l; ins; inv INS.
49   apply IHL in INS0.
50   etransitivity.
51   { by apply perm_swap. }
52   | by constructor.
53 Qed.
54
55 Lemma insert_sorted a l (SORT : is_sorted l) :
56 {l' | is_inserted a l l' & is_sorted l'}.
57 (* Hint: le_gt_dec *)
58 Proof.
59   induction l; eauto with myconstr.
60   edestruct IHL as [l'].
61   { clear -SORT. inv SORT. auto with myconstr. }
62   destruct (le_gt_dec a a0).
63   { exists (a::a0::l); auto with myconstr.
64     apply sorted_cons; auto.
65     eapply smallest_head; eauto.
66     inv SORT. auto with myconstr. }
67   exists (a0::l'); auto. constructor; auto.
68
69   clear -SORT i i0 g.
70   induction i; auto.

```



Proof

Main 1 Shelved 0 Given up 1

Goal 1

```

a : nat
l, l' : list nat
p : Permutation l l'
i : is_sorted l'

```

(1/1) _____

is_sorted l'

Messages

```
s.v M X ...
src > s.v
92
93 Definition sort l : {l' | Permutation l l' & is_sorted l'}.
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHL as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105 Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108 ✨
109
110
111
112
113
114
115
116
117
118
119
120
```

```


```

Coq Goals X

Proof

Not in proof mode

Messages

```
= let (x, _, _) := sort [3; 2; 4; 1] in x
: list nat
```

vscoq-language-server 2.1.0, coq 8.19.0 Go Live Record Story (beta)

```
s.v M X ...
src > s.v
92
93 Definition sort l : {l' | Permutation l l' & is_sorted l'}.
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHL as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105 Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108 ✨
109
110
111
112
113
114
115
116
117
118
119
120
```

```


```

Coq Goals X

Proof

Not in proof mode

Messages

```
= let (x, _, _) := sort [3; 2; 4; 1] in x
: list nat
```

vscoq-language-server 2.1.0, coq 8.19.0 Go Live Record Story (beta)


```

src > s.v
89 Definition p1 {l} (x : {l' | Permutation l l' & is_so
90   destruct x as [l']. exact l.
91 Defined.
92
93 Definition sort l : {l' | Permutation l l' & is_sorte
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHl as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105 Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108
109
110
111
112
113
114
115
116
117

```



Proof

Main 1 Shelved 0 Given up 1

Goal 1

a : nat

l, l' : list nat

p : Permutation l l'

i : is_sorted l'

(1/1)

is_sorted l'

Messages

```

src > s.v
89 Definition p1 {l} (x : {l' | Permutation l l' & is_so
90   destruct x as [l']. exact l.
91 Defined.
92
93 Definition sort l : {l' | Permutation l l' & is_sorte
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHl as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105 Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108
109
110
111
112
113
114
115
116
117

```



Proof
 Main 1 Shelved 0 Given up 1

Goal 1
 a : nat
 l, l' : list nat
 p : Permutation l l'
 i : is_sorted l'
 (1/1) _____
 is_sorted l'

Messages

Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.

Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.



...



Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.



...



Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.



...



Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.

Checker

Lemma unsort_sorted
Proof 1.
...
Defined.

Lemma unsort_sorted
Proof 2.
...
Defined.



...



Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.

Checker

Lemma unsort_sorted
Proof 1.
...
Defined.

Lemma unsort_sorted
Proof 2.
...
Defined.



...



Few-shot prompt

Lemma insert_sorted ...
Proof.
...
Defined.

...

Lemma is_inserted_perm ...
Proof.
...
Defined.

Query

Lemma unsort_sorted ...
Proof.
???
Admitted.



...



Checker

Lemma unsort_sorted
Proof 1.
...
Defined.

Lemma unsort_sorted
Proof 2.
...
Defined.

src > s.v

```

92
93 Definition sort l : {l' | Permutation l l' & is_sorte
94 Proof.
95   induction l.
96   { admit. }
97   destruct IHL as [l'].
98   edestruct (insert_sorted a l') as [l''].
99   { admit. }
100  exists l''.
101  2: { admit. }
102  transitivity (a::l').
103  { admit. }
104  admit.
105  Admitted.
106
107 Eval compute in (p1 (sort [3;2;4;1])).
108
109
110
111
112
113
114
115
116
117
118

```

