



Red Hat

# NetConfEval: Can LLMs Facilitate Network Configuration?

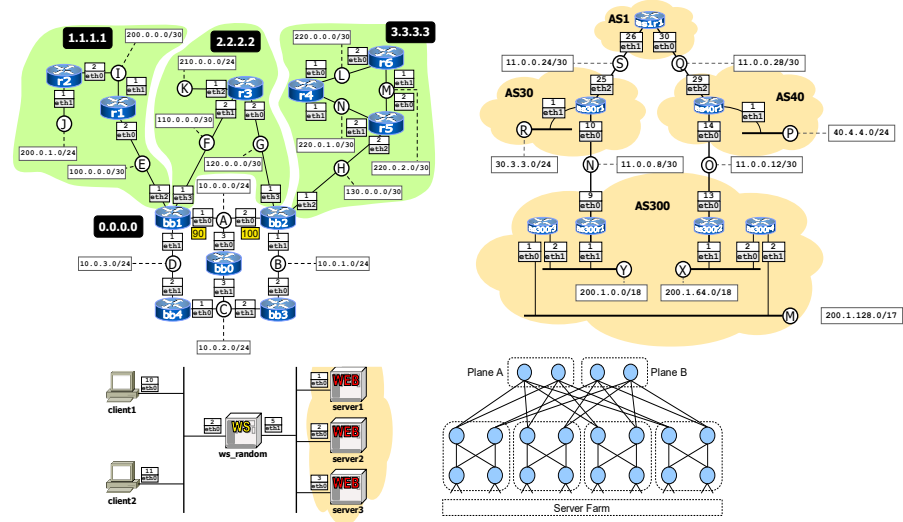
**Changjie Wang**<sup>1</sup>, Mariano Scazzariello<sup>1 2</sup>, Alireza Farshin<sup>2</sup>, Simone Ferlin<sup>3</sup>, Dejan Kostic<sup>1 2</sup>, Marco Chiesa<sup>1</sup>

KTH Royal Institute of Technology <sup>1</sup> - RISE Research Institutes of Sweden <sup>2</sup> - Red Hat <sup>3</sup>

# Orchestrating networks is a **cumbersome** task



network  
operator



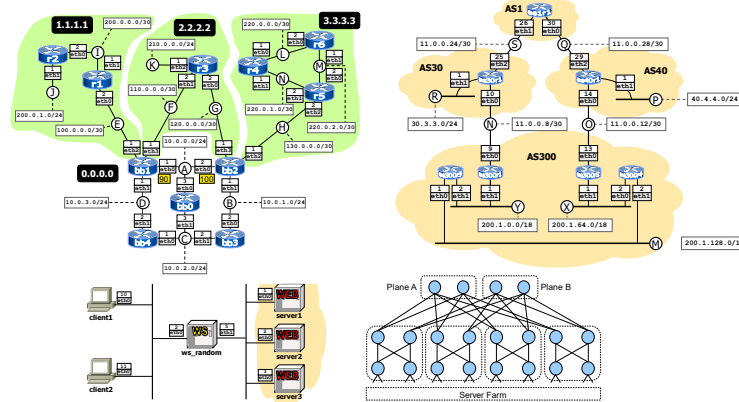
# Orchestrating networks is a **cumbersome** task



network  
operator



configure



large software  
documentation



complex  
network/application  
failures



cyber-attacks  
or vulnerabilities



misunderstanding  
within the team

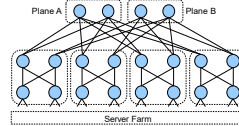
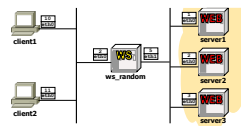
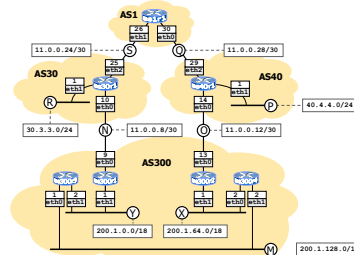
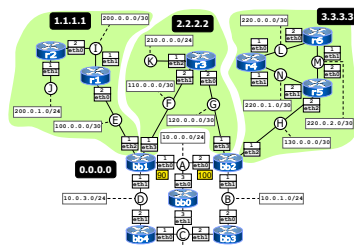
# Orchestrating networks is a **cumbersome** task



network  
operator



configure



Manual configuration is  
costly, difficult, and  
susceptible to **human error**

minor error



**Flight  
disruption**



**Outage**



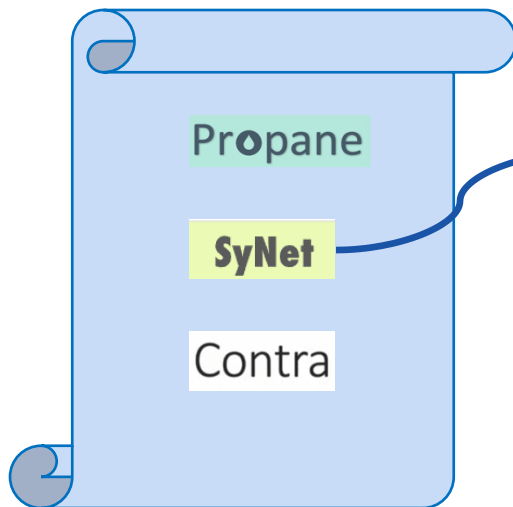
# Network Synthesis Tools **Help** but ...

Network Synthesizer



Generate low-level configurations from intents

But employ a self-defined specification or language to describe the network intents



$$\begin{array}{lll}
 \text{(Program)} \ P ::= \bar{r} & \text{(Literal)} \ l ::= a \mid \neg a & \text{(Variables)} \ X, Y \in Vars \\
 \text{(Rule)} \ r ::= a \leftarrow \bar{l} & \text{(Predicates)} \ p, q \in Preds & \text{(Values)} \ v \in Vals \\
 \text{(Atom)} \ a ::= p(\bar{t}) & \text{(Term)} \ t ::= X \mid v & 
 \end{array}$$

Syntax for Stratified Datalog



network  
operator

It's hard to learn and extend!

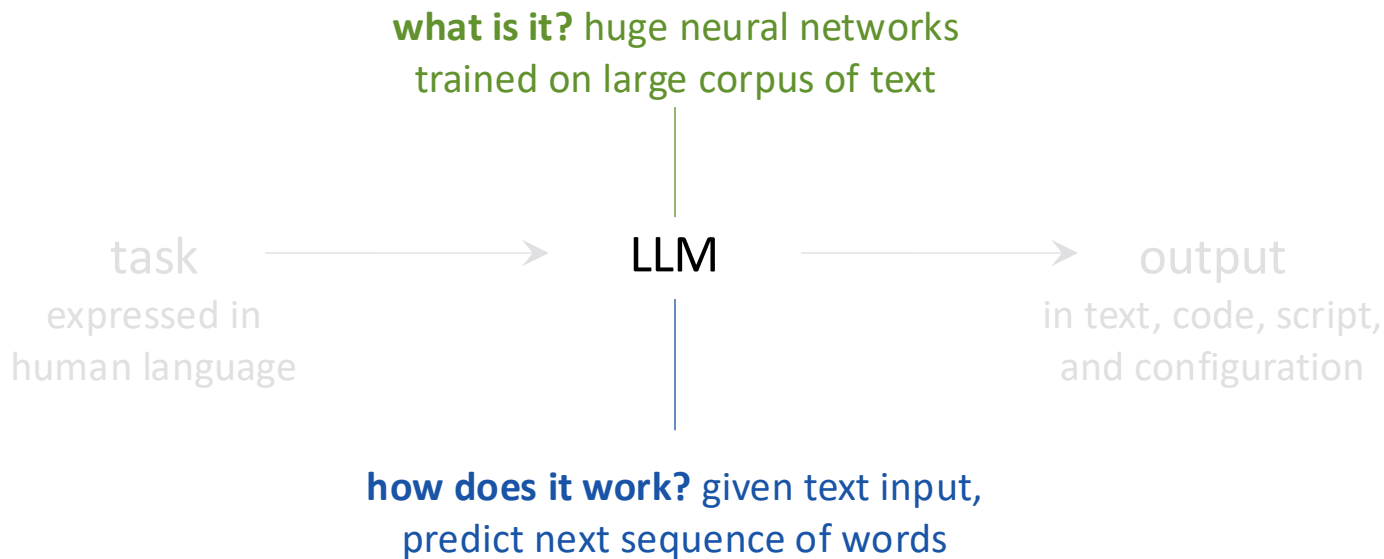


# Can **Large Language Models (LLMs)** help?





# Can **Large Language Models (LLMs)** help?





# Can **Large Language Models (LLMs)** help?



# Can **Large Language Models** (LLMs) help?



OpenAI  
ChatGPT **4.0**

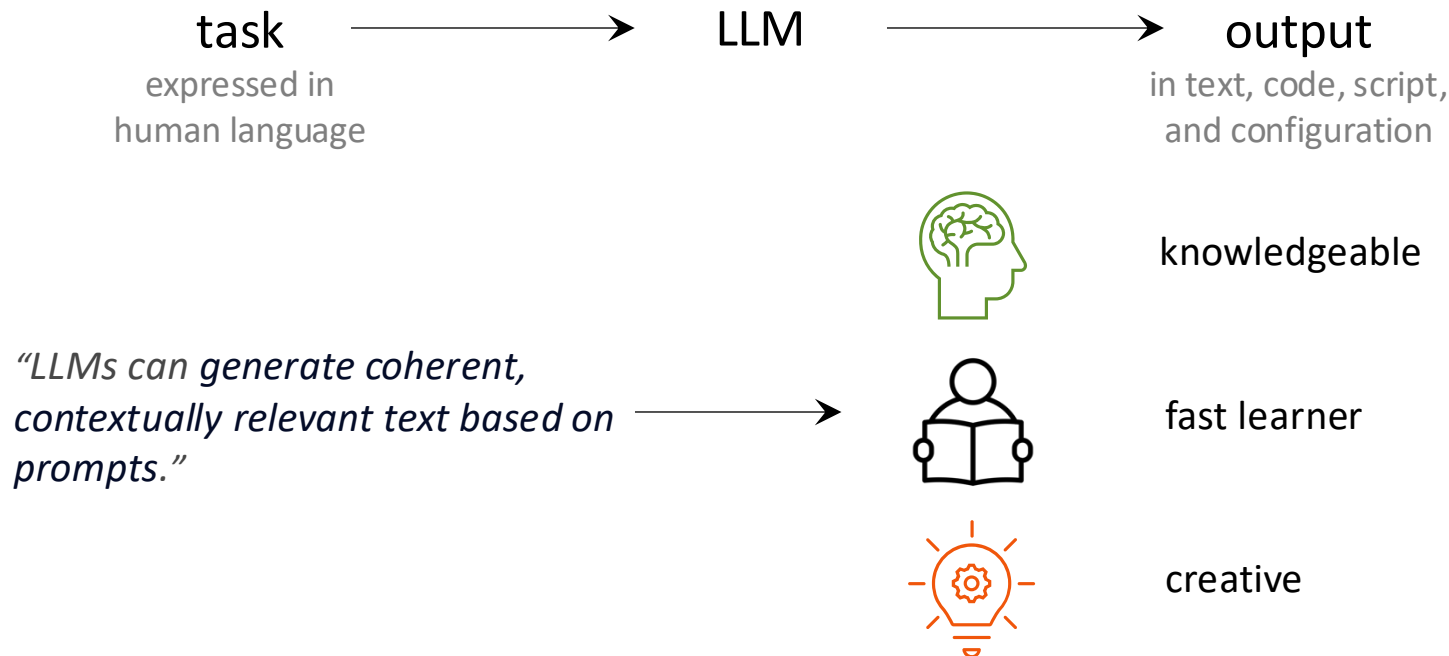


Claude 3

**LLaMA**  
by  Meta

Example of LLMs

# Can Large Language Models (LLMs) help?



# Can Large Language Models (LLMs) help?



## Can LLMs help in network orchestration?

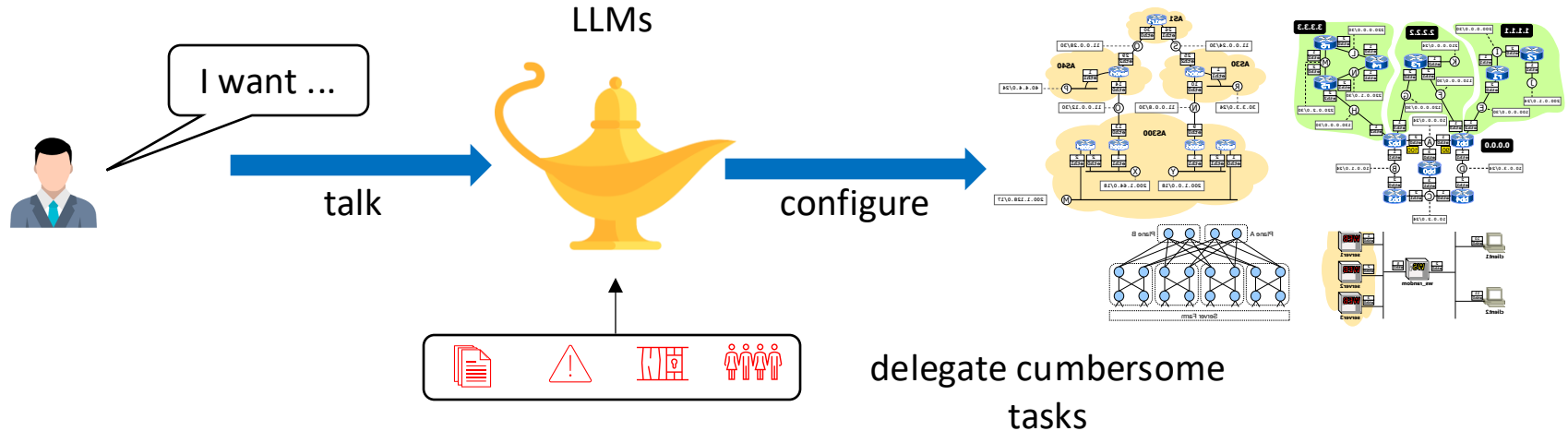
*"LLMs can generate coherent,  
contextually relevant text based on  
prompts."*



