

---

# Divide and Translate: Compositional First-Order Logic Translation and Verification for Complex Logical Reasoning



Hyun Ryu



Gyeongman Kim



Hyemin S. Lee



Eunho Yang

*ryuhyun1905@kaist.ac.kr*

# Logical Reasoning

---

- What is a Logical reasoning?
  - A structured process of reaching conclusions from premises
  - One of the most challenging metrics to measure intelligence
- Large language models (LLMs) unlock the ability of machine to reason
- Chain-of-thought (CoT) prompting improves logical reasoning ability of LLMs
- However, CoT falls short in *complex* logical reasoning tasks which require long sequence of reasoning

# Neurosymbolic Approaches

- Logic-LM, SatLM, LINC, etc. utilize an LLM with a symbolic solver
  1. LLM translates NL problem into satisfiability (SAT) problem that consists of FOL formulas
  2. Symbolic solver returns a mathematically correct solution
- By considering LLM as a semantic parser, it can avoid errors in reasoning steps
- Example

Each of five students—Hubert, Lori, Paul, Regina, and Sharon—will visit exactly one of three cities—Montreal, Toronto, or Vancouver—for the month of March, according to the following conditions: Sharon visits a different city than Paul. Hubert visits the same city as Regina. Lori visits Montreal or else Toronto. If Paul visits Vancouver, Hubert visits Vancouver with him. Each student visits one of the cities with at least one of the other four students.

Question: Which one of the following must be true?

(A) If any of the students visits Montreal, Lori visits Montreal. (B) ...



LLM as a semantic parser (NL-to-FOL)

```
students = [Hubert, Lori, Paul, Regina, Sharon]
cities = [Montreal, Toronto, Vancouver]
visits = Function(students -> cities)
```

```
visits(Sharon) != visits(Paul)
```

```
Or(visits(Lori) == Montreal, visits(Lori) == Toronto)
```

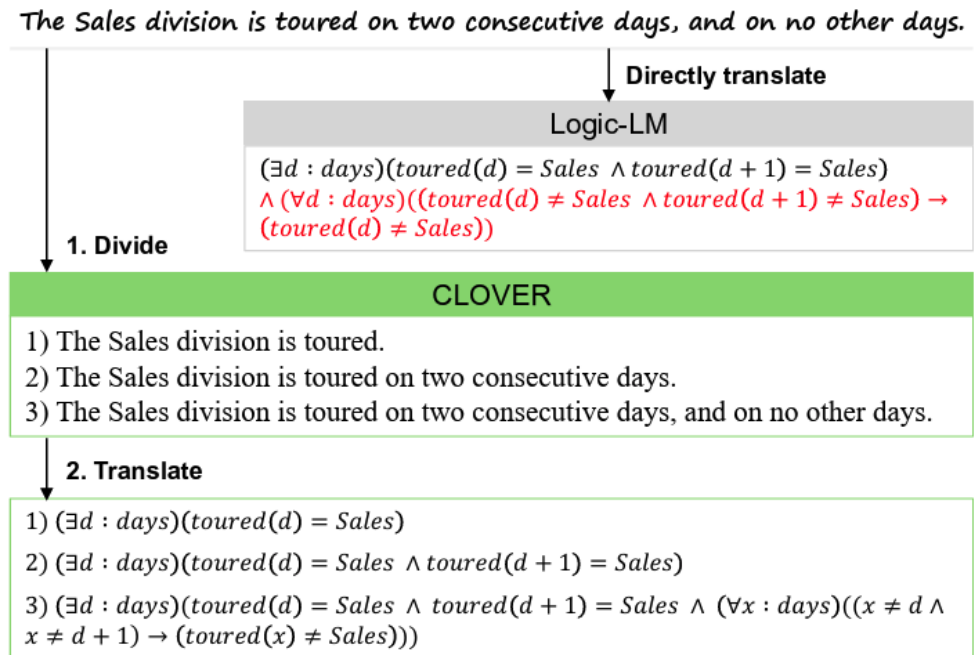
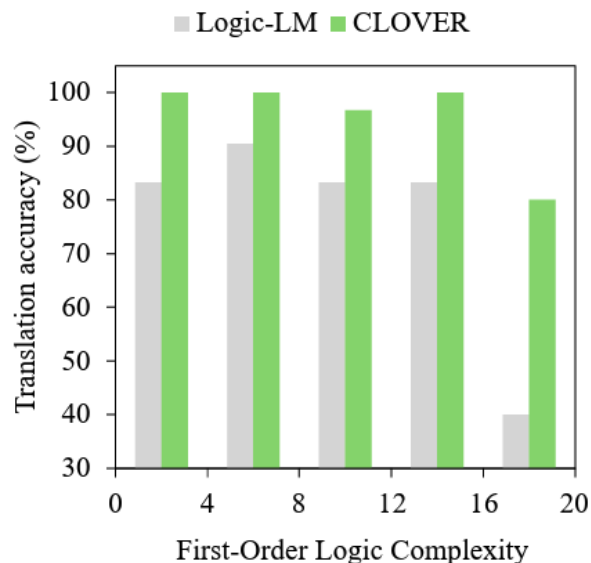
```
ForAll([s1], Exists([s2], And(s2 != s1, visits(s1) == visits(s2))))
```

```
...
```

```
solve(Implies(Exists([s], visits(s) == Montreal), visits(Lori) == Montreal))
```

