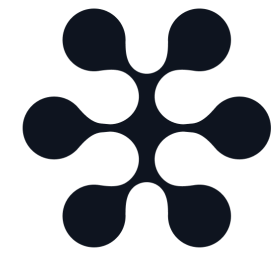




Massachusetts
Institute of
Technology



Subconscious

Beyond Context Limits

Subconscious Threads for Long-horizon Reasoning

Hongyin Luo, August 21

hongyin@subconscious.dev

Outline

- Motivation
- Working memory, thread, and TIM model
- TIMRUN inference engine
- Experiments
- Limitations

Motivation

Challenges around agent building

- Complexity in agent frameworks
- Redundant computation and repeated payments
- Decreasing throughput and performance

Motivation

Challenges around agent building

- Complexity in agent frameworks
- Redundant computation and repeated payments
- Decreasing throughput and performance

Language models have their context limits.

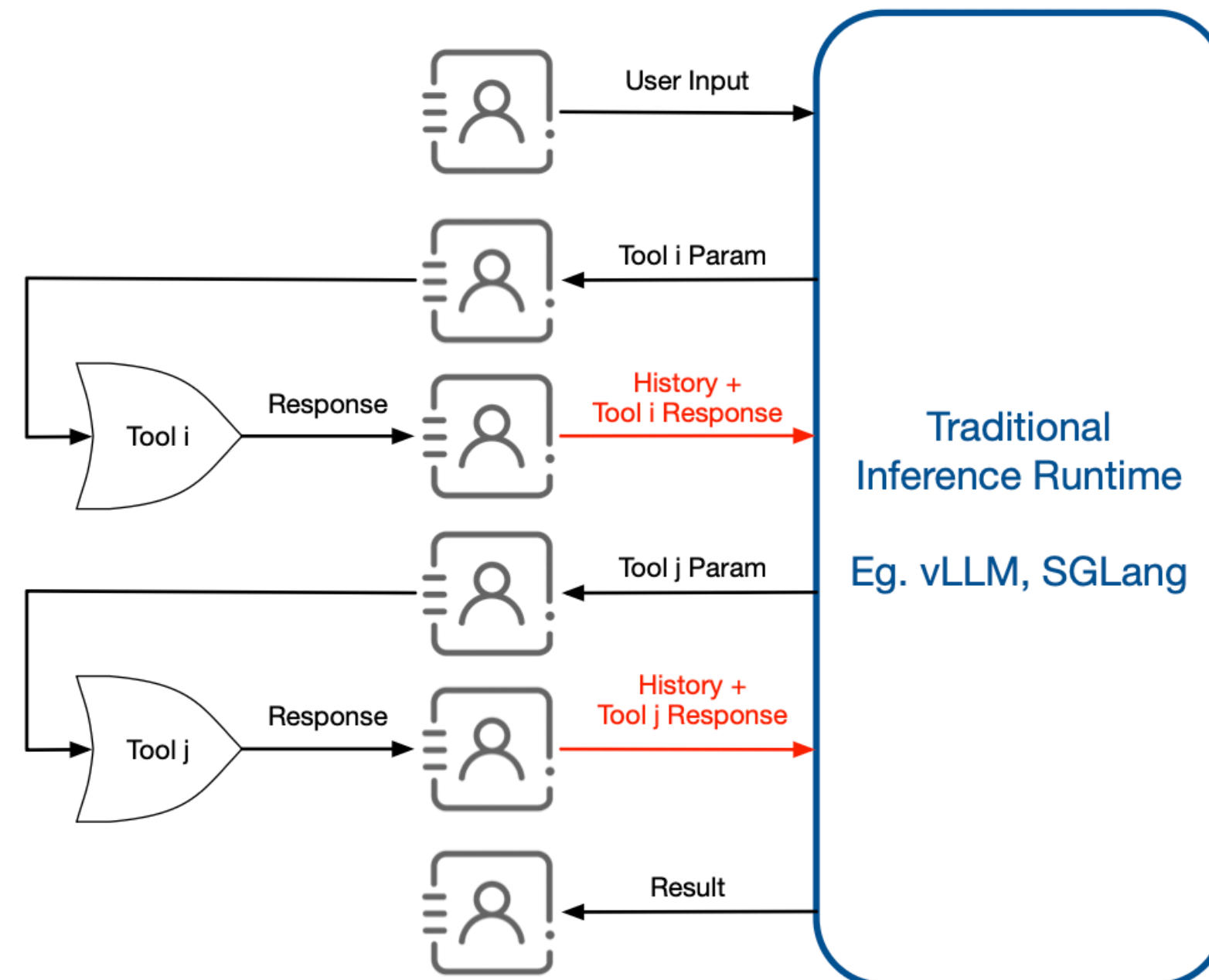
Engineering in chaining BERT (2019) -> LLMs (2024) together for complex tasks

“We tried a bunch of different multi-agent frameworks and landed on our own.”

Motivation

Challenges around agent building

- Complexity in agent frameworks
- Redundant computation and repeated payments
- Decreasing throughput and performance



Token Economics: Why Conversational Agents Don't Scale



Motivation

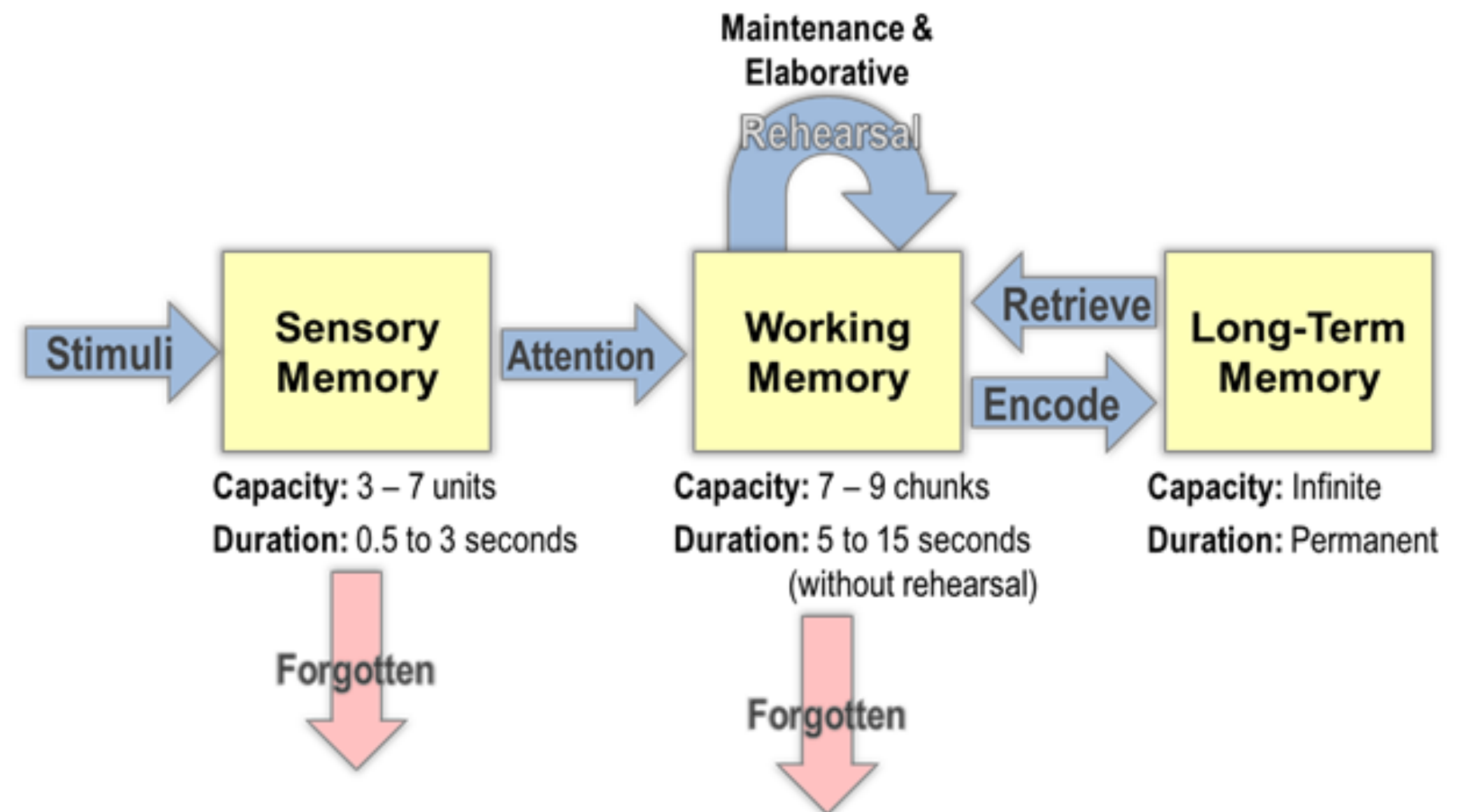
Challenges around agent building

- Complexity in agent frameworks
- Redundant computation and repeated payments
- Decreasing throughput and performance
 - Long context hallucination
 - Compute complexity

Working Memory

Basic concept behind Thread

- Subset of global memory
- Most careful computation only applied to working memory
- Interface with other modules



When we take a phone call, we don't remember how our muscles were moving to pick up the phone.

Background: Thread

Think longer and deeper

- Task-specific prompting (>4k tokens)
- No context sharing
- Information loss
- Complicated special token config
- Difficult to deploy



TIM LLM

Thread Inference Model

- Simple instruction
- Continuous state management
- Recursive **JSON** generation
- Straightforward structure and tool parameter extraction

	Thread	TIM
Structure Control	Special Tokens	JSON
State Management	No	Working Memory
Task General	No	Yes
Tool call	Developer	Server

JSON Generation via Constrained decoding

- Define a Pydantic class
- Send via “text_format”
- LLM returns valid JSON dictionaries follows the defined format.

