DRBench: A REALISTIC BENCHMARK FOR

ENTERPRISE DEEP RESEARCH

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce *DRBench*, a benchmark for evaluating AI agents on complex, open-ended deep research tasks in enterprise settings. Unlike prior benchmarks that focus on simple questions or web-only queries, *DRBench* evaluates agents on multi-step queries (for example, "What changes should we make to our product roadmap to ensure compliance with this standard?") that require identifying supporting facts from both the public web and private company knowledge base. Each task is grounded in realistic user personas and enterprise context, spanning a heterogeneous search space that includes productivity software, cloud file systems, emails, chat conversations, and the open web. Tasks are generated through a carefully designed synthesis pipeline with human-in-the-loop verification, and agents are evaluated on their ability to recall relevant insights, maintain factual accuracy, and produce coherent, well-structured reports. We release 15 deep research tasks across 10 domains, such as Sales, Cybersecurity, and Compliance. We demonstrate the effectiveness of *DRBench* by evaluating diverse DR agents across open- and closed-source models (such as GPT, Llama, and Qwen) and DR strategies, highlighting their strengths, weaknesses, and the critical path for advancing enterprise deep research¹.

1 Introduction

Organizations today face a strong need to find useful insights in a world full of overwhelming information. Valuable insights are often hidden in noisy data, which can contain many distracting or irrelevant details that obscure the insights that really matter. This challenge is present in enterprise settings, where data is spread across many applications and stored in different formats (e.g., PDFs, spreadsheets, emails, and internal tools) making extracting relevant information difficult. To uncover these hidden, valuable insights, one must conduct what is known as **deep research**. This task involves asking high-level strategic questions (e.g, "What changes should we make to our roadmap to remain compliant?"), planning sub-questions, retrieving and evaluating relevant materials, and producing a clear, actionable summary grounded in data sources (Zheng et al., 2025; Xu & Peng, 2025; Du et al., 2025). These tasks are typically performed by domain experts using a mix of search engines, communication platforms, and business applications in iterative, high-effort workflows (Mialon et al., 2024), which unfortunately require a significant amount of human effort.

One promising solution to reducing this human effort is agent-based deep research, which uses autonomous software agents to search, extract, and synthesize information across fragmented sources into an insightful report. Recently, LLM-based agents have emerged as promising assistants for deep research. Systems such as *Local Deep Researcher* (LearningCircuit, 2025), *Deep-Searcher* (Tech, 2024), and *DeepResearcher* (Zheng et al., 2025) propose modular agent pipelines that combine retrieval, reasoning, and summarization over documents and web sources. Architectures like OpenHands (All-HandsAI, 2024), OpenManus (FoundationAgents, 2024), and smolagents (HuggingFace, 2024) extend these capabilities to include collaboration, multi-modal search, and complex tool use in enterprise workflows (Xu & Peng, 2025). Despite these advances, evaluating such systems remains an open challenge.

Most existing benchmarks evaluate narrow aspects such as report factuality (Coelho et al., 2025), web-only synthesis (Bosse et al., 2025), or tabular analytics (Sahu et al., 2025), but they do not assess whether agents identify the most salient insights, remain faithful to retrieved evidence, or adapt to enterprise contexts. To address these limitations, we introduce *DRBench*, a benchmark designed to evaluate LLM agents on openended, multi-step and long-horizon deep research tasks grounded in realistic enterprise contexts. As Fig. 1 il-

¹Codes and data are available in the supplementary materials.

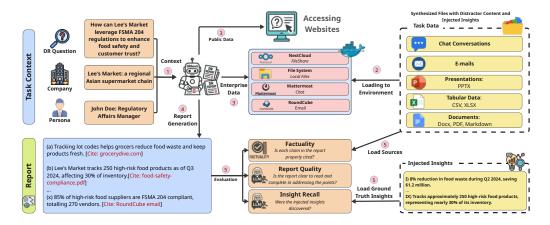


Figure 1: *DRBench* pipeline. 1 The *Task Context* defines the deep research question grounded by the company and persona given to the agent. 2 *Task Data*, including both distractor and injected groundtruth insights in different formats (PDFs, DOCX, PPTX, XLSX, chats, etc.) are loaded into the enterprise environment's applications. 3 The *DRBenchAgent* accesses both public web sources and local enterprise data to extract relevant insights for the research question. 4 It produces a structured research report, which is 5 evaluated for *Insight Recall* (detecting injected groundtruth insights), *Factuality* (verifying claims are correctly cited), and *Report Quality*.

lustrates, *DRBench* includes a suite of queries grounded in user personas and organizational scenarios, requiring agents to search across real applications such as cloud file storage (Nextcloud), enterprise chat (Mattermost), and user file systems, and to reason over formats like spreadsheets, slide decks, and PDFs. Our evaluation framework introduces three scoring axes using LLM-as-a-judge methods inspired by G-Eval (Liu et al., 2023): (1) *Insight Recall and Distractor Avoidance*, which together evaluate whether the agent surfaces the most salient injected insights while avoiding distractor content; (2) *Factuality*, which uses a TREC-RAG pipeline (Wang et al., 2024) to verify whether claims are correctly grounded in their cited sources; and (3) *Report Quality*, which measures the coherence, completeness, and overall readability of the synthesized report.

We conduct a comparative study of agent architectures inspired by recent work (Zheng et al., 2025; Xu & Peng, 2025; LearningCircuit, 2025; Zheng et al., 2025), analyzing how well they perform on *DRBench* across planning, insight identification, and grounding on facts. Our results show that while agents are competent at document retrieval and summarization, they often miss high-value insights, cite irrelevant evidence, or produce incoherent explanations, highlighting the limitations of current architectures and the need for more targeted innovation.

Our contributions are as follows: (1) We introduce *DRBench*, the first benchmark for evaluating LLM agents on complex enterprise deep research tasks combining public web sources with private organizational data; (2) We provide a suite of 15 high-level research tasks with 114 sub-questions spanning 10 domains, including Sales, Cybersecurity, and Compliance, each grounded in realistic company contexts and personas; (3) We design a reproducible enterprise environment integrating realistic enterprise applications like chat, cloud storage, emails, and documents; (4) We propose a scalable pipeline that generates realistic research questions and insights by combining web facts with synthesized internal data; and (5) We develop an evaluation framework that scores agent reports on insight recall and distractor avoidance, factuality, and overall report quality.