```

[... Coq proof state output ...]

```

Proof

Main 1 Shelved 0 Given up 0

Goal 1

l : list nat

(1/1)

 {l' : list nat | Permutation l l' & is_sorted l'}

Messages

Proof repair

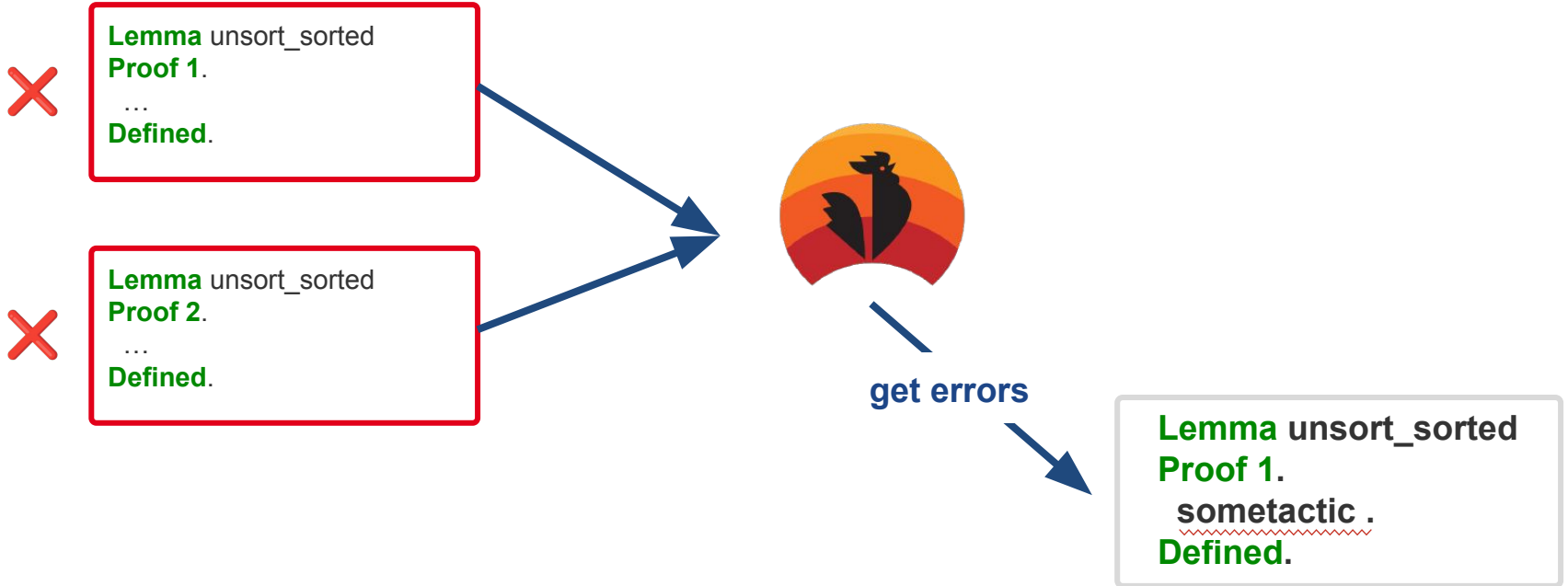


Lemma unsort_sorted
Proof 1.
...
Defined.

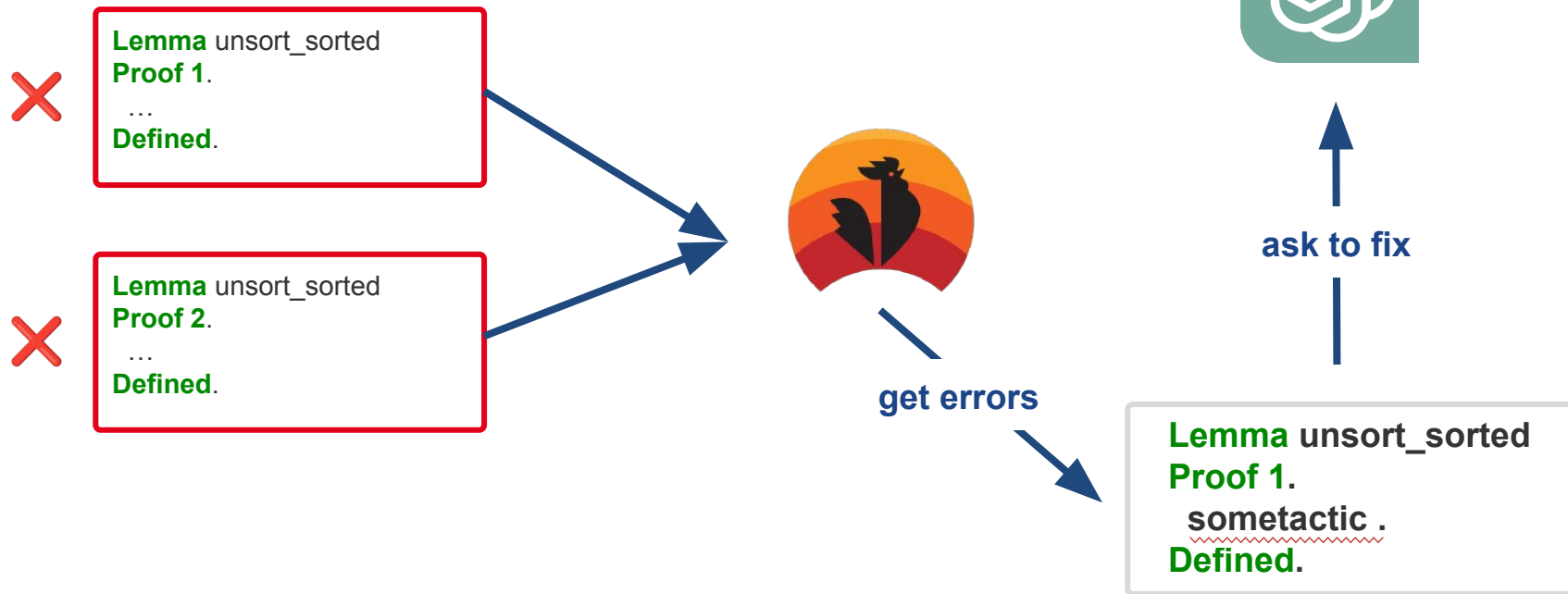


Lemma unsort_sorted
Proof 2.
...
Defined.

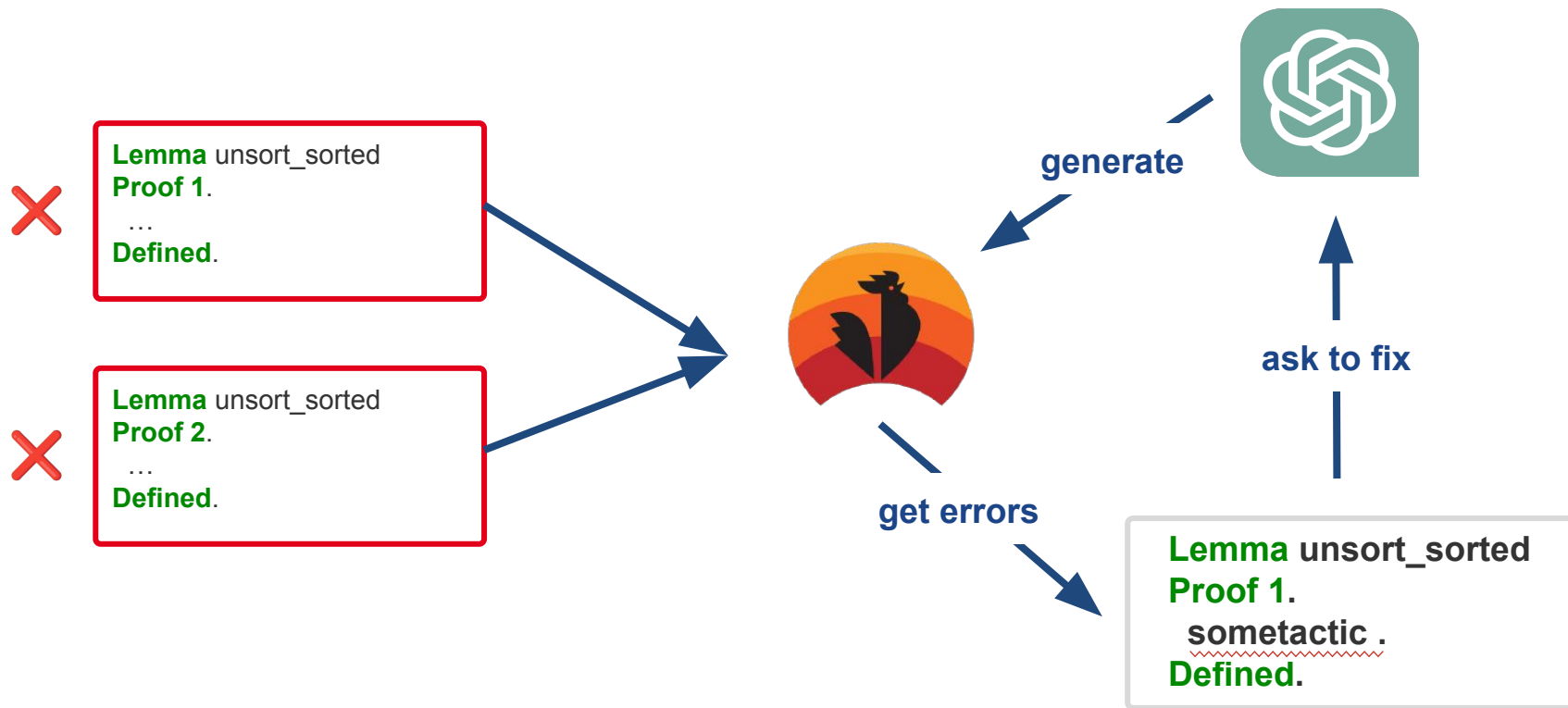
Proof repair



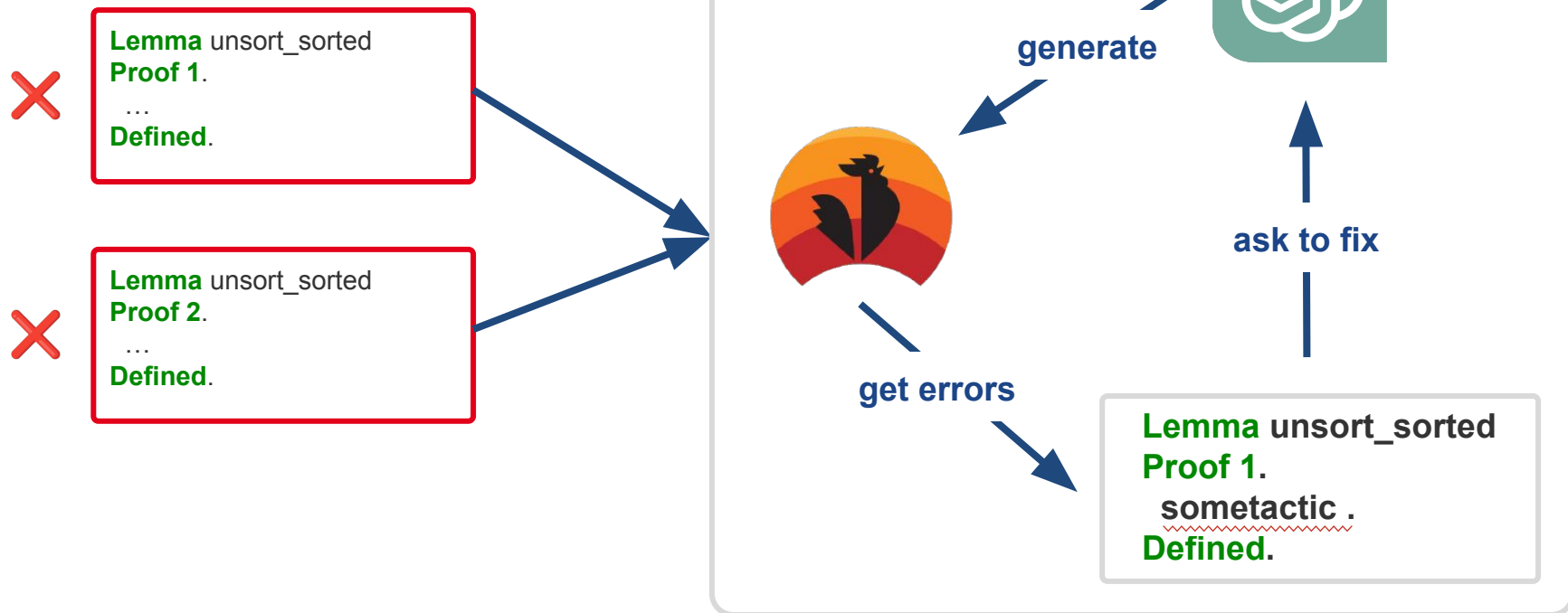
Proof repair



Proof repair



Proof repair



Research Questions

Research Questions