Getting a structured response

python ↕

```
1 from openai import OpenAI
2 from pydantic import BaseModel
3
4 client = OpenAI()
5
6 class CalendarEvent(BaseModel):
7     name: str
8     date: str
9     participants: list[str]
10
11 response = client.responses.parse(
12     model="gpt-4o-2024-08-06",
13     input=[
14         {"role": "system", "content": "Extract the event information."},
15         {
16             "role": "user",
17             "content": "Alice and Bob are going to a science fair on Friday.",
18         },
19     ],
20     text_format=CalendarEvent,
21 )
22
23 event = response.output_parsed
```

Generating Recursive JSON

JSON mode handles this, too!

```
class CalendarEvent(BaseModel):  
    name: str  
    date: str  
    participants: List[str]  
    related_events: List['CalendarEvent']
```

```
CalendarEvent.model_rebuild()
```

```
{  
  "name": "AI Symposium 2025",  
  "date": "2025-09-15",  
  "participants": ["Dr. Alice Lee", "Prof. Mark Brown"],  
  "related_events": [  
    {  
      "name": "AI Workshop: Large Language Models",  
      "date": "2025-09-14",  
      "participants": ["Dr. Alice Lee", "Chris Wong"],  
      "related_events": [  
        {  
          "name": "LLM Paper Reading Group",  
          "date": "2025-09-13",  
          "participants": ["Chris Wong", "Samantha Lin"],  
          "related_events": []  
        }  
      ]  
    }  
  ],  
  {  
    "name": "AI Symposium Closing Ceremony",  
    "date": "2025-09-15",  
    "participants": ["Prof. Mark Brown", "Julia Yu"],  
    "related_events": []  
  }  
]  
}
```

Thread in JSON

More flexibility, less prompting

Recursion Hierarchy

```
class TimResponse(BaseModel):  
    reasoning: List[Task]  
    answer: str  
  
class Task(BaseModel):  
    thought: str  
    tooluse: Optional[ToolUse] = None  
    subtasks: Optional[List['Task']] = None  
    conclusion: str
```

Reasoning process

Thread in JSON

More flexibility, less prompting

Recursion Hierarchy

```
class TimResponse(BaseModel):  
    reasoning: List[Task]  
    answer: str  
  
class Task(BaseModel):  
    thought: str  
    tooluse: Optional[ToolUse] = None  
    subtasks: Optional[List['Task']] = None  
    conclusion: str
```

Step-level reasoning

Thread in JSON

More flexibility, less prompting

Recursion Hierarchy

```
class TimResponse(BaseModel):  
    reasoning: List[Task]  
    answer: str  
  
class Task(BaseModel):  
    thought: str  
    tooluse: Optional[ToolUse] = None  
    subtasks: Optional[List['Task']] = None  
    conclusion: str
```

→ Subconscious Threads

Thread in JSON

More flexibility, less prompting

Recursion Hierarchy

```
class TimResponse(BaseModel):  
    reasoning: List[Task]  
    answer: str  
  
class Task(BaseModel):  
    thought: str  
    tooluse: Optional[ToolUse] = None  
    subtasks: Optional[List['Task']] = None  
    conclusion: str
```

→ Tool use

Thread in JSON

More flexibility, less prompting

Recursion Hierarchy

```
class TimResponse(BaseModel):  
    reasoning: List[Task]  
    answer: str  
  
class Task(BaseModel):  
    thought: str  
    tooluse: Optional[ToolUse] = None  
    subtasks: Optional[List['Task']] = None  
    conclusion: str
```

Tool Use Schema

```
class ToolUse(BaseModel):  
    tool_name: Literal[ "SearchTool", "WebReaderTool" ]  
    parameters: Union[ SearchTool, WebReaderTool ]  
    tool_result: dict  
  
class SearchTool(BaseModel):  
    query: str  
  
class WebReaderTool(BaseModel):  
    goal: str  
    url: str
```

Thread in JSON

More flexibility, less prompting

Tool Use Schema

Tools are called
right after this
field is generated

```
class ToolUse(BaseModel):  
    tool_name: Literal[ "SearchTool", "WebReaderTool" ]  
    parameters: Union[ SearchTool, WebReaderTool ]  
    tool_result: dict  
  
class SearchTool(BaseModel):  
    query: str  
  
class WebReaderTool(BaseModel):  
    goal: str  
    url: str
```


Thread in JSON

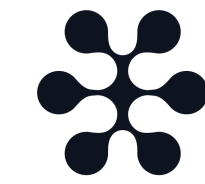
More flexibility, less prompting

JSON grammar engine
accepts any valid tool
responses as dictionaries

Tool Use Schema

```
class ToolUse(BaseModel):  
    tool_name: Literal[ "SearchTool", "WebReaderTool" ]  
    parameters: Union[ SearchTool, WebReaderTool ]  
    tool_result: dict  
  
class SearchTool(BaseModel):  
    query: str  
  
class WebReaderTool(BaseModel):  
    goal: str  
    url: str
```





TIM Outputs

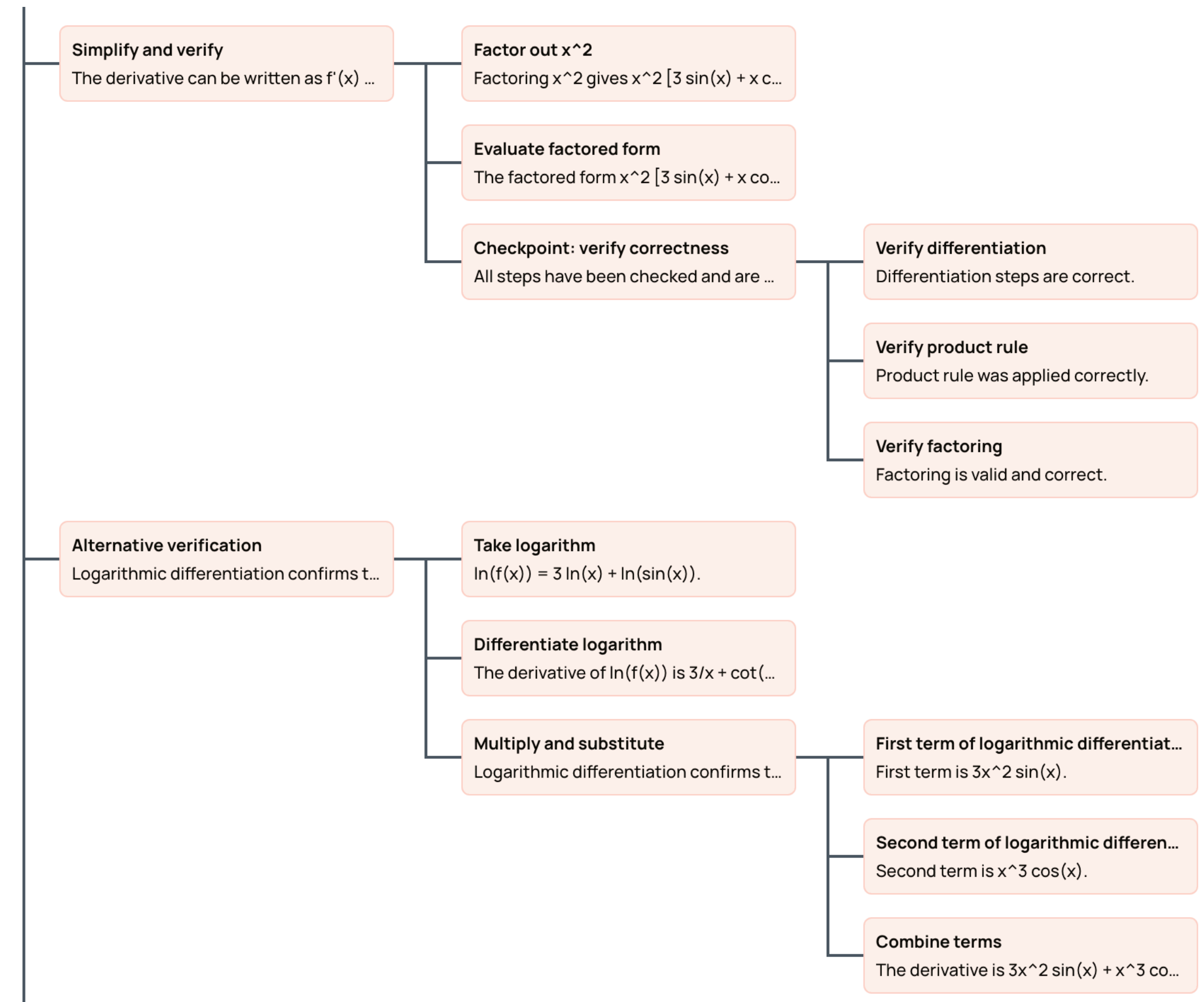
More flexibility, less prompting

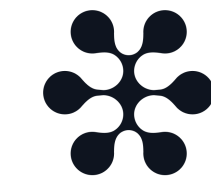
```
{
  "reasoning": [
    {
      "thought": "First, I need to identify the type of function given:  $f(x) = x^3 * \sin(x)$ . This is a product of two functions, so I should use the product rule for differentiation.",
      "subtasks": [
        {
          "thought": "Recall the product rule: if  $f(x) = u(x) * v(x)$ , then  $f'(x) = u'(x) * v(x) + u(x) * v'(x)$ . I need to choose  $u(x)$  and  $v(x)$  for this problem.",
          "subtasks": [
            {
              "thought": "Let  $u(x) = x^3$  and  $v(x) = \sin(x)$ . This splits the function into two parts suitable for the product rule.",
              "subtasks": [
                {
                  "thought": "Now, I need to compute the derivative of  $u(x) = x^3$ . Using the power rule, the derivative of  $x^n$  is  $n*x^{n-1}$ .",
                  "subtasks": [
                    {
                      "thought": "Apply the power rule to  $x^3$ : derivative is  $3x^2$ .",
                      "conclusion": " $u'(x) = 3x^2$ .",
                      "title": "Derivative of  $x^3$ "
                    },
                    {
                      "thought": "Check if there are any additional terms or constants to consider in  $u(x)$ . Since  $u(x) = x^3$ , there are none.",
                      "conclusion": "No additional terms in  $u(x)$ .",
                      "title": "Check  $u(x)$  for extra terms"
                    }
                  ],
                },
                {
                  "thought": "Summarize the result for  $u'(x)$ :  $u'(x) = 3x^2$ .",
                  "conclusion": " $u'(x) = 3x^2$ .",
                  "title": "Summarize  $u'(x)$ "
                }
              ],
            }
          ],
        }
      ],
    }
  ],
}
```


