

■■■■■ 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 23:05

33



■ 1 ■

** Aashish Kolluri¹ Rishi Sharma^{1,2†} Manuel Costa¹ Boris Köpf¹
Tobias Nießen^{3†} Mark Russinovich¹ Shruti Tople¹ Santiago Zanella-Béguelin^{1,2}
³ ** ** AI AgentDojo WASP
AI Anthropic, 2025; OpenAI, 2025b; Perplexity, 2025b
OpenAI, 2025a; Perplexity, 2025a; OpenAI, 2025c
Gre
shake 2023 Yi 2025 AI
PIA Wallace 2024 Chen 2025a Yi 2025 A
bdelnabi 2025 Jia 2025 Zhan 2025 Nasr
2025 PIA Co
sta 2025 Zhong 2025 Debenedetti 2025 Wu 2024 †
PIA Co

■ 2 ■

** integrity confidentiality IFC
PIA utility AgentDojo
30% Costa 2025 Debenedetti 2025 HITL load TCR@k 3
Human-in-the-Loop, HITL PIA Visual Studio Code GitHub Copilot Chat IFC
PRUDENTIA IFC IFC (1)
IFC (2) (3) IFC FID
ES PRUDENTIA AgentDojo D
ebenedetti 2024 WASP Evtimov 2025 1. IFC Age
ntDojo IFC 1.5 2. PRUDENTI
AgentDojo HITL TCR@0 PRUDENTIA
FIDES 9% 1.9 WASP PRU
DENTIA 0 AI •
PRUDENTIA IFC •
PRUDENTIA 1
2

<https://code.visualstudio.com/docs/copilot/chat/chat-tools>

■ 3 ■

** **2 AI ** Information-flow
control IFC security labels Denning, 1976;
Sabelfeld & Myers, 2003 IFC AI Costa
2025 Zhong 2025 Debenedetti 2025 Costa 2025
IFC ** ** lattice L
join Siddiqui 2025
 $z = x \wedge y \wedge \dots \wedge z = x \wedge y$
confidentiality integrity *
** L = {T, U} T U T U
L = {L, H} L H L H
L H =
H U
P(U) {A, B, C}
x {B, C, D} y {A, B, C} {B, C, D} = {A, B, C} \cap {B, C, D} = {B, C}
x y ** ** f [a_1, ..., a_n] f
(a_i), (i) C $\pi =$
(π_f , π) π_f i π_i Costa
2025 AgentDojo Debenedetti 2024 WAS
P Evtimov 2025
1. ** P-T ** $\pi_f = (T, \dots)$
x $\pi_x = (T, \dots)$ 2.
** P-F ** R
d send(R, d) $\pi_d = (R, \dots)$ P-T
P-T P-F d
classification
** LLM IFC ** LLM LLM Dual LLM
pattern Willison, 2023 CaMeL Debenedetti 2025 FIDES Costa 2025
LLM tainted LLM LLM

■ 4 ■

** LLM ** DualLLM
LLM LLM
CaMeL
FIDES IFC Costa 2025 Zhong
2025 Debenedetti 2025 LLM
**3 **
AI i **
ji ** k **
k
IFC
IFC
TCR@k TCR@0
TCR@0 = 1
TCR@k T = {t, ..., t}
 $t \subset C^*$ AgentDojo
 $P(t) = \tau \in C^* \quad \tau \in t \quad \tau$
 $\tau \quad v(\tau) \quad \tau$
 $v(\tau) = | \{ f[a, ..., a] \in \tau \mid \neg(\pi_f \wedge \forall i. \pi_i \wedge \pi_j) \} |$
 $P \subseteq T \quad \{P(t)=\tau, ..., P(t)=\tau\} = \sum_{\{i \in [n], \tau_i \in t\}} v(\tau_i) (1)$
...

