# Representing Prompting Patterns with PDL: Compliance Agent Case Study

Mandana Vaziri   Louis Mandel   Yuji Watanabe   Hirokuni Kitahara   Martin Hirzel   Anca Sailer

## Prompt engineering is hard
## How does Prompt Declaration Language (PDL) help?

### Prompts at the forefront

PDL written in YAML

Single declarative language with control structures, and functions for pattern reuse
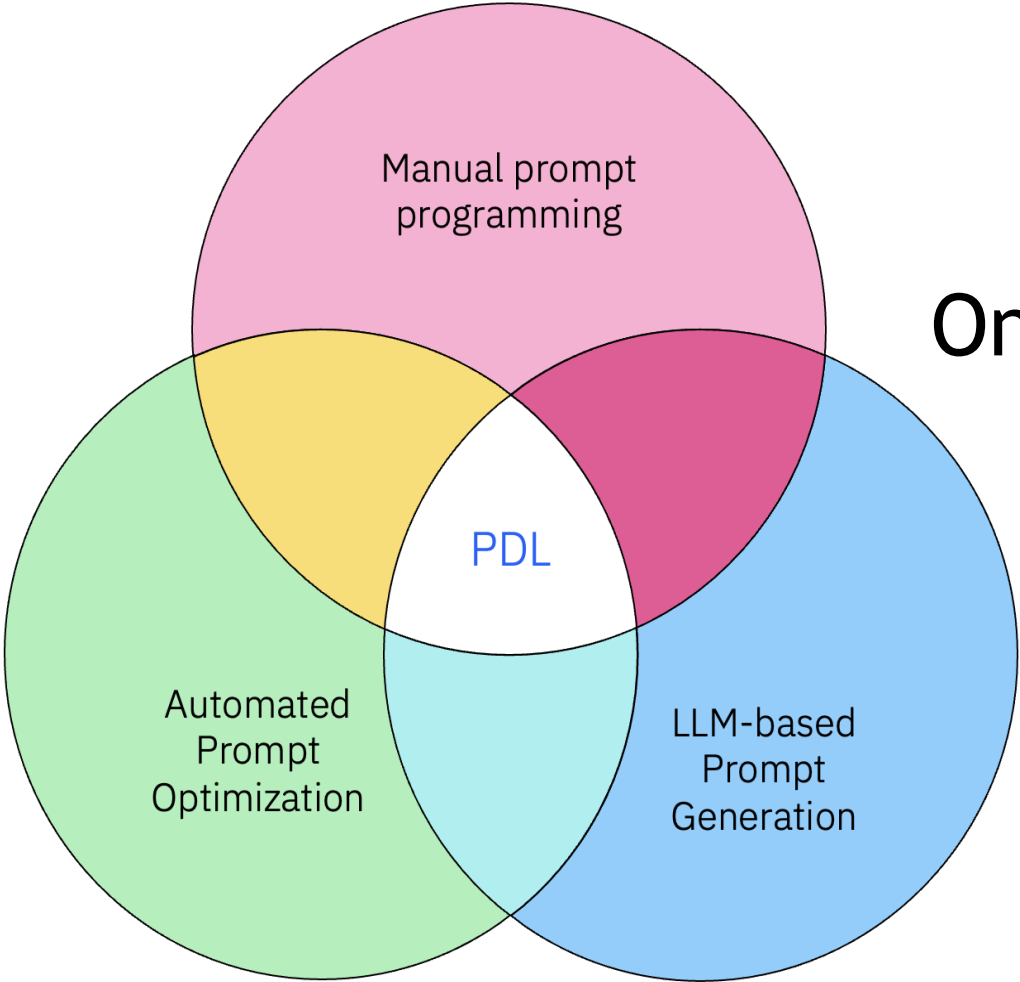
Few orthogonal features

### Intrinsics

PDL is based on granite-io and supports the following intrinsics:
thinking, hallucinations answerability, certainty citations, query-rewrite

### Composition of LLMs & code

PDL abstracts away the plumbing necessary for such compositions

Supports a wide variety of model providers and models, based on LiteLLM

### Automated parallelization

In PDL, all model calls are asynchronous and will be executed in parallel in the absence of data dependencies

### Implicit accumulation of messages

PDL keeps track of the context (list of messages) implicitly, making programs much less verbose