TIM Outputs - Tree View

More flexibility, less prompting

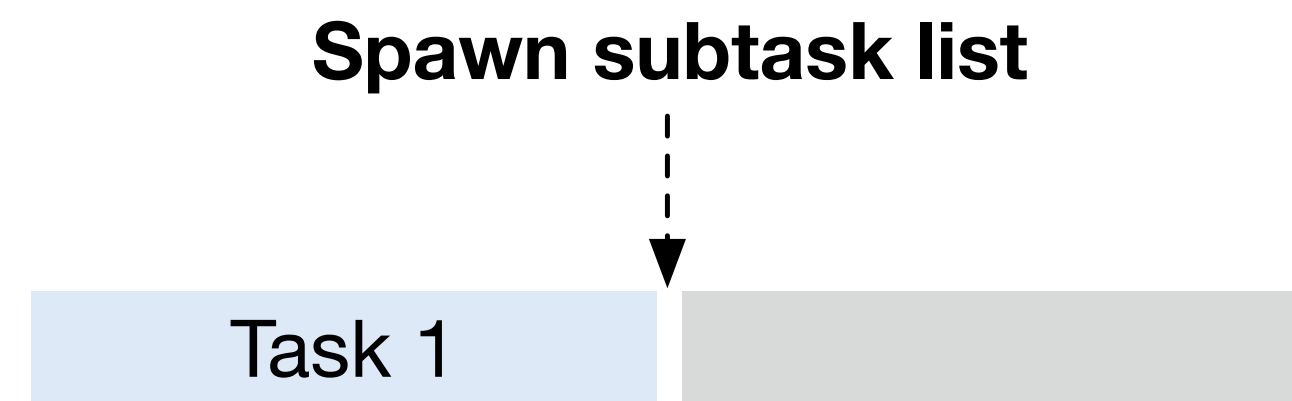
- Automatic workflow generation
- Critical path tracking
- Enabling subtask pruning