- **RQ1:** How well general purpose LLMs can write Coq proofs?

Research Questions

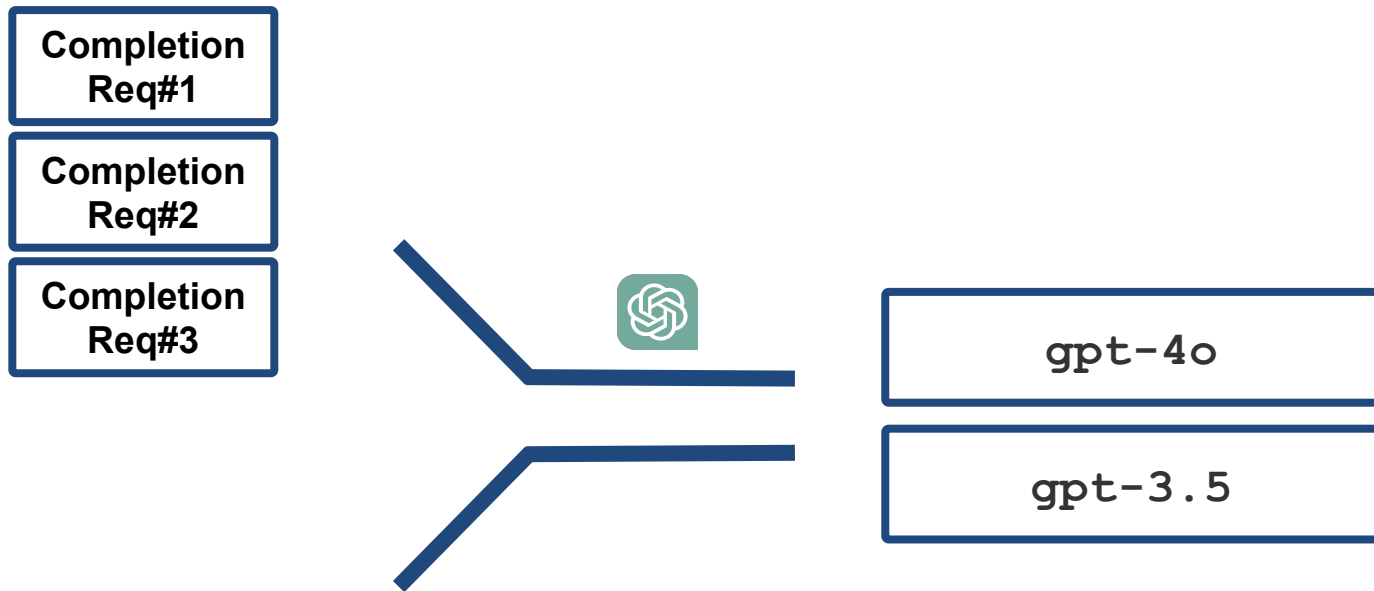
- **RQ1:** How well general purpose LLMs can write Coq proofs?
- **RQ2:** To which extent does CoqPilot improve the LLM approach to Coq generation?

Research Questions

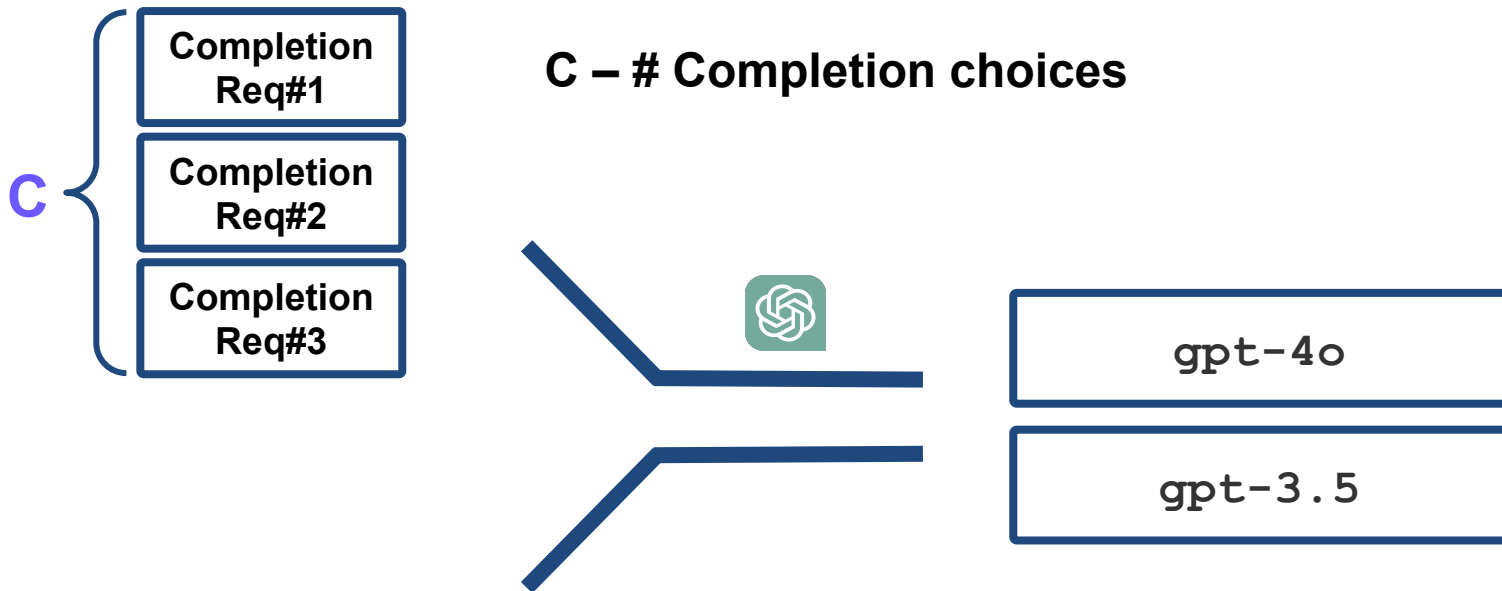
- **RQ1:** How well general purpose LLMs can write Coq proofs?
- **RQ2:** To which extent does CoqPilot improve the LLM approach to Coq generation?