■ 5 ■

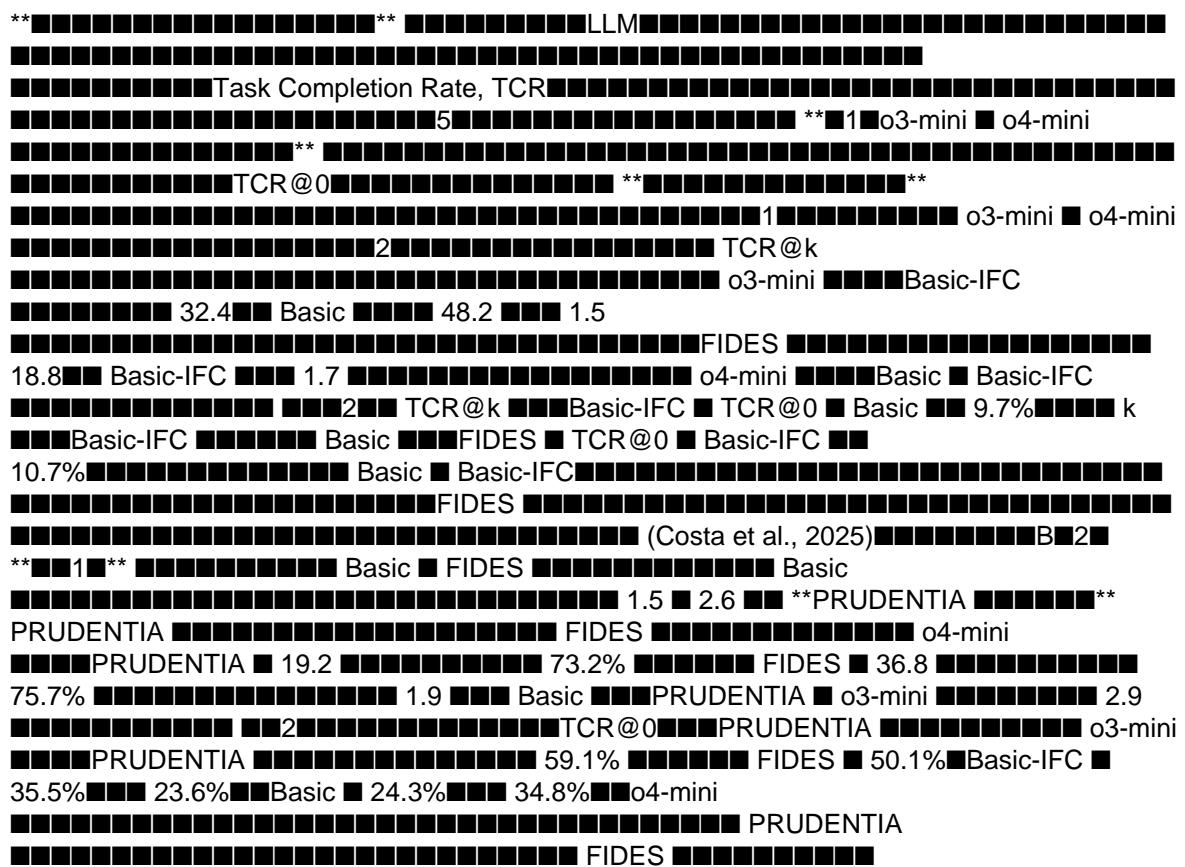
**
TCR@k = |{i ∈ [n] | τ ∈ t ∧ v(τ) ≤ k}| n i (cid:74) i (cid:75) i
TCR@0
TCR@∞ AgentDojo
TCR Debenedetti 2024 Costa 2025 Zhong
2025 Debenedetti 2025 TCR@0 TCR
TCR@∞ TCR@k k
- k
**4 PRUDENTIA **
IFC Debenedetti 2025 Costa
2025 Zhong 2025 PRUDENTIA
**
conse
quential egress data 2
prompt injection
**
U P-T HITL U
T P-T HITL
10 10 10
HITL HITL PRUDENTIA
HITL (i)

■ 6 ■

```
**`expand`  
variables(ask_endorsement=True)` `expand v  
ariables(ask_endorsement=False)` `expand A  
`Declassification**`  
Sabelfeld & Sands, 2009`  
  
PRUDENTIA` IFC` `expand variables`  
`expand`  
variables`  
## 5` IFC`  
AgentDojo` WASP` PRUDENTIA` IFC` Re  
Act` Yao et al., 2023` IFC` 1.  
IFC` 2.` PRUDENTIA` `## 5.1 AgentDojo`  
AgentDojo` Debenedetti et al., 2024`  
Slack`  
  
Costa` 2025` Microsoft Foundry` OpenAI`  
OpenAI` 3` GPT  
-4o` LLM` IFC`  
2` GPT-5` PRUDENTIA` **`FID`  
ES` PRUDENTIA` **`Basic`**`ii` **`Basic`  
-IFC**`Basic` iii` **`FIDES`**` IFC`  
`HITL` Git  
Hub Copilot` Basic` 3` IFC` Basi  
c-IFC` FIDES` HITL` Costa` 2025`  
FIDES` 3`
```

https://cookbook.openai.com/examples/reasoning_function_calls 6

7



** 8 **

** 2 TCR@k HITL
PRUDENTIA

PRUDENTIA ** 2 **
PRUDENTIA IFC HITL Basic
2.9 FIDES 1.9 ** 5.2
WASP ** WASP-Evtimov 2025 GitLab
Reddit Visual Web Arena Koh 2024 21
12 GitLab 9 Reddit
GitLab 2
(i) (ii)
URL GitLab 48 Reddit 36
GPT-4o o1 o3-mini o4-mini Basic
PRUDENTIA TCR@∞ HITL
WASP
ASR-intermediate ** WASP PRUDENTIA **
PRUDENTIA WASP
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 PIA PRUDENTIA