TIM Reasoning Breakdown

Recursive task decomposition and aggregation



Thinking

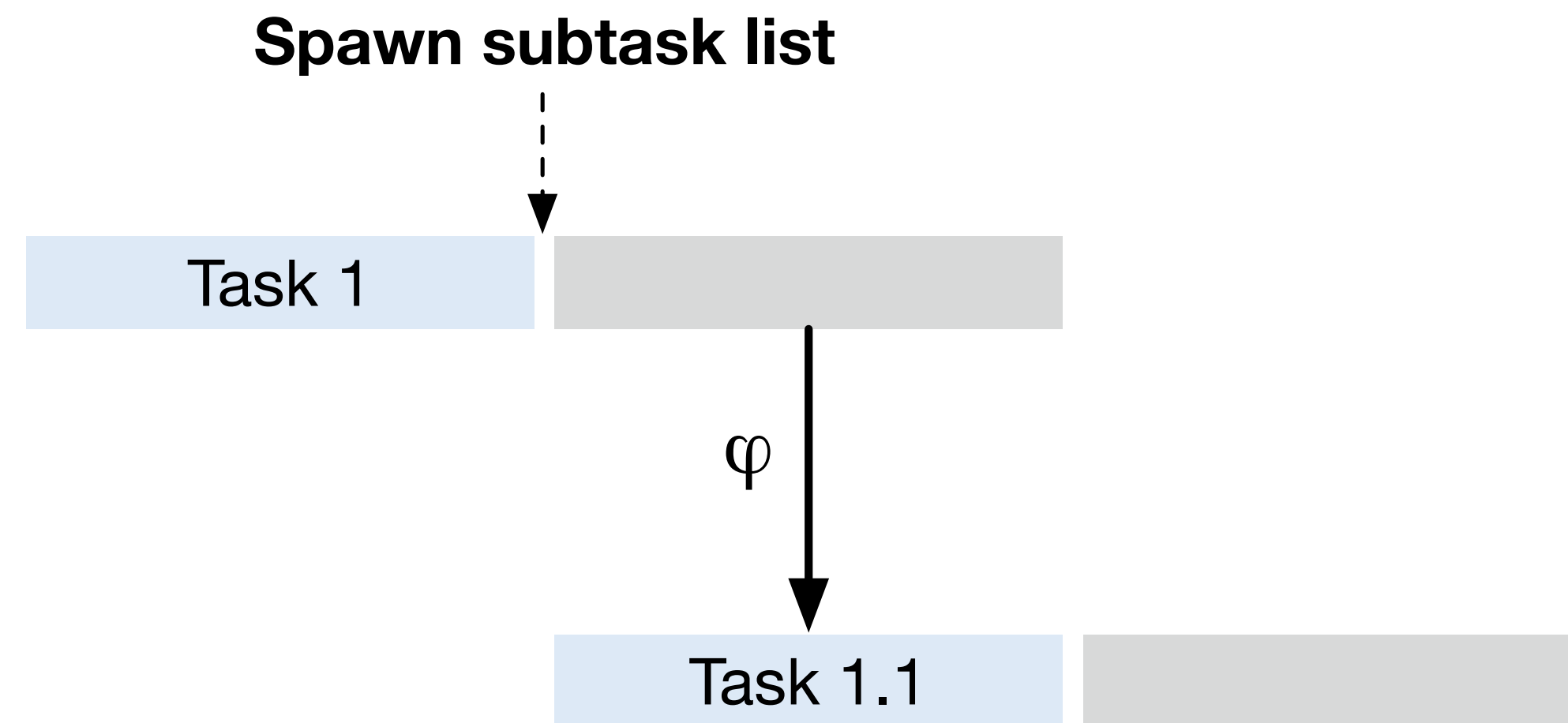
ToolUse

Subtasks

Conclusion

TIM Reasoning Breakdown

Recursive task decomposition and aggregation



Thinking

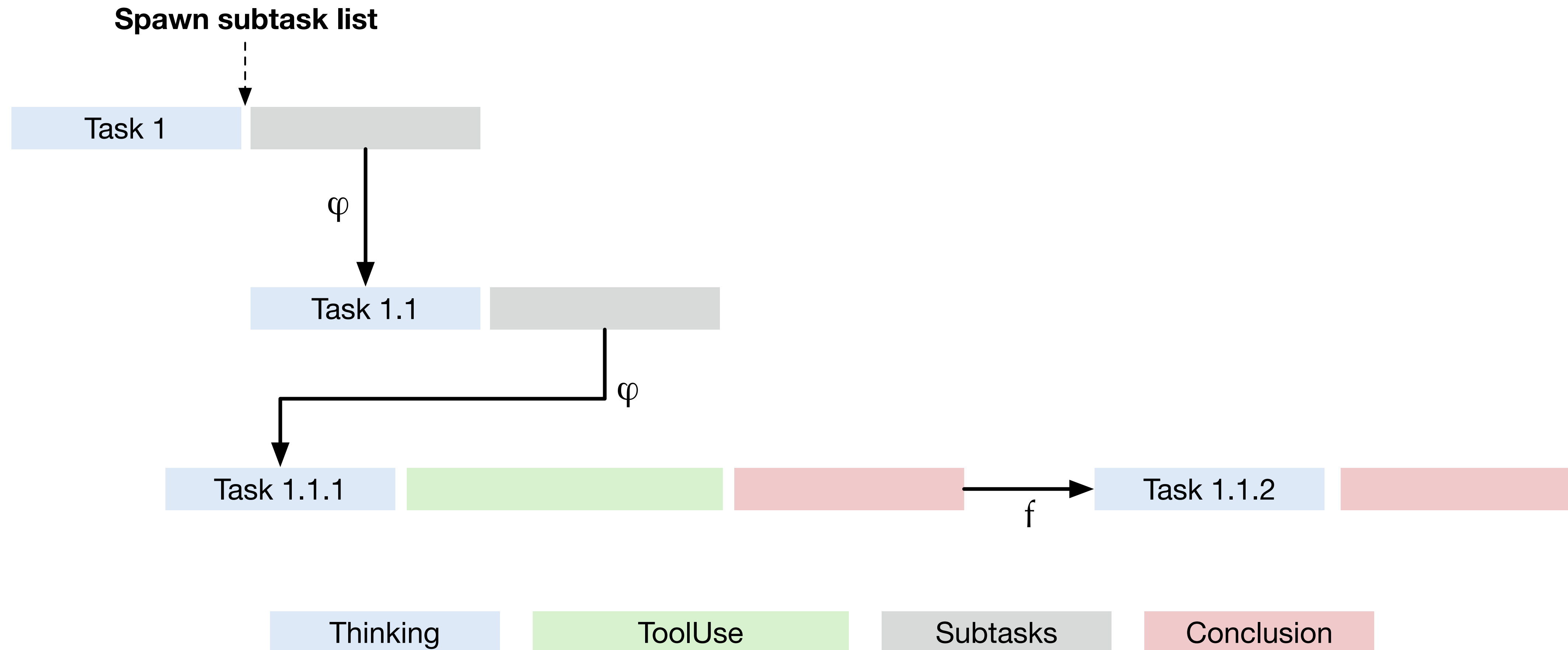
ToolUse

Subtasks

Conclusion

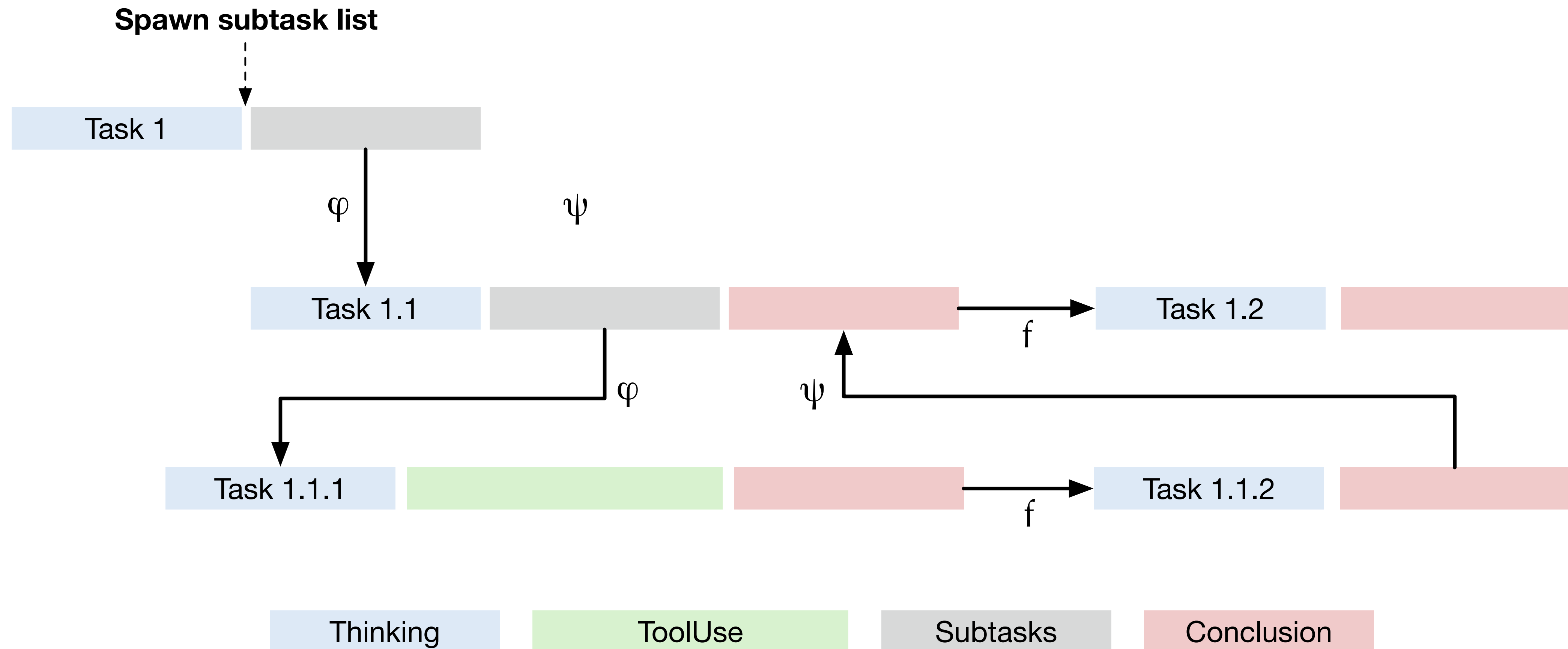
TIM Reasoning Breakdown

Recursive task decomposition and aggregation



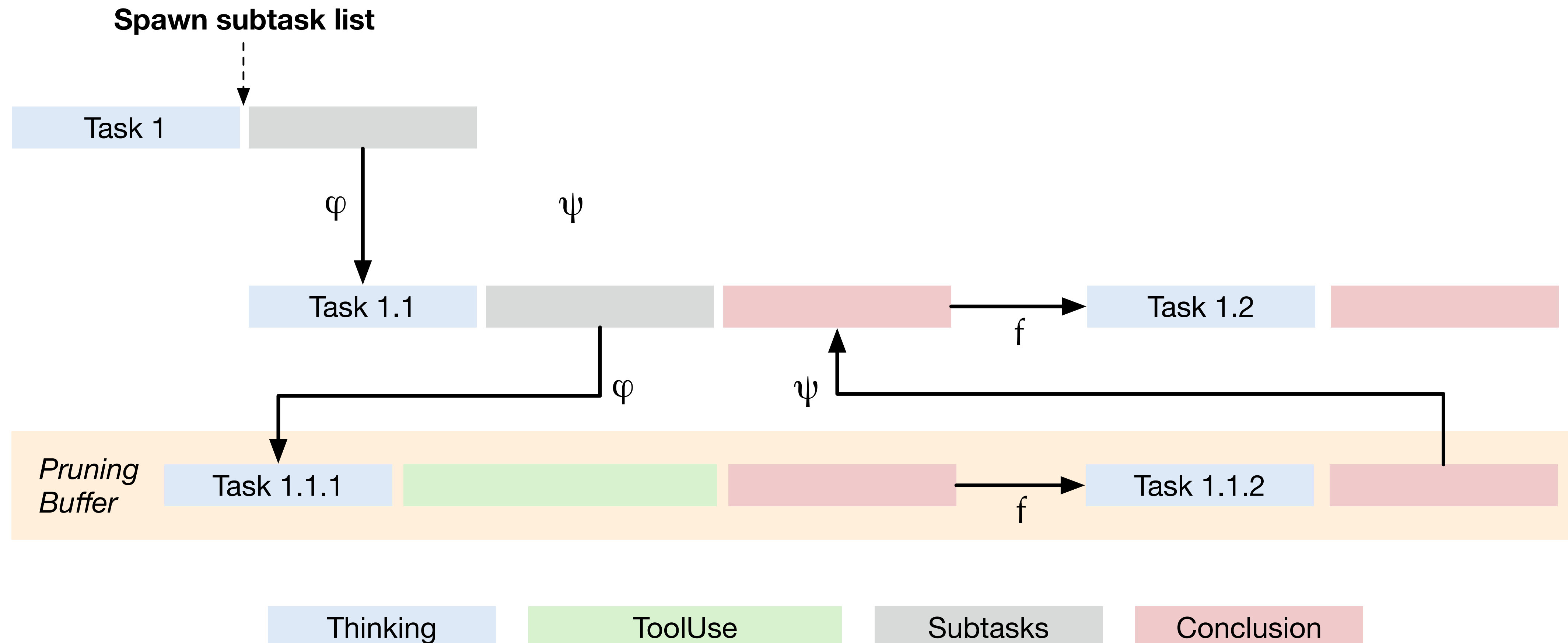