- **RQ3:** What is the additional value CoqPilot contributes to other Coq automation tools such as CoqHammer and Tactician?

Benchmarking Framework

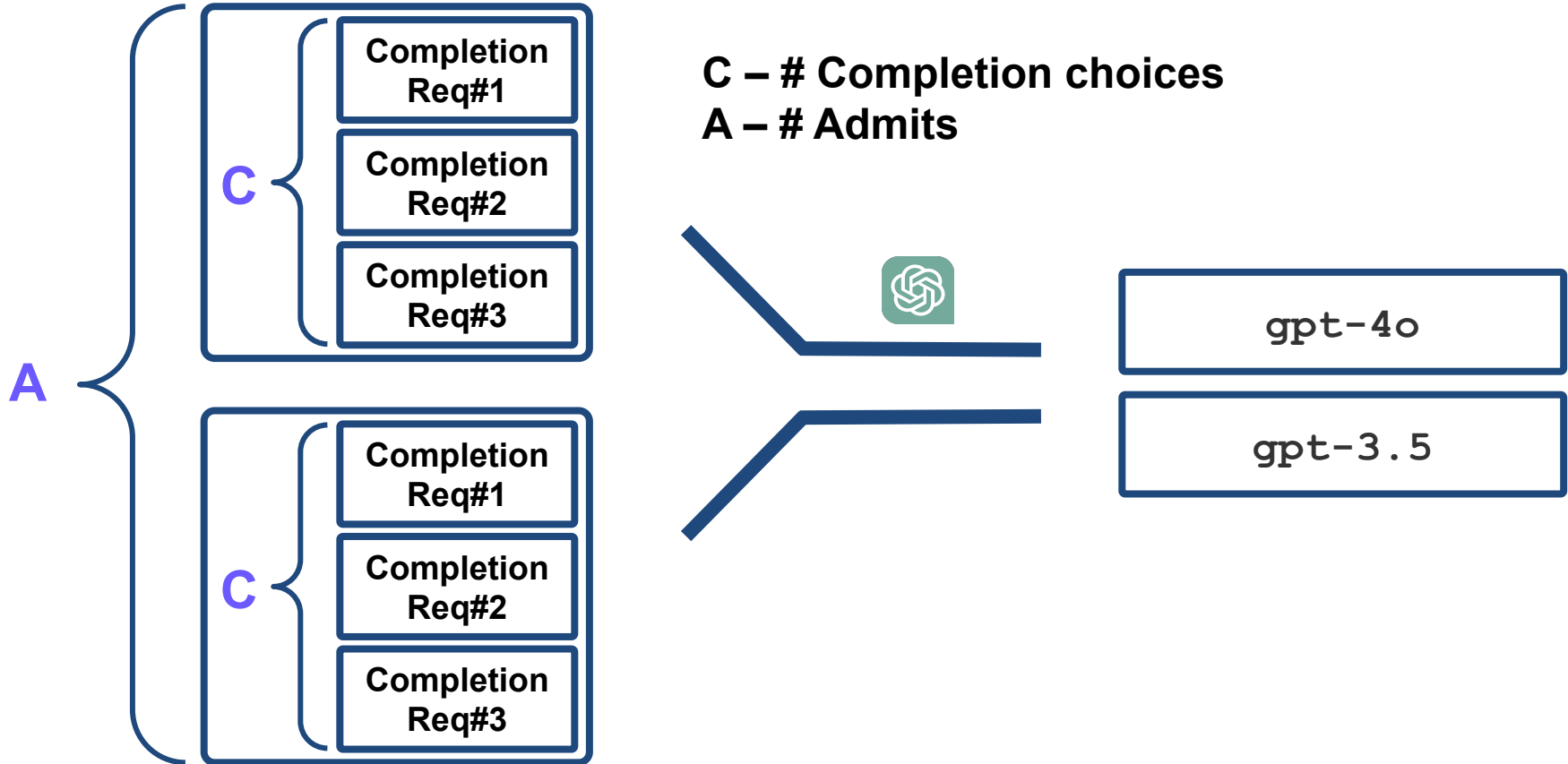
Benchmarking Framework



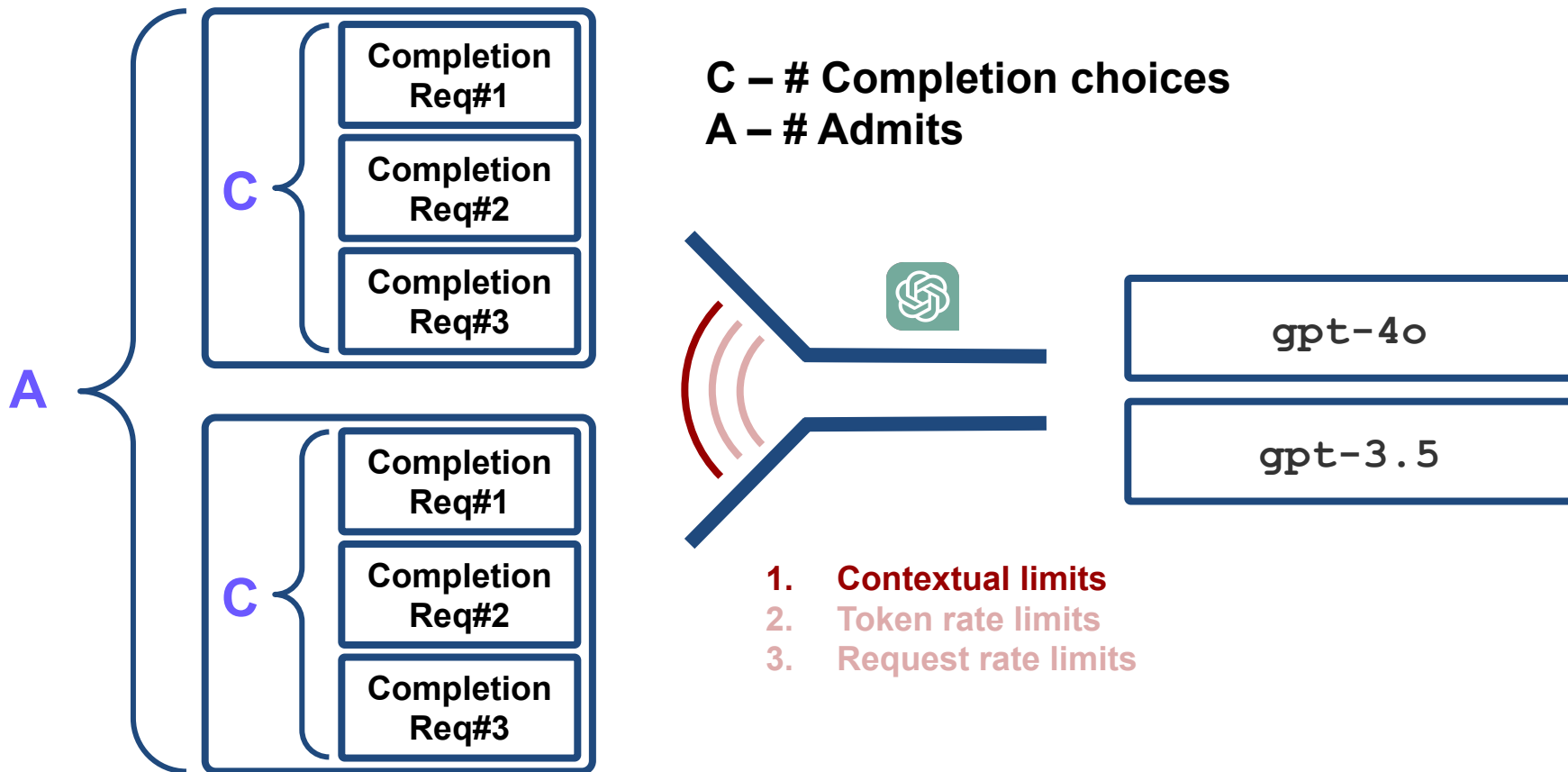
Benchmarking Framework



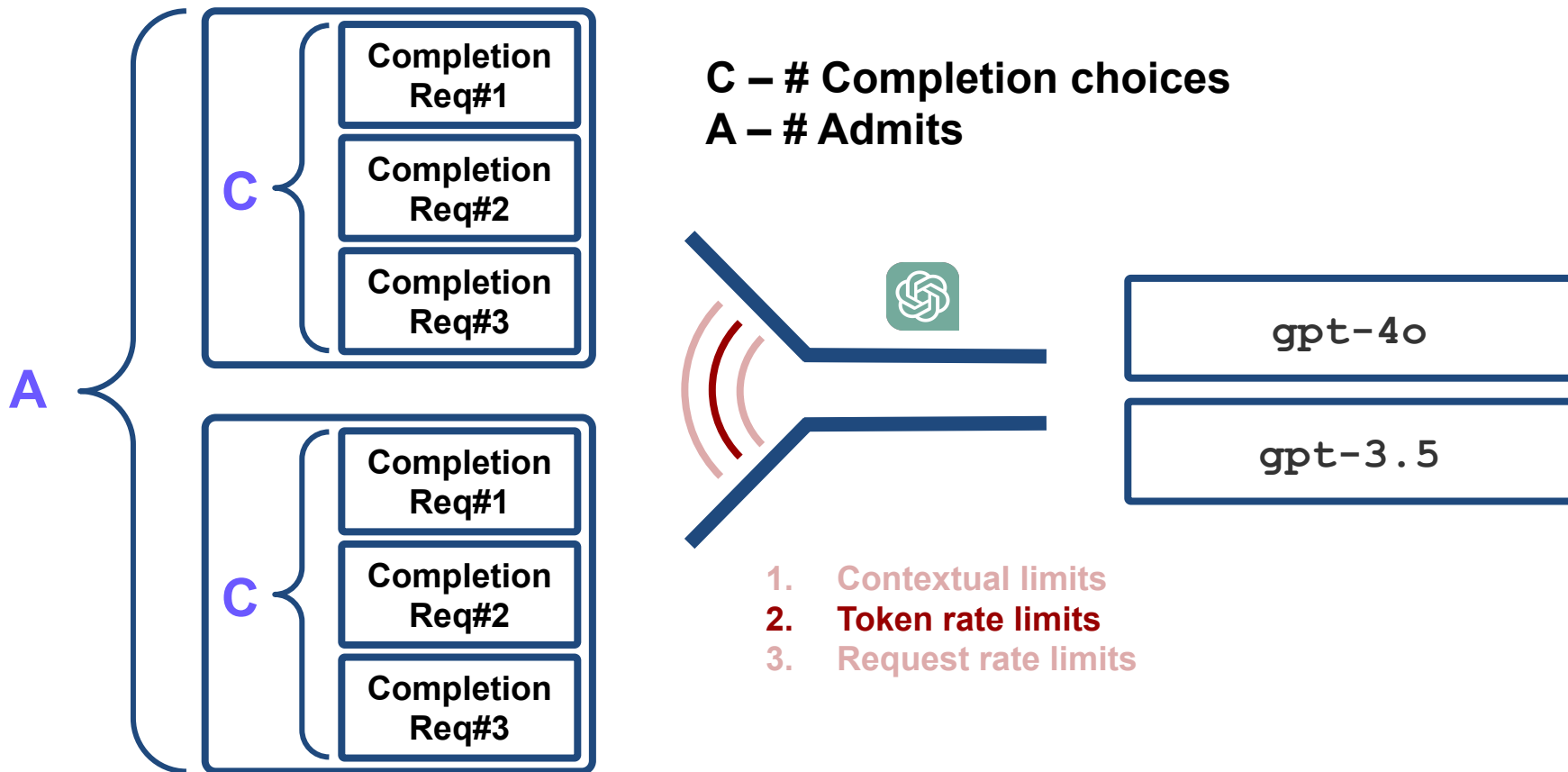
Benchmarking Framework



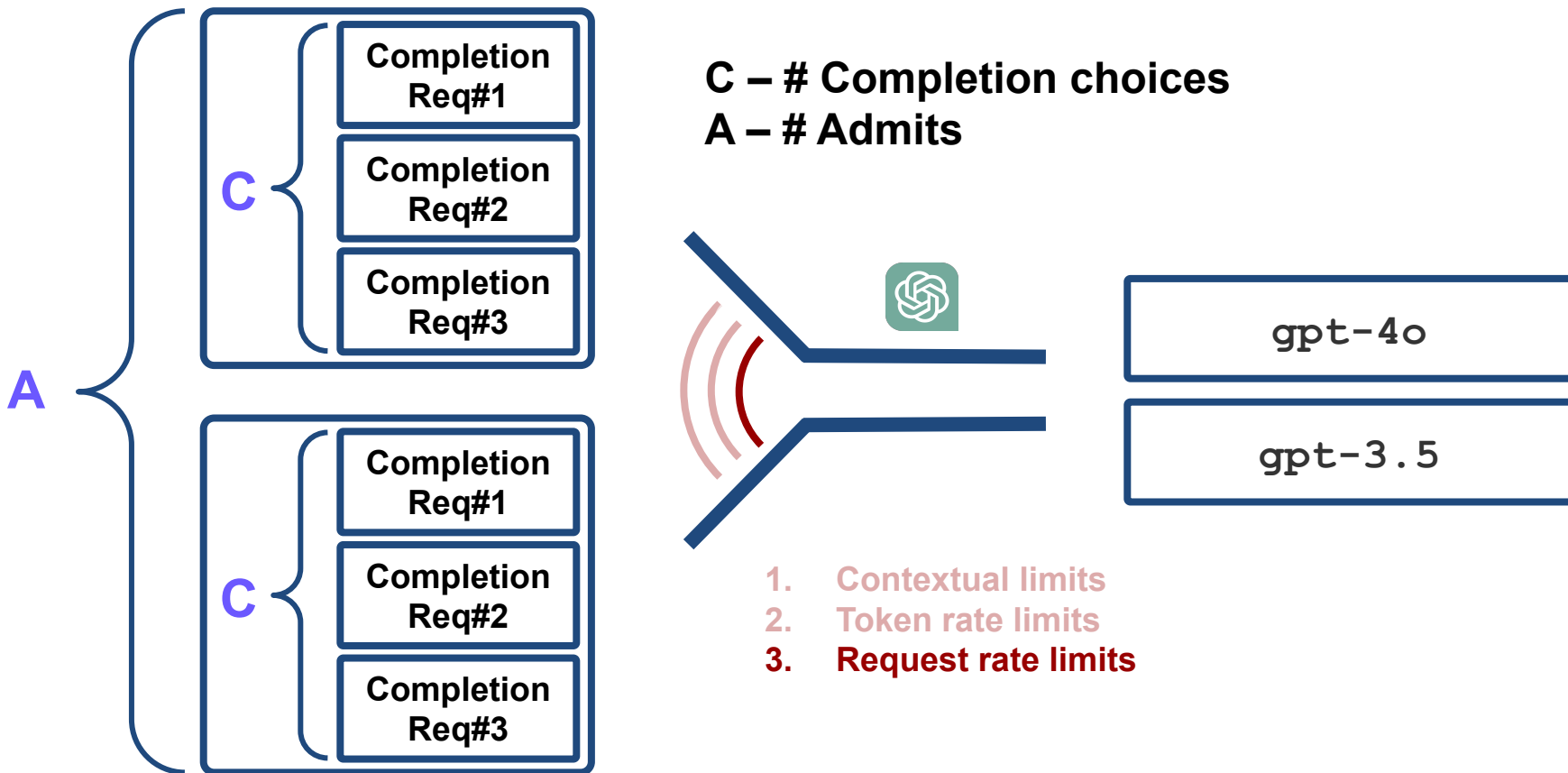
Benchmarking Framework



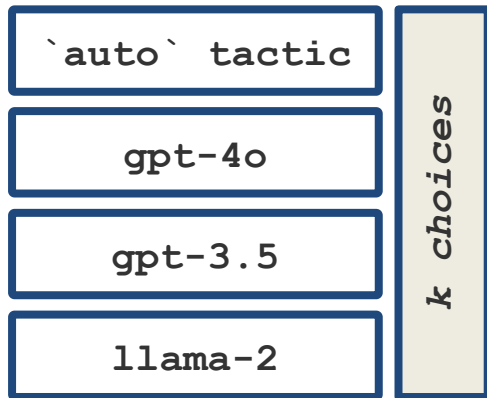
Benchmarking Framework



Benchmarking Framework



Benchmarking Framework



Results

In the table below you can find the results of our experiments. For each of the theorems we show its group and the methods which proved the theorem during our experiments.

File	Theorem Name	Predefined tactic	OpenAI GPT-4	OpenAI GPT-3.5	LLaMA-2 13B Chat
basic/Execution_eco.v	rf_rmw_ct_in_co	×	✓	×	×
imm/imm_hb.v	cr_hb_hb	×	✓	✓	×
basic/FinExecution.v	fin_exec_same_events	✓	✓	✓	×
basic/Execution.v	sb_trans	×	×	×	×
basic/Execution.v	sb_same_loc_W_trans	×	✓	×	×
basic/Events.v	restr_eq_rel_same_loc	✓	✓	✓	✓
basic/Events.v	same_loc_loceq	✓	✓	✓	×

Benchmarking Framework