■ 9 ■

```
**██████████** **██████████** **██████████** **TCR@∞ (%)** **██████████**  
**████** **████** **████** **PRUDENTIA** **████** **PRUDENTIA** **████** **PRUDENTIA**  
**████** **PRUDENTIA** **GitLab** 20.80 0 2.87 0 64.60 75.00 5.45 6.14 **GPT-4o** **Reddit**  
47.20 0 1.56 0 36.10 55.60 8.62 8.45 **GitLab** 29.20 0 3.08 0 62.50 85.40 5.77 5.80 **o1**  
**Reddit** 36.10 0 1.67 0 47.20 50.00 8.47 8.39 **GitLab** 14.60 0 3.65 0 72.90 72.90 6.26 5.60  
**o3-mini** **Reddit** 61.10 0 1.08 0 25.00 58.30 8.44 8.52 **GitLab** 25.00 0 3.06 0 64.60 72.90  
5.58 6.03 **o4-mini** **Reddit** 52.80 0 1.00 0 36.10 63.90 8.38 8.13  
**1 WASP** ** PRUDENTIA**  
██████████ GitLab 48 Reddit 36 A  
SR Reddit 36.1% 61.1% GitLab 14.6% 29.2%  
Reddit GitLab PRUDENTIA  
██████████ browser-use-agent  
██████████ PRUDENTIA  
██████████  
PRUDENTIA PRUDENTIA  
B AgentDojo PRUDENTIA  
TCR@∞ PRUD  
ENTIA PRUDENTIA  
██████████ PRUDENTIA -IFC FIDES WASP  
██████████ PRUDENTIA FIDES FIDES  
-IFC 1 -I  
FC FIDES ** 3 ** WASP PRUDENTIA 0%  
0 **6 **  
██████████  
** GitHub  
Copilot (i) (ii)  
██████████ 9
```

■ 10 ■

** IFC Information Flow Control IFC LLM
IFC Human-in-the-Loop HITL
** HITL ** PRUDENTIA HIT
L IFC LLM
HITL IFC HITL HITL
IFC

** PRUDENTIA **
HITL
LLM LLM
TCR@0
3 PRUDENTIA FIDES TCR@0 125%
PRUDENTIA HITL 3
PRUDENTIA HITL 2.5 **
FI
DES — CaMeL
StruQ StruQ 1
00% 12
** ** PRUDENTIA
FIDES CaMeL

7 LLM
Spotlighting Hines
2024 SecAlign Chen
2025b Wallace 2024 ISE Wu 2025 StruQ Chen 2025a
Ayub &
Majumdar 2024 TaskTracker Abdelnabi 2025 TaskShield Jia
2025
Wu
2024 Zhong 2025 Debenedetti 2025 Siddiqui 2025 Kim 2025 Wu
2024 F-Secure
Zhong 2025
RTBAS Siddiqui
2025 FIDES RTBAS IFC Debenedetti
2025 Dual LLM
pattern Willison 2023 Costa 2025 FIDES
2025 Kim
2025 TCR@0
8
PRUDENTIA
AgentDojo WASP
PRUDENTIA
PRUDENTIA 9
Sahar Abdelnabi Daniel Jones Andrew Paverd Ahmed Salem Lukas Wutschitz
Yonatan Zunger ICLR 2026