fast learner

creative

# Can LLMs help in network orchestration?





# Can LLMs help in network orchestration?

Opportunities come with **challenges**

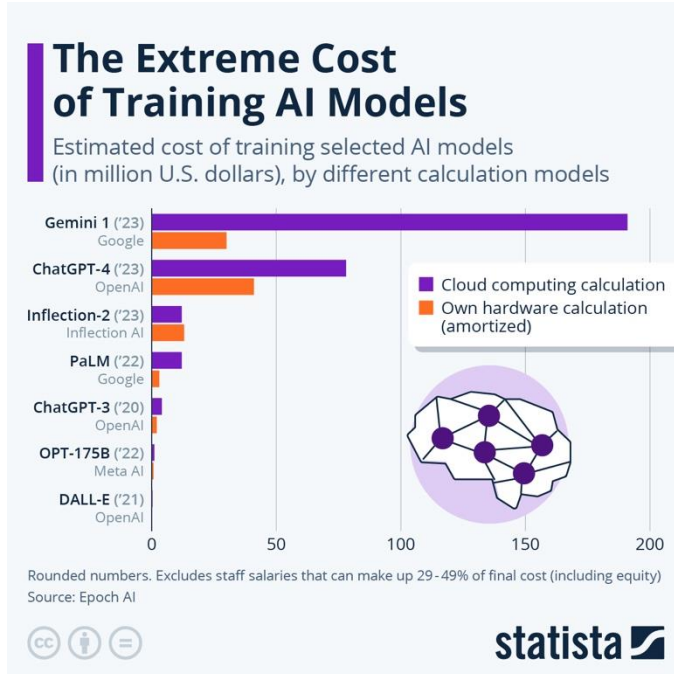


unreliability



costs

# Can LLMs help in network orchestration?

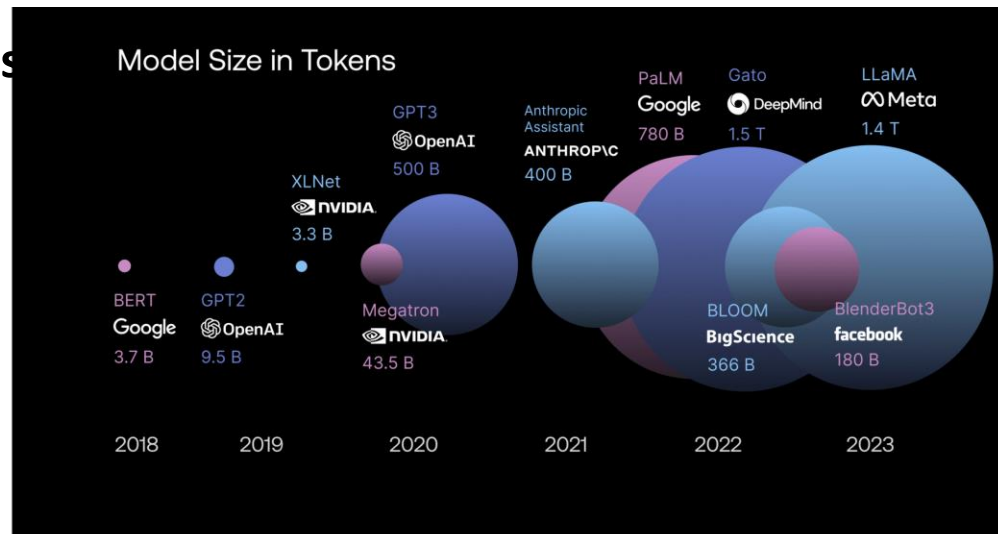
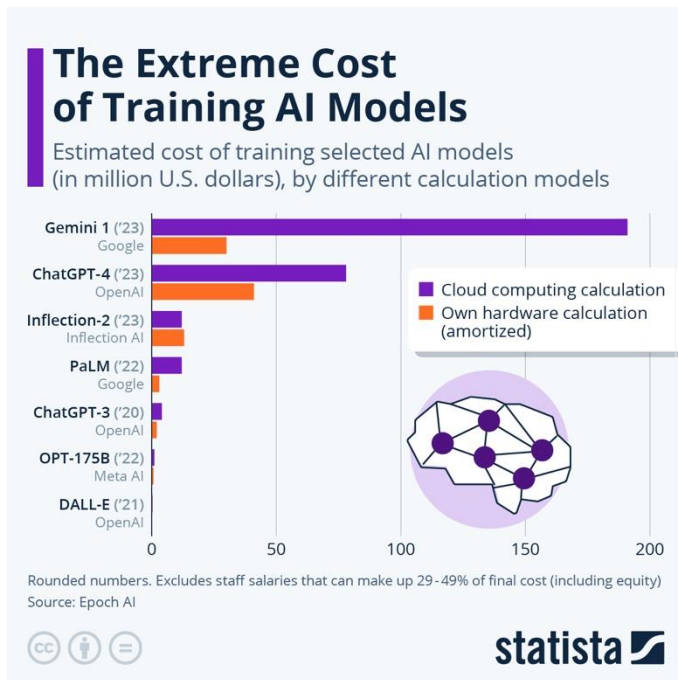


es come with **challenges**



**costs**

# Can LLMs help in network orchestration?





# Assess LLMs in today's networking tasks

1. Design a set of **benchmarks** (NetConfEval) to evaluate LLMs for networking
2. Formulate **takeaways** based on our benchmarking experiment
3. Present **prototypes** for LLM-based networking systems



# We'll focus on **three** tasks in **orchestrating** networks

1. Translating high-level requirements to a formal specification format
2. Adapting code to new requirements
3. Generating low-level configurations



# We'll focus on three tasks in orchestrating networks

1. Translating high-level requirements to a formal specification format
2. Adapting code to new requirements
3. Generating low-level configurations

# Translating **high-level requirements** to a **formal specification format**

*"traffic from **Rome** to **Milan**  
must **traverse** a **firewall**"*



network  
operator

LLM

dictionary

```
{ "reachability": {  
  "rome": [ "milan" ] },  
  "waypoint": {  
    [ "rome", "milan" ] :  
      [ "fw1", "fw2" ] },  
  "avoidance": { }  
}
```

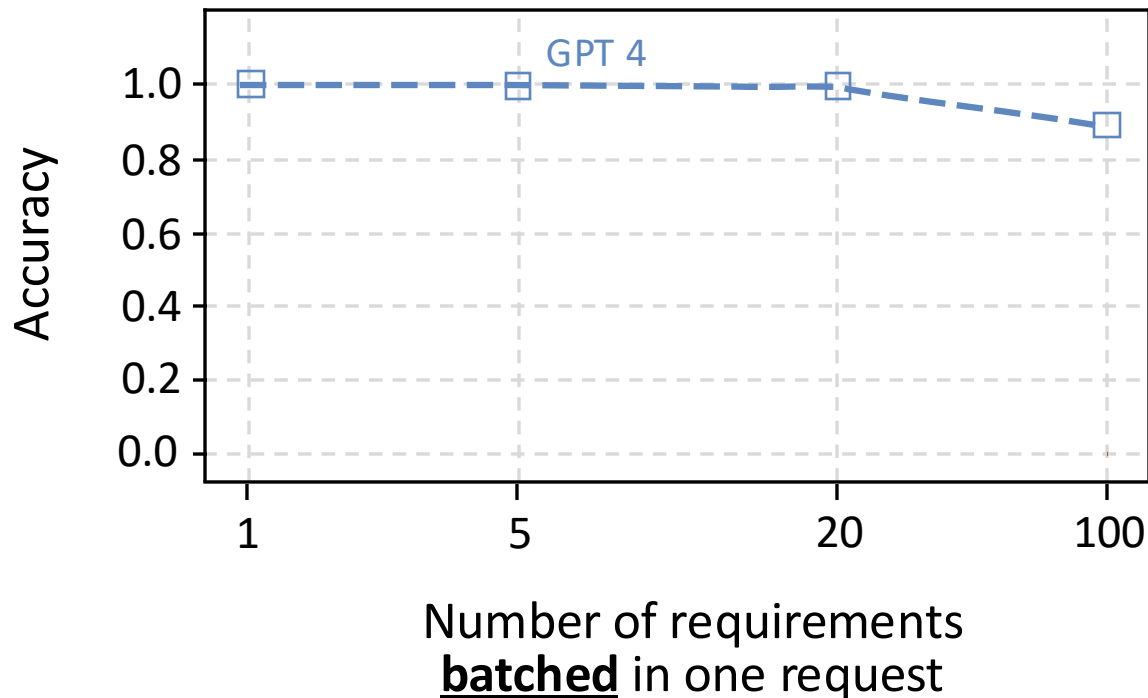


# Translating **high-level requirements** to a **formal specification format**

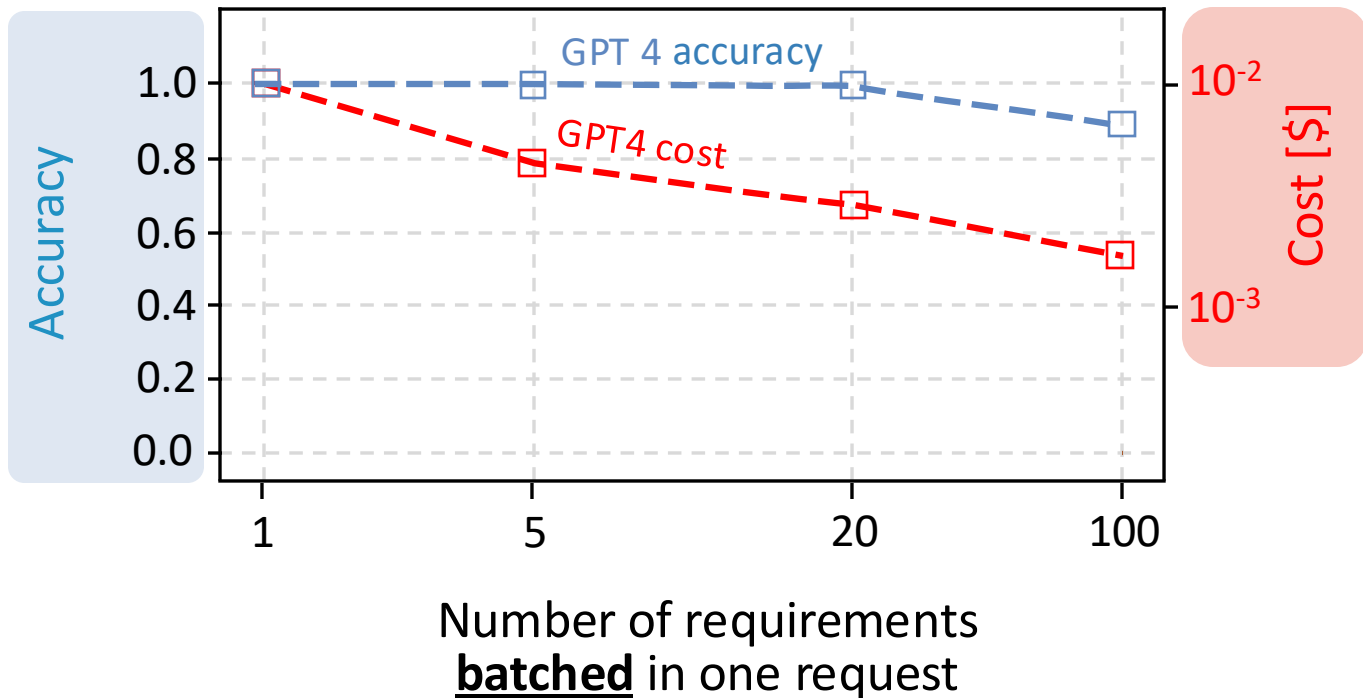
1. Generate 3200 network requirements focusing on reachability, waypoint, and load-balancing using Config2Spec<sup>1</sup>
2. Pick a certain number of requirements and sliced them with various batch sizes
3. Transform them to natural language based on predefined templates
4. Evaluate the efficiency of different LLMs in translation



# GPT4 translates **accurately** requirements



# GPT4 translates **accurately** requirements at a **cost**



# Issues with very large language models

## MegaScale: Scaling Large Language Model Training to More Than 10,000 GPUs

### Authors:

Ziheng Jiang and Haibin Lin, *ByteDance*; Yinmin Zhong, *Peking University*; Qi Huang, Yangrui Chen, Zhi Zhang, Yanghua Peng, Xiang Li, Cong Xie, Shibiao Nong, Yulu Jia, Sun He, Hongmin Chen, Zhihao Bai, Qi Hou, Shipeng Yan, Ding Zhou, Yiyao Sheng, Zhuo Jiang, Haohan Xu, Haoran Wei, Zhang Zhang, Pengfei Nie, Leqi Zou, Sida Zhao, Liang Xiang, Zherui Liu, Zhe Li, Xiaoying Jia, and Jianxi Ye, *ByteDance*; Xin Jin, *Peking University*; Xin Liu, *ByteDance*

# Problems with very large language models

## Large language models:

- 1000B parameters
- slow inferences
- resource intensive
- hard to deploy

## MegaScale: Scaling Large Language Model Training to More Than 10,000 GPUs

### Authors:

Ziheng Jiang and Haibin Lin, *ByteDance*; Yinmin Zhong, *Peking University*; Qi Huang, Yangrui Chen, Zhi Zhang, Yanghua Peng, Xiang Li, Cong Xie, Shibiao Nong, Yulu Jia, Sun He, Hongmin Chen, Zhihao Bai, Qi Hou, Shipeng Yan, Ding Zhou, Yiyao Sheng, Zhuo Jiang, Haohan Xu, Haoran Wei, Zhang Zhang, Pengfei Nie, Leqi Zou, Sida Zhao, Liang Xiang, Zherui Liu, Zhe Li, Xiaoying Jia, and Jianxi Ye, *ByteDance*; Xin Jin, *Peking University*; Xin Liu, *ByteDance*