```
interface Benchmark {  
    name: string;  
    items: DatasetItem[];  
    inputModelsParams: InputModelsParams;  
    requireAllAdmitsCompleted: Boolean;  
    benchmarkFullTheorems: Boolean;  
    benchmarkAdmits: Boolean;  
    timeoutMinutes: number;  
    groupName: string;  
    additionalImports?: AdditionalFileImport[];  
    maximumUsedPremisesAmount?: number;  
    perProofTimeoutMillis: number;  
}
```

Evaluation

Evaluation

Reference proof length	≤ 4	5 – 8	9 – 20	Total
Group size	131	98	71	300
firstorder auto with *	11%	2%	1%	6%
OpenAI GPT-3.5	29%	17%	6%	20%
OpenAI GPT-4o	50%	26%	15%	34%
LLaMA-2 13B Chat	2%	0%	0%	0.5%
Anthropic Claude	21%	7%	7%	13%
All models together	57%	32%	18%	39%
Tactician	45%	23%	10%	29%
CoqHammer	23%	4%	0%	11%
All methods together	71%	45%	23%	51%

IMM Project: <https://github.com/weakmemory/imm>

CoqPilot ASE'24 Tool Demo paper on: <https://podkopaev.net/>

Evaluation

Reference proof length	≤ 4	5 – 8	9 – 20	Total
Group size	131	98	71	300
firstorder auto with *	11%	2%	1%	6%
OpenAI GPT-3.5	29%	17%	6%	20%
OpenAI GPT-4o	50%	26%	15%	34%
LLaMA-2 13B Chat	2%	0%	0%	0.5%
Anthropic Claude	21%	7%	7%	13%
All models together	57%	32%	18%	39%
Tactician	45%	23%	10%	29%
CoqHammer	23%	4%	0%	11%
All methods together	71%	45%	23%	51%

IMM Project: <https://github.com/weakmemory/imm>

CoqPilot ASE'24 Tool Demo paper on: <https://podkopaev.net/>

Evaluation

Reference proof length	≤ 4	5 – 8	9 – 20	Total
Group size	131	98	71	300
firstorder auto with *	11%	2%	1%	6%
OpenAI GPT-3.5	29%	17%	6%	20%
OpenAI GPT-4o	50%	26%	15%	34%
LLaMA-2 13B Chat	2%	0%	0%	0.5%
Anthropic Claude	21%	7%	7%	13%
All models together	57%	32%	18%	39%
Tactician	45%	23%	10%	29%
CoqHammer	23%	4%	0%	11%
All methods together	71%	45%	23%	51%

IMM Project: <https://github.com/weakmemory/imm>

CoqPilot ASE'24 Tool Demo paper on: <https://podkopaev.net/>

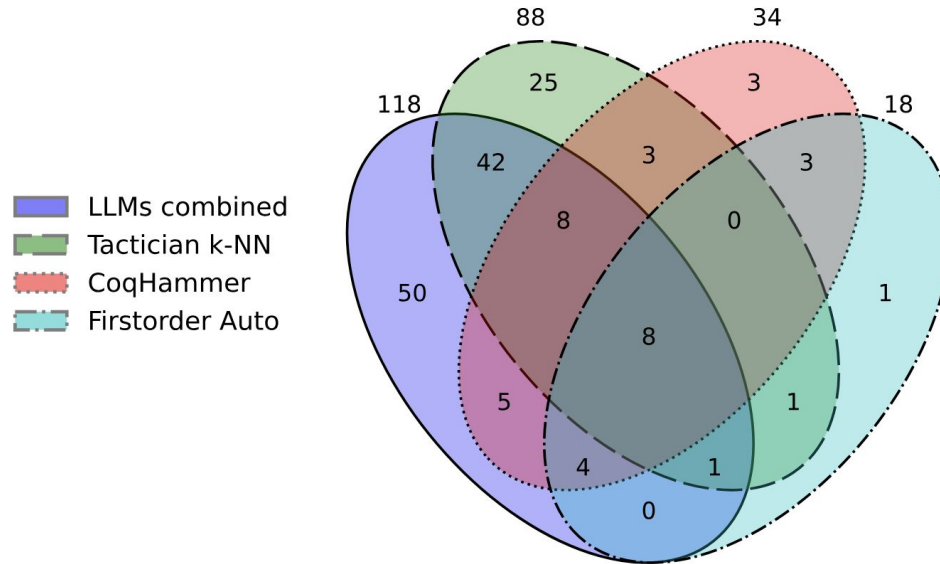
Evaluation

Reference proof length	≤ 4	5 – 8	9 – 20	Total
Group size	131	98	71	300
firstorder auto with *	11%	2%	1%	6%
OpenAI GPT-3.5	29%	17%	6%	20%
OpenAI GPT-4o	50%	26%	15%	34%
LLaMA-2 13B Chat	2%	0%	0%	0.5%
Anthropic Claude	21%	7%	7%	13%
All models together	57%	32%	18%	39%
Tactician	45%	23%	10%	29%
CoqHammer	23%	4%	0%	11%
All methods together	71%	45%	23%	51%

IMM Project: <https://github.com/weakmemory/imm>

CoqPilot ASE'24 Tool Demo paper on: <https://podkopaev.net/>

Evaluation



IMM Project: <https://github.com/weakmemory/imm>

CoqPilot ASE'24 Tool Demo paper on: <https://podkopaev.net/>

Technical challenges: Coq-LSP

Technical challenges: Coq-LSP

Theorem `plus` : forall n:nat, 1 + n = S n.

Proof.

`admit.`

Admitted.

Technical challenges: Coq-LSP

Theorem plus : forall n:nat, 1 + n = S n.

Proof.

admit.

Admitted.

File prefix

Theorem plus : forall n:nat, 1 + n = S n.

Proof.

Proof Candidates

1. reflexivity. 

2. simpl. 

3. constructor. 

Technical challenges: Coq-LSP