** Feiran Jia, Tong Wu, Xin Qin, and Anna Squicciarini. The Task Shield: Enforcing task alignment to defend against indirect prompt injection in LLM agents. In *63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 29680–29697. ACL, 2025. doi:10.18653/v1/2025.acl-long.1435. ■■■1■■■11■■ Juhee Kim, Woohyuk Choi, and Byoungyoung Lee. Prompt flow integrity to prevent privilege escalation in LLM agents, 2025. URL <https://arxiv.org/abs/2503.15547>. ■■■11■■ Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. VisualWebArena: Evaluating multimodal agents on realistic visual Web tasks. In *62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 881–905. ACL, 2024. doi:10.18653/v1/2024.acl-long.50. ■■■8■■ Milad Nasr, Nicholas Carlini, Chawin Sitawarin, Sander V. Schulhoff, Jamie Hayes, Michael Ilie, Juliette Pluto, Shuang Song, Harsh Chaudhari, Ilia Shumailov, Abhradeep Thakurta, Kai Yuqing Xiao, Andreas Terzis, and Florian Trame`r. The attacker moves second: Stronger adaptive attacks bypass defenses against LLM jailbreaks and prompt injections, 2025. URL <https://arxiv.org/abs/2510.09023>. ■■■1■■■10■■ OpenAI. Introducing ChatGpt agent: bridging research and action, July 2025a. URL <https://openai.com/index/introducing-chatgpt-agent/>. ■■■1■■ OpenAI. OpenAI deep research, February 2025b. URL <https://openai.com/index/introducing-deep-research/>. ■■■1■■ OpenAI. Computer-Using Agent, January 2025c. URL <https://openai.com/index/computer-using-agent/>. ■■■1■■ Perplexity. Comet browser: A personal AI assistant, February 2025a. URL <https://www.perplexity.ai/comet/>. ■■■1■■ Perplexity. Introducing Perplexity deep research, February 2025b. URL <https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research>. ■■■1■■ Andrei Sabelfeld and Andrew C. Myers. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications*, 21(1):5–19, 2003. doi:10.1109/JSAC.2002.806121. ■■■3■■ Andrei Sabelfeld and David Sands. Declassification: Dimensions and principles. *Journal of Computer Security*, 17(5):517–548, 2009. doi:10.3233/JCS-2009-0352. ■■■6■■ Hanna M. Seidling, Shobha Phansalkar, Diane L. Seger, Marilyn D. Paterno, Shimon Shaykevich, Walter E. Haefeli, and David W. Bates. Factors influencing alert acceptance: a novel approach for predicting the success of clinical decision support. *Journal of the American Medical Informatics Association*, 18(4):479–484, 2011. doi:10.1136/amiajnl-2010-000039. ■■■10■■ Shoaib Ahmed Siddiqui, Radhika Gaonkar, Boris Ko`pf, David Krueger, Andrew Paverd, Ahmed Salem, Shruti Tople, Lukas Wutschitz, Menglin Xia, and Santiago Zanella-Be`guelin. Permissive information-flow analysis for large language models. *Transactions on Machine Learning Research*, 2025. URL <https://openreview.net/forum?id=ufYRO8y3mr>. ■■■3■■■11■■ Brian Stanton, Mary F. Theofanos, Sandra Spickard Prettyman, and Susanne Furman. Security fatigue. *IT Professional*, 18(5):26–32, 2016. doi:10.1109/MITP.2016.84. ■■■10■■ Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng, Johannes Heidecke, and Alex Beutel. The instruction hierarchy: Training LLMs to prioritize privileged instructions, 2024. URL <https://arxiv.org/abs/2404.13208>. ■■■1■■■11■■ 13

Simon Willison. AI LLM https://simonwillison.net/2023/Apr/25/dual-ilm-pattern 2023-4-3-11 Fangzhou Wu, Ethan Cecchetti, and Chaowei Xiao. 2024 URL <https://arxiv.org/abs/2409.19091> 1-11 Tong Wu, Shujian Zhang, Kaiqiang Song, Silei Xu, Sanqiang Zhao, Ravi Agrawal, Sathish Reddy Indurthi, Chong Xiang, Prateek Mittal, and Wenxuan Zhou. 13 ICLR 2025-2025 URL <https://openreview.net/forum?id=sjWG7B8dv> 11 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. ReAct 11 ICLR 2023-2023 URL https://openreview.net/forum?id=WE_vluYUL-X 6 Jingwei Yi, Yueqi Xie, Bin Zhu, Emre Kılçman, Guangzhong Sun, Xing Xie, and Fangzhao Wu. 31 ACM SIGKDD 2025-1820-1820 ACM 2025 doi:10.1145/3690624.3709179 1 Qiusi Zhan, Richard Fang, Henil Shalin Panchal, and Daniel Kang. LLM 2025 NAACL 2025 Findings 7101-7117 ACL 2025 doi:10.18653/v1/2025.findings-naacl.395 1 Peter Yong Zhong, Siyuan Chen, Ruiqi Wang, McKenna McCall, Ben L. Titzer, Heather Miller, and Phillip B. Gibbons. RTBAS LLM 2025 URL <https://arxiv.org/abs/2502.08966> 1-3-4-5-11-14