TIM Reasoning Breakdown

Recursive task decomposition and aggregation



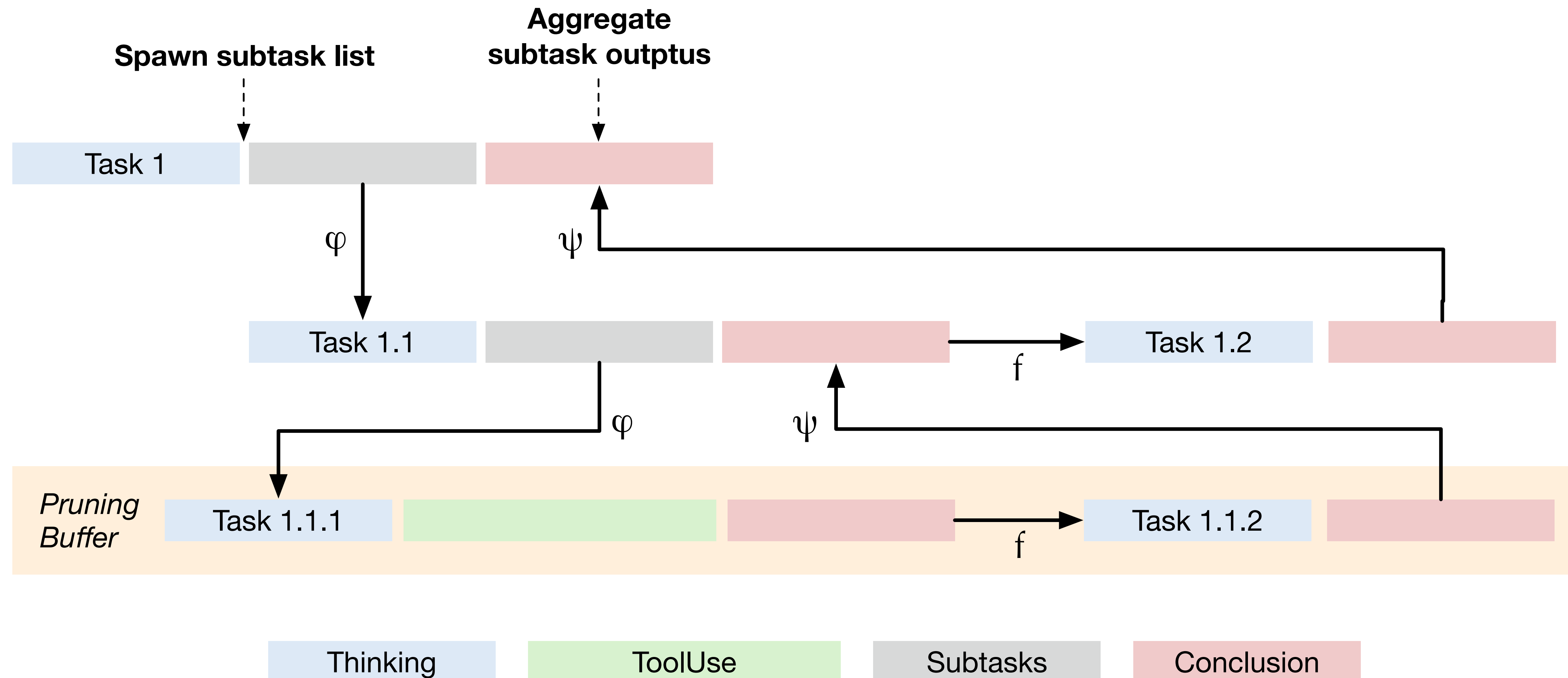
TIM Reasoning Breakdown

Recursive task decomposition and aggregation



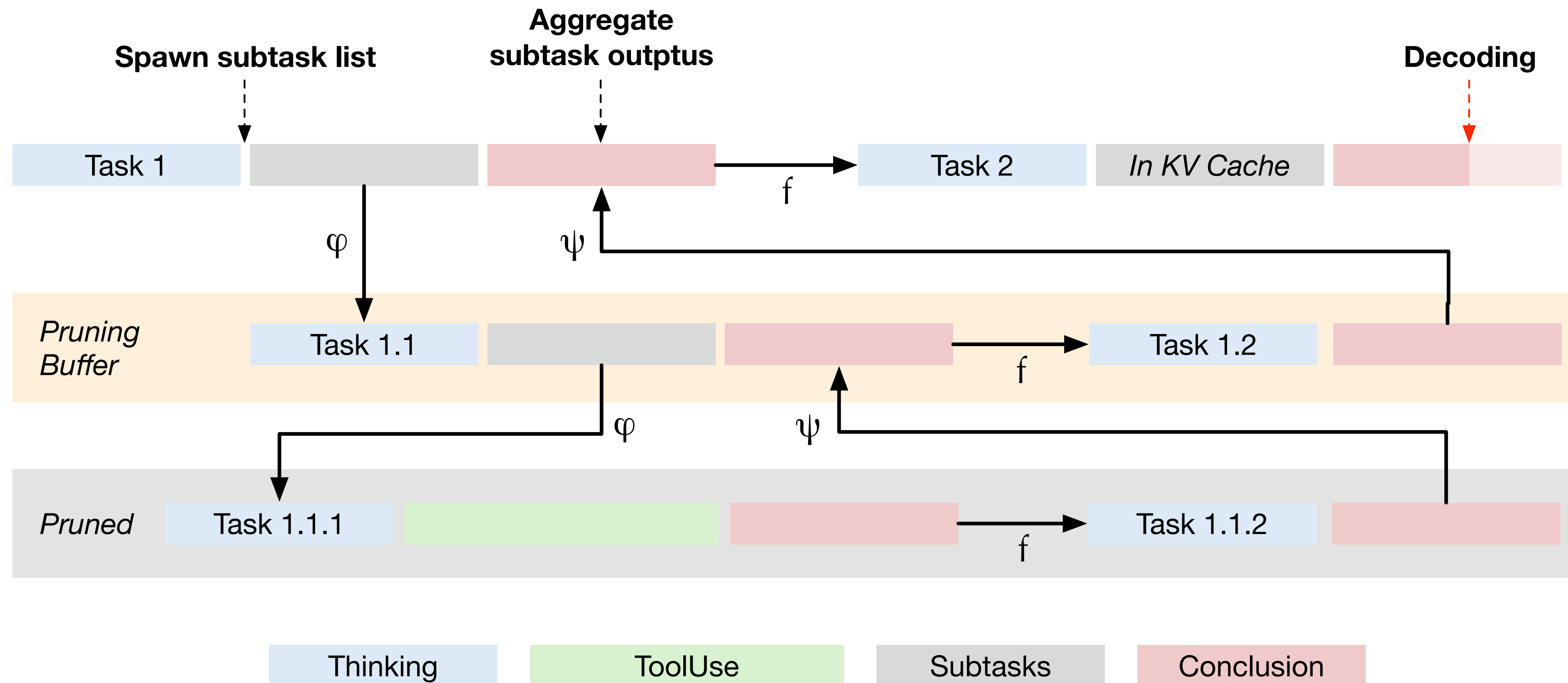
TIM Reasoning Breakdown

Recursive task decomposition and aggregation



Subtask Pruning

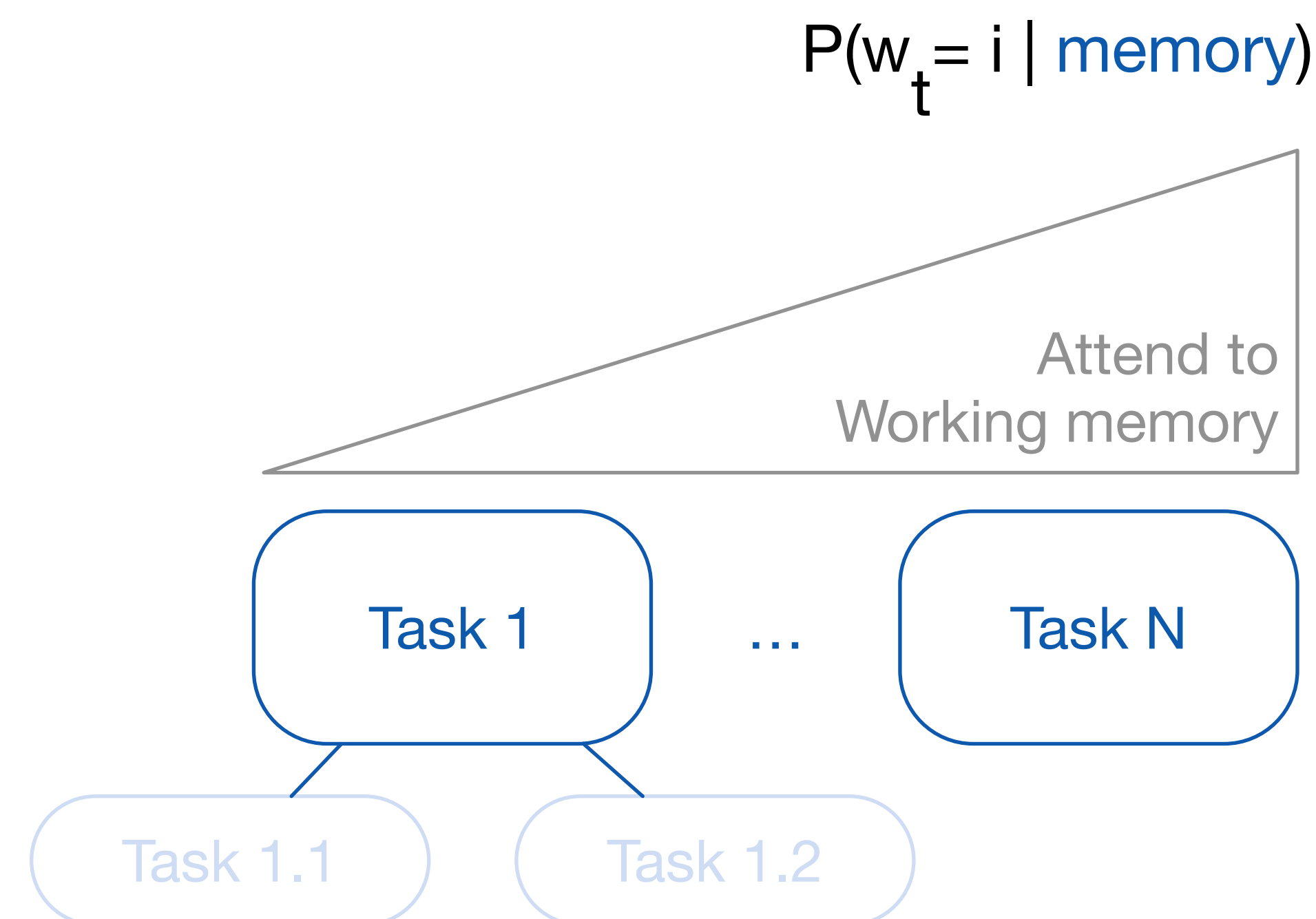
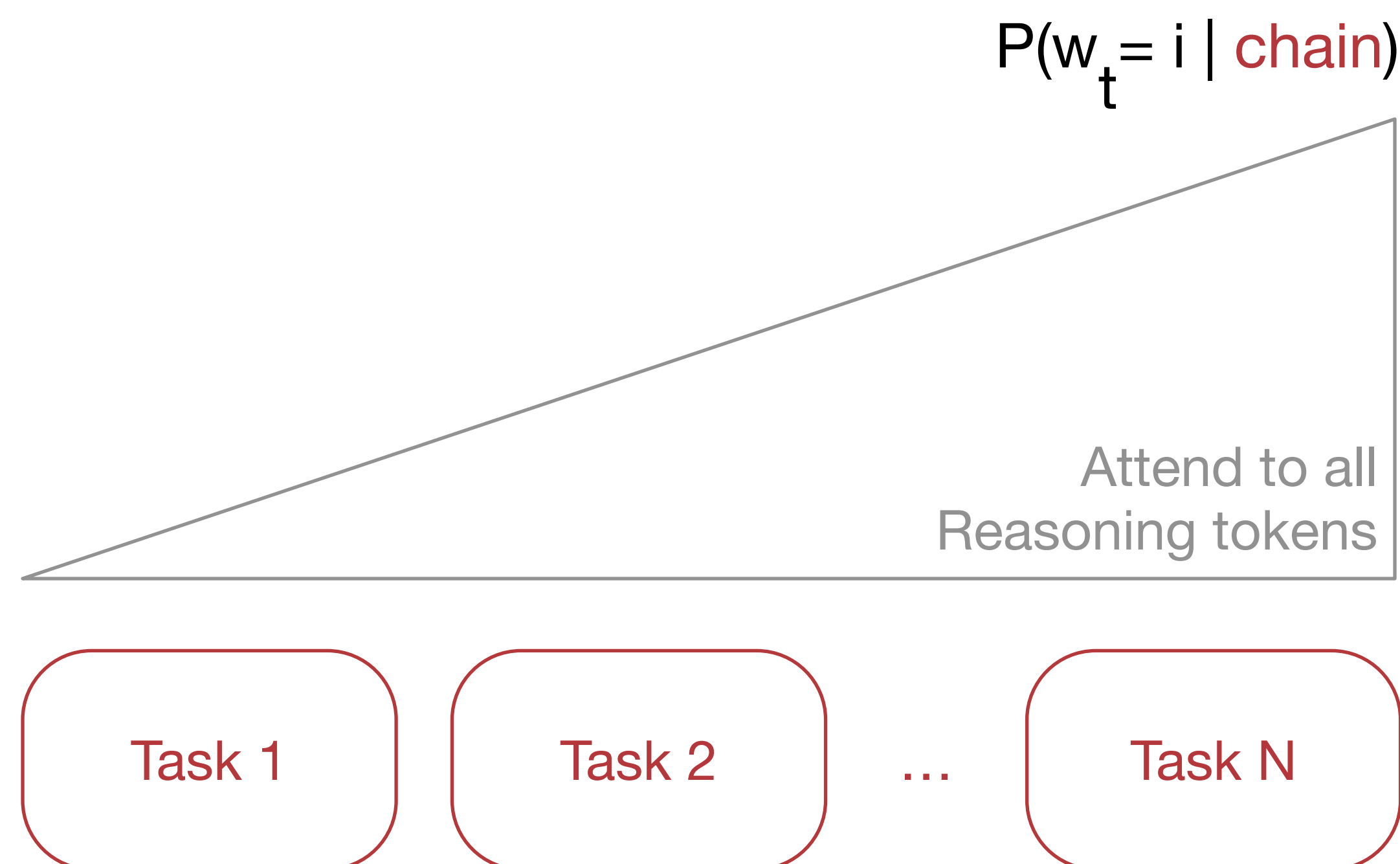
A rule-based approach