# The quest towards **smaller (cheaper!)** models

## **Large** language models:

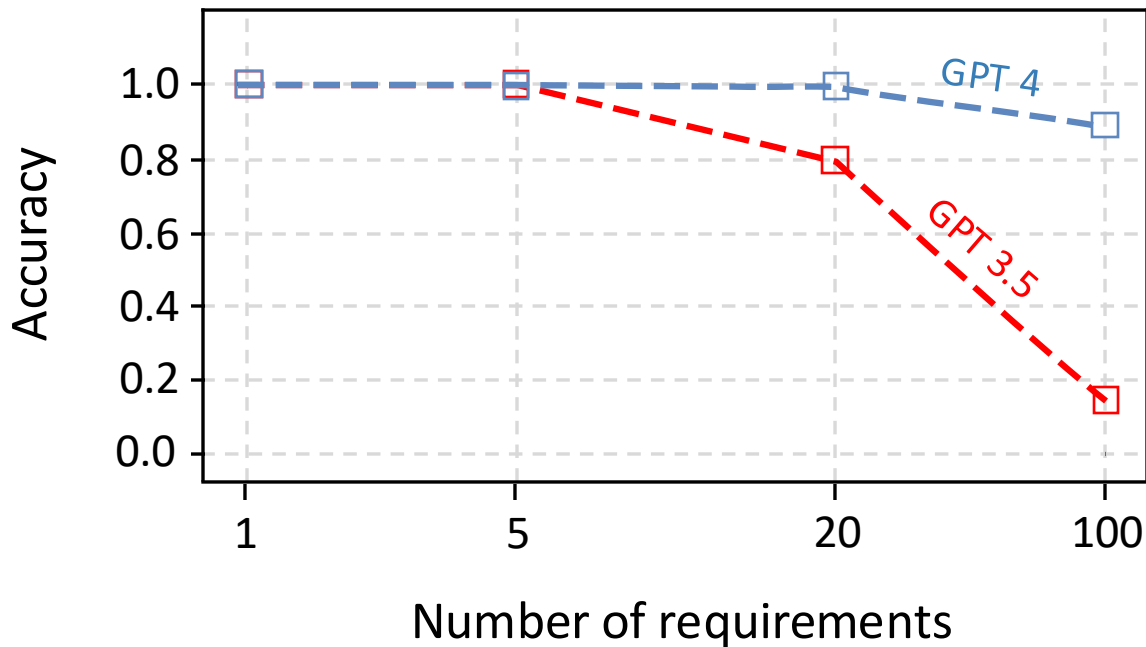
- 1000B parameters
- slow inferences
- resource intensive
- hard to deploy

## **Small/Medium** language models:

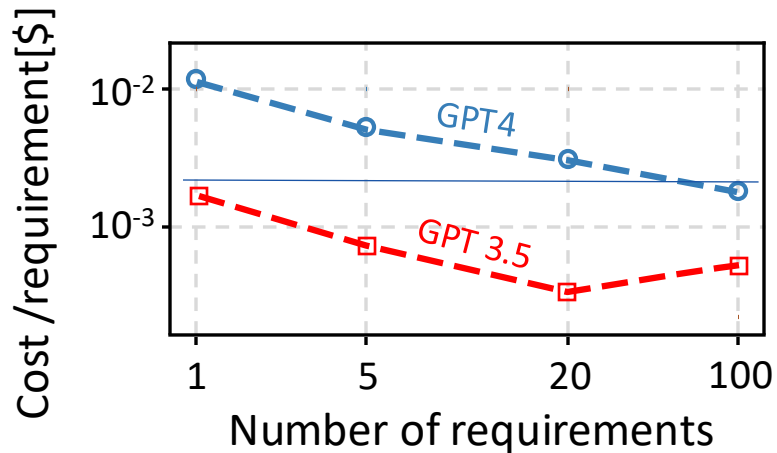
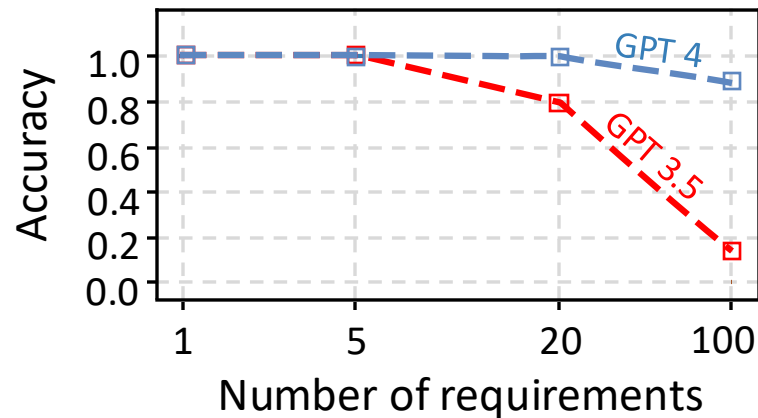
- 1B-100B parameters
- faster inferences
- deployable on a few GPUs
  - sometimes even on a laptop

How do **smaller** models perform?

Sadly, **smaller** models perform **worse**



Sadly, **smaller** models perform **worse**, yet **cost less**





# Can one **specialized** language models for **one task**?

## General-purpose models:

- trained on any text
- know **everything** (almost)
- but may fail in **something**

how many r's are in the phrase "network orchestration"?



The phrase "network orchestration" contains one 'r'.





# Can one **specialized** language models for **one task**?

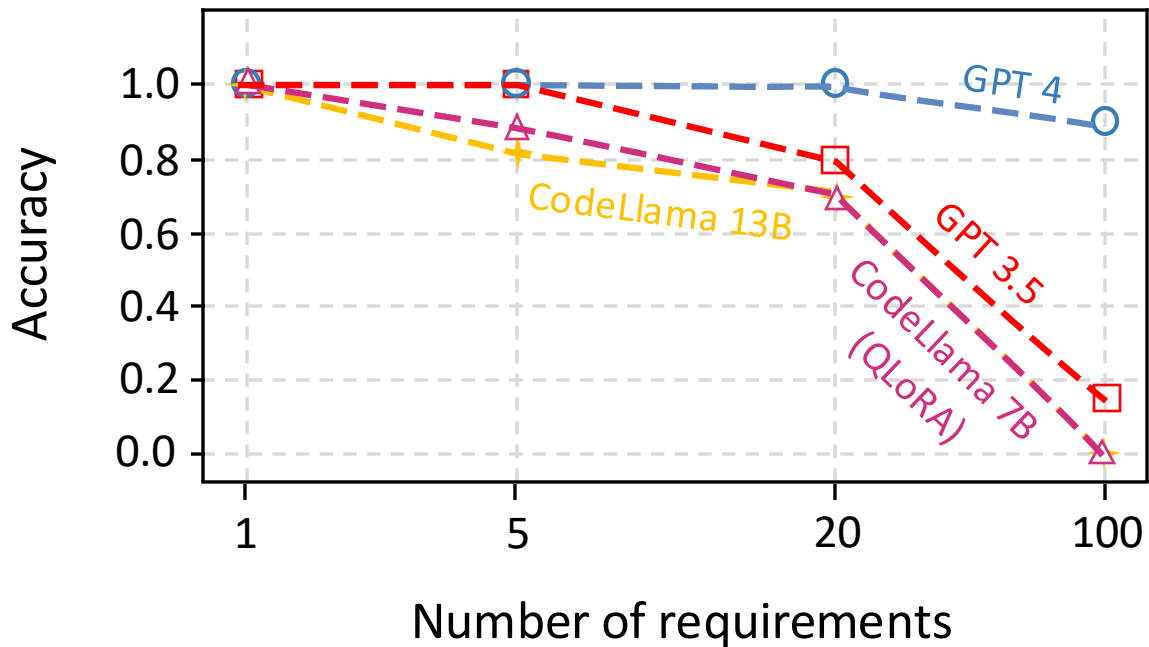
## **General-purpose** models:

- trained on any text
- know **everything** (almost)
- but may fail in **something**

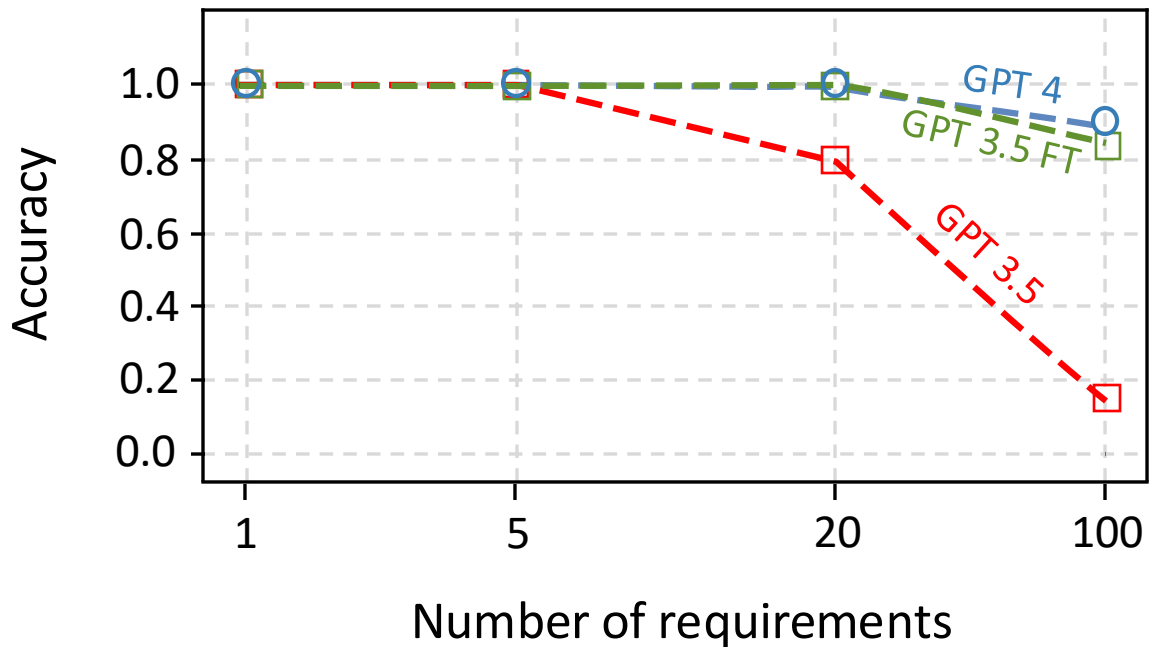
## **Specialized** models:

- pre-trained on specific tasks, or
- fine-tuned from general-purpose

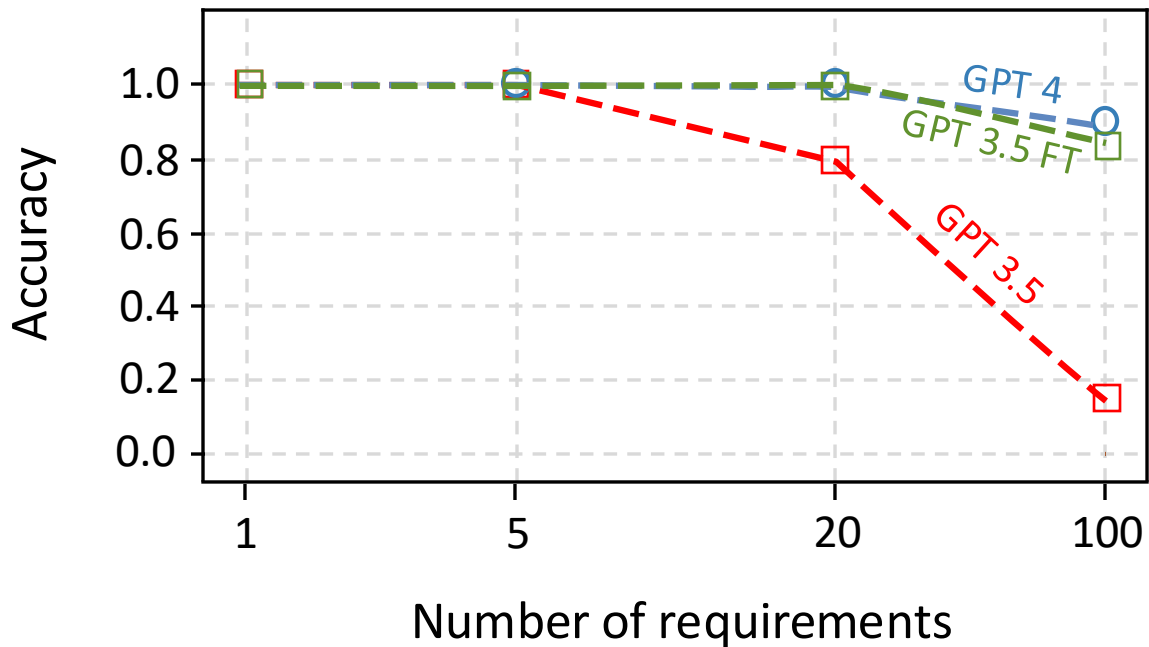
# Small Specialized models perform **poorly**, but **better**



# Larger Specialized models perform better

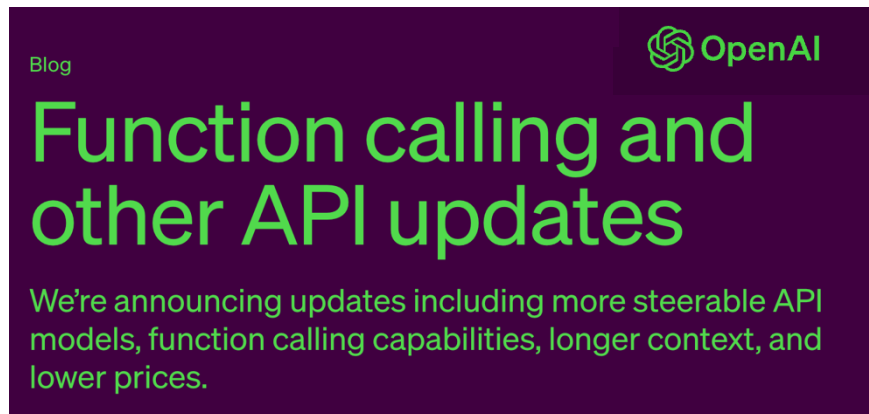


# Larger Specialized models perform better





# Can LLMs call API functions?



June 2023 (one function call) and November 2023 (parallel function calls)

# Translating high-level requirements to a formal specification format

*"traffic from **Rome** to **Milan**  
must **traverse** a **firewall**"*



network  
operator

LLM

dictionary

```
{ "reachability": {  
  "rome": [ "milan" ] },  
  "waypoint": {  
    [ "rome", "milan" ] :  
      [ "fw1", "fw2" ] },  
  "avoidance": { }  
}
```

function calling

```
add_reachability("rome", "milan");  
add_waypoint("rome", "milan", [ "fw1", "fw2" ] );
```