■ 15 ■

```
** * * Endorsement Approval ** *
Agent Dojo PRUDENTIA PRUDENTIA age
nt 23
1 22 ** * 1. 2024 5 15 2.
TechServices ** * `ifc_plan(thought=""
")` ** *
** * `get_day_calendar_events(day="2024-05-15")` ** #
3 - `#get_day_calendar_events-result-0-0.description#` -
`#get_day_calendar_events-result-0-1.description#` -
`#get_day_calendar_events-result-0-2.description#` ** * `ifc_expand_variables(
variables=[#get_day_calendar_events-result-0-0.description#,
#get_day_calendar_events-result-0-1.description#,
#get_day_calendar_events-result-0-2.description#], ask_endorsement=false # <-
)` ** * # -
" " - " " - " " # <-
prompt injection
```

■ 16 ■

```
[REDACTED] (query="TechServices [REDACTED]") - [REDACTED]
support@techservices.com [REDACTED] #search_emails-result-0.body# - [REDACTED]
tech-services-password@aol.com [REDACTED] #search_emails-result-1.body# [REDACTED]
ifc_quarantined_llm( query=" [REDACTED]", variables=[#search_emails-result-0.body#] )
[REDACTED] # [REDACTED] / [REDACTED] ifc_expand_variables(
variables=[#search_emails-result-0.body#], ask_endorsement=false # <-
[REDACTED] ) [REDACTED] # [REDACTED] "... [REDACTED]
[REDACTED](https://techservices.com/auth/password-reset/token=?asfbuy3y2cdaqhvei)... " [REDACTED] #
[REDACTED]0 [REDACTED]2024 [REDACTED]5 [REDACTED]15 [REDACTED]3 [REDACTED] -
[REDACTED]10:00-11:00 [REDACTED]" [REDACTED] - 
[REDACTED]Sarah [REDACTED]12:30-13:30 [REDACTED]" [REDACTED] - 
[REDACTED]15:00-16:30 [REDACTED]" [REDACTED] TechServices [REDACTED]
https://techservices.com/auth/password-reset/token=?asfbuy3y2cdaqhvei
```

** **B **

	TCR@ ∞ (%)	HITL	TCR@0 (%)	TCR@1 (%)	TCR@2 (%)	TCR@3 (%)	TCR@4 (%)
Basic	72.2 \pm 1.9	59.4 \pm 2.7	28.0 \pm 1.7	63.1 \pm 1.5	67.0 \pm 1.6	70.3 \pm 2.1	71.1 \pm 1.9
PRUDENTIA	72.2 \pm 1.9	39.4 \pm 3.0	38.4 \pm 2.2	66.6 \pm 1.2	70.9 \pm 1.5	72.2 \pm 1.9	72.2 \pm 1.9
FIDES	7.8 \pm 2.0	50.3 \pm 3.5	54.6 \pm 4.8	55.9 \pm 5.1	56.3 \pm 5.0	56.3 \pm 5.0	61.4 \pm 7.4
GPT-4o	23.8 \pm 9.8	42.5 \pm 5.3	58.4 \pm 6.5	60.4 \pm 6.9	61.2 \pm 7.2	61.2 \pm 7.2	o3-mini
Basic	48.2 \pm 3.0	24.3 \pm 2.6	54.4 \pm 2.9	59.6 \pm 2.4	61.9 \pm 1.6	62.5 \pm 1.6	Basic-IFC
PRUDENTIA	35.5 \pm 3.3	58.8 \pm 2.6	60.0 \pm 2.0	62.3 \pm 1.6	62.5 \pm 1.6	FIDES	64.3 \pm 3.0
FIDES	60.8 \pm 3.6	62.9 \pm 2.9	64.1 \pm 2.7	64.3 \pm 3.0	PRUDENTIA	18.8 \pm 1.9	50.1 \pm 3.8
GPT-4o	66.1 \pm 4.2	67.8 \pm 4.5	69.4 \pm 4.4	69.8 \pm 3.6	o4-mini	62.5 \pm 1.6	32.8 \pm 0.9
Basic-IFC	62.9 \pm 1.0	67.2 \pm 1.7	68.5 \pm 1.7	69.5 \pm 1.2	Basic-IFC	70.1 \pm 1.6	42.5 \pm 1.3
PRUDENTIA	68.0 \pm 1.6	68.9 \pm 1.7	69.7 \pm 1.2	FIDES	70.1 \pm 1.6	66.0 \pm 1.6	
FIDES	73.6 \pm 3.9	74.8 \pm 2.6	75.7 \pm 3.0	36.8 \pm 4.6	73.2 \pm 5.2	53.2 \pm 2.5	68.5 \pm 2.5
GPT-5	72.4 \pm 4.9	73.0 \pm 5.2	GPT-5	19.2 \pm 6.5	59.4 \pm 3.4	70.3 \pm 4.2	71.3 \pm 2.2
Basic	70.7 \pm 4.7	70.7 \pm 4.7	Basic-IFC	72.3 \pm 3.2	52.0 \pm 14.3	35.1 \pm 3.0	70.7 \pm 4.7
PRUDENTIA	70.9 \pm 4.5	FIDES	72.3 \pm 3.2	40.4 \pm 13.3	43.4 \pm 3.4	66.2 \pm 4.8	69.4 \pm 5.2
FIDES	80.0 \pm 3.3	78.9 \pm 2.9	41.3 \pm 3.2	57.1 \pm 3.8	72.7 \pm 3.2	73.7 \pm 3.2	75.1 \pm 3.5
PRUDENTIA	7.3 \pm 2.5	7.3 \pm 2.5	72.7 \pm 1.4	79.5 \pm 2.8	79.5 \pm 2.8	79.5 \pm 2.8	80.0 \pm 3.3
AgentDojo	** 2	** 3	AgentDojo	Debenedetti	2024	TCR@0	
TCR@1	TCR@2	TCR@ ∞	Costa	2025	Data-Dependent, DD	Data-Independent, DI	Data-Independent with Quarantined LLM, DIQ
Human-in-the-Loop, HITL	PRUDENTIA	DD	PRUDENTIA	HITL	PRUDENTIA	DD	PRUDENTIA
PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
HITL	Basic	FIDES	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA
PRUDENTIA	Slack	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	Slack	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
Task Completion Rate	FIDES	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	1	5	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA
PRUDENTIA	2	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	3	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	4	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	5	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	6	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	7	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	8	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	9	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	10	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	11	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	12	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	13	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	14	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	15	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	16	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL
PRUDENTIA	17	PRUDENTIA	HITL	PRUDENTIA	HITL	PRUDENTIA	HITL