# Limitations of Neurosymbolic Approaches

- LLM still cannot translate complex sentences
  - It falls short beyond a certain degree of FOL complexity
  - It correctly translates NL with simple logic, but fails to translate NL with more complex logic
- Our Idea: Divide and Translate!
  - First translate an atomic subsentence, and then translate more complex ones sequentially



# Proposed Method: CLOVER

- But, how to divide the sentences?
- *Logical Dependency Parsing*
  - New parsing method for NL that represents FOL
  - Categorize sentence components into the following three:
    - Logic Units (U), Logic Couplers (C), and Logic Dependents (D)
  - Example
    - Target sentence: “Hamadi cannot be appointed to the same court as Perkins.”
    - Logical Dependency Structure:

## Components:

$U_1$  = “Hamadi is appointed to a court”,  
 $U_2$  = “Perkins is appointed to a court”,  
 $D_1$  = “cannot”,  $C_1$  = “the same ... as”

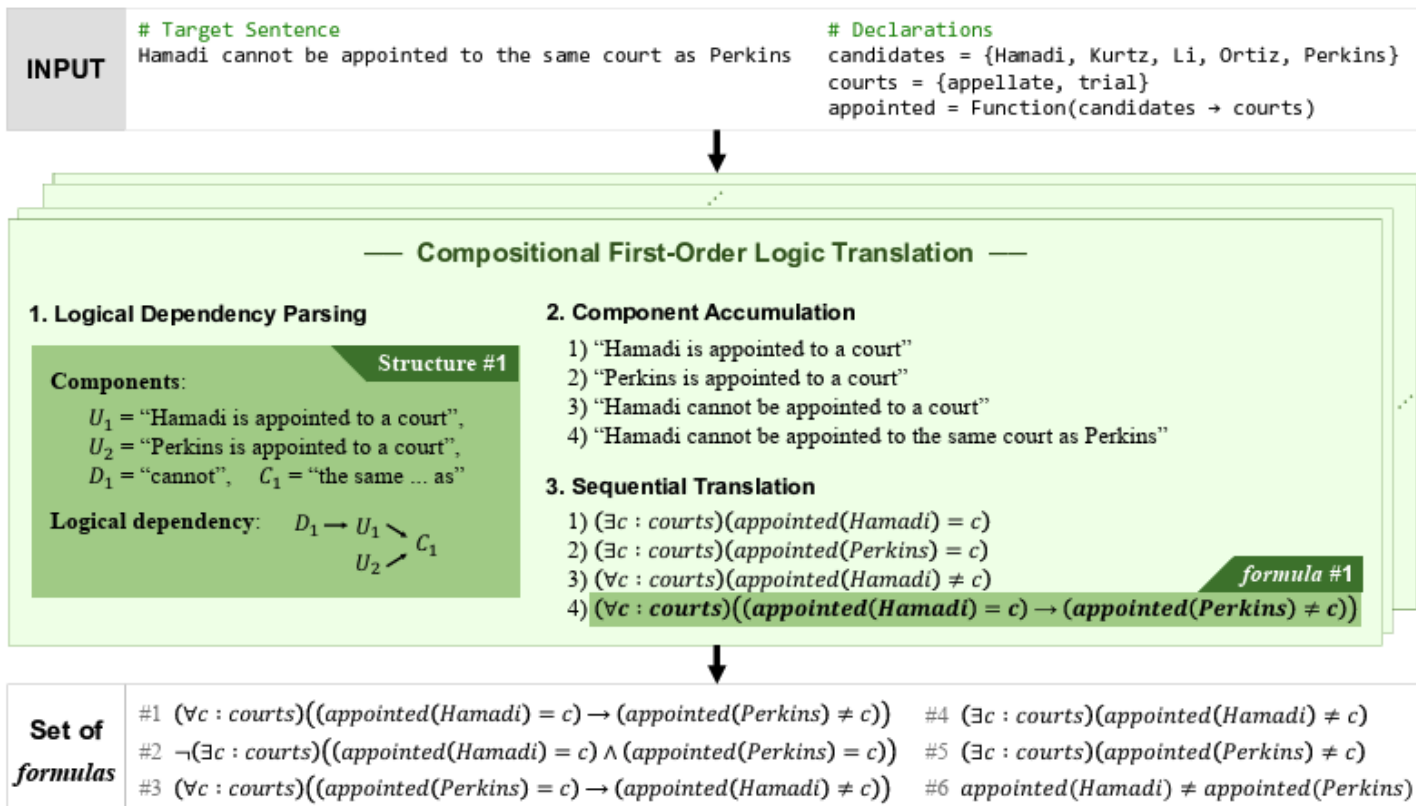
Logical dependency:  $D_1 \rightarrow U_1 \searrow C_1$   
 $U_2 \nearrow C_1$