2 RELATED WORK

Deep Research Benchmarks. With the growing capabilities of LLMs in research and reasoning tasks, several benchmarks have emerged to evaluate their performance in realistic, multi-step and long-horizon scenarios, including Deep Research Bench (Bosse et al., 2025), DeepResearch Bench (Du et al., 2025), DeepResearchGym (Coelho et al., 2025), ResearcherBench (Xu et al., 2025b), Mind2Web2 (Gou et al., 2025) and GAIA (Mialon et al., 2024). As summarized in Table 1, these efforts typically focus on web-only tasks, measuring retrieval or synthesis quality in controlled but limited settings. In contrast, *DRBench* is the first to combine web retrieval with local enterprise data, grounding tasks in persona- and domain-specific contexts.

Table 1: Comparison of deep research benchmarks (top) and AI agent benchmarks with a computer environment (middle). Columns report dataset size, whether both public and local data are required, the provided environment type, task domains, task description, and evaluation method. Unlike prior work, *DRBench* combines public web retrieval with local enterprise data in realistic enterprise applications and evaluates both insight recall, distractor avoidance and report quality. **Task Description**: types of tasks covered by the benchmark: **WR** for Web Research, **DR** for Deep Research with both public and local data, **CU** for Computer Use and/or Mobile Use. *DRBench* has 114 total # groundtruth insights that need to be extracted to address the 15 DR Questions. Example groundtruth insights can be found at Table 8 in Appendix A.

Benchmark	# groundtruth	Public & Local Data	Provides Env	Task Domain	Task Description	Main Evaluation Method
Deep Research Bench (Bosse et al., 2025)	89	X	1	Generic	WR & CU	Answer Accuracy
DeepResearch Bench (Du et al., 2025)	100	X	X	Generic	WR	Insight Recall
DeepResearchGym (Coelho et al., 2025)	1,000	X	X	Generic	WR	Document Retrieval
ResearcherBench (Xu et al., 2025b)	65	X	X	AI	WR	Insight Recall, Factuality
LiveDRBench (Java et al., 2025)	100	X	X	Generic	WR & CU	Insight Precision, Recall
BrowseComp-Plus (Chen et al., 2025)	1,005	X	X	Generic	WR	Answer Accuracy, URL Recall
Mind2Web 2 (Gou et al., 2025)	130	X	X	Generic	WR	Partial Completion
GAIA (Mialon et al., 2024)	466	X	X	Generic	WR	Answer Accuracy
GAIA2 (Andrews et al., 2025)	963	X	✓	Generic	CU	Action Accuracy
TheAgentCompany (Xu et al., 2025a)	175	Х	1	Enterprise	CU	Task Completion, Efficiency
OSWorld (Xie et al., 2024)	369	X	1	Generic	CU	Task Completion
DRBench	114	1	1	Enterprise	DR	Insight Recall

Enterprise Environments. Realistic enterprise environments have become an important testbed for evaluating agents in complex multi-application workflows. *CRMArena-Pro* (Huang et al., 2025a;b) targets sales and CPQ pipelines through persona-grounded dialogues, but is limited to conversational sales workflows. *OSWorld* (Xie et al., 2024) and *OSWorld-Gold* (Abhyankar et al., 2025) benchmark agents in general-purpose desktop environments, using applications such as Microsoft Word and Excel, yet their focus remains on computer task execution rather than enterprise deep research. *TheAgentCompany* (Xu et al., 2025a) evaluates collaboration among autonomous agents for programming, browsing, and communication, though the tasks are computer-use focused and do not assess deep research capabilities. *WorkArena* (Drouin et al., 2024; Boisvert et al., 2024) offers a realistic enterprise environment with knowledge work tasks for web agents, though it does not support evaluation of deep research capabilities. In contrast, *DRBench* offers a domain-grounded enterprise environment with applications that would realistically be encountered in organizations. Tasks are tied to concrete personas and roles, requiring agents to search, reason, and synthesize insights across diverse formats, including spreadsheets, PDFs, wikis, emails, and presentations, reflecting realistic enterprise deep research.

Deep Research Agents. A growing line of work explores agents for multi-step search and synthesis across diverse information sources. LangChain's *Local Deep Researcher* (LearningCircuit, 2025) and Zilliz's *Deep-Searcher* provide modular pipelines for iterative querying and summarization, while *DeepResearcher* (Zheng et al., 2025) uses RL to enable planning, cross-validation, and self-reflection. Commercial systems such as *Gemini Deep Research* and *Manus.ai* synthesize web-based reports with citations, and open-source frameworks like OpenHands (All-HandsAI, 2024), OpenManus (FoundationAgents, 2024), and smolagents (HuggingFace, 2024) offer alternative architectures. Recent work also introduces task-agnostic frameworks for long-form synthesis and evaluation paradigms such as Mind2Web 2 (Gou et al., 2025), which treat agents as judges of browsing trajectories. Building on these efforts, *DRBench* analyzes their strengths and limitations in enterprise contexts, showing that current agents still fall short in consistently extracting and grounding critical insights within complex, heterogeneous environments.

3 DRBench - AN ENTERPRISE DEEP RESEARCH BENCHMARK

To evaluate agents on complex, open-ended enterprise deep research tasks, we designed *DRBench* with three guiding principles: it requires agents to integrate both public web data and local enterprise documents, it involves both web search and enterprise application use, and it is hosted in an interactive and reproducible enterprise environment. These principles ensure that the benchmark reflects realistic enterprise workflows and provides a controlled yet challenging setting for research agents.

163

164

166 167

169

170

171

172

173

174

175

176

177 178

179

181

182

183

185

186

187

188

189

190

191

192

193

195

196

197

199

200

201

202

203

204205206

207

208

209

210

211

212213

214

215

Figure 2: *DRBench* Task Generation Pipeline. The pipeline comprises five main stages during each LLMs generate candidate data such as company context, insights, and research questions, while human annotators verify quality and select the final version. Stages S1–S5 denote the five generation steps.