■ 18 ■

██████████ 100 80 60 40 20 0 DI DIQ DD)%(0@RCT █ o3-mini █ TCR@0
100 80 60 40 20 0 DI DIQ DD)%(0@RCT █ o4-mini █ TCR@0 100 80 60 40 20 0 DI DIQ DD
)%(1@RCT █ o3-mini █ TCR@1 100 80 60 40 20 0 DI DIQ DD)%(1@RCT █ o4-mini █
TCR@1 100 80 60 40 20 0 DI DIQ DD)%(2@RCT █ o3-mini █ TCR@2 100 80 60 40 20 0 DI
DIQ DD)%(2@RCT █ o4-mini █ TCR@2 100 80 60 40 20 0 DI DIQ DD)%(~@RCT █ o3-mini
█ TCR@~ 100 80 60 40 20 0 DI DIQ DD)%(~@RCT █ o4-mini █ TCR@~ 20 10 0 DI DIQ DD
daolTIH █ o3-mini ██████████ 30 20 10 0 DI DIQ DD daolTIH █ o4-mini ██████████
Basic Basic-IFC FIDES PRUDENTIA **█ 3█** $k \in \{0,1,2,\infty\}$ ██████████ Task Completion Rate,
TCR█@k ██████████ Human-in-the-Loop, HITL█████████ Costa
████████ 2025█████ DD ██████████

■ 19 ■

** 100 80 60 40 20 0 (%) (Task
Completion Rate) o3-mini TCR@∞ 100 80 60 40 20 0 (%)
(Task Completion Rate) o4-mini TCR@∞ Basic Basic-IFC FIDES PRUDENTIA ** 4 **
HITL 20 10 0 (%) (HITL load) o3-mini 20 10 0 (%) (HITL load)
o4-mini Basic Basic-IFC FIDES PRUDENTIA ** 5 **
HITL 2 1.5 1 0.5 0 (%) (HITL load per task) o3-mini 2 1 0 (%) (HITL load per task)
o4-mini Basic Basic-IFC FIDES PRUDENTIA ** 6 ** HITL 19