```
Theorem plus : forall n:nat, 1 + n = S n.
```

```
Proof.
```

```
  admit.
```

```
Admitted.
```

File prefix

```
Theorem plus : forall n:nat, 1 + n = S n.
```

```
Proof.
```

Proof Candidates

1. reflexivity. 

2. simpl. 

3. constructor. 

```
createTempFile()
```

```
file.write(prefix)
```

```
lsp.openDocument()
```

```
for proof in proofCandidates:
```

```
  file.append(proof)
```

```
  lsp.checkFile()
```

```
  file.remove(proof)
```

Technical challenges: Coq-LSP

```
Theorem plus : forall n:nat, 1 + n = S n.
```

```
Proof.
```

```
  admit.
```

```
Admitted.
```

File prefix

```
Theorem plus : forall n:nat, 1 + n = S n.
```

```
Proof.
```

Proof Candidates

1. reflexivity. 

2. simpl. 

3. constructor. 

```
createTempFile()
```

```
file.write(prefix)
```

```
lsp.openDocument()
```

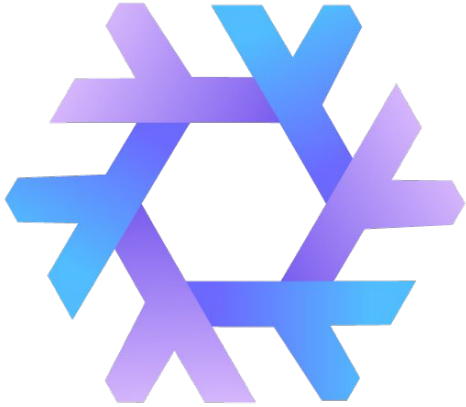
```
for proof in proofCandidates:
```

```
  file.append(proof)
```

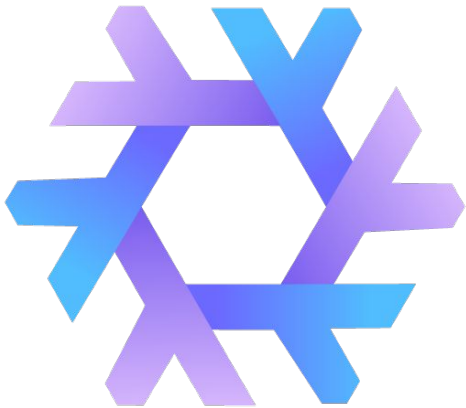
```
  lsp.checkFile()
```

```
  file.remove(proof)
```

Technical challenges: Nix (1/2)



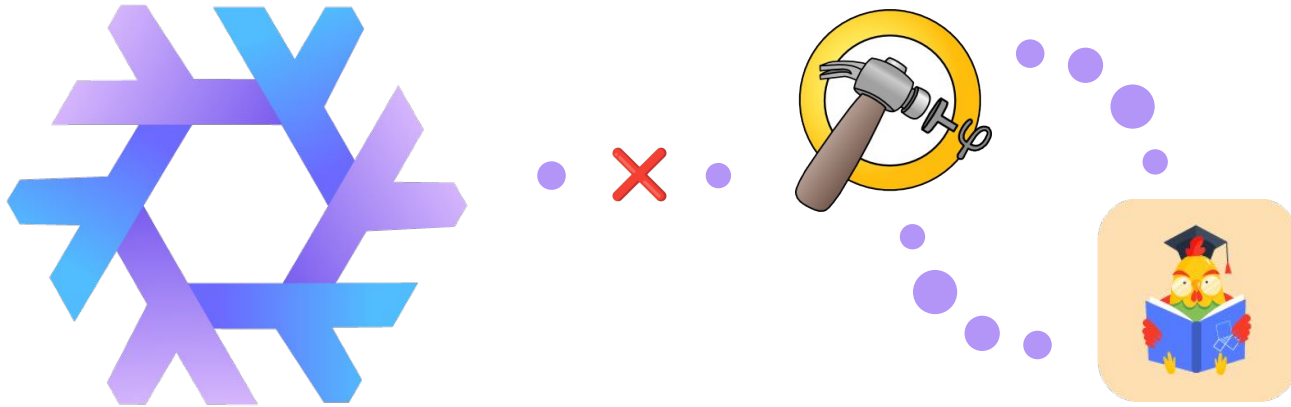
Technical challenges: Nix (1/2)