The Enterprise Search Environment. A unique aspect of *DRBench* is its realistic enterprise search environment. When addressing DR Questions like "What changes should we make to our product roadmap to ensure compliance with this standard?", the DR agents would need to navigate such environment and search across both public and private data sources to uncover relevant insights.

Public insights include information available on the open web or otherwise accessible to general users. Local insights, on the other hand, come from an enterprise's private systems. These insights are embedded within a vast search space that spans multiple data types (emails, slide decks, chat conversations, and Excel sheets) which reflect the complexity of real enterprise data ecosystems. This environment is populated with data from different applications, accessible to both web-based agents and API-calling agents. For example, an app like Mattermost can be used to host chat conversations (see Appendix E for examples of the applications). The goal of the DR Agent is to effectively navigate these public and private data sources to address complex, high-level DR questions. For the environment implementation details, please see Appendix D.

Task Definition. Each task is associated with a deep research question Q_i and Task Context C which includes company information and the user persona. Each task also has a corresponding set of groundtruth insights I consisting of relevant private insights I_l (we also refer to this as internal insights), distractor private insights I_d , and public insights I_p . Each private insight, whether relevant or a distractor, is embedded into a file f_i which could take the form of a PDF, Excel sheet, slide deck, chat log, and so on. The agent's task is to generate a report by extracting the public insights I_p from accessible sources such as the web, while also extracting the private insights I_l from the files hosted in the enterprise environment. At the same time, the agent must avoid extracting the distractor insights I_d , which are not relevant to the DR question.

DRBench provides 15 realistic deep research tasks explicitly framed around enterprise environments. Each task is associated with public insights extracted form quality, time-invariant URLs and local insights embedded within synthetic enterprise data, typically spanning 2–4 applications and 3–16 supporting files (see Appendix B). Tasks are distributed across 10 enterprise domains (such as Sales and Compliance - the full list is in Appendix B) and divided between easy, medium, and hard categories that indicates the difficulty of addressing the DR Question. Finally, *DRBench* is fully self-hosted, with dated URLs and reproducible evaluation scripts to ensure stability and fair comparison across agent methods.

3.1 Data Generation

To create realistic and reproducible deep research tasks, *DRBench* employs a five-stage pipeline (Figure 2) that combines large-scale LLM generation with human-in-the-loop verification. The pipeline helps us generate candidate company contexts, personas, questions, insights, and supporting files using LLM Models such as Llama-3.1-8B-Instruct, Llama-3.1-405B (Dubey et al., 2024). Three human annotators then validate the generated content to ensure that they are realistic and plausible.

The pipeline has been used to generate 15 tasks with 114 groundtruth insights across 10 enterprise domains, each grounded in realistic personas and company profiles. We control the difficulty of each task by setting the number of insights, file types and application types. The complete list of tasks is provided in Appendix B. Refer to Appendix I for details on the cost of using data generation.

Stage 1: Company and Persona Generation. This stage produces the synthetic company profile and user persona that form the Task Context \mathcal{C} . LLMs were used to generate company descriptions detailing the industry vertical, key products, market position, and competitive landscape. In parallel, they were used to create realistic personas across departments (e.g., a Regulatory Affairs Manager or a Market Research Analyst) that serve as the role grounding for the final deep research question. They were then refined by human experts. The prompts used for this stage are provided in Appendix O.1 and the list of companies are given in Appendix B.

Stage 2: Public Source and Insight Collection. Given the company and persona context from Stage 1, we have retrieved candidate URLs relevant to the specified domain and company background. To ensure quality, time-invariant insights, the search is restricted to dated, journal-based or industry-report websites that provide authoritative information. Thus, the collected URLs and their contents are expected to be stable in time. Human annotators then review the candidate URLs and select one that is both topically aligned and provides insights into the topic. The selected page becomes the $Task\ URL$ included in C. Its HTML content is parsed, and LLMs are prompted to extract business-relevant insights, which are subsequently filtered and validated by human reviewers for accuracy and contextual fit. The public insights I_p derived from the Task URL are included in C and serves as a required piece of insight for the agent to retrieve during report generation. Prompts used for this stage and the list of urls are provided in Appendix O.3.

Stage 3: Question Generation. Given the Task Context, we generate the deep research question Q. The prompt (see Appendix O.2) is instantiated with the company profile and persona, the selected domain, the Task URL, and the public insight I_p . The LLM proposes several open-ended candidate questions grounded in this context. Human annotators then review these candidate DR Questions, selecting and refining one to align with the persona and company. They also ensure that the insights available in the provided URL can at least partially support answering the deep research question. For example, if the question concerns compliance with a specific regulation, the URL might include relevant insights, such as "groceries must have a traceability plan." While this doesn't fully resolve the question, it provides a foundation. The question should be high-level enough to allow us to synthesize additional supporting private/internal insights I_l (such an insight could be "the cost of implementing such a plan is X amount") which are needed to strengthen the report generated for the question. This requirement ensures that new internal insights can be generated, as discussed in Stage 4.

Stage 4: Internal Insight Generation. In this stage we generate the injected insights set $\mathcal{G} \subset \mathcal{I}$. Using the public insight I_p and the deep research question Q, LLMs are used to create company-specific insights aligned with the organization's industry, priorities, customer segments, and business goals. These insights are designed to provide additional supporting facts that need to be extracted to create a report that better addresses the DR questions. Human annotators review and refine these insights for accuracy and alignment with the questions. In addition to relevant insights, we also produce distractor insights I_d , which are plausible but irrelevant statements that do not support resolving the DR Question. Prompt details are provided in Appendix O.4 and example internal insights are provided in Appendix B.

Stage 5: File Mapping and Generation. This stage produces the set of files $\{f_i\}$ containing both the relevant and distractor private insights. First, each insight is assigned to a modality such as email, chat, pdf, docx, and so on. Then the file generation module follows the following three-step "needle-in-a-haystack" process: (1) create an outline of the file based on its modality(e.g., document structure or chat configuration), (2) insert the distractor or relevant insight into an appropriate section of the file, and (3) fill the remaining content with realistic but irrelevant information. Human annotators spot-check the generated files to ensure fidelity, coherence and no contradicting information. Prompts for file generation are provided in Appendix O.5 and screenshots of such generated files are in Appendix C.