# Which one would be best?

## dictionary

```
{ "reachability": {  
    "rome": ["milan"] },  
  "waypoint": {  
    ["rome", "milan"]:  
      ["fw1", "fw2"] },  
  "avoidance": {}  
}
```

- + compact
- rearranging items

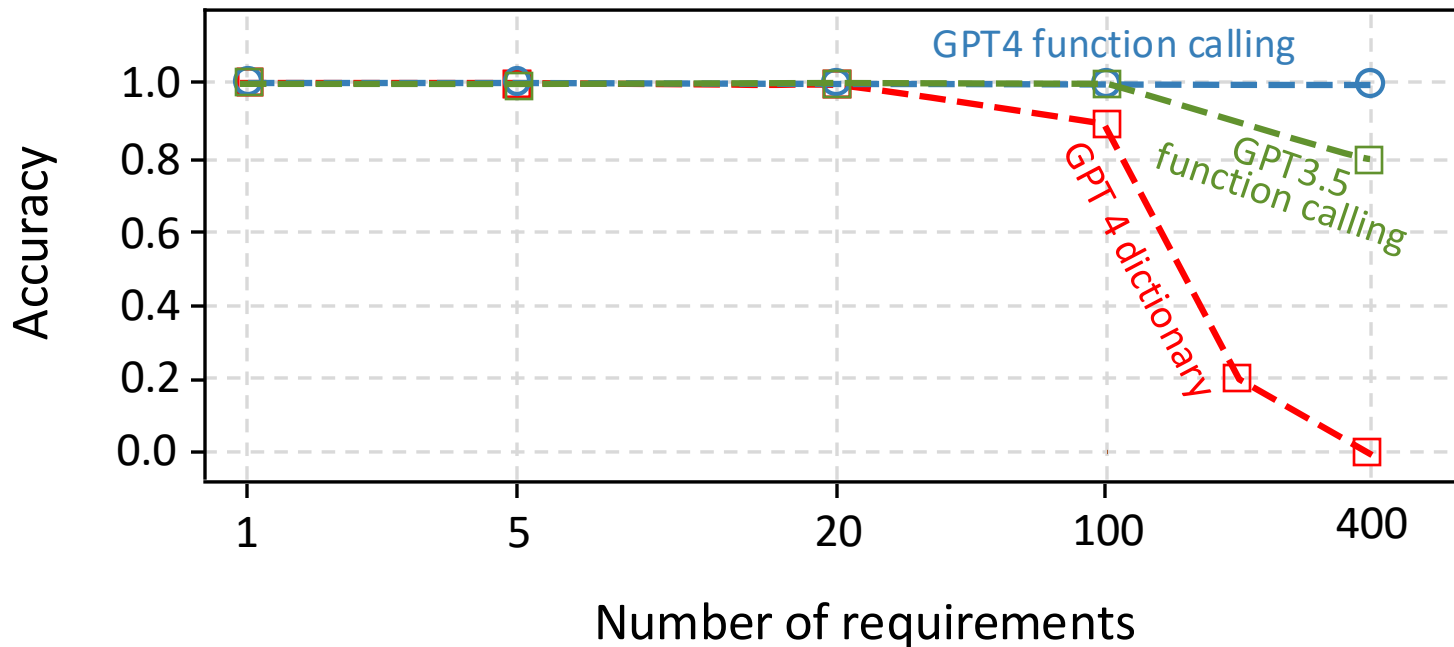
## function calling

```
add_reachability("rome", "milan");  
add_waypoint("rome", "milan", ["fw1", "fw2"]);
```

- less compact
- + no re-ordering

# Function calling versus **dictionary** data structure

LLMs are good at 1:1 translations







# We'll focus on three tasks in orchestrating networks

1. Translating high-level requirements to a formal specification format
2. **Adapting code to new requirements**
3. Generating low-level configurations



# Adapting code to new **requirements**. Why?

Developing modern software is **difficult**

- fast-paced due to **rapid** technological changes
- higher performance, resilience, and security guarantees

Developing modern software is **expensive**

- hire developers with a deep understanding of **numerous** systems, protocols, etc.
- development process becomes **time-consuming**, error-prone, and cumbersome

# Adapting code to new requirements

*“Create a **function** that takes as input [...] and produces **waypoint paths** as output”*



network  
operator

LLM

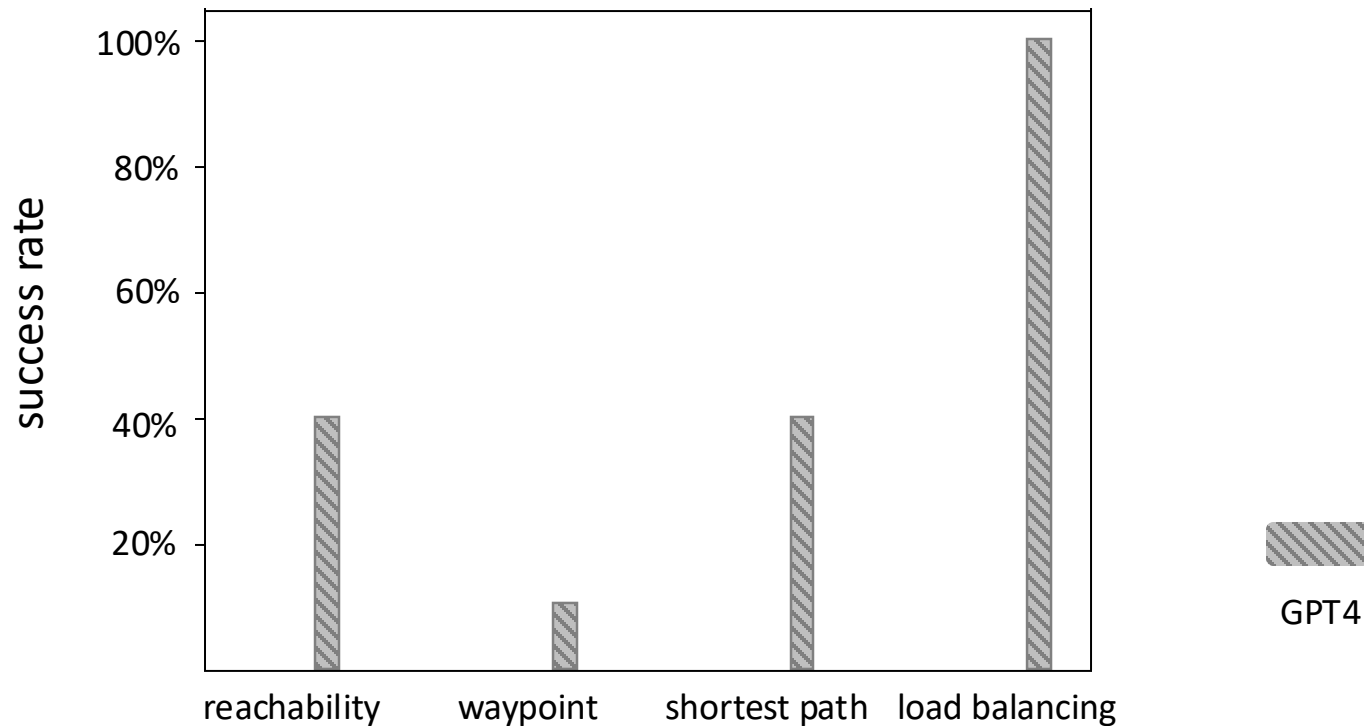
code generation

```
void Dijkstra(int source, const vector<vector<pair<int, int>>> &graph, int n) {
    int n = graph.size();
    dist.assign(n, INF);
    set<pair<int, int>> active_vertices;
    dist[source] = 0;
    active_vertices.insert({0, source});

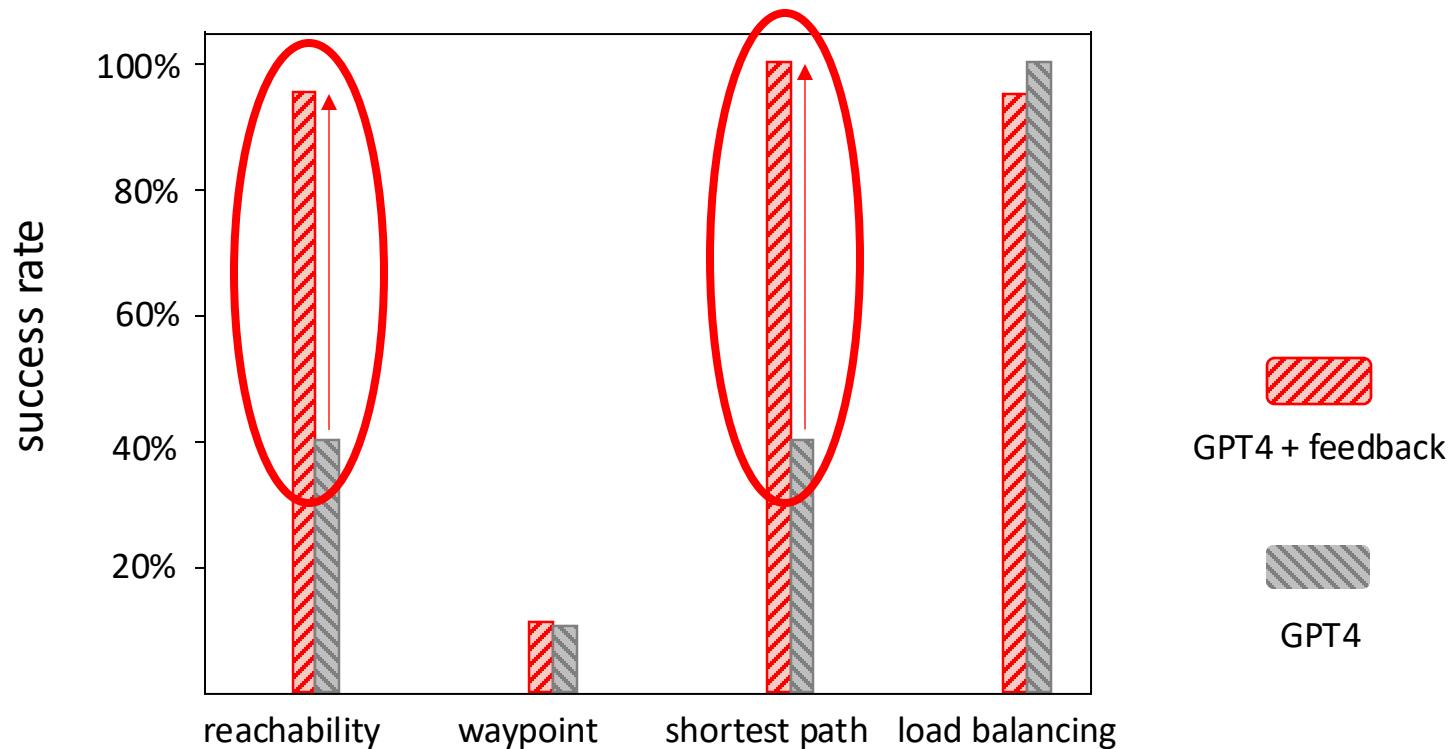
    while (!active_vertices.empty()) {
        int vertex = active_vertices.begin()->second;
        active_vertices.erase(active_vertices.begin());

        for (auto edge : graph[vertex]) {
            int neighbor = edge.first;
            int weight = edge.second;
```

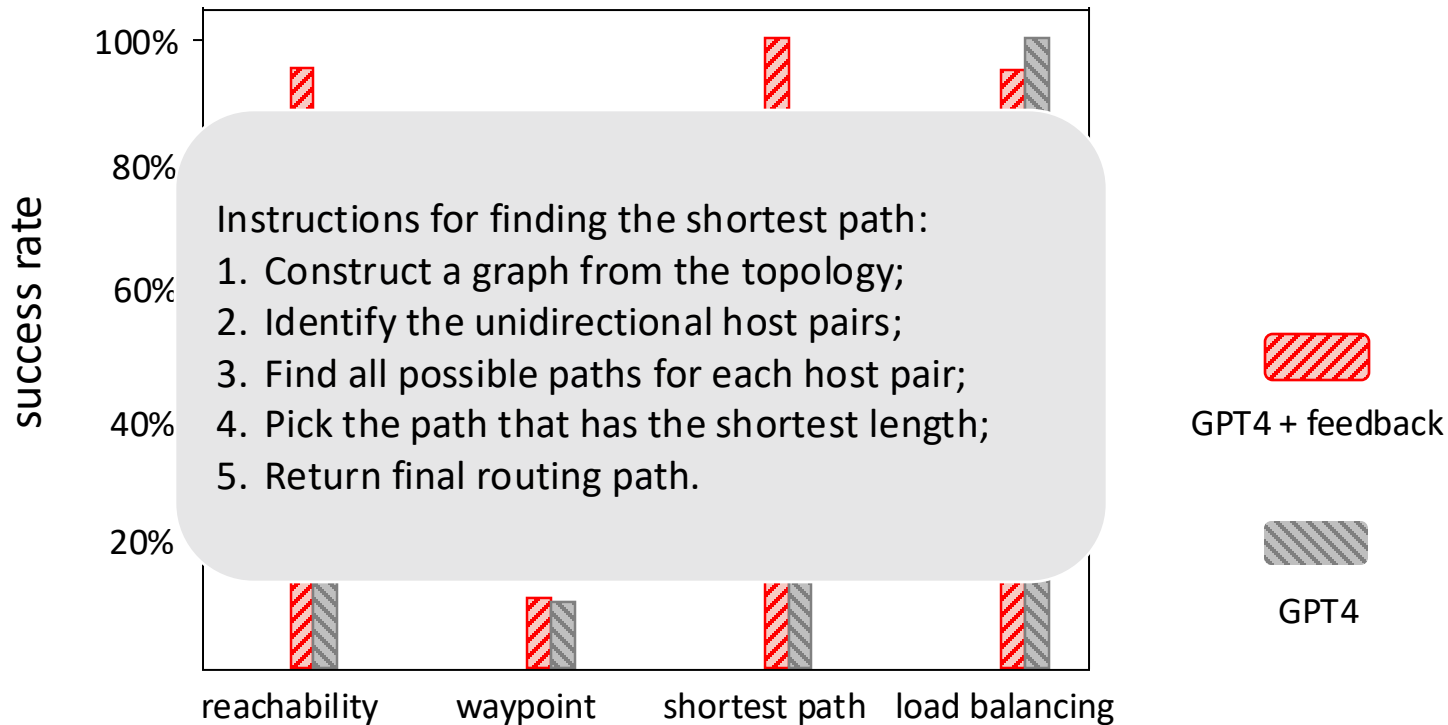
## Poor performance even for simple tasks