## Structure #1

# Proposed Method: CLOVER

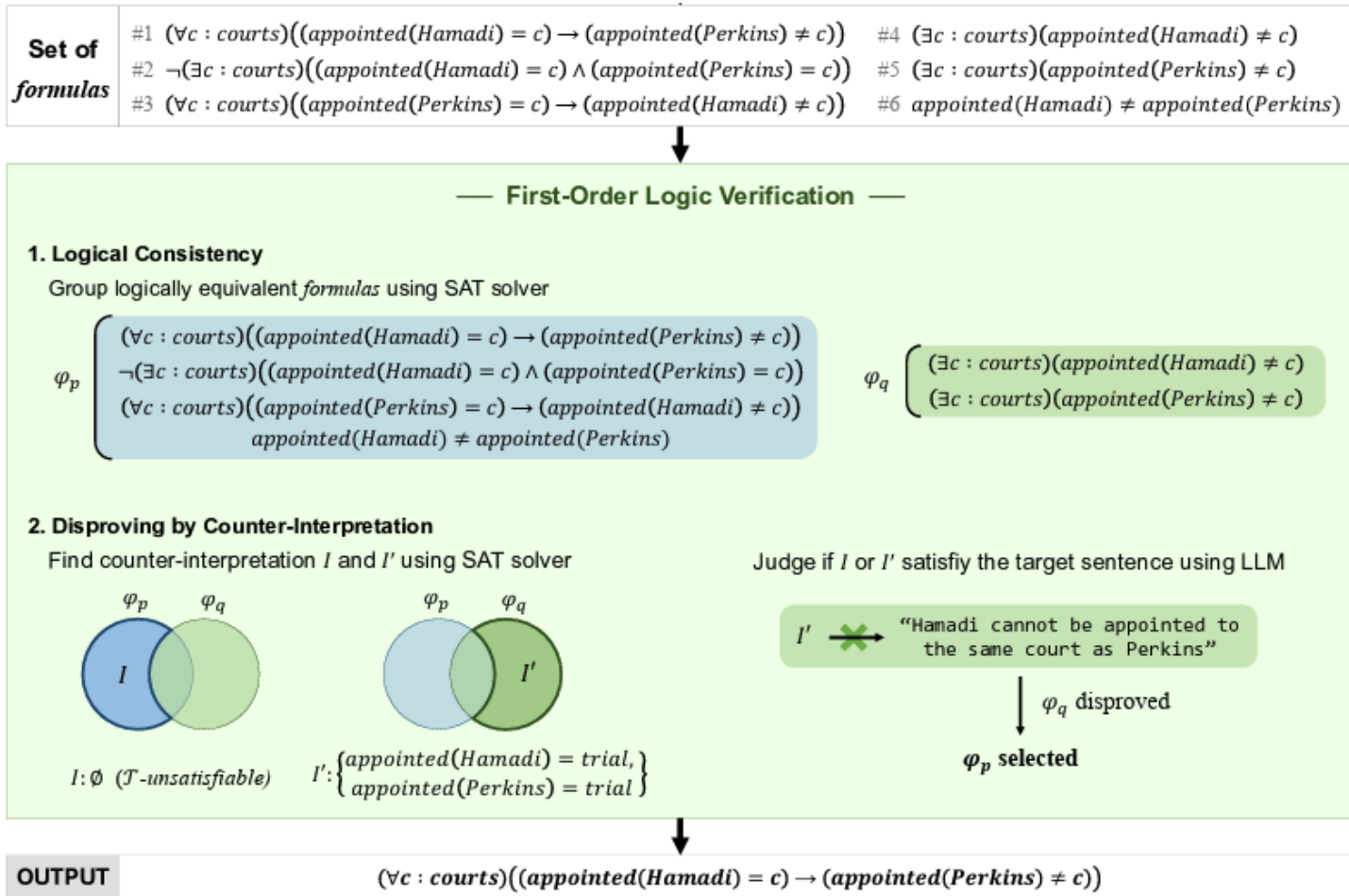
- Compositional First-Order Logic Translation

- Consists of three steps
- There could be multiple output formulas, so we need to select the most reliable formula



# Proposed Method: CLOVER

- *First-Order Logic Verification*
  - Two verification algorithms



# Experiments

- Setup
  - Language models: gpt-4o, gpt-4o-mini, etc.
  - Baselines: Standard/CoT prompting and Neurosymbolic approaches
- CLOVER outperforms all the baselines!

	AR-LSAT	ZebraLogic	Puzzle	Symbol	Deduction	FOLIO	ProofWriter
Standard	30.3	0.4	63.0	74.7	84.7	70.9	53.7