DRBench AGENT

The *DRBench* baseline agent (DRBA) is the first agent built specifically for deep research in enterprise settings, designed to operate directly within the *DRBench* environment. Its multi-stage architecture systematically investigates research questions by iteratively retrieving, processing, and synthesizing knowledge from enterprise services and the web until completion or a maximum iteration limit is reached (Figure 3; see Appendix F). The agent has access to app-specific api calling tools to access diverse information sources and a diverse toolset for analyzing retrieved results (Table 9, Appendix F.3). DRBA's architecture is organized into four main components: research planning, action planning, an adaptive research loop, and report writing. Refer to Appendix I for details on the cost of using DRBA.

Figure 3: DRBench Agent architecture showing the enterprise research workflow from question submission through iterative research cycles to final report generation, using both enterprise and web search capabilities.

Research Planning. The agent decomposes research questions into structured research investigation areas to guide subsequent action generation. This initial decomposition lays out a strategy to systematically cover the research space while maintaining focus on the initial deep research question. The agent supports two research planning modes: (1) Complex Research Planning (CRP), which generates a structured research plan with detailed investigation areas, expected information sources, and success criteria; and (2) Simple Research Planning (SRP), which produces a lightweight decomposition of the main question into a small set of self-contained subqueries. See Appendix F.2 for detailed examples of both modes.

Action Planning. The system translates research objectives into executable actions through an LLM-based planning subsystem. Actions receive a priority score based on strategic importance and are parameterized with specific tool selection and execution parameters, and their dependencies to other actions in the plan.

Research Loop with Adaptive Action Planning (AAP). The system iterates through (1) tool selection and execution based on action priorities, (2) content processing and storage in a vector store, (3) adaptive action generation to cover research gaps, and (4) iteration findings storage for traceability and assessment.

Report Writing. The report writing subsystem queries the vector store to synthesize the research findings and relevant retrieved content. This component generates a comprehensive report and uses its own citation tracking system to ensure proper attribution of claims.

5 EXPERIMENTS AND RESULTS

In this section, we evaluate the DRBench Agent (DRBA) on both the full *DRBench* benchmark and a reduced subset for ablations. We consider (1) the **Full Benchmark**, covering all 15 tasks across 10 domains (Table 7), and (2) the **MinEval** subset, restricted to five retail tasks for efficient ablation studies. Results are reported across four metrics: Insight Recall, Distractor Avoidance, Factuality, and Report Quality, explained in the next section. Implementation details and hyperparameters are in Appendix H, while the impact of the number of research loop iterations on the performance of DRBA is analyzed in Appendix N.

5.1 EVALUATION METRICS

Insight Recall. We first decompose each report into atomic insights with their associated citations using an LLM (see Prompt 14). Each extracted insight is then compared against the set of groundtruth insights embedded in the task files using an LLM Judge with the Prompt 15. If a match is found, the insight is marked as *detected* and contributes to the insight recall score; otherwise, it is ignored. This metric thus measures recall rather than precision, since judging whether an unmatched insight is nonetheless *useful* for answering the deep research question is inherently subjective and difficult to automate. To prevent agents from trivially achieving 100% recall by copying all content into the generated report, the LLM Judge evaluates only the first *k* insights, where *k* equals the number of ground-truth insights plus five. This buffer ensures that reports are not penalized for including seemingly relevant insights that are not part of the groundtruth insight set.

Table 2: DRBA performance with different planning configurations on *DRBench*. We compare the base agent with variants using Simple Research Planning (SRP), Complex Research Planning (CRP), Adaptive Action Planning (AAP), and their combinations. See Appendix K for the standard error across 3 runs. Note that higher numbers correspond to better scores, and the best result on each metric is bolded.

Configuration	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
Base DRBA	16.92	64.06	97.94	91.08	41.71
+ SRP	17.00	69.38	98.89	92.46	42.48
+ CRP	16.94	66.78	99.52	90.72	42.07
+ AAP	20.56	67.44	98.89	93.00	47.43
+ SRP + AAP	19.20	61.72	98.97	91.86	44.80
+ CRP + AAP	18.69	57.20	98.89	90.12	43.39

Distractor Avoidance. To measure precision, we track whether the agent's report includes distractor insights that are irrelevant to the research question. We compute *distractor recall* analogously to insight recall, and define *distractor avoidance* as 1— distractor recall.

Factuality. Using the same set of extracted insights (that we used for Insight Recall), we follow the methodology of FactScore (Min et al., 2023). If an insight lacks a citation or references a non-existent source, it is labeled unfactual. Otherwise, we apply a retrieval-augmented system based on text-embedding-3-large (OpenAI, 2024) to fetch the top-5 most relevant chunks from the cited document (Appendix H). The LLM Judge with Prompt 16 then determines whether the cited evidence supports the claim. We also store justifications and model confidence scores for interpretability.

Report Quality. Inspired by prior work (Coelho et al., 2025; Abaskohi et al., 2025), we query the LLM Judge with Prompt 17 to assign a 1–10 rating across six dimensions: (1) depth and quality of analysis, (2) relevance to the research question, (3) persona consistency, (4) coherence and conciseness, (5) absence of contradictions, and (6) completeness and coverage. The final report quality score is obtained by averaging these six ratings.

5.2 MAIN RESULTS

We first evaluate our DRBA agent (Section 4) using GPT-40 as the backbone model, a maximum of 15 research loop iterations, and different combinations of planning modules: Simple Research Planning (SRP), Complex Research Planning (CRP), and Adaptive Action Planning (AAP). The results are reported in Table 2. Overall, the agent demonstrates moderate ability to ground its answers in factual evidence but struggles to consistently surface the main injected insights necessary for answering the deep research questions. In many cases, the agent relies on prior knowledge or external web content rather than integrating the crucial enterprise-specific information available in the files. By contrast, it is consistently strong in avoiding distractors, showing that the agent is robust against misleading or irrelevant information but less effective at prioritizing decision-critical insights. Note that our LLM Judge backbone is GPT-40.