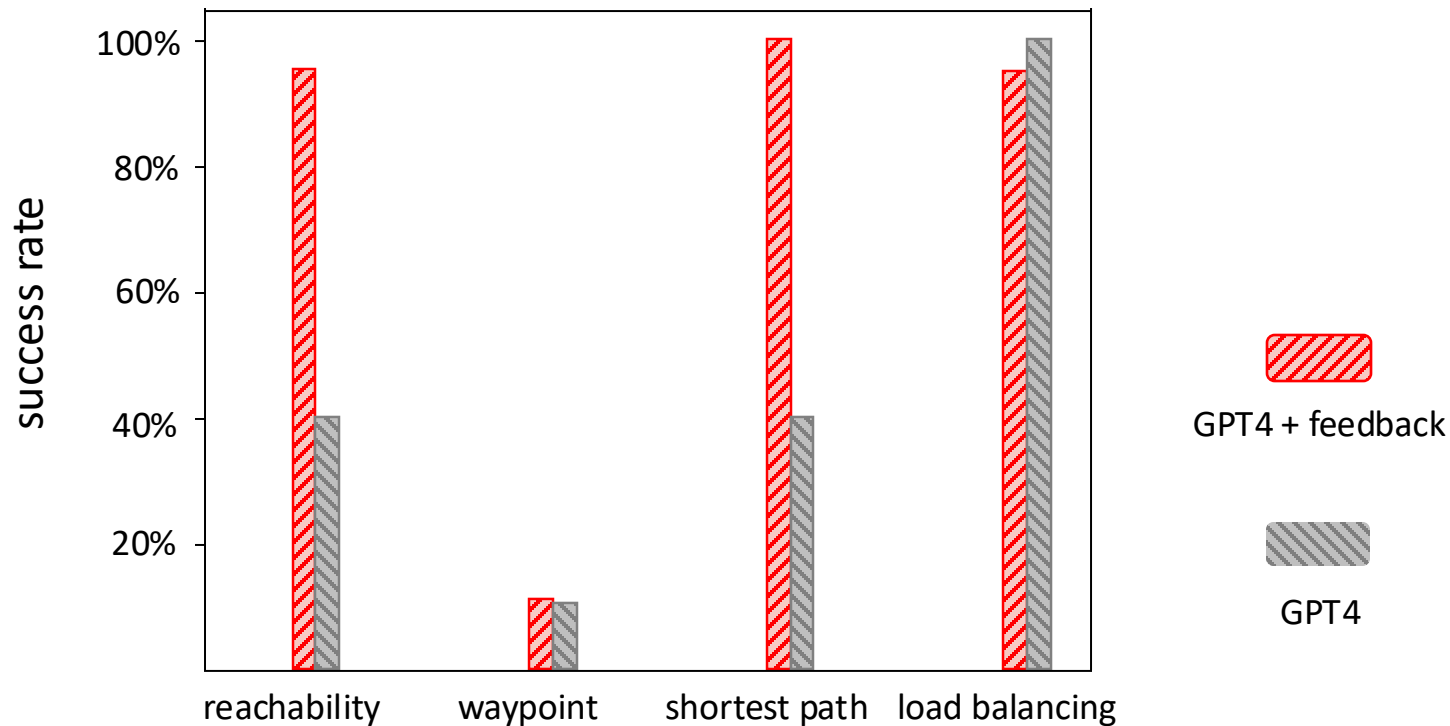
# What if we give **feedback**?



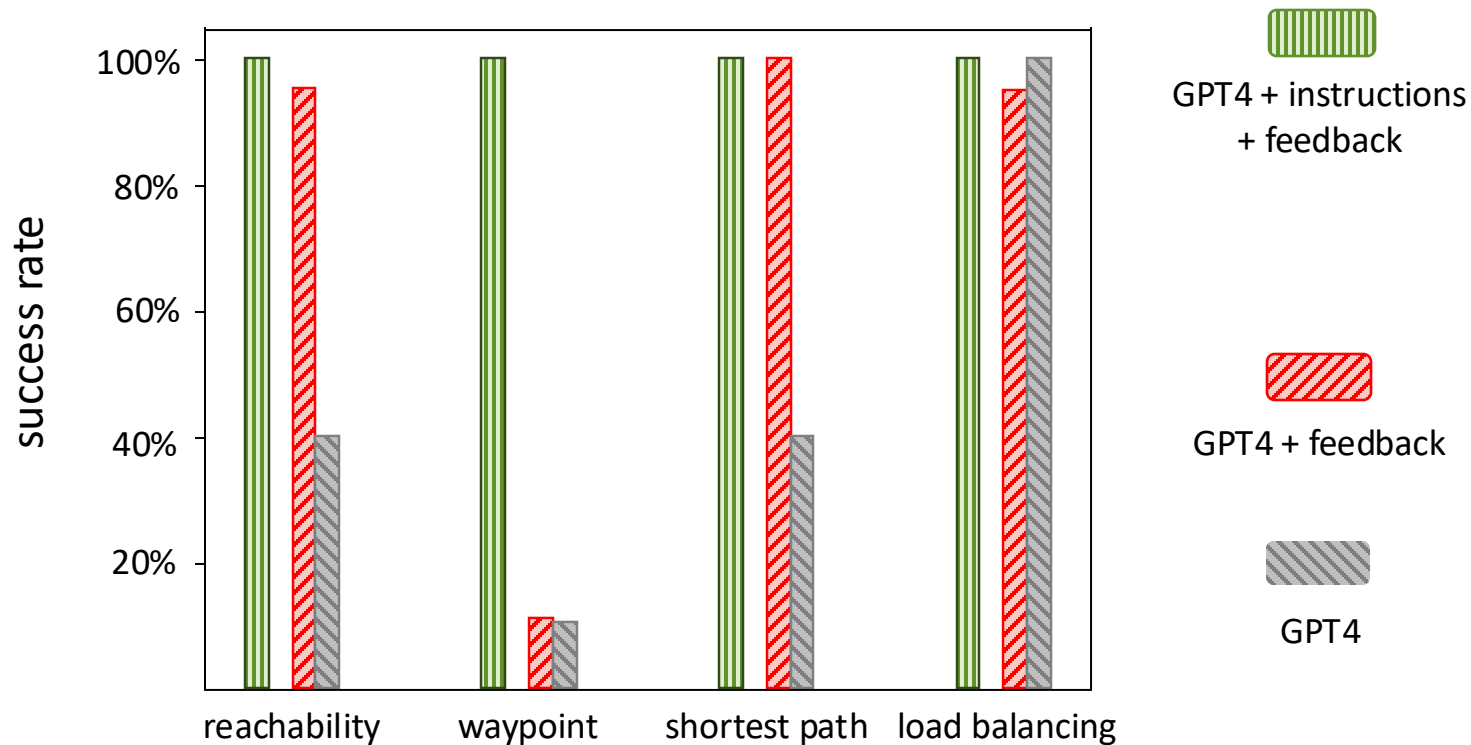
# What if we also provide some **algorithmic help**?