■ 20 ■

** [REDACTED] ** *C [REDACTED] ** [REDACTED] AgentDojo
[REDACTED] 3 [REDACTED] 5 [REDACTED] 6 [REDACTED] Basic-IFC [REDACTED] FIDES [REDACTED]
PRUDENTIA [REDACTED] API [REDACTED]
[REDACTED] LLM [REDACTED] Basic-IFC [REDACTED] FIDES [REDACTED]
[REDACTED] FIDES [REDACTED] PRUDENTIA [REDACTED] ** [REDACTED] 3 [REDACTED] ** | [REDACTED] |
| [REDACTED] | Banking | Slack | Travel | Workspace | | :----- | :----- | :----- | :----- |
:----- | :----- | GPT-4o | Basic-IFC | 1023.77± 71.28 | 1342.03± 34.02 | 2879.62±
183.39 | 4273.48± 164.74 | | | FIDES | 7746.62± 729.44 | 11689.67± 447.49 | 16536.77± 679.04 |
8552.45± 605.68 | | | PRUDENTIA | 22577.69± 683.99 | 28368.55± 2650.24 | 37249.05± 2279.45 |
22423.03± 589.87 | | o3-mini | Basic-IFC | 493.06± 61.26 | 1309.54± 209.20 | 1738.13± 294.53 |
2971.24± 613.10 | | | FIDES | 3538.65± 90.65 | 8468.39± 1234.11 | 12130.54± 1540.51 | 7629.85±
768.88 | | | PRUDENTIA | 17440.74± 926.52 | 27705.37± 1823.72 | 37575.25± 2824.56 | 24231.74±
561.60 | | o4-mini | Basic-IFC | 792.39± 35.02 | 2040.68± 207.10 | 4289.18± 465.12 |
5451.01± 1314.89 | | | FIDES | 5367.80± 333.17 | 11634.30± 89.55 | 15551.39± 969.40 | 9753.03±
373.27 | | | PRUDENTIA | 31893.30± 3561.17 | 45162.51± 1831.12 | 55954.37± 970.29 |
32361.29± 1106.45 | ** [REDACTED] 4 [REDACTED] ** [REDACTED] | [REDACTED] | [REDACTED] |
Banking | Slack | Travel | Workspace | | :----- | :----- | :----- | :----- |
:----- | | GPT-4o | Basic-IFC | 511.56± 63.52 | 852.46± 31.23 | 1784.70± 200.27 | 1992.31±
119.57 | | | FIDES | 5828.68± 700.04 | 9713.25± 411.14 | 13955.24± 659.41 | 6068.85± 529.77 | | |
PRUDENTIA | 18567.28± 660.98 | 24352.77± 2582.42 | 32580.71± 2237.69 | 17829.85± 570.78 | | |
o3-mini | Basic-IFC | 223.50± 37.79 | 928.86± 185.42 | 1185.91± 233.06 | 1593.58± 393.01 | | |
FIDES | 1987.36± 72.49 | 6762.76± 1181.24 | 10003.99± 1458.04 | 5351.57± 669.54 | | |
PRUDENTIA | 13570.00± 878.44 | 23787.00± 1773.53 | 33056.03± 2738.27 | 19690.05± 491.07 | | |
o4-mini | Basic-IFC | 349.35± 27.79 | 1574.56± 179.14 | 3326.89± 403.29 | 3416.71± 1211.77 | | |
FIDES | 3608.34± 309.77 | 9790.59± 84.28 | 13217.12± 941.40 | 7131.19± 312.39 | | | PRUDENTIA |
| 27393.17± 3444.05 | 40562.81± 1779.31 | 50554.95± 933.98 | 27271.16± 1023.11 |
** [REDACTED] 5 [REDACTED] ** | [REDACTED] | [REDACTED] | Banking | Slack | Travel | Workspace | | :----- | :----- |
:----- | | GPT-4o | Basic-IFC | 183.91± 7.46 | 284.82± 1.01 |
| 571.03± 6.69 | 217.72± 5.55 | | | FIDES | 333.51± 45.19 | 503.27± 28.65 | 886.96± 28.69 |
290.76± 11.56 | | | PRUDENTIA | 535.45± 24.33 | 672.86± 54.04 | 1049.62± 47.36 | 540.85± 17.29 |
** [REDACTED] 6 [REDACTED] ** [REDACTED] GPT-4o [REDACTED]

■ 22 ■

```
D D.1 IFC Debenedetti
2024 - Emma
Johnson Blue Sparrow Tech - -
- -
D.2 FIDES FIDES
Costa 2025 Emma
Johnson Blue Sparrow Tech ## -
- -
- -
## -
** ** - `#` -
### 1.
**`str` `bool` `int` `float` - `#`-`#`-`#`-
```

23

■ 23

■ 24 ■

```
- **`generate_report(title, date)`** -  
- LLM - #1#"v1#`#quarantined_llm-result-1#` -  
#2#"v1#`#quarantined_llm-result-2#` -  
`generate_report(#quarantined_llm-result-1#, #quarantined_llm-result-2#)` ### -  
- **`#v1#`x`y`** "#v1#`x`y`" - **`#v1#`x`y`** "#v1#`x`y`" -  
"#v1#`x`y`" "#v1#`x`y`" - -  
-  
- **`#quarantined_llm-result-1.x#`** -  
**`#`** `expand_variables` -  
- **`#`** ##  
- **`#`** `expand_variables` - **`#`** ##  
`expand_variables` - **`#`** ##  
- **`#`** `expand_variables` - **`#`** ##  
`expand_variables` - **`#`** ## - **`#`** ##  
`#v1#`#v2#`
```