Comparison Across Planning Strategies. Here we see that SRP tends to produce more factually grounded answers, while CRP excels at filtering out distractors through structured decomposition. AAP, on the other hand, provides the largest improvements in both insight recall and report quality, suggesting that dynamically adapting the plan during execution helps the agent recover missed evidence and refine its use of sources. However, combining CRP or SRP with AAP does not yield clear gains, and in some cases reduces factuality, likely because overlapping strategies create redundant or unstable planning behavior. These findings indicate that adaptive mechanisms are key for improving coverage of injected insights, while lightweight planning is more effective for maintaining factual grounding, and that carefully balancing the two remains an open challenge. See Appendix M for detailed results for each task.

5.3 ABLATION: EFFECT OF BACKBONE LANGUAGE MODEL ON DRBA

We evaluate the impact of backbone language models on DRBA using the MinEval subset for controlled comparison (Table 3). GPT-5 achieves the best balance of factual grounding, insight recall, and report quality. Open-source models show mixed results: Llama-3.1-405B excels in factuality but lags in recall, DeepSeek-V3.1 delivers balanced performance through targeted fine-tuning, and Qwen-2.5-72B is reliable but trails GPT-5. These results underline the importance of backbone choice; larger and more advanced

Table 3: Performance of DRBA on the MinEval subset using different backbone language models and planning strategies. Note that higher numbers correspond to better scores, and the best result on each metric is bolded. The full table with more models is given in Appendix L.

Model	Planning	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
GPT-5	None	38.33	74.52	95.14	94.56	79.80
GPT-5	Simple	37.86	72.09	97.14	95.34	78.92
GPT-5	Complex	39.63	65.17	92.86	93.42	77.74
Llama-3.1-405B-Instruct	None	17.37	78.91	100.00	90.48	49.78
Llama-3.1-405B-Instruct	Simple	16.97	79.27	98.10	92.34	48.87
Llama-3.1-405B-Instruct	Complex	20.16	69.75	97.90	91.26	53.86
DeepSeek-V3.1	None	25.15	72.66	97.43	86.52	62.59
DeepSeek-V3.1	Simple	25.56	73.45	96.67	87.36	63.29
DeepSeek-V3.1	Complex	30.26	70.27	96.67	86.88	69.28
Qwen-2.5-72B-Instruct	Complex	26.82	58.35	97.65	89.64	61.75
Qwen-2.5-72B-Instruct	None	25.55	69.39	98.10	90.24	62.64
Qwen-2.5-72B-Instruct	Simple	23.20	67.23	98.10	88.14	58.58

models generally yield stronger overall performance, though some open-source options are competitive in specific metrics. In addition, our experiments also revealed a significant limitation in agents retrieving critical insights from the open web. As shown in Figure 8 in Appendix L, no agent managed to successfully source external knowledge, highlighting the difficulty of extracting relevant information for deep research applications within an unboundedly large search space.

Table 4: Insights Recall Improvement Areas (Task DR0002). We highlight in bold where each model was able to accurately find details relevant to the groundtruth insight. We also show the corresponding score where 1.0 is considered a successful recall and 0.0 an unsuccessful recall. The full table with all groundtruth insights and predicted insights is given in Appendix G.

Groundtruth Insight	Insight Predicted by Llama 3.1 405B	Insight Predicted by GPT-5
	45% of Lee's Market online customers engage with personalized product recommendations, resulting in a $25%$ increase in average order value. (Score = $1.0)$	00 1
	85% of transactions are linked to loyalty accounts at Lee's Market, providing a solid foundation for personalized marketing and improving customer engagement. (Score = 0.0)	

5.4 QUALITATIVE ANALYSIS

In Table 4 we show a sample of three groundtruth insights as well as the predicted insights from using both Llama 3.1 405B and GPT-5. We see that for the first insight, both models are able to effectively recover the groundtruth insight. For the second insight GPT-5 can extract the relevant time of year, where as Llama 3.1 405B fails to do so. This possibly suggests that GPT-5 may be better at extracting fine details.

5.5 PERFORMANCE OF WEB AGENTS ON DRBench

We evaluated Generic WebAgents from AgentLab in a browser-only setting (without API access). The GPT-4.1-powered agent achieved only 1.11% insight recall, 6.67% factuality, and 33.07% report quality. While the reports appeared well-structured, they lacked grounded insights, with most trajectories degenerating into repetitive clicks on irrelevant files or windows. This shows that browser-only agents are currently far from effective for deep research tasks. Further trajectory examples are shown in Appendix J.

5.6 APP-BASED ENVIRONMENT VS LOCAL ENVIRONMENT

In Table 5, we compare results across two settings in *DRBench*: (1) **local**, where all the task files (e.g., PDFs, PPTX, DOCX, XLSX, chats) are placed in a local folder that the agent can access, and (2) **app-based**, where the same files must be retrieved through our standard enterprise environment and its apps,

Table 5: Model Performance Comparison Across Local or App-based Environments. Note that higher numbers correspond to better scores, and the best result on each metric is bolded.

Model	Env	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
DRBA (GPT-5)	App	38.33	74.52	95.14	94.56	66.01
DRBA (GPT-5) + CRP	App	39.63	65.17	92.86	93.42	64.46
DRBA (DeepSeek-V3.1)	App	25.15	72.66	97.43	86.52	53.09
DRBA (DeepSeek-V3.1) + CRP	App	30.26	70.27	96.67	86.88	57.86
DRBA (GPT-5)	Local	41.25	82.43	98.62	91.08	69.57
DRBA (GPT-5) + CRP	Local	42.18	83.91	98.45	92.46	70.67
DRBA (DeepSeek-V3.1)	Local	35.62	77.12	97.95	89.16	64.03
DRBA (DeepSeek-V3.1) + CRP	Local	36.54	78.35	97.62	90.12	65.07
Perplexity	Local	39.14	81.06	98.84	90.36	67.72
OpenAI Deep Research (GPT-5)	Local	44.78	87.53	99.12	94.92	73.56
Gemini	Local	43.92	85.68	98.97	93.24	72.37

introducing additional interaction complexity. We find that OpenAI's Deep Research (GPT-5) achieves the highest scores across all metrics. Our agent with GPT-5 and DeepSeek backbones achieves similar performance to Perplexity in the local-only setting, but lags behind OpenAI and Gemini. In the app-based setting, performance declines across both backbones, highlighting the added difficulty of navigating multi-application environments. This gap underscores that the environment in *DRBench* is intentionally challenging, enabling a more realistic evaluation of model capabilities in enterprise research scenarios.