# What if we also provide some **algorithmic help**?

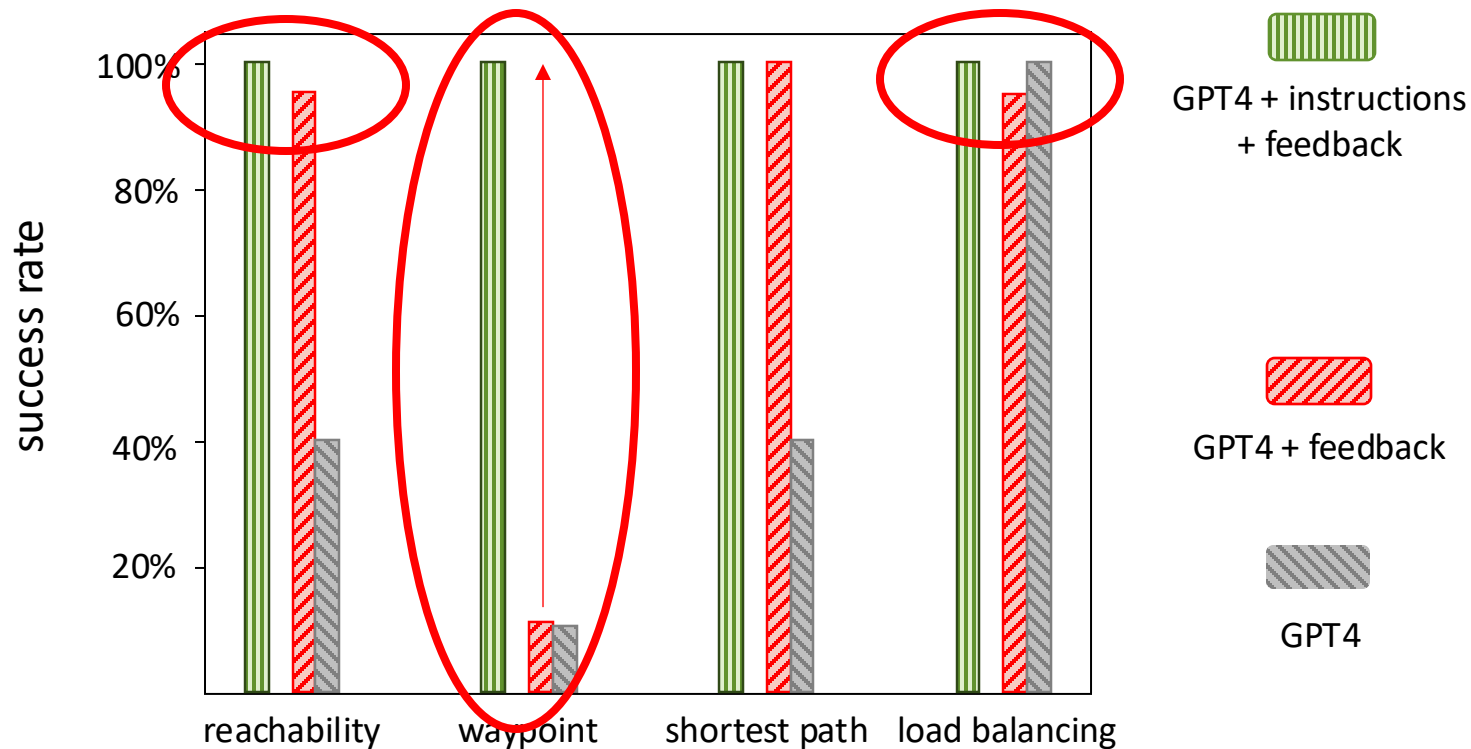


# What if we also provide some **algorithmic help**?

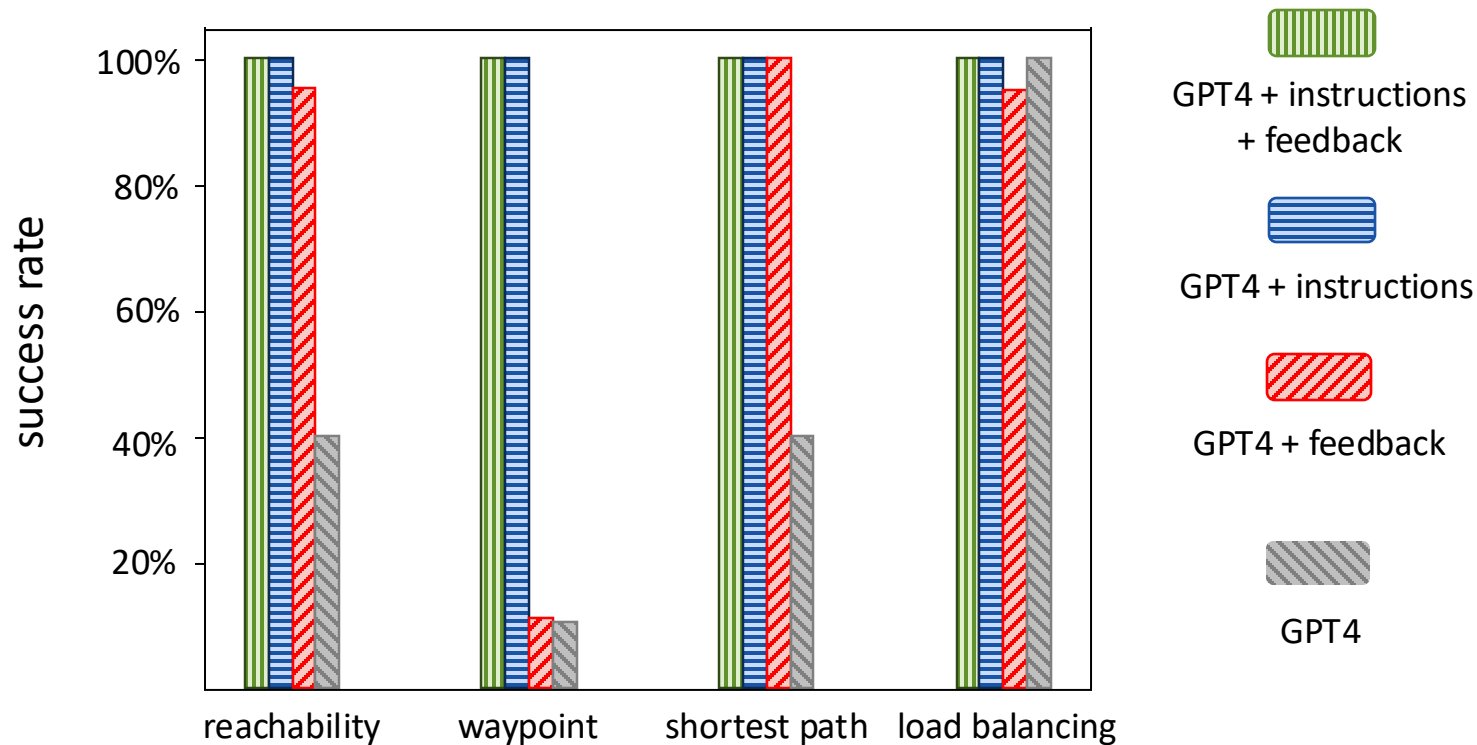




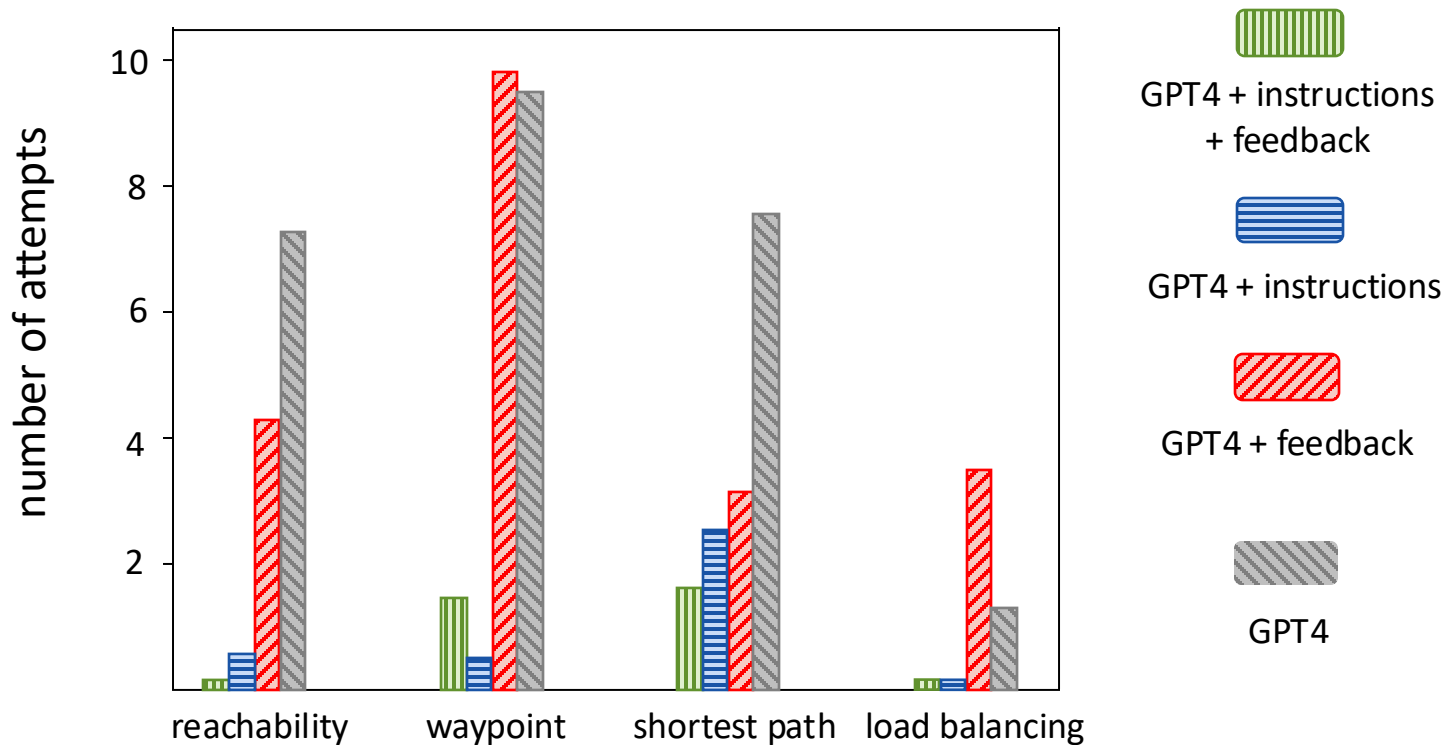
# What if we also provide some **algorithmic help**?



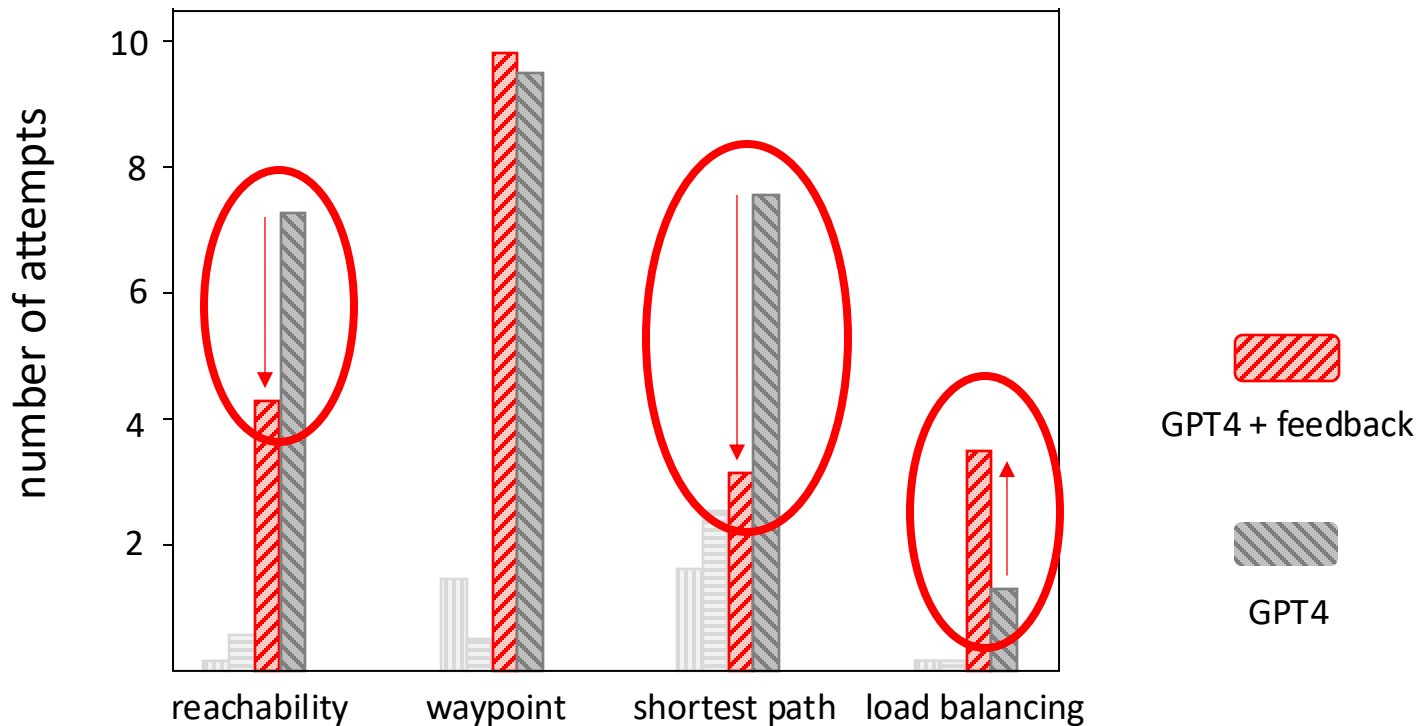
# What if we provide **algorithmic help without feedback**?



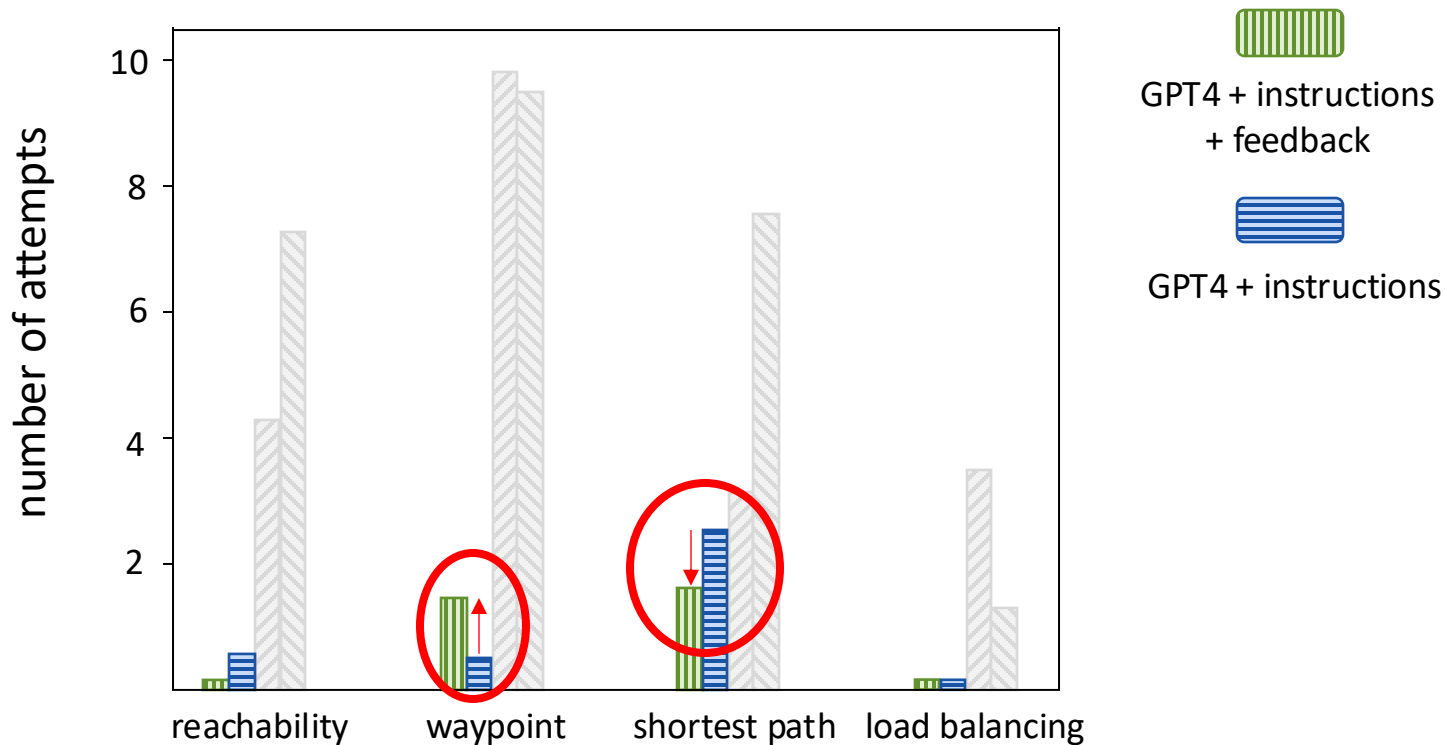
# Does providing precise feedback **always** help?



# Does providing precise feedback **always** help?



# Does providing precise feedback **always** help?





# Smaller models could **not** produce **meaningful** code

We tested a few additional models:

- **phy** (specialized in Python)
- **mistral**
- **codellama** 7B, 13B, 34B (with 4-bit quantization)
- **GPT 3.5**

None of these models generated correct code

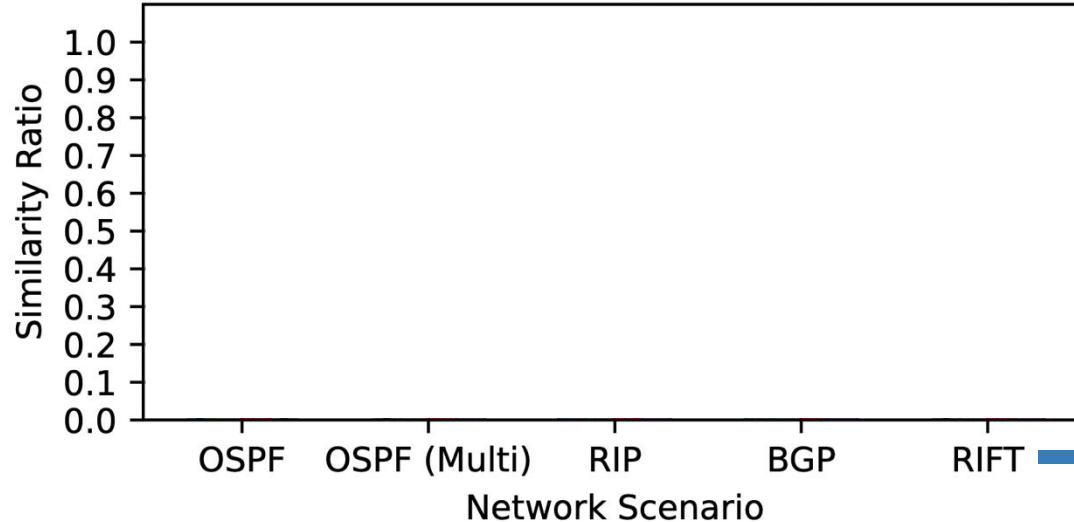
- from basic syntax errors to wrong semantic of data structures, logic, ...