TIMRUN

Co-designed inference engine for TIM

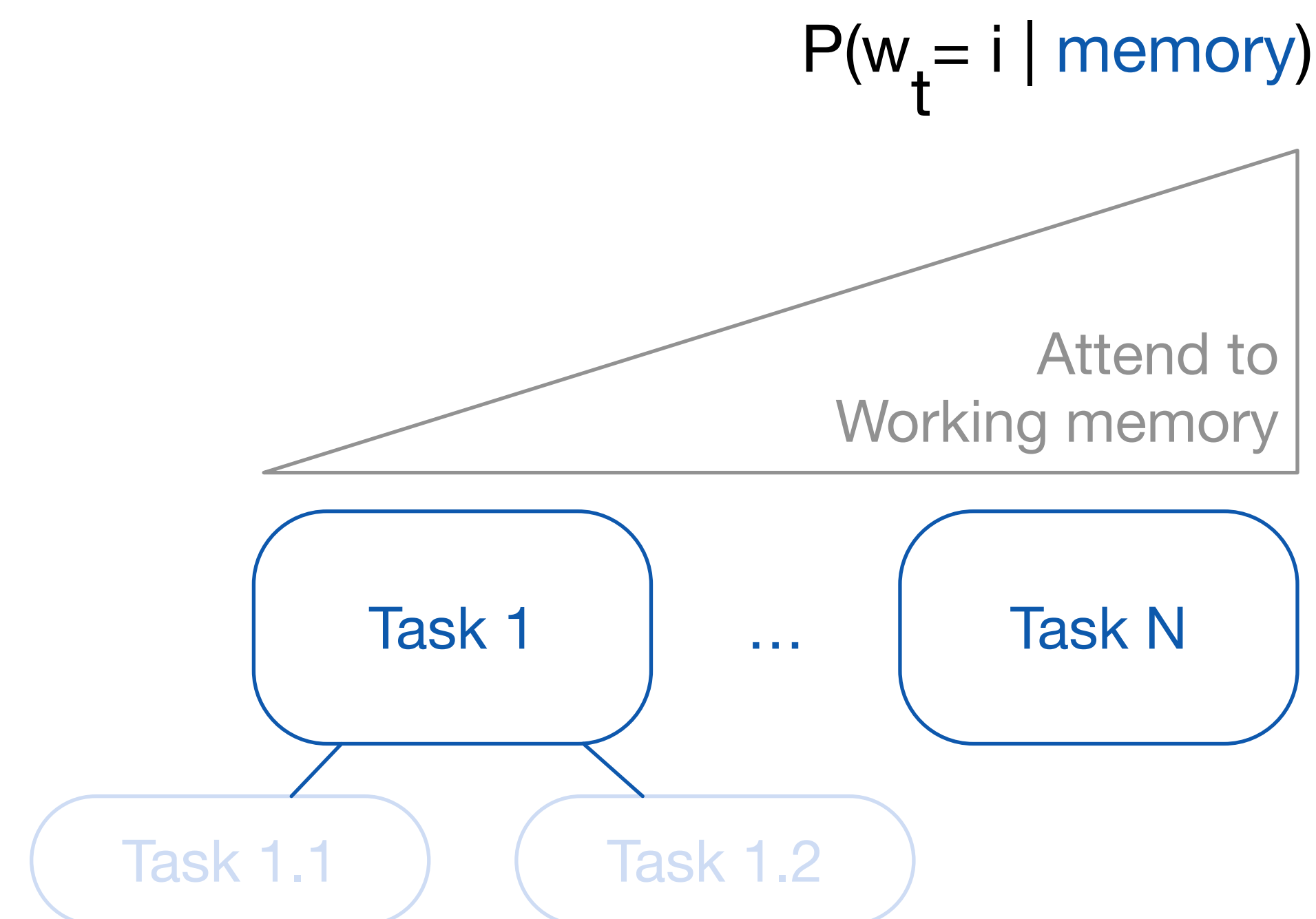
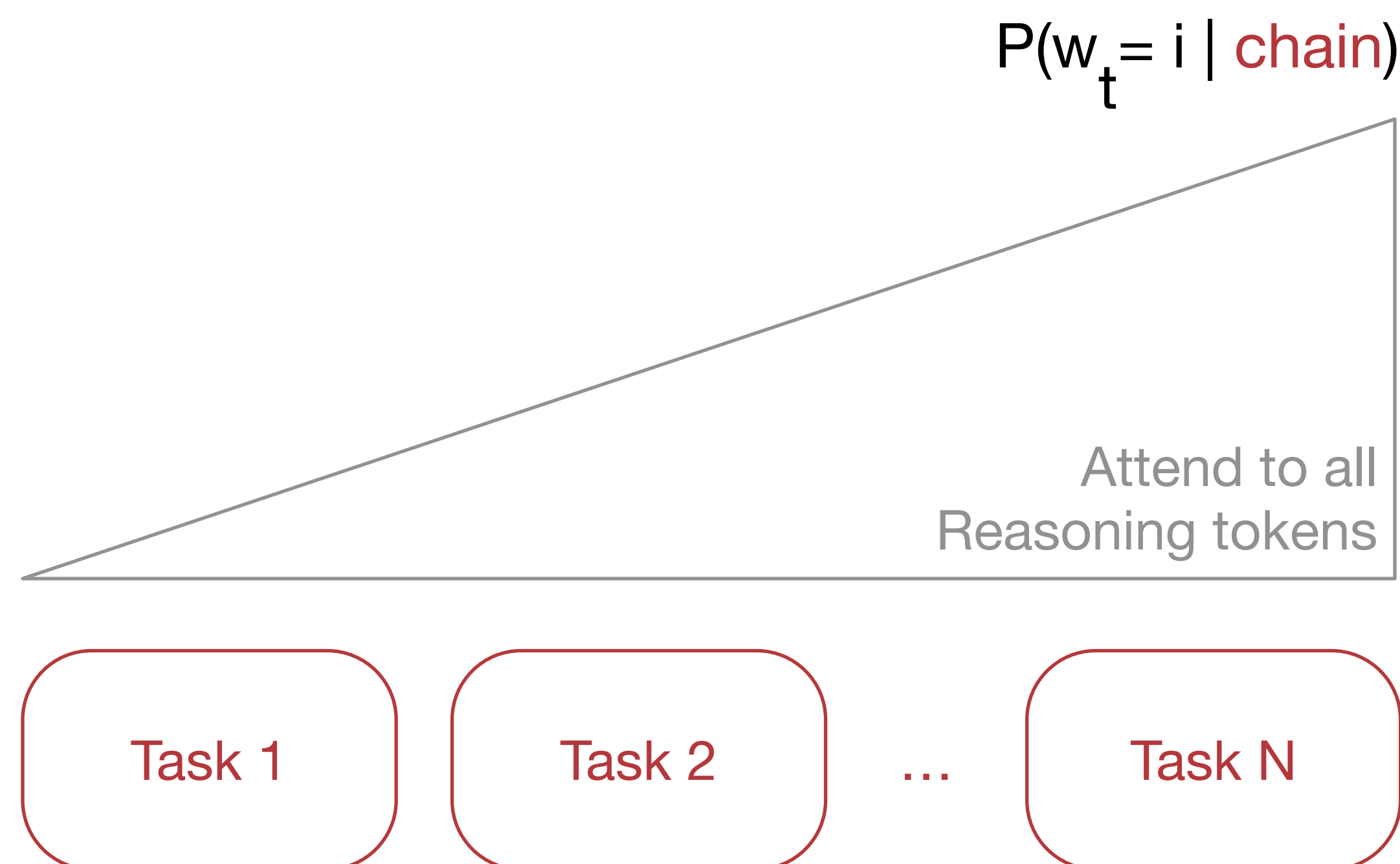
1. Efficient subtask pruning;
2. Working memory attention;
3. Tool calling



