

■■■■■ Optimizing Agent Planning for Security and Autonomy

A horizontal row of 15 small black squares, likely representing a binary sequence or a set of data points.

input_papers\test.pdf

2026-02-14 22:37

33

■ 1 ■

** Aashish Kolluri¹ Rishi Sharma^{1,2†} Manuel Costa¹ Boris Köpf¹ Tobias Nießen^{3†} Mark Russinovich¹ Shruti Tople¹ Santiago Zanella-Béguelin¹ ^{1,2} [REDACTED]

³ [REDACTED] ** [REDACTED] AI [REDACTED]
 [REDACTED] ** [REDACTED] Human-in-the-Loop, HITL [REDACTED]
 [REDACTED] HITL [REDACTED] Information-Flow Control, IFC [REDACTED]
 [REDACTED] AgentDojo WASP [REDACTED]

**1 ** AI [REDACTED] Anthropic, 2025; OpenAI, 2025b; Perplexity, 2025b [REDACTED] OpenAI, 2025a; Perplexity, 2025a; OpenAI, 2025c [REDACTED] Prompt Injection Attacks, PIAs [REDACTED] Greshake [REDACTED] 2023 Yi [REDACTED] 2025
 AI [REDACTED] PIAs [REDACTED] Wallace [REDACTED] 2024 Chen [REDACTED] 2025a
 Yi [REDACTED] 2025 Abdelhabib [REDACTED] 2025 Jia [REDACTED] 2025
 Zhan [REDACTED] 2025 Nasr [REDACTED] 2025 PIAs [REDACTED]
 Information-Flow Control, IFC [REDACTED] PIAs Costa [REDACTED]
 2025 Zhong [REDACTED] 2025 Debenedetti [REDACTED] 2025 Wu [REDACTED] 2024
 integrity [REDACTED] confidentiality [REDACTED]
 IFC [REDACTED] PIAs [REDACTED]

AgentDojo [REDACTED] 30% Costa [REDACTED] 2025 Debenedetti [REDACTED] 2025
 HITL [REDACTED] TCR [REDACTED] @k [REDACTED] 3 [REDACTED] Hum
 an-in-the-Loop, HITL [REDACTED] PIAs [REDACTED] Visual Studio Code [REDACTED] GitHub Copilot Chat [REDACTED]

IFC [REDACTED] IFC [REDACTED]
 IFC [REDACTED] IFC [REDACTED]
 1 [REDACTED] 2 [REDACTED]
 3 [REDACTED] endorse [REDACTED]
 IFC [REDACTED] FIDES [REDACTED] PRUDENTIA [REDACTED] Agent
 Dojo Debenedetti [REDACTED] 2024 WASP Evtimov [REDACTED] 2025
 1. ** IFC [REDACTED] ** IFC [REDACTED] HITL [REDACTED] 1.5
 AgentDojo [REDACTED] IFC [REDACTED] HITL [REDACTED] 1.5
 2. **PRUDENTIA [REDACTED] ** AgentDojo [REDACTED] HITL
 TCR@0 PRUDENTIA [REDACTED] FIDES [REDACTED] 9% HITL [REDACTED] 1.9
 WASP PRUDENTIA [REDACTED] HITL [REDACTED] 0 [REDACTED] *
 AI [REDACTED] *

PRUDENTIA [REDACTED] IFC [REDACTED] HITL [REDACTED] *
 PRUDENTIA [REDACTED] ** [REDACTED] **2
 AI [REDACTED] ** [REDACTED] Denning, 1976;
 Sabelfeld & Myers, 2003 IFC [REDACTED] AI [REDACTED] Costa [REDACTED] 20
 25 Zhong [REDACTED] 2025 Debenedetti [REDACTED] 2025 IFC [REDACTED] Costa [REDACTED] Costa
 2025 ** [REDACTED] ** [REDACTED] lattice L [REDACTED]
 join [REDACTED] join [REDACTED] Siddiqui [REDACTED] 2025 z [REDACTED]
 x ly [REDACTED] join [REDACTED] _z = x _y [REDACTED]

* ** [REDACTED] L_integrity = {T, U}

2025 **PRUDENTIA**
** ** 2
** plan tool
** Endorsement Approval **
U T
P-T HITL
U 10 TODO P-T
HITL 10
10 HITL
HITL PRUDENTIA
HITL i
expand_variables(ask_endorsement=True) iii ex
expand_variables(ask_endorsement=False) A
** Declassification **
Sabelfeld & Sands, 2009