# We'll focus on three tasks in orchestrating networks

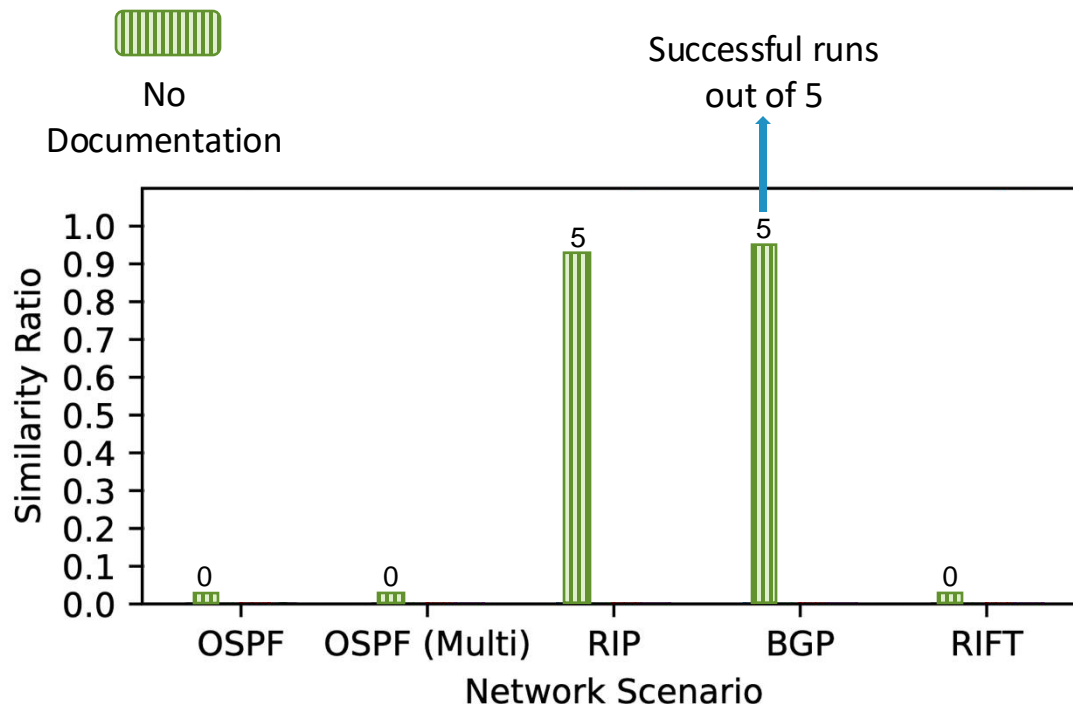
1. Translating high-level requirements to a formal specification format
2. Adapting code to new requirements
3. **Generating low-level configurations**

# Generating low-level configurations





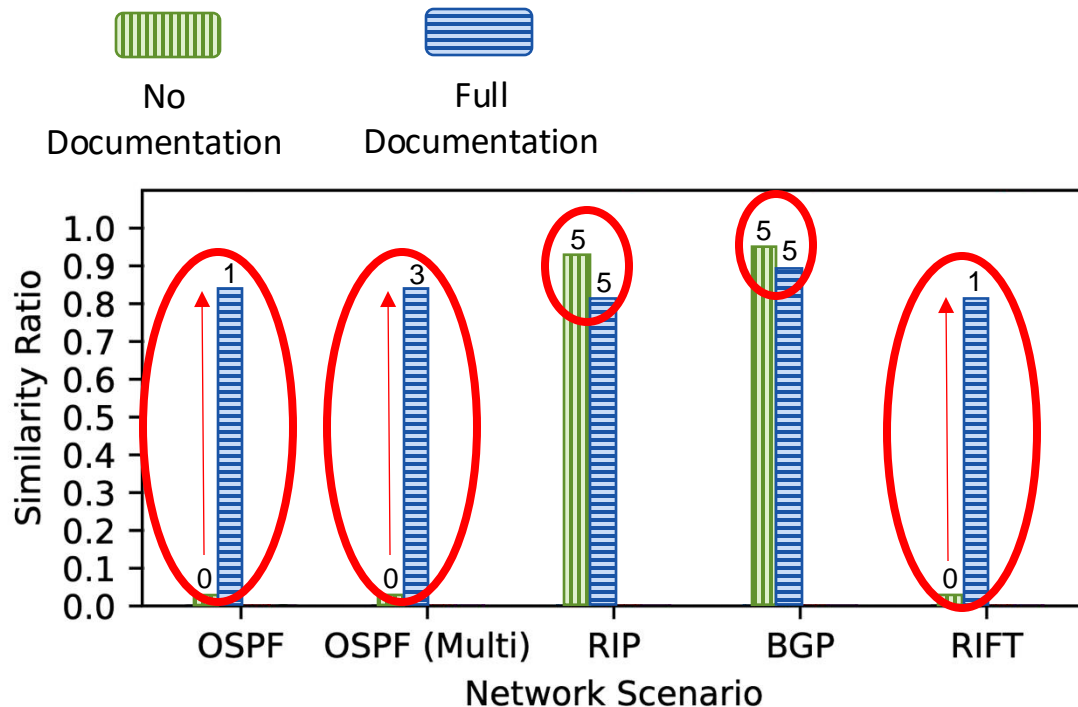
# Generating low-level configurations



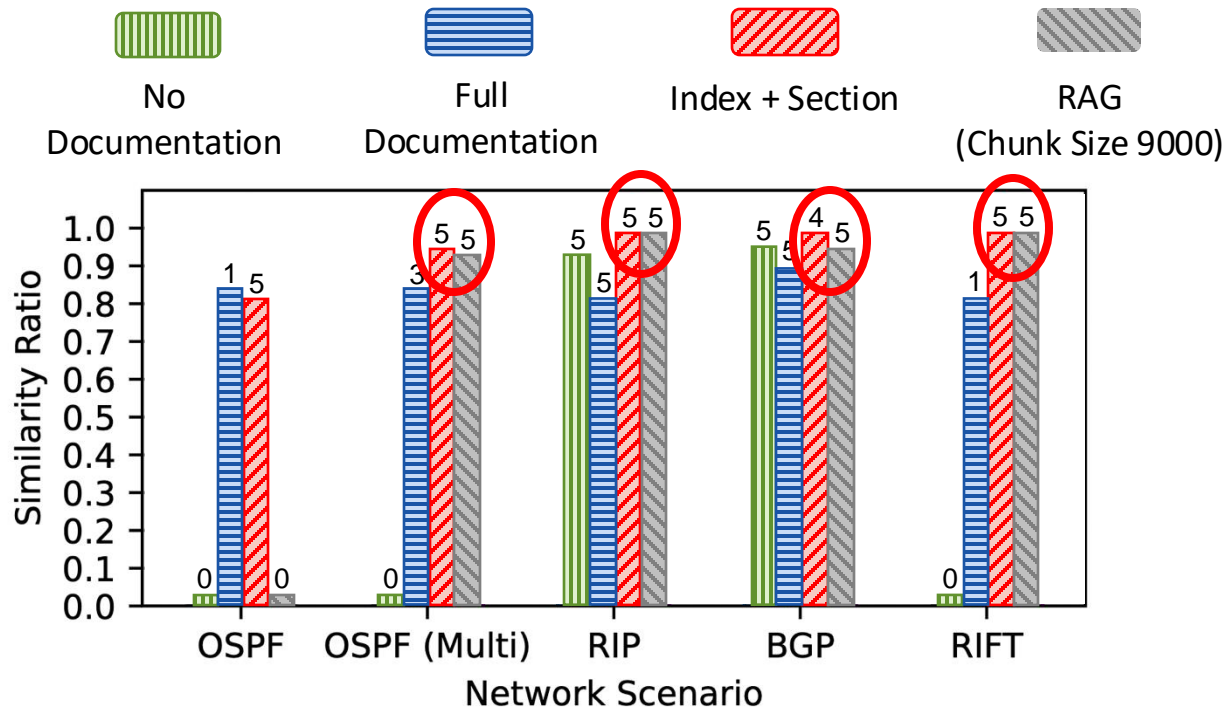
LLMs **know** something already!



# Generating low-level configurations



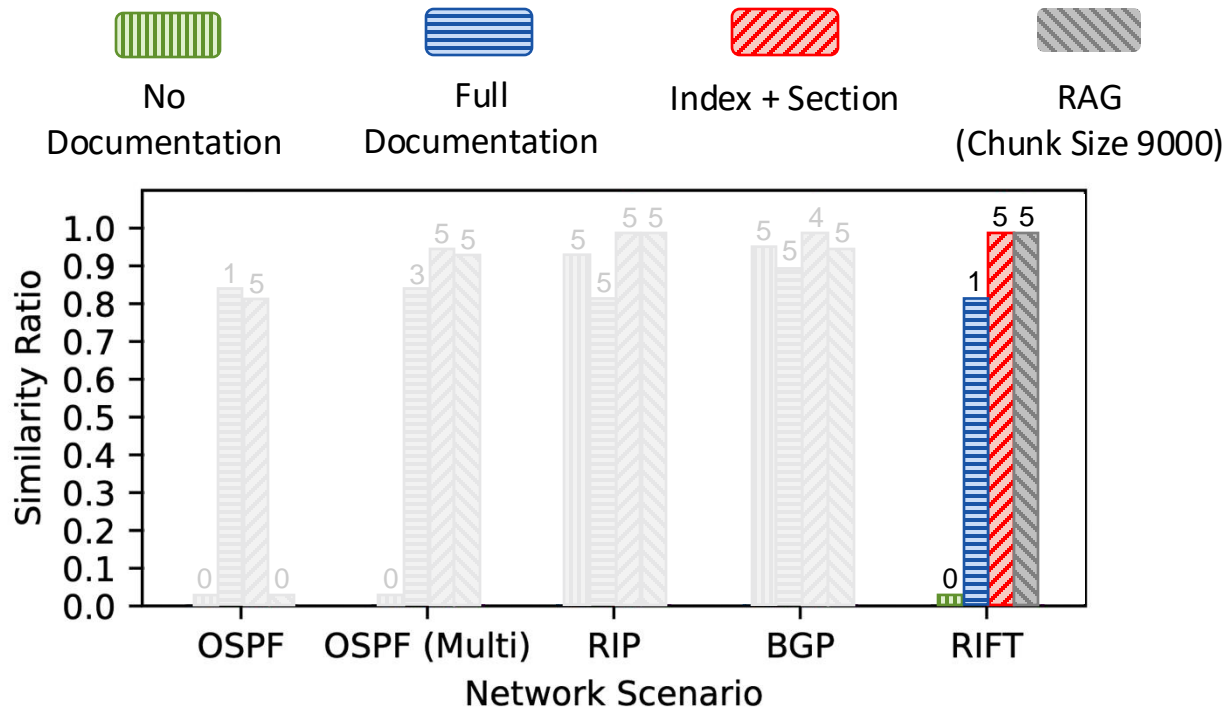
# Generating low-level configurations



**Index + Section and RAG reduce the context size.**



# Generating low-level configurations



FRROUTING

**LLMs can take advantage of knowledge without fine-tuning.**



# Building LLM-based system for networks

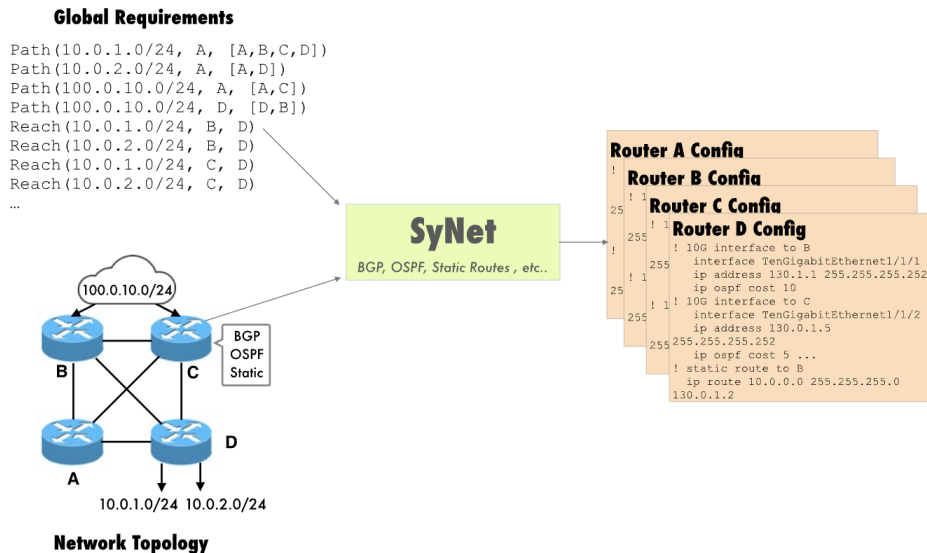
1. Split complex tasks into smaller subtasks
2. Support task-specific verifiers
3. Keep humans still in the loop



# Prototypes

1. LLMs in action with network synthesizers
2. LLMs from intents to low-level configuration

# LLMs in action with network synthesizers



# LLMs in action with network synthesizers

network operator



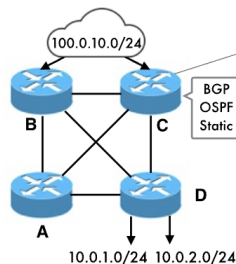
*“Traffic classified as **media** should be routed from Rome to Milan exclusively using the **OSPF** protocol”*



Does GPT know how to write SyNet code?  
**Not at all!**

## Global Requirements

```
Path(10.0.1.0/24, A, [A,B,C,D])
Path(10.0.2.0/24, A, [A,D])
Path(100.0.10.0/24, A, [A,C])
Path(100.0.10.0/24, D, [D,B])
Reach(10.0.1.0/24, B, D)
Reach(10.0.2.0/24, B, D)
Reach(10.0.1.0/24, C, D)
Reach(10.0.2.0/24, C, D)
...
```



Network Topology

## SyNet

BGP, OSPF, Static Routes, etc..