TIMRUN

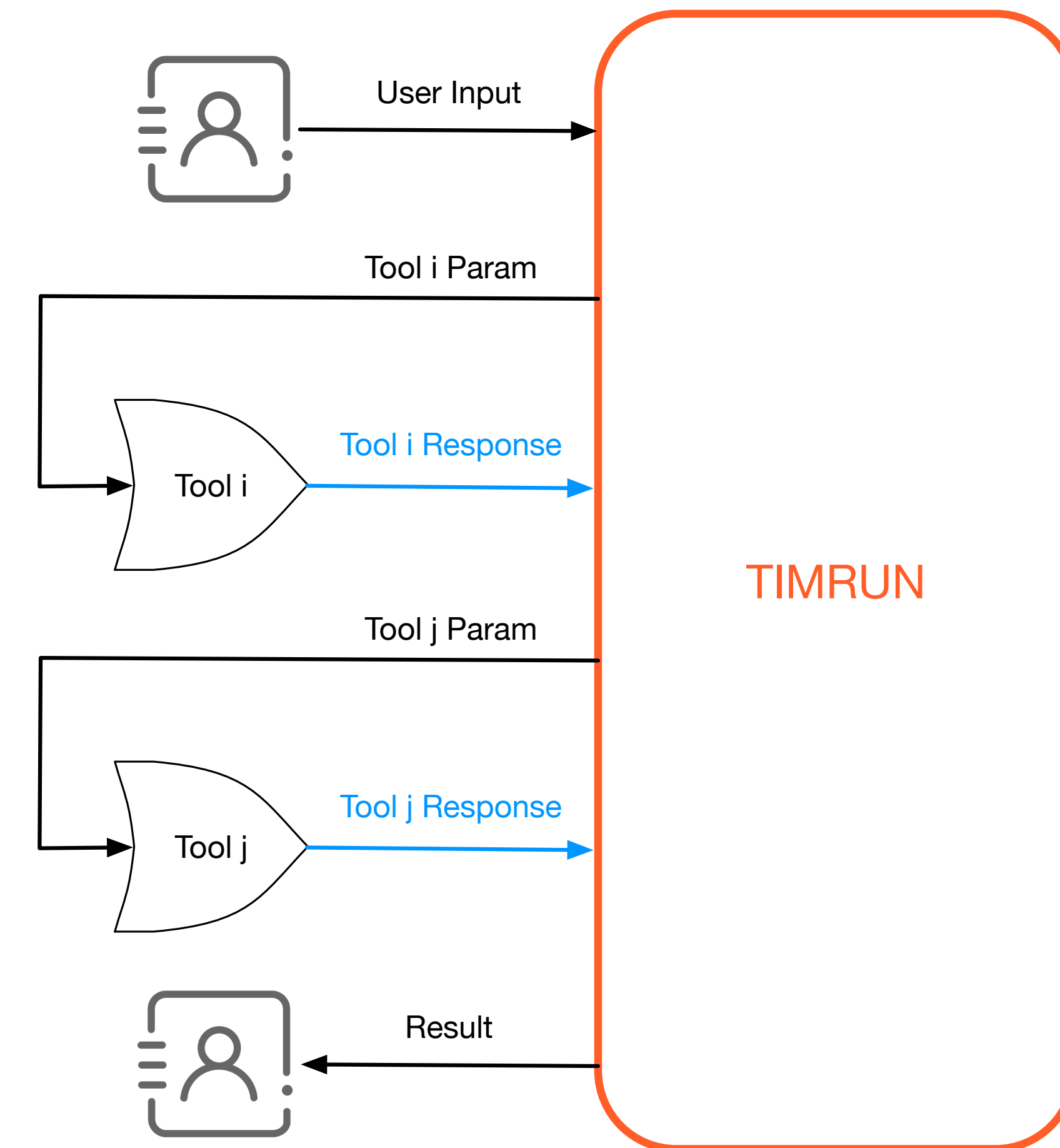
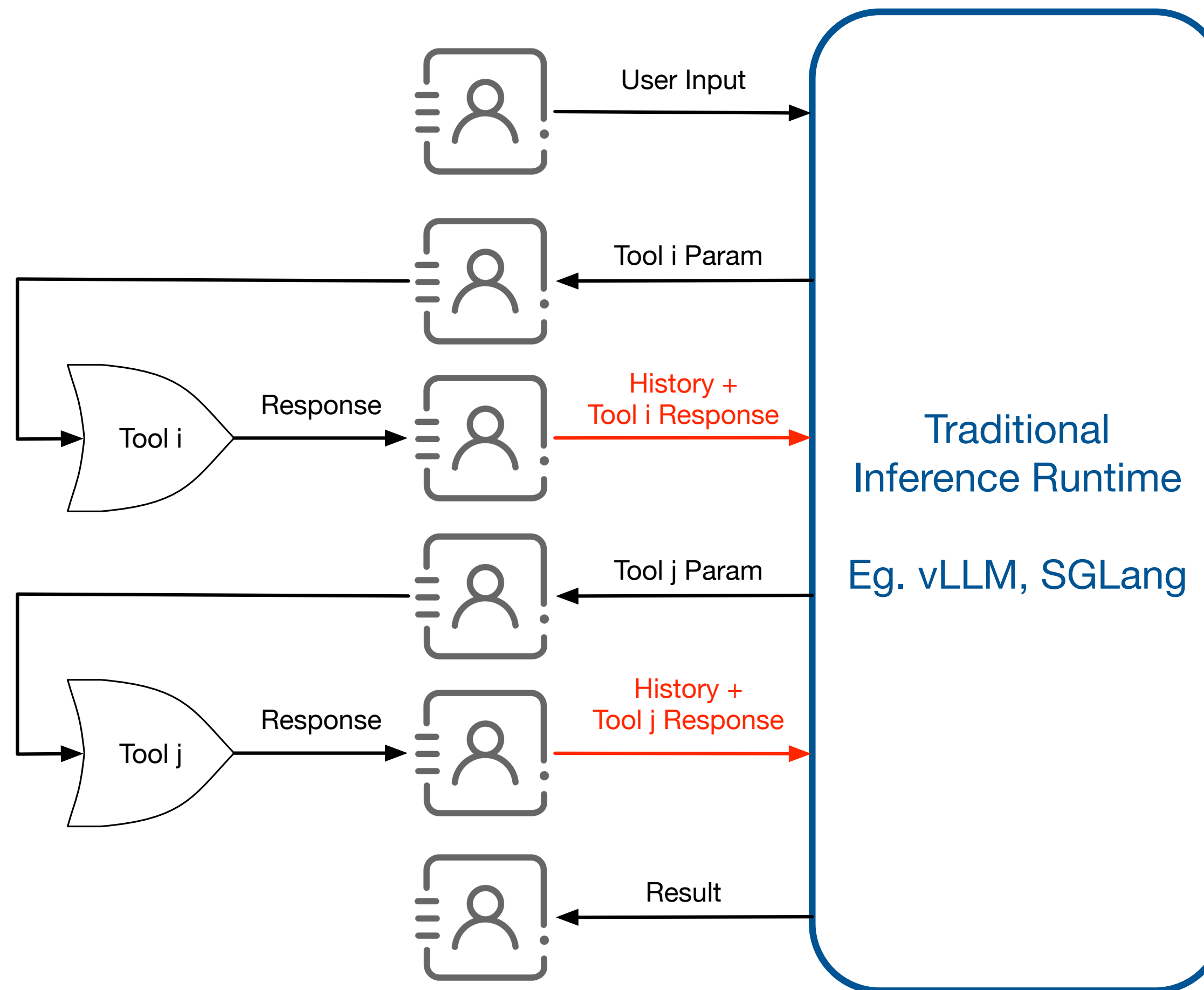
Co-designed inference engine for TIM

1. Efficient subtask pruning;
2. Working memory attention;
3. Tool calling



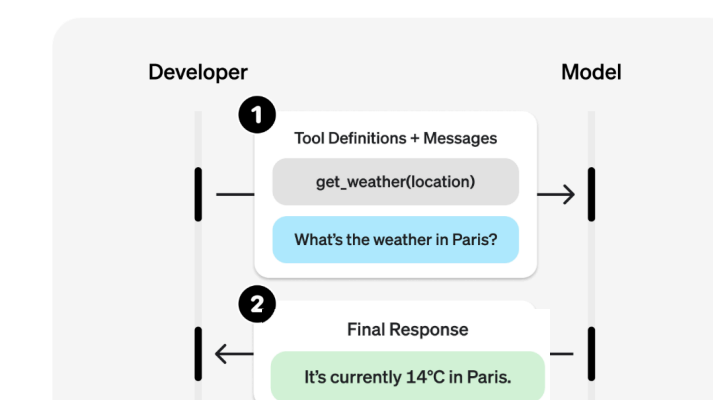
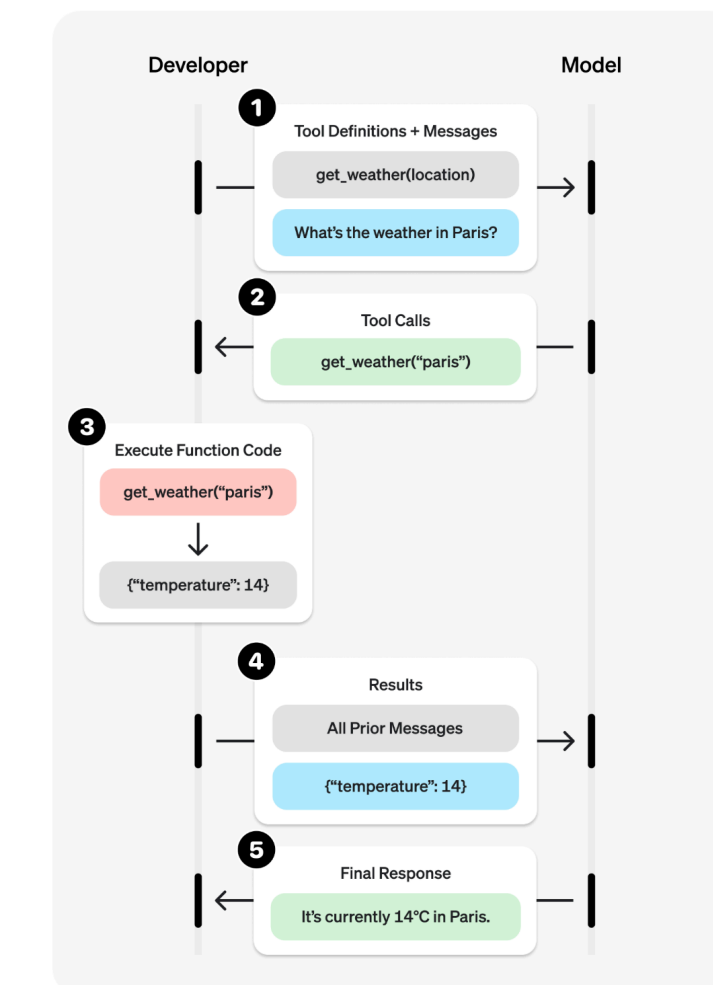
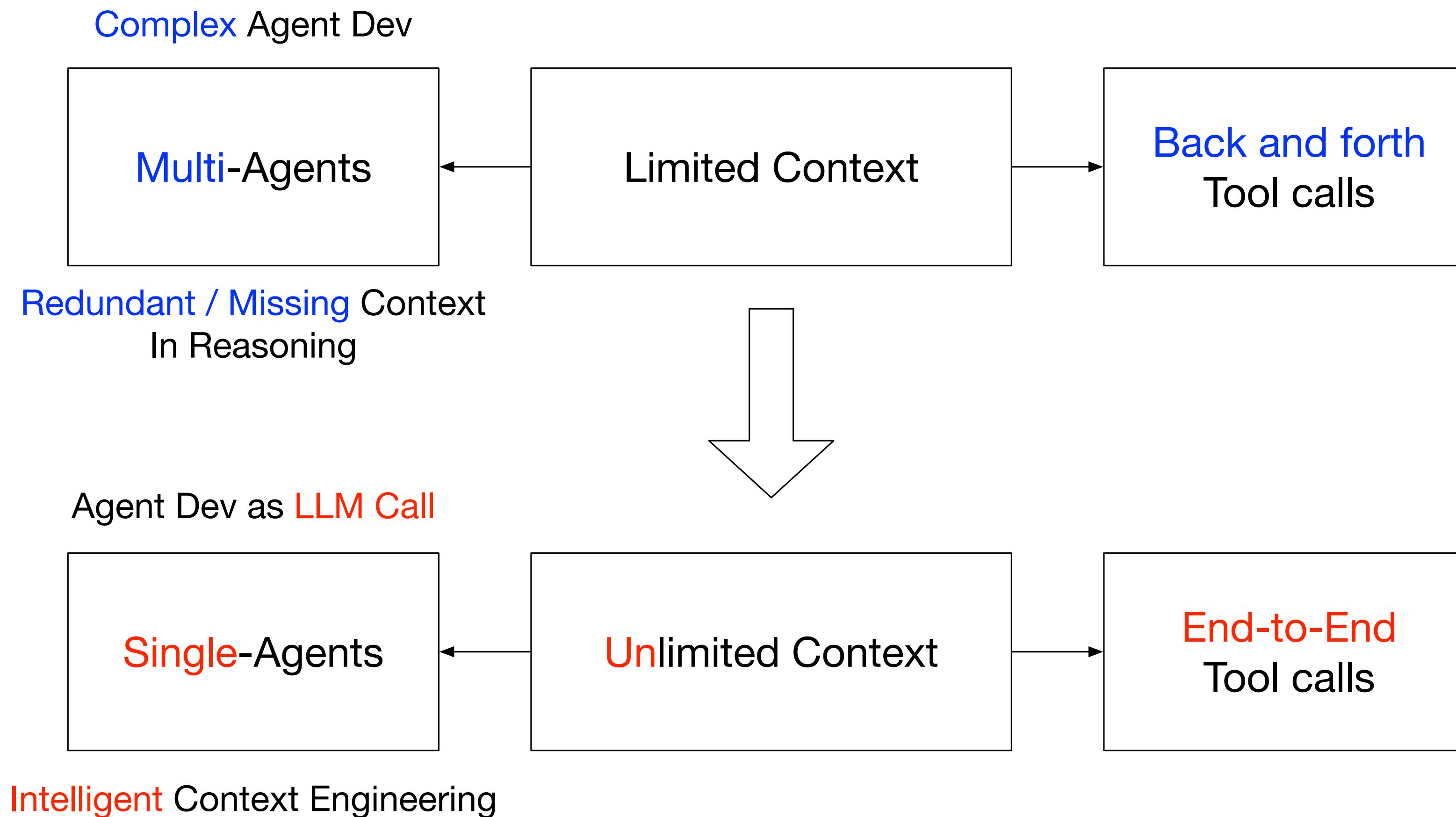
TIMRUN vs Others