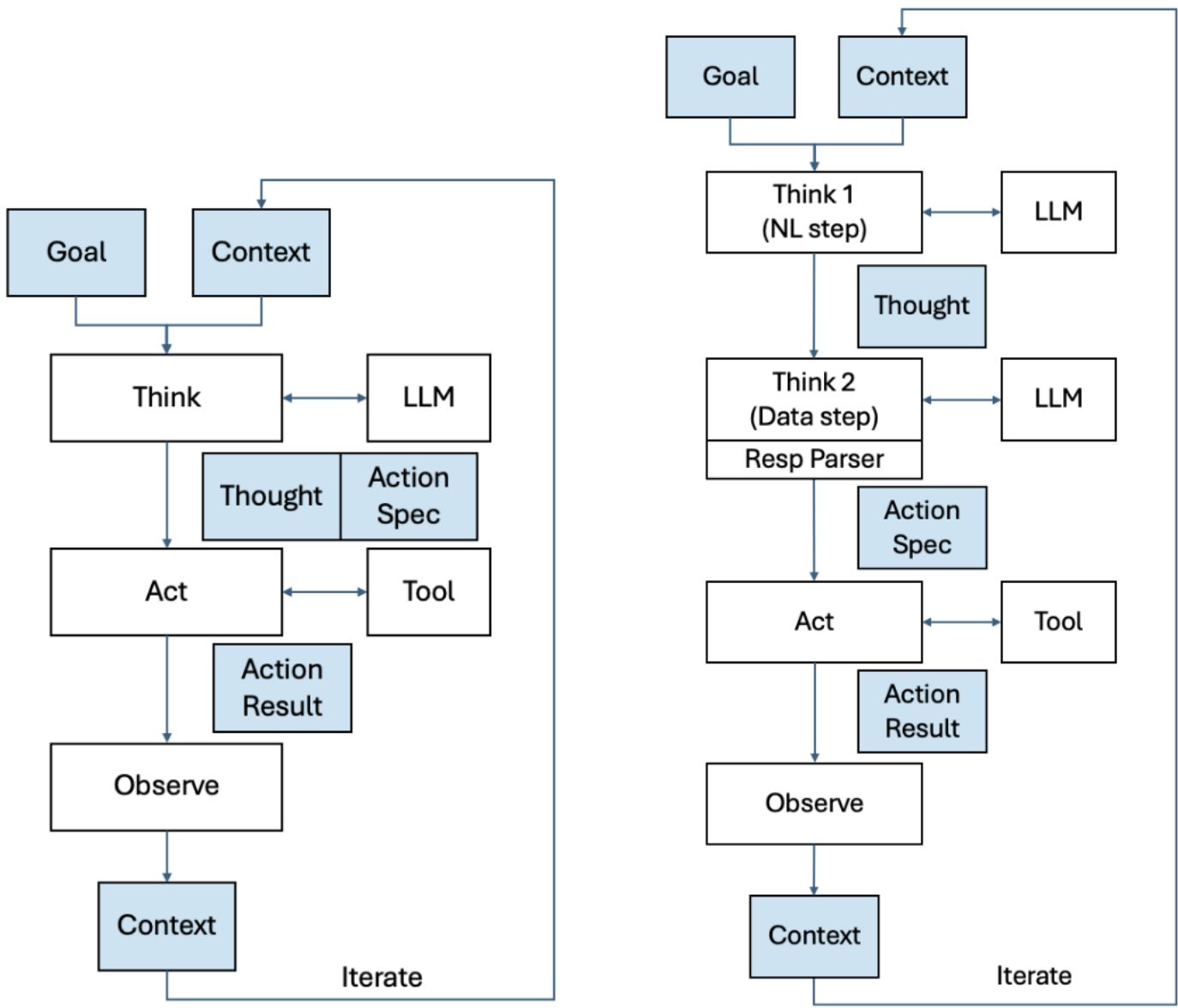Support for chat APIs

### Automated Prompt Optimization

AutoPDL starts with a PDL program with variables, a domain specification, and a dataset. It optimizes prompting patterns and few-shots. Paper at AutoML'25
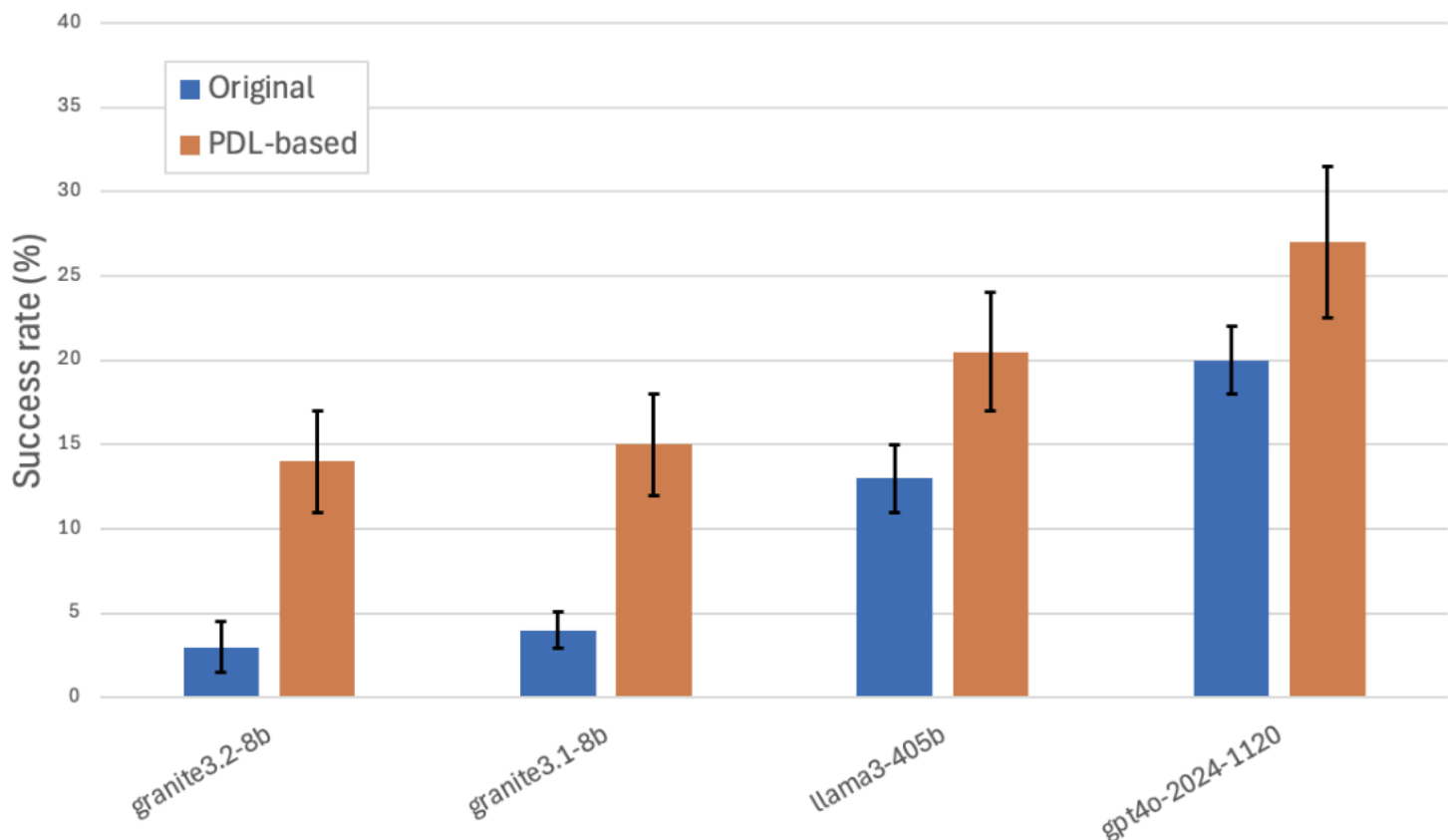
```
1   description: tool use
2 ∨ defs:
3 ∨   search:
4       description: Wikipedia search      # Tool definition
5       function:
6 ∨       topic:
7           type: string
8           description: Topic to search
9       return:
10        lang: python                     # Python code block
11 ∨      code: |
12          import warnings, wikipedia
13          warnings.simplefilter("ignore")
14          try:
15            result = wikipedia.summary("${ topic }")
16          except wikipedia.WikipediaException as e:
17            result = str(e)
18 ∨ text:
19 ∨ - role: system                        # System prompt
20 ∨   content: |
21        You are a helpful AI assistant with access to the
22        following tools. If a tool does not exist in the
23        provided list of tools, notify the user that you
24        do not have the ability to fulfill the request.
25      contribute: [context]
26 ∨ - role: tools                         # Tools prompt
27 ∨   content:
28        text: ${ [ search.signature ] }  # Using the function signature
29      contribute: [context]
30   - "What is the circumference of planet Earth?\n"   # Query
31 ∨ - def: actions
32      model: ollama_chat/granite3.3:8b
33      parser: json
34      spec: [{ name: string, arguments: { topic: string }}]   # Type checking
35   - "\n"
36 ∨ - if: ${ actions[0].name == "search" }   # Conditional
37 ∨   then:
38        call: ${ search }                # Tool calling
39        args:
40          topic: ${ actions[0].arguments.topic }
```

### Type checking

PDL provides type checking of both input and output of models. Types feed seamlessly into constrained decoding



One representation
Many uses

## IT Compliance

IT compliance tasks require specialized expertise because of complex standards and internal organizational policies. Automation tools exist for routine operations, but policy assessment remains largely a manual task.