6 Human Evaluation

Quality of Deep Research Questions. We evaluated the quality of the deep research questions in *DRBench* through a human study with five expert annotators across all 15 tasks. Each task was judged on three criteria: (1) grounding in the external website, (2) relevance to the domain and company context, and (3) alignment with associated insights. Annotators provided binary ratings plus optional feedback. Results show strong quality: 12 tasks received unanimous approval, while only three (tasks DR1, DR11, and DR13) received a single negative vote due to minor issues with specificity or distractor difficulty. This corresponds to a 96% approval rate (72/75 votes).

Correlation of Used Metrics with Human Preference. We collected human preference on a subset of 11 tasks². Each annotator was shown a golden insight with aligning insights from two models³ and asked to choose which they preferred, or label both as good/bad. Missing alignments were shown as empty strings. We compared agents with AAP no RP against GPT-5 and Llama-3.1-405B-Instruct. The Fleiss κ (Fleiss, 1971) across five annotators was 0.67. Most outputs were judged *both bad* due to missing alignments, but when preferences were expressed, GPT-5 was favored 61.1% over Llama-405B-Instruct, consistent with our metric-based findings in Section 5.3. Additional analyses are in Appendix Q.

7 CONCLUSION & FUTURE WORK

In this work, we introduced *DRBench*, a benchmark for evaluating AI agents on complex, open-ended enterprise deep research tasks that require reasoning over both public and private data. Unlike prior benchmarks focused on surface-level or web-only queries, *DRBench* offers 15 persona-grounded tasks situated in realistic enterprise contexts and evaluated through an environment that integrates real-world enterprise applications and heterogeneous data formats. We also presented the DRBench Agent (DRBA) as a strong baseline and analyzed its behavior across planning strategies and backbone models. Our results show that while agents are generally effective at avoiding distractors and capable of producing structured reports, they still struggle to consistently extract decision-critical insights. Adaptive planning improves recall of injected insights, while lightweight strategies tend to preserve factual accuracy, underscoring the difficulty of balancing exploration with reliability. Looking ahead, we plan to extend *DRBench* with tasks requiring cross-file integration, reasoning across modalities such as PDFs and chats, and richer distractors. We also aim to add multimodal sources like images and video, as well as privacy-sensitive tasks to assess data protection. Together, these extensions will move research agents closer to enterprise readiness and provide a stronger foundation for studying deep research in realistic organizational settings.

²We selected tasks with fewer than 8 insights for a reasonable amount of manual work.

³The gold-prediction alignment is provided by the insight recall metric.

ETHICS STATEMENT

This work raises important considerations around data privacy, fairness, and potential misuse. Although *DRBench* simulates enterprise research environments with private data, all datasets are synthetically generated or drawn from public, time-invariant web sources. No personal or sensitive user data is included. The synthetic personas and companies are fictional, designed to prevent any risk of harm or re-identification. We highlight that agents evaluated on *DRBench* must handle sensitive-like contexts (e.g., healthcare, compliance, cybersecurity), which underscores the importance of designing systems that prioritize data protection and avoid exposing private enterprise content. Human annotators were involved in validating task quality; they were compensated at fair rates and gave informed consent.

Large Language Models (LLMs) were used solely to assist with polishing the writing of this paper, such as improving readability and clarity of exposition. All ideas, experimental designs, implementations, analyses, and conclusions are original contributions of the authors.

REPRODUCIBILITY STATEMENT

We have taken multiple steps to ensure reproducibility. The *DRBench* benchmark, including all generated tasks, data generation scripts, supporting files, and evaluation scripts, will be released under a permissive license. Each task is fully self-contained with dated URLs for public insights and synthetic enterprise files for private insights, ensuring stability over time. Detailed descriptions of the task generation pipeline, environment implementation, evaluation prompts, and cost considerations are included in the supplementary materials. We provide open-source code for running agents in the *DRBench* environment and for reproducing all reported results. Hyperparameter settings, backbone models, and planning strategies are documented. Together, these design choices make our benchmark transparent, reproducible, and extensible for future research.

REFERENCES

- Amirhossein Abaskohi, Amrutha Varshini Ramesh, Shailesh Nanisetty, Chirag Goel, David Vazquez, Christopher Pal, Spandana Gella, Giuseppe Carenini, and Issam H. Laradji. AgentAda: Skill-adaptive data analytics for tailored insight discovery, 2025. URL https://arxiv.org/abs/2504.07421.
- Reyna Abhyankar, Qi Qi, and Yiying Zhang. OSWorld-Gold: Benchmarking the efficiency of computer-use agents. In *ICML 2025 Workshop on Computer Use Agents*, 2025. URL https://openreview.net/forum?id=sV3n6mYy7J.
- All-HandsAI. OpenHands. https://github.com/All-Hands-AI/OpenHands, 2024. Accessed: 2024-06-01.
 - Pierre Andrews, Amine Benhalloum, Gerard Moreno-Torres Bertran, Matteo Bettini, Amar Budhiraja, Ricardo Silveira Cabral, Virginie Do, Romain Froger, Emilien Garreau, Jean-Baptiste Gaya, Hugo Laurençon, Maxime Lecanu, Kunal Malkan, Dheeraj Mekala, Pierre Ménard, Grégoire Mialon, Ulyana Piterbarg, Mikhail Plekhanov, Mathieu Rita, Andrey Rusakov, Thomas Scialom, Vladislav Vorotilov, Mengjue Wang, and Ian Yu. ARE: Scaling up agent environments and evaluations, 2025. URL https://arxiv.org/abs/2509.17158.
 - Léo Boisvert, Megh Thakkar, Maxime Gasse, Massimo Caccia, Thibault Le Sellier de Chezelles, Quentin Cappart, Nicolas Chapados, Alexandre Lacoste, and Alexandre Drouin. WorkArena++: Towards compositional planning and reasoning-based common knowledge work tasks. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=PCjK8dqrWW.
 - Nikos I Bosse, Jon Evans, Robert G Gambee, Daniel Hnyk, Peter Mühlbacher, Lawrence Phillips, Dan Schwarz, Jack Wildman, et al. Deep Research Bench: Evaluating AI Web Research Agents. *arXiv* preprint arXiv:2506.06287, 2025.
 - Zijian Chen, Xueguang Ma, Shengyao Zhuang, Ping Nie, Kai Zou, Andrew Liu, Joshua Green, Kshama Patel, Ruoxi Meng, Mingyi Su, Sahel Sharifymoghaddam, Yanxi Li, Haoran Hong, Xinyu Shi, Xuye Liu, Nandan Thakur, Crystina Zhang, Luyu Gao, Wenhu Chen, and Jimmy Lin. BrowseComp-Plus: A more fair and transparent evaluation benchmark of Deep-Research agent, 2025. URL https://arxiv.org/abs/2508.06600.
 - Thibault Le Sellier De Chezelles, Maxime Gasse, Alexandre Drouin, Massimo Caccia, Léo Boisvert, Megh Thakkar, Tom Marty, Rim Assouel, Sahar Omidi Shayegan, Lawrence Keunho Jang, Xing Han Lù, Ori Yoran, Dehan Kong, Frank F. Xu, Siva Reddy, Quentin Cappart, Graham Neubig, Ruslan Salakhutdinov, Nicolas Chapados, and Alexandre Lacoste. The BrowserGym ecosystem for web agent research, 2025. URL https://arxiv.org/abs/2412.05467.
 - João Coelho, Jingjie Ning, Jingyuan He, Kangrui Mao, Abhijay Paladugu, Pranav Setlur, Jiahe Jin, Jamie Callan, João Magalhães, Bruno Martins, et al. Deepresearchgym: A free, transparent, and reproducible evaluation sandbox for deep research. *arXiv preprint arXiv:2505.19253*, 2025.
- Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H. Laradji, Manuel Del Verme, Tom Marty, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. WorkArena: How capable are web agents at solving common knowledge work tasks? 2024.
- Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. DeepResearch Bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Joseph L Fleiss. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76 (5):378, 1971.
- FoundationAgents. OpenManus. https://github.com/FoundationAgents/OpenManus, 2024. Accessed: 2024-06-01.