End-to-end, multi-hop tool use



Impact on Agent dev

Agent **inference** instead of agent loop



TIM Inference vs Agent Loop

Deep research challenges

BrowseComp	Model	Deepseek-R1	GPT-4o	TIM-large	TIM-8b
	Paradigm	ReACT	Browsing	Browsing	Browsing
	Success (%)	9.5	1.9	7.8	2.3

DataCommons	Method	Reflection	NLEP+ReACT	Decomp	THREAD	TIM
	Accuracy	24.3	27.1	57.9	67.9	67.9

4k task-specific
prompt tokens

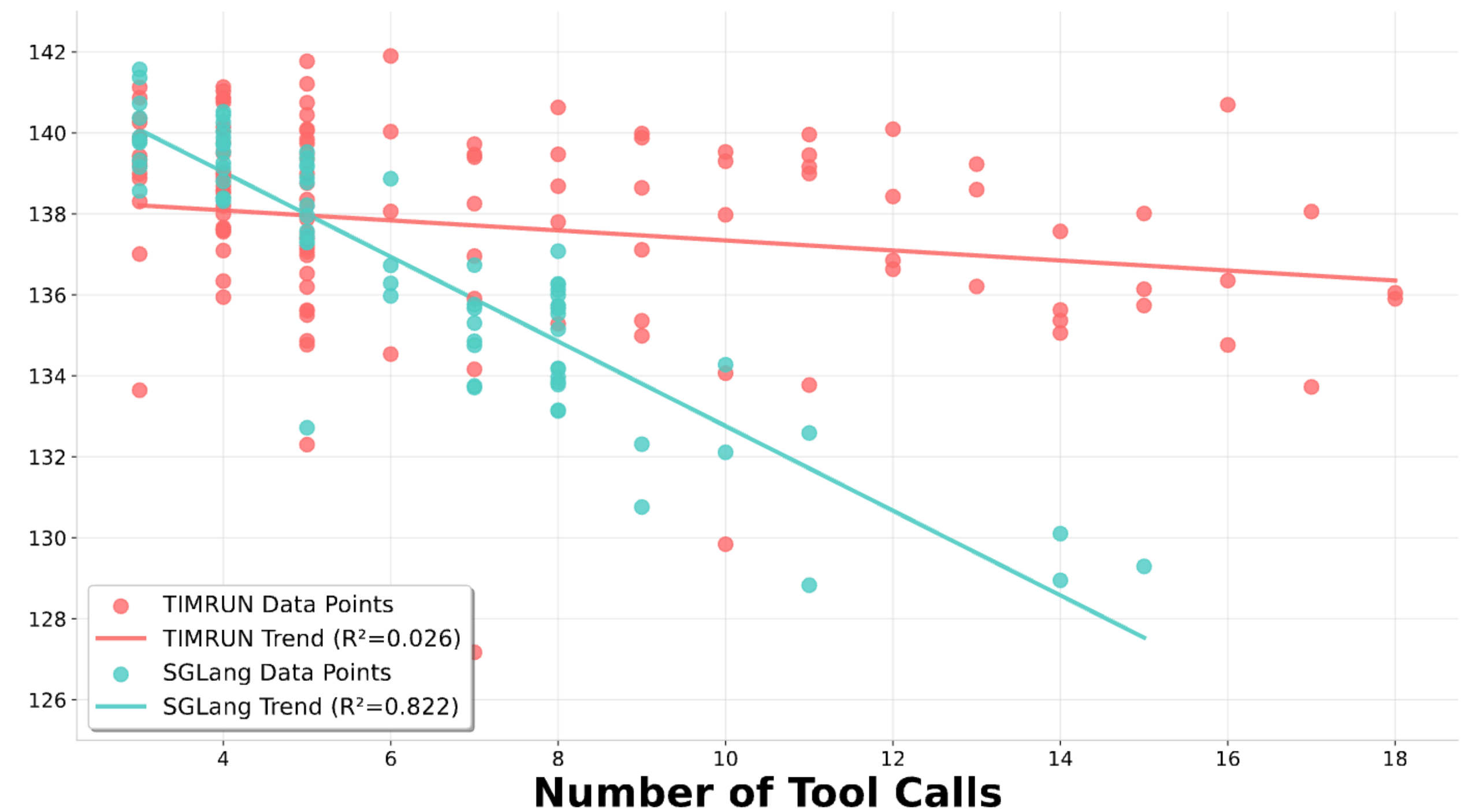
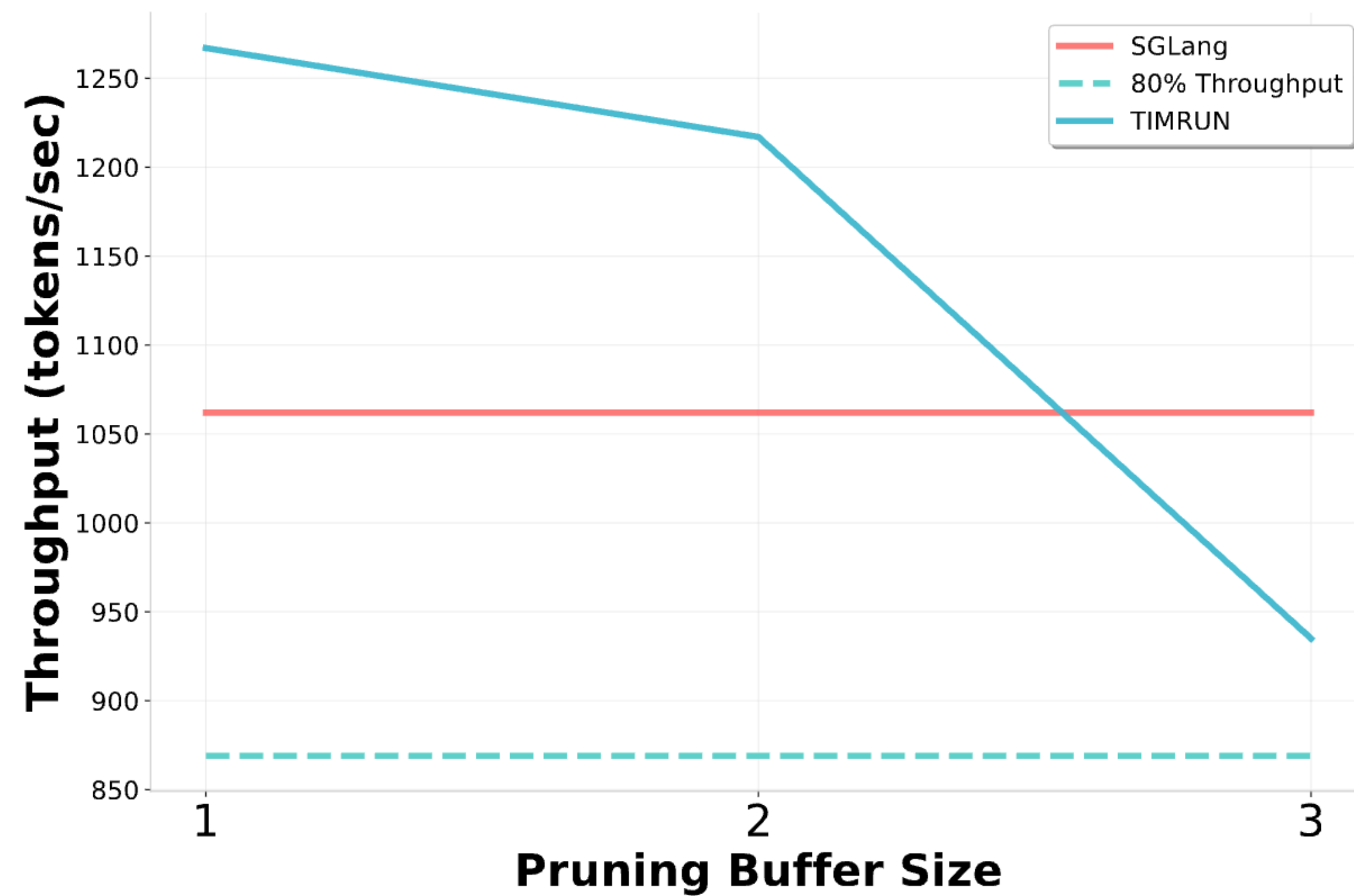
Effect of Context Pruning

Not harming -> improving performance

Task	TIM-8b + SGLang		TIM-8b + TIMRUN		
	Accuracy	Accuracy	Max Cache	Output Len.	KV Pruned (%)
MATH500	69.6	69.0	1569.2	3362.2	53.3
MMLUSTEM500	88.4	87.6	1330.9	2747.0	51.6
AMC 2022	60.5	60.5	2203.9	5131.7	57.1
AMC 2023	80.0	80.0	1876.5	4547.4	58.7
AIME 2024	40.0	46.7	3218.6	8974.7	64.1
GPQADiamond	44.9	48.5	1712.9	3742.6	54.2

Throughputs

Subtasks pruning improves efficiency



Limitations

In reasoning and tool using

- Repeat generations
- Sensitive to prompts and sampling parameters on benchmarks
 - On AIME 2024
 - Temperature = 0.9, top_p = 0.95 -> consistent > 43% acc.
 - Temperature = 0.8 -> 23% acc
- If tool responses are too long, still triggers context limit



TIM-8b-preview is available on Hugging Face at
[SubconsciousDev/TIM-8b-preview](https://huggingface.co/SubconsciousDev/TIM-8b-preview)



Try TIMRUN at Subconscious.dev



Thank you!

<https://brand.mit.edu/logos-marks/tim-beaver>