```
[nix-shell:~/imm]$  
[nix-shell:~/compcert]$  
[nix-shell:~/math-comp]$  
[nix-shell:~/topology]$
```

The image shows a sequence of four terminal prompts, each on a new line. The prompts are: `[nix-shell:~/imm]$`, `[nix-shell:~/compcert]$`, `[nix-shell:~/math-comp]$`, and `[nix-shell:~/topology]$`. The text is in a green monospace font. Four curved purple arrows point from the first prompt to the second, the second to the third, and the third to the fourth, indicating a sequence of operations or a flow of execution.

Technical challenges: Nix (2/2)

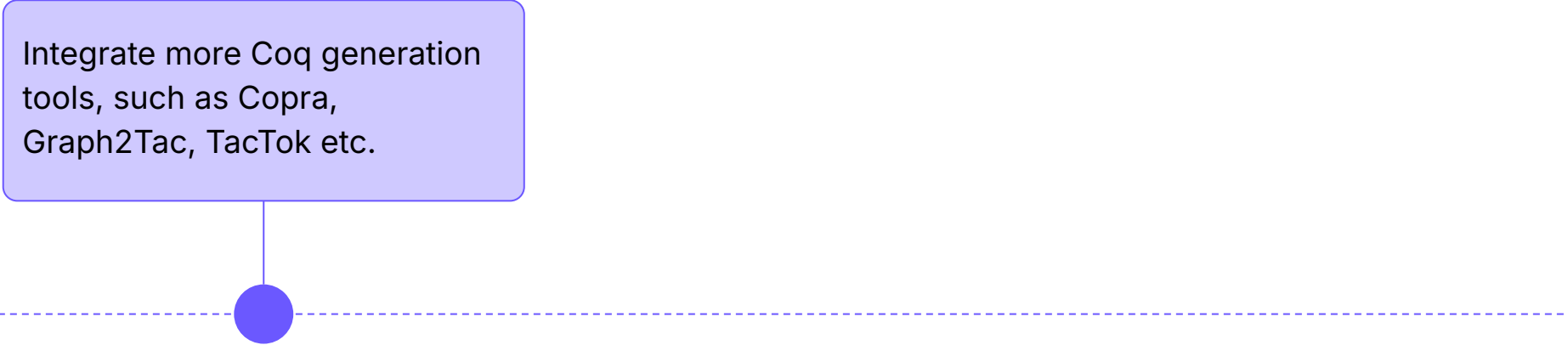


Improvement directions



Improvement directions

Integrate more Coq generation tools, such as Copra, Graph2Tac, TacTok etc.



Improvement directions

Improve premise selection and
add new selection techniques

Integrate more Coq generation
tools, such as Copra,
Graph2Tac, TacTok etc.



Improvement directions

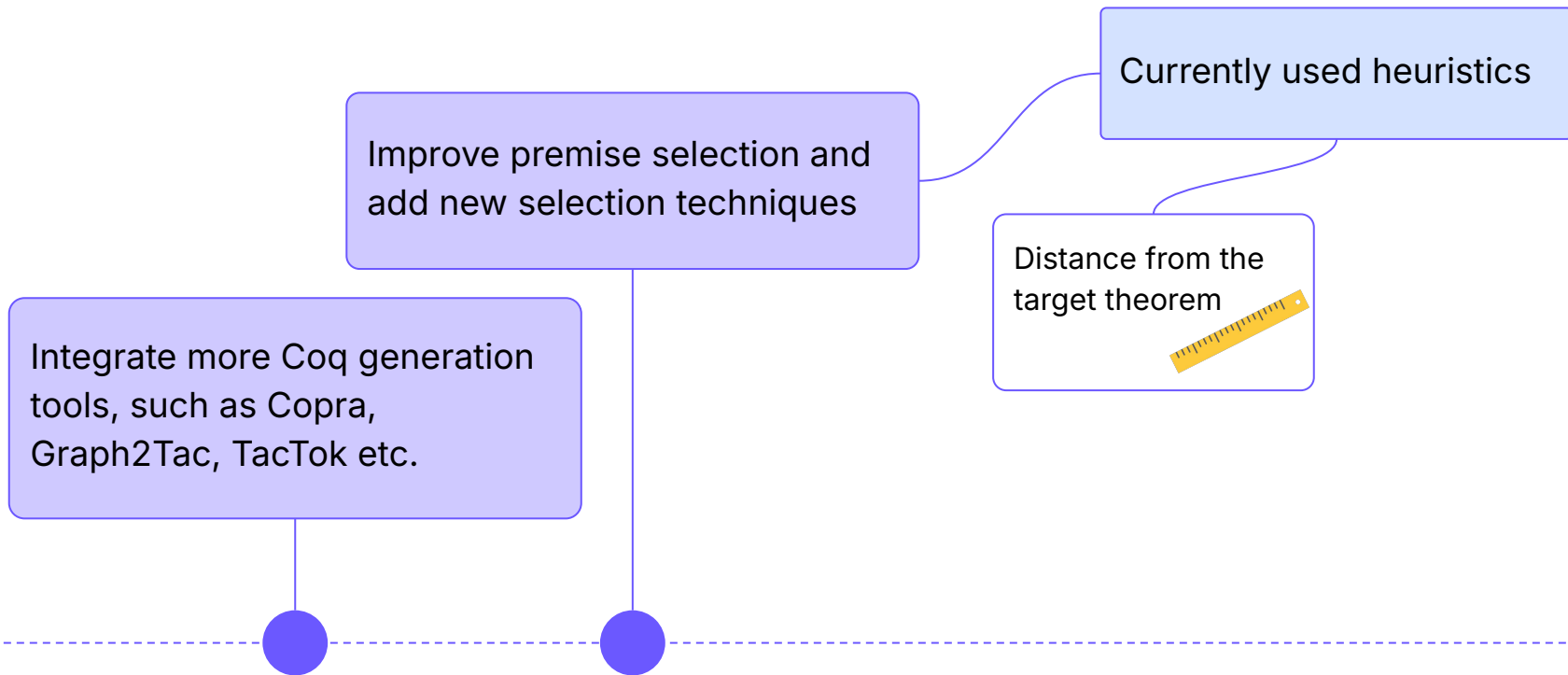
Currently used heuristics

Improve premise selection and
add new selection techniques

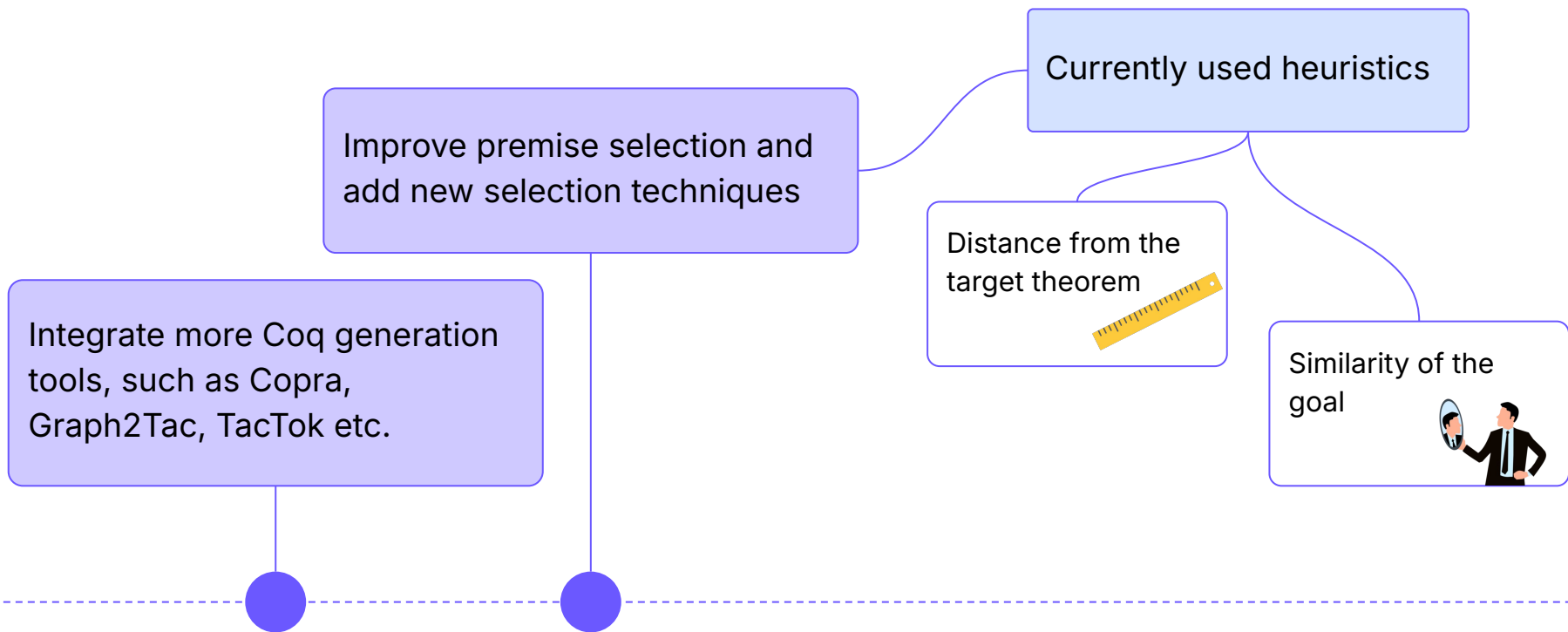
Integrate more Coq generation
tools, such as Copra,
Graph2Tac, TacTok etc.



Improvement directions



Improvement directions



Improvement directions

Improve premise selection and add new selection techniques

Integrate more Coq generation tools, such as Copra, Graph2Tac, TacTok etc.

Explore and improve locally available models in order to make inference cheaper and preserve privacy

Improvement directions

Improve premise selection and
add new selection techniques

Integrate more Coq generation
tools, such as Copra,
Graph2Tac, TacTok etc.

**Please talk to us if you
have ideas!**





CoqPilot: a plugin for LLM-based generation of proofs



JetBrains-Research/**coqpilot**



extension: **coqpilot**



{andrei.kozyrev, gleb.solovev, nikita.khramov, anton.podkopaev}@jetbrains.com

Programming Languages and Program Analysis Lab (PLAN), JetBrains Research