CoT	36.8	0.4	51.0	80.8	94.0	73.9	78.0
SymbCoT	34.2	0.8	66.5	55.6	90.7	76.9	80.2
Logic-LM	42.4	45.4	64.0	81.8	95.3	75.4	95.3
CLOVER	<b>62.8</b>	<b>75.4</b>	<b>83.5</b>	<b>89.9</b>	<b>99.3</b>	<b>78.8</b>	<b>96.7</b>

	Program Acc		Execution Rate		Execution Acc	
	Logic-LM	CLOVER	Logic-LM	CLOVER	Logic-LM	CLOVER
AR-LSAT	17.3	<b>46.8</b>	33.8	<b>59.7</b>	51.3	<b>78.3</b>
Puzzle	60.0	<b>79.0</b>	79.5	<b>80.0</b>	75.5	<b>98.8</b>
Symbol	49.5	<b>76.8</b>	52.5	<b>82.8</b>	<b>94.2</b>	92.7
Deduction	92.7	<b>99.0</b>	97.3	<b>99.7</b>	95.2	<b>99.3</b>
FOLIO	51.2	<b>62.6</b>	65.5	<b>74.9</b>	78.2	<b>83.6</b>
ProofWrtier	94.2	<b>96.5</b>	96.8	<b>99.2</b>	97.2	<b>97.3</b>