### Router A Config

### Router B Config

### Router C Config

### Router D Config

```
!
!
! 10G interface to B
! interface TenGigabitEthernet1/1/1
! ip address 130.1.1 255.255.255.252
! ip ospf cost 10
!
! 10G interface to C
! interface TenGigabitEthernet1/1/2
! ip address 130.0.1.5
! ip ospf cost 5 ...
! static route to B
! ip route 10.0.0.0 255.255.255.0
! 130.0.1.2
```



# LLMs in action with network synthesizers

network operator



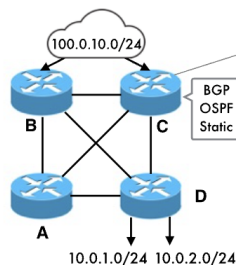
*"Traffic classified as **media** should be routed from Rome to Milan exclusively using the **OSPF** protocol"*



SyNET  
paper

Global Requirements

Fwd(media, rome, milan, ospf)



Network Topology

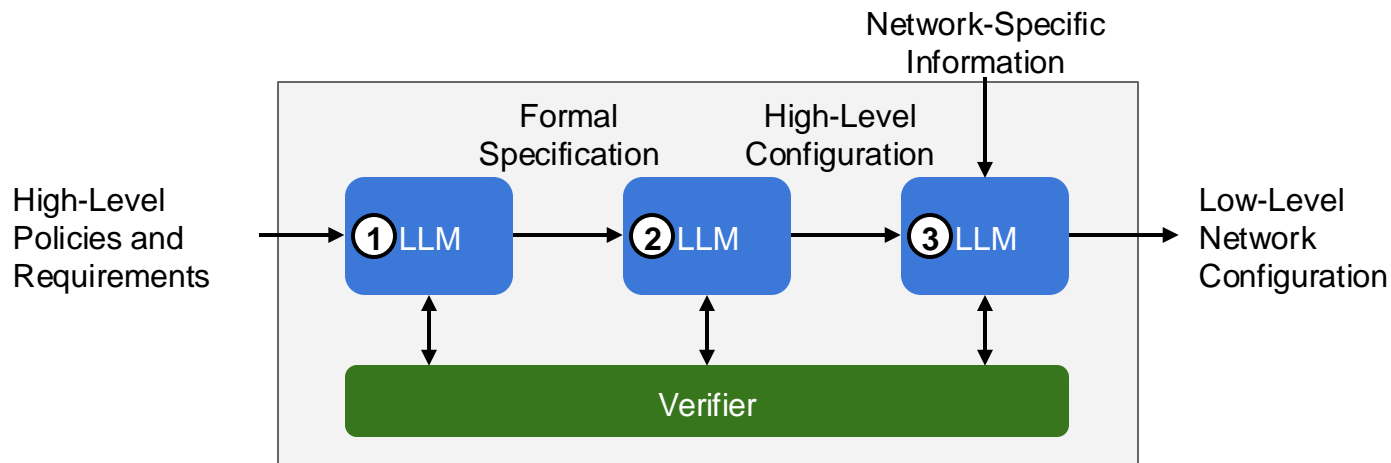
**SyNet**

BGP, OSPF, Static Routes, etc..

```

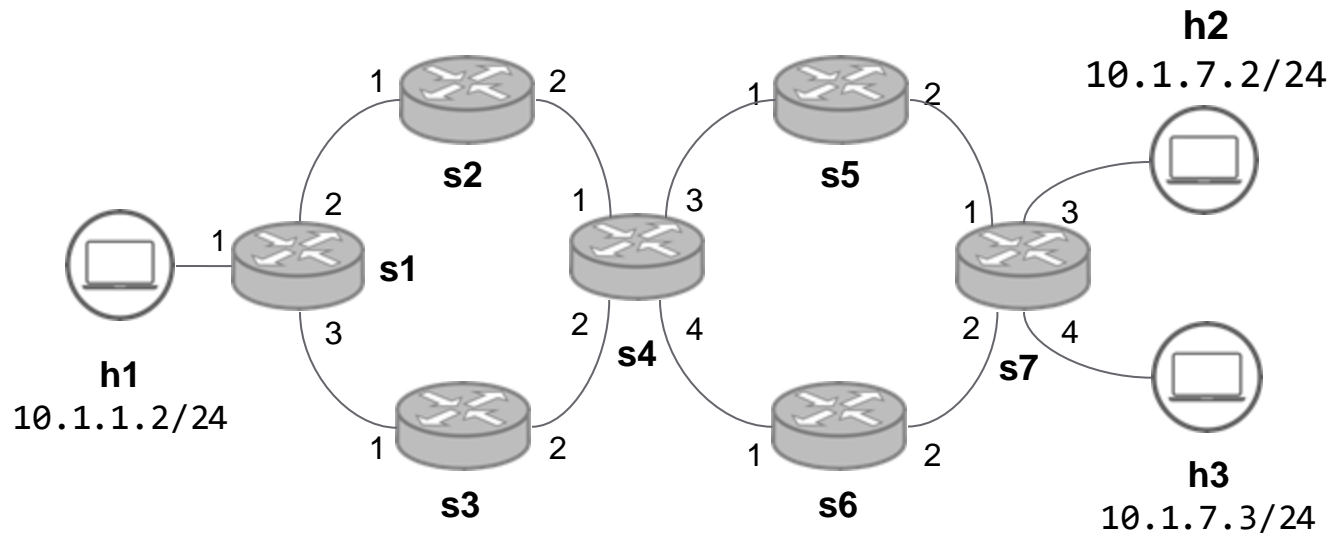
Router A Config
Router B Config
Router C Config
Router D Config
!
! 10G interface to B
! interface TenGigabitEthernet1/1/1
! ip address 130.1.1 255.255.255.252
! ip ospf cost 10
!
! 10G interface to C
! interface TenGigabitEthernet1/1/2
! ip address 130.0.1.5
! ip ospf cost 5 ...
! static route to B
! ip route 10.0.0.0 255.255.255.0
! 130.0.1.2
    
```

# Netbuddy: LLMs from intents to configuration



# Evaluated Topology

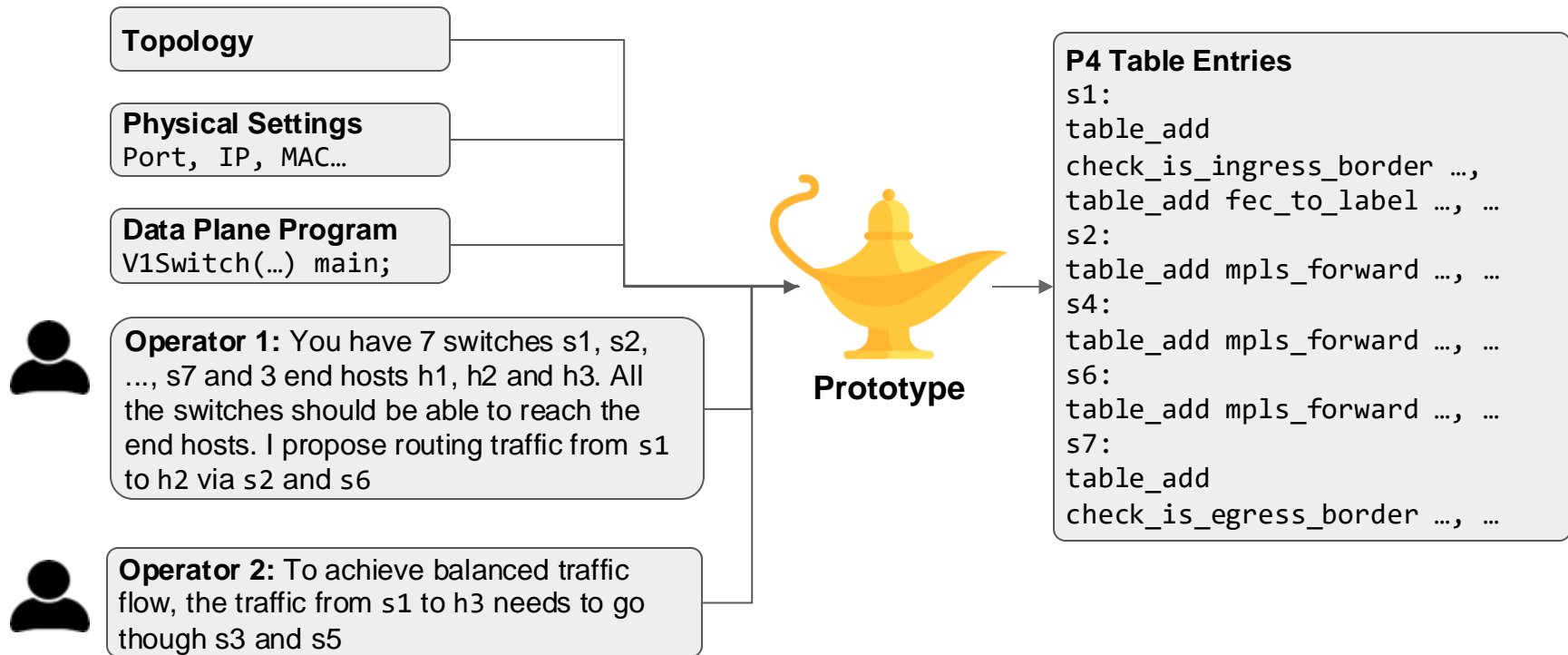
prototype  
video



Emulated using  Kathará

# From requirements to P4 code

prototype  
video



# Controlling the network

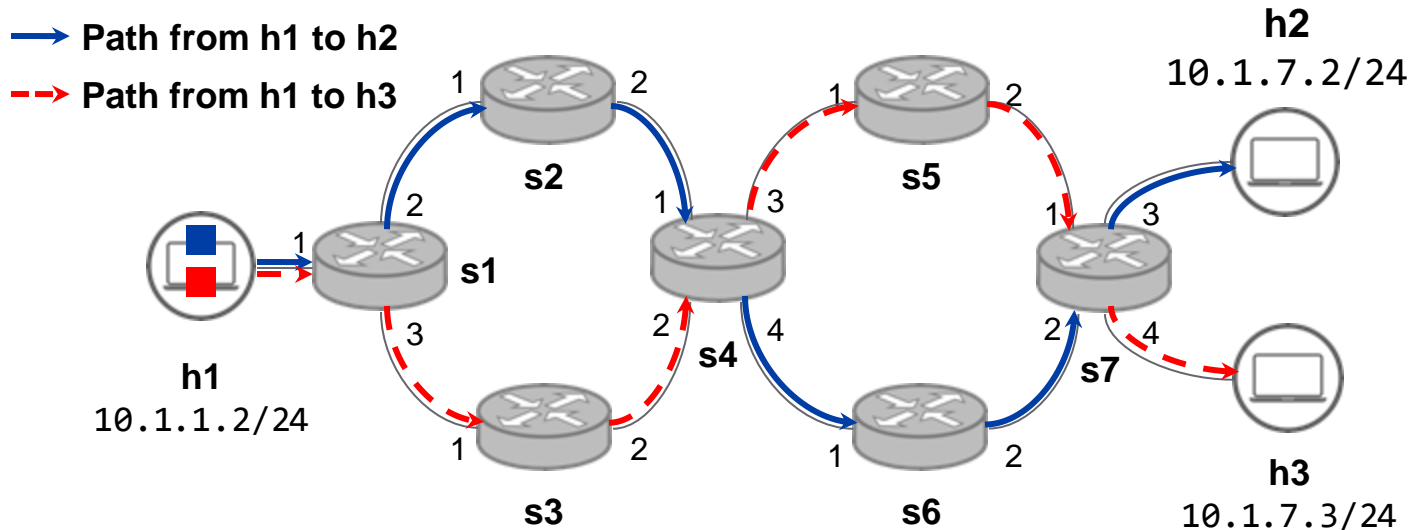
prototype  
video



**Operator 1:** You have 7 switches s1, s2, ..., s7 and 3 end hosts h1, h2 and h3. All the switches should be able to reach the end hosts. I propose routing traffic from s1 to h2 via s2 and s6



**Operator 2:** To achieve balanced traffic flow, the traffic from s1 to h3 needs to go through s3 and s5



Activities

Test Editor

373%

NetBuddy-design.jpg

require.txt

Open

Save

Ln 1, Col 20

INS

Operator 1: You have 7 switches s1, s2, ..., s7 and 3 end hosts h1, h2 and h3. All the switches should be able to reach the end hosts. I propose routing traffic from s1 to h2 via s2 and s3.

Operator 2: To achieve balanced traffic flow, the traffic from s1 to h3 needs to go through s3 and s5.

Terminal

Local

python3 NetBuddy\_once.py

h1

10.1.1.1.2/24

s1

s2

s3

s4



# More application of LLMs

Applications	Related work
Network Configuration and Management	LLMs to Synthesize Router Configurations? (HotNets '23) Network Management Using Code Generated by LLMs (HotNets '23)
Protocol Specification Extraction	Extracting Protocol Specifications Using LLMs (HotNets '23)
More Specific Tasks	NetLLM: Adapting LLMs for Networking (SIGCOMM '24)
Network Research	Reproducing Network Research Results Using LLMs (HotNets '23)
Network for LLMs	vLLM (SIGOPS '23), Orca (OSDI '22), FlexGen (ICML '23), ...

# Reliable LLM applications

Write a Python function to calculate the maximum number in a set.



Reliability

Test units

```
def find_max_in_set(numbers_set):  
    """  
    Find the maximum number in a set of numbers.  
  
    Parameters:  
    numbers_set (set): A set of numbers.  
  
    Returns:  
    int or float: The maximum number in the set.  
    """  
    # Check if the set is empty  
    if not numbers_set:  
        raise ValueError("The set is empty, cannot determine the maximum.")  
  
    # Initialize the maximum number  
    max_number = float('-inf')  
  
    # Iterate through each number in the set  
    for number in numbers_set:  
        # Update max_number if the current number is greater  
        if number > max_number:  
            max_number = number  
  
    return max_number  
  
# Example usage  
example_set = {1, 3, 5, 7, 9}  
print("The maximum number in the set is:", find_max_in_set(example_set))
```

Can we trust?



# Reliable LLM applications

Write a Python function to calculate the maximum number in a set.



Reliability

vulnerability

```
def find_max_in_set(numbers_set):  
    """  
    Find the maximum number in a set of numbers.  
  
    Parameters:  
    numbers_set (set): A set of numbers.  
  
    Returns:  
    int or float: The maximum number in the set.  
    """  
    # Check if the set is empty  
    if not numbers_set:  
        raise ValueError("The set is empty, cannot determine the maximum.")  
  
    # Initialize the maximum number  
    max_number = float('-inf')  
  
    # Iterate through each number in the set  
    for number in numbers_set:  
        # Update max_number if the current number is greater  
        if number > max_number:  
            max_number = number  
  
    return max_number  
  
# Example usage  
example_set = 1  
print("The maximum number in the set is:", find_max_in_set(example_set))
```

Can we trust?

# Reliable LLM applications by formal verification

Write me a Dafny code to calculate.



Reliability



```
1 method findMax(arr: array<int>) returns (max: int)
2   requires arr.Length > 0 // The array must not be empty
3   ensures forall k :: 0 <= k < arr.Length ==> arr[k] <= max // Ensure the result is the maximum
4 {
5   max := arr[0]; // Initialize max with the first element of the array
6
7   // Loop through the array to find the maximum
8   var i := 1;
9   while i < arr.Length
10    ...
11    invariant 0 < i <= arr.Length // Loop invariant to help Dafny verify the loop
12    invariant forall j :: 0 <= j < i ==> arr[j] <= max // All elements before i are less than or equal to max
13    {
14      if arr[i] > max {
15        max := arr[i];
16      }
17      i := i + 1;
18    }
19 }
```

We can trust the code, as long as

- the **verified property** is **correct**, and
- **verification passes**



# Conclusions

We

- present the first **benchmark** for network orchestration and share experience
- automate common workflows with LLM-powered **prototypes**

**Opportunities:** LLMs can dramatically **simplify** and **automate** complex network orchestration tasks.

**Challenges:** unreliability, cost, ...

## Thank you!