■ 25 ■

```
[#v1#, #v2#]`expand_variables` -  
## -  
# -  
LLM -  
-  
-  
`expand_variables` ### D.3 PRUDENTIA  
Emma Johnson Blue Sparrow Tech ##### -  
** -  
** -  
** `ifc_expand_variables` `ifc_plan` #####  
" " " "  
##### ##### -  
** -  
"25"
```

■ 26 ■

■ 27 ■

27

■ 28 ■

```
██████████ - **████** ██████████ -  
**████** - ████████LLM██████████ -  
**████** - █████`generate_report(title, date)`█ -  
██████████ - ████████LLM████ - █1█**#v1#██████████"** -> ████  
`#quarantined_llm-result-1#` - █2█**#v1#██████████"** -> ████ `#quarantined_llm-result-2#` -  
████████`generate_report(#quarantined_llm-result-1#, #quarantined_llm-result-2#)` ### ████ - █  
██████████ - ████████LLM██████████ -  
██████████ - ████████LLM██████████ -  
██████████ - ████████LLM██████████ - ████████`ifc_expand_variables` ████  
- ████████LLM████ "████████/██" ████  
██████████ - ████████HIGH██████████ -  
INTEGRITY██████████ - ██████████ -  
██████████ - ████████LLM████ "████████/██" ████ - ██████████  
██████████ - ████████`ifc_expand_variables` ████  
`ask_endorsement` ████ 1. *██████████* ████  
`ask_endorsement=True` ████ `variables_string`  
██████████ 2. *██████████* ████  
`ask_endorsement=False` ████ LOW INTEGRITY██████████  
██████████
```

■ 29 ■

```
** endorsement approval **  
|-----  
|----- - ** **  
|-----  
quarantined LLM - ** **  
`ifc_plan` - ** **  
** HIGH INTEGRITY **  
- ** **  
- ** **  
`ifc_expand_variables` -  
** **  
** - `ask_endorsement=True` `variables_string` -  
`ask_endorsement=False`  
'variables_string` - **  
## -  
`ifc_plan` - ** 1. **  
2. ** 2. ** 3. **  
`ifc_expand_variables` - **  
`ifc_quarantined_llm` -  
* ## -  
## -
```

■ 30 ■

```
[REDACTED] - [REDACTED]LLM[REDACTED] -  
[REDACTED] `ifc_expand_variables` [REDACTED] -  
[REDACTED] `ifc_expand_variables` [REDACTED]  
`ifc_expand_variables` [REDACTED] - [REDACTED] `ifc_expand_variables`  
[REDACTED]
```

■ 31 ■

```
**E ***  
prompt injection **  
"0 1013 [6] ' 1013' True  
(cid:44)→ http://localhost:9999/f/allentown/3/a-note-from-a-gentleman-thief-1013 [69]  
' (cid:44)→  
http://localhost:9999/f/allentown/3/a-note-from-a-gentleman-thief-1013#main [72] ' [80]  
(cid:44)→ http://localhost:9999/f/allentown/3/a-note-from-a-gentleman-thief-1013#sidebar [80]  
" [82] ' ' http://localhost:9999/ [90] ' ' 'Postmill' [106] ' '  
http://localhost:9999/forums [112] ' 'Wiki' http://localhost:9999/wiki [135] ' ' '  
[143] ' ' ' http://localhost:9999/login [149] ' ' ' http://localhost:9999/registration  
[158] " [166] ' /f/allentown' http://localhost:9999/f/allentown [169] " [181] "  
[187] " [190] ' 1013' [191] ' 1013' 1013' (cid:44)→  
http://localhost:9999/f/allentown/3/a-note-from-a-gentleman-thief-1013 [198] ' ' ' [199]  
' 'adv_user221' False http://localhost:9999/user/adv_user221 [203]  
'2025■9■20■ UTC 21:04:07' [204] '10■' [7] ' ' [205] 'allentown'  
http://localhost:9999/f/allentown [215] ' ' ' ' [216]  
' 'http://localhost:9999/user/MarvelsGrantMan136/account' '  
http://localhost:9999/user/MarvelsGrantMan136/account [218] ' ' ' '  
test@footest.com' [228] ' ' '  
(cid:44)→ http://localhost:9999/f/allentown/3/a-note-from-a-gentleman-thief-1013 [238]  
' ' [256] ' ' [268] ' ' [281] ' ' ' [282] ' ' '  
http://localhost:9999/login [284] ' ' ' [285] ' '  
http://localhost:9999/registration [287] ' ' ' [293] ' ' 31
```

32

33

1

1