- Boyu Gou, Zanming Huang, Yuting Ning, Yu Gu, Michael Lin, Botao Yu, Andrei Kopanev, Weijian Qi, Yiheng Shu, Jiaman Wu, Chan Hee Song, Bernal Jimenez Gutierrez, Yifei Li, Zeyi Liao, Hanane Nour Moussa, TIANSHU ZHANG, Jian Xie, Tianci Xue, Shijie Chen, Boyuan Zheng, Kai Zhang, Zhaowei Cai, Viktor Rozgic, Morteza Ziyadi, Huan Sun, and Yu Su. Mind2web 2: Evaluating agentic search with agent-as-a-judge. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025. URL https://openreview.net/forum?id=AUaW6DS9si.
- Kung-Hsiang Huang, Akshara Prabhakar, Sidharth Dhawan, Yixin Mao, Huan Wang, Silvio Savarese, Caiming Xiong, Philippe Laban, and Chien-Sheng Wu. CRMArena: Understanding the capacity of LLM agents to perform professional CRM tasks in realistic environments. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 2025a.
- Kung-Hsiang Huang, Akshara Prabhakar, Onkar Thorat, Divyansh Agarwal, Prafulla Kumar Choubey, Yixin Mao, Silvio Savarese, Caiming Xiong, and Chien-Sheng Wu. CRMArena-Pro: Holistic assessment of LLM agents across diverse business scenarios and interactions. *arXiv preprint arXiv:2505.18878*, 2025b.
- HuggingFace. smolagents. https://github.com/huggingface/smolagents, 2024. Accessed: 2024-06-01.
- Abhinav Java, Ashmit Khandelwal, Sukruta Midigeshi, Aaron Halfaker, Amit Deshpande, Navin Goyal, Ankur Gupta, Nagarajan Natarajan, and Amit Sharma. Characterizing deep research: A benchmark and formal definition, 2025. URL https://arxiv.org/abs/2508.04183.
- LearningCircuit. Local deep research. https://github.com/LearningCircuit/local-deep-research, 2025.
- Yao Liu, Deming Ye, Shuohang Wang, Furu Wei, Yujia Ma, and Minlie Huang. G-Eval: NLG evaluation using GPT-4 with better human alignment. *EMNLP*, 2023.
- Xing Han Lù, Zdeněk Kasner, and Siva Reddy. WebLINX: Real-world website navigation with multi-turn dialogue, 2024. URL https://arxiv.org/abs/2402.05930.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA: a benchmark for general AI assistants. *ICLR*, 2024.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12076–12100, 2023.
- OpenAI. text-embedding-3-large. https://platform.openai.com/docs/guides/embeddings, 2024. [Large language model]. Accessed September 2025.
- Gaurav Sahu, Abhay Puri, Juan A Rodriguez, Amirhossein Abaskohi, Mohammad Chegini, Alexandre Drouin, Perouz Taslakian, Valentina Zantedeschi, Alexandre Lacoste, David Vazquez, et al. InsightBench: Evaluating insight extraction for business analytics agents. *ICLR*, 2025.
- Zilliz Tech. Deep-Searcher. https://github.com/zilliztech/deep-searcher, 2024. Accessed: 2024-06-01.
- Yuhao Wang, Ruiyang Ren, Junyi Li, Xin Zhao, Jing Liu, and Ji-Rong Wen. REAR: A relevance-aware retrieval-augmented framework for open-domain question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 5613–5626, 2024.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *NeurIPS*, 2024.

Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z. Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su, Leander Maben, Raj Mehta, Wayne Chi, Lawrence Jang, Yiqing Xie, Shuyan Zhou, and Graham Neubig. The Agent Company: Benchmarking LLM agents on consequential real world tasks, 2025a. URL https://arxiv.org/abs/2412.14161.

- Renjun Xu and Jingwen Peng. A comprehensive survey of deep research: Systems, methodologies, and applications. *arXiv preprint arXiv:2506.12594*, 2025.
- Tianze Xu, Pengrui Lu, Lyumanshan Ye, Xiangkun Hu, and Pengfei Liu. ResearcherBench: Evaluating deep AI research systems on the frontiers of scientific inquiry, 2025b. URL https://arxiv.org/abs/2507.16280.
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. DeepResearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv*, 2025.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. WebArena: A Realistic Web Environment for Building Autonomous Agents, April 2024. URL http://arxiv.org/abs/2307.13854.arXiv:2307.13854 [cs].