## What is CISO Agent?

The CISO Agent provides automated support to IT teams. When it receives new regulatory requirements the agents: analyzes their content, identifies target systems, generates and deploys any scripts, validates outcomes, and provides comprehensive posture reporting.

## Chief Information Security Officer (CISO) Compliance Agent Case Study



(a) Original ReAct         (b) PDL-based

## Prompt pattern customization with PDL

We compare two implementations of the CISO Agent: with CrewAI (a) and with PDL (b). As a framework CrewAI does not provide sufficient flexibility to customize the prompting pattern. With PDL, the Think step in the ReAct loop is split in 2 steps: first generate a thought and then generate an action spec. PDL also allowed to easily add a specialized parser for the second step.

## 4x better performance

The PDL implementation demonstrates consistent improvements across all models, with particularly dramatic gains in smaller models like granite3.2-8b, achieving 4 times better performance. The tool call success rate improved significantly.
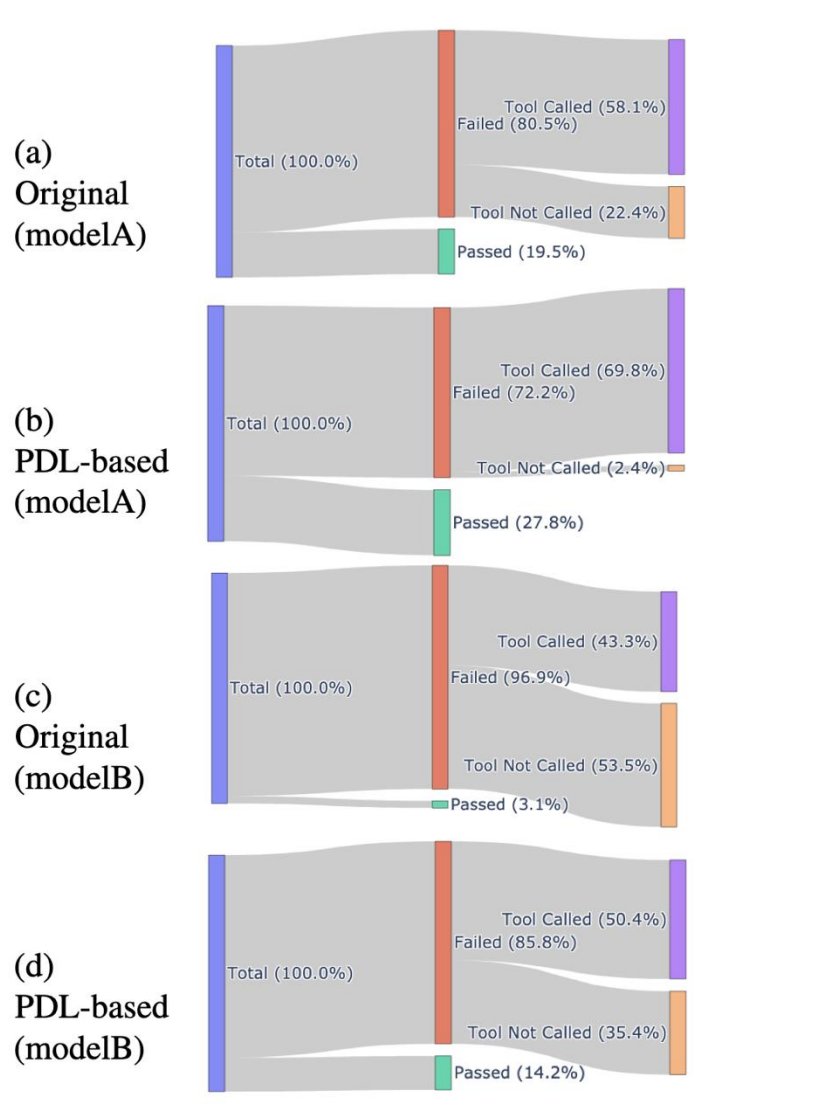


Performance evaluation results using IT-Bench comparing CrewAI (blue) and PDL-based (orange) implementations



Tool Call Success Rate Comparison
model A: gpt4o-2024-11-20
model B: granite3.2-8b-instruct