context-engineering PRUDENTIA IFC
expand_variables IFC
expand_variables IFC
**5 ** IFC
AgentDojo WASP PRUDENTIA IFC Re
Act Yao 2023 IFC 1.
IFC 2. PRUDENTIA **5.1 AgentDojo ** AgentDojo
Debenedetti 2024 Slack Costa
2025 Microsoft Foundry OpenAI GPT-4o
OpenAI LLM IFC
GPT-5 PRUDENTIA ** ** FIDES PRUDEN
TIA i **Basic** iii **Basic-IFC**
Basic iii **FIDES** IFC
HITL GitHub Copilot
Basic 3 IFC Basic-IFC FIDES
HITL Costa 2025 FIDES
LLM TCR
HITL 5
**IFC ** IFC 1
o3-mini o4-mini HITL 2 TCR@k
HITL o3-mini Basic-IFC HITL 32.4 Basic 48.2
1.5 IFC FIDES 18.8 HITL
Basic-IFC 1.7 o4-mini Basic Basic-IFC 2
TCR@k Basic-IFC TCR@0 Basic 9.7% k Basic FIDES TCR@
0 Basic-IFC 10.7% Basic Basic-IFC IFC
FIDES HITL Costa 2025 B 2
** 1 ** Basic FIDES IFC Basic
HITL 1.5-2.6 **PRUDENTIA ** PRUDENTIA
FIDES o4-mini PRUDENTIA 19.2 HITL 73.2%
FIDES 36.8 HITL 75.7% 1.9 Basic PRUDENTIA
o3-mini HITL 2.9 2 PRUDENTIA TCR@0

o3-mini PRUDENTIA 59.1% FIDES 50.1% Bas
ic-IFC 35.5% 23.6% Basic 24.3% 34.8% o4-mini
PRUDENTIA FIDES IF
C PRUDENTIA ** 2 ** PRUDE
NTIA IFC HITL Basic 2.9 FIDES 1.9
**5.2 WASP **
WASP Evtimov 2025 Visual Web Arena Koh 2024
GitLab Reddit 21 GitL
ab 12 Reddit 9 GitLab 2 URL
i GitLab 48 Reddit 36 GPT-4o o1 o
3-mini o4-mini Basic PRUDENTIA TCR@∞ HITL
turns WASP ASR-inter
mediate ** WASP PRUDENTIA **
PRUDENTIA WASP accessibility
tree Chromium, 2021
E 12
click type press goto tab_focus go_back go_forward
P-T hover scroll new_tab close_tab stop
Reddit GitLab
** ** 1 PRUDENTIA Basic WASP Basic PIAs P
RUDENTIA
Basic ASR Reddit 36.1%–61.1% GitLab
14.6%–29.2% Reddit GitLab Basic HITL P
RUDENTIA Basic HITL Basic
PRUDENTIA HITL HITL PRUDENTIA
PRUDENTIA HITL B AgentDojo H
ITL Basic PRUDENTIA TCR@∞
Basic PRUDENTIA PRUDENTIA
PRUDENTIA turns Basic
PRUDENTIA Basic-IFC FIDES WASP PRUDENTIA FIDE
S FIDES Basic-IFC Basic
HITL 1 Basic-IFC FIDES ** 3 ** WASP
PRUDENTIA 0% ASR HITL 0 Basic
**6 **
** HITL ** GitHub Copilot
HITL HITL i ii
confirmation fatigue
Stanton 2016 Seidling 2011 IFC
LLM IFC
HITL IFC
** HITL ** PRUDENTIA HITL IFC
LLM HITL IFC
HITL HITL IFC
** PRUDENTIA **
HIT LLM LLM
TCR@0
3 PRUDENTIA TCR@0 FI
DES 25% PRUDENTIA
HITL 3 PRUDENTIA HITL 2.5
** **
FIDES CaMeL Debenedetti 2025

■ 2 ■

[■■■■■ 2]

■ 3 ■

[■■■■■ 3]

■ 4 ■

[■■■■■ 4]

■ 5 ■

[■■■■■ 5]

■ 6 ■

[■■■■■ 6]

■ 7 ■

[■■■■■ 7]

■ 8 ■

[■■■■■ 8]

■ 9 ■

[■■■■■ 9]

■ 10 ■

[■■■■■ 10]

■ 11 ■

[■■■■■ 11]

■ 12 ■

[■■■■■ 12]

■ 13 ■

[■■■■■ 13]

■ 14 ■

[■■■■■ 14]

■ 15 ■

[■■■■■ 15]

■ 16 ■

[■■■■■ 16]

■ 17 ■

[■■■■■ 17]

■ 18 ■

[■■■■■ 18]

■ 19 ■

[■■■■■ 19]

■ 20 ■

[■■■■■ 20]

■ 21 ■

[■■■■■ 21]

■ 22 ■

[■■■■■ 22]

■ 23 ■

[■■■■■ 23]

■ 24 ■

[■■■■■ 24]

■ 25 ■

[■■■■■ 25]

■ 26 ■

[■■■■■ 26]

■ 27 ■

[■■■■■ 27]

■ 28 ■

[■■■■■ 28]

■ 29 ■

[■■■■■ 29]

■ 30 ■

[■■■■■ 30]

■ 31 ■

[■■■■■ 31]

■ 32 ■

[■■■■■ 32]

■ 33 ■

[■■■■■ 33]

1

1