Table 6: Comparison of Deep Research Tasks of Different Benchmarks.

Benchmark	Sample Question
DeepResearchGym (Coelho et al., 2025)	Is the COVID vaccine dangerous
Deep Research Bench (Bosse et al., 2025)	Find a reliable, known number on the internet. The total number of FDA Class II Product Recalls of medical devices.
DeepResearch Bench (Du et al., 2025)	While the market features diverse quantitative strategies like multi-factor and high-frequency trading, it lacks a single, standardized benchmark for assessing their performance across multiple dimensions such as returns, risk, and adaptability to market conditions. Could we develop a general yet rigorous evaluation framework to enable accurate comparison and analysis of various advanced quant strategies?
ResearcherBench (Xu et al., 2025b)	Compare the Transformer and Mamba model architectures, analyzing their performance and technical characteristics in different application scenarios. Based on the latest research, discuss the advantages and disadvantages of both models and their applicable scenarios.
LiveDRBench (Java et al., 2025)	For complex reasoning tasks (e.g., tasks involving multiple citations or extended reasoning chains), what are the strengths of current agent technologies, and what are their limitations? Please analyze this in the context of research since June 2024.
BrowseComp-Plus (Chen et al., 2025)	Identify the title of a research publication published before June 2023, that mentions Cultural traditions, scientific processes, and culinary innovations. It is co-authored by three individuals: one of them was an assistant professor in West Bengal and another one holds a Ph.D.
GAIA2 (Andrews et al., 2025)	Update all my contacts aged 24 or younger to be one year older than they are currently.
DRBench	How can Lee's Market leverage FSMA 204 regulations to enhance food safety and customer trust?

A COMPARISON OF DEEP RESEARCH BENCHMARKS AND AI AGENT BENCHMARKS WITH A COMPUTER ENVIRONMENT

In Table 1, we compare existing deep research benchmarks and AI agent benchmarks that provide a computer environment with *DRBench*. While the questions in existing benchmarks focus on public interest topics and require generic web search and computer use, *DRBench* provides realistic questions that real personas in organizations need to resolve.

B DRBench TASKS

As shown in Table 7, *DRBench* contains 15 tasks in total, covering 3 industries (retail, healthcare and electric vehicles), 10 task domains (compliance, sales, customer relationship management, market analysis, customer service management, IT service management, cyber security, marketing, quality assurance, and research), and 3 difficulty levels (easy, medium, hard). In addition, we generate the following 3 companies (one for each industry type): (1) a supermarket chain called Lee's Market, (2) a virtual healthcare company called MediConn Solutions, and (3) an electric vehicle company called Elexion Automotive.

Table 8 presents a deep research question from *DRBench* and its supporting groundtruth insights. We also visualize the DR Question and all QA pairs by embedding them with OpenAI's text-embedding-3-large model and projecting into 2D using t-SNE in Figure 4. The plot shows that injected supporting insights lie closer to the DR Question, while distractors appear farther away, confirming that our injected insights are semantically aligned with the research objective.

C DRBench Examples of Injected Insights

As shown in Figure 2, supporting documents are generated with enterprise insights injected. In Figure 5, we show two examples of a generated files (PPTX and Mattermost chat) with their embedded insights.

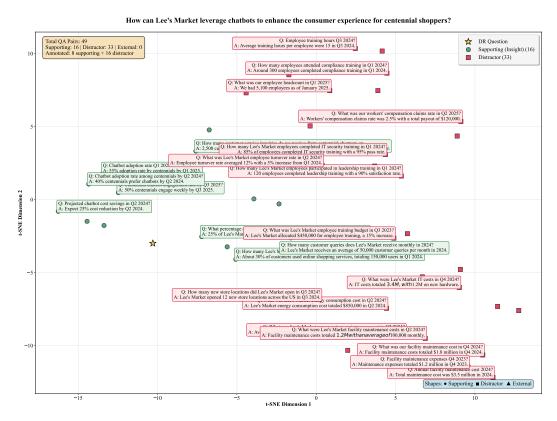


Figure 4: t-SNE visualization of QA pairs for the DR Question in Task DR0005. The plot shows the distribution of annotated pairs across **Supporting Insights** (green), **Distractors** (red), and the central **Deep Research** (**DR**) **Question** (gold star). Out of 49 pairs, 16 correspond to supporting insights and 33 are distractors. The visualization illustrates how relevant insights cluster separately from distractors, highlighting the challenge of retrieving salient information in a distractor-heavy environment.



(a) Example supporting file named food-safety-regulatory-compliance.pdf with an injected insight "Lee's Market reduced food waste by 8% in Q2 2024, saving \$1.2M."



(b) Example supporting email in the Sent mailbox with an injected insight "85% of high-risk food suppliers are FSMA 204 compliant, totaling 270 vendors."

Figure 5: Example files with injected insights in *DRBench*.

Table 7: *DRBench* Questions and Statistics. **Industry**: target industry of the deep research question. **Domain**: domain of the deep research task. **DR Question**: question of the deep research task. **Difficulty**: difficulty of the task defined based on the rubric mentioned in 3. **# Applications**: the number of total applications in the task environment. **# Insights**: the number of relevant insights to the deep research question. **# Distractors**: the number of non-supportings documents that do not contain relevant insights.

Industry	Domain	DR Question	Difficulty	# Applications	# Insights	# Distractors
Retail	Compliance	How can Lee's Market leverage FSMA 204 reg- ulations to enhance food safety and customer trust?	easy	2	3	7
Retail	Sales	How can personalization drive sales in the retail industry and what strategies can be used for Lee's Market in action?	easy	2	4	10
Retail	CRM	How can we leverage data-driven loyalty programs to enhance customer engagement?	medium	4	7	15
Retail	Market Analysis	What are the new trends in the grocery retail market and what strategies can Lee's Market adopt to remain competitive?	medium	3	6	14
Retail	CSM	How can Lee's Market leverage chatbots to enhance the consumer experience for centennial shoppers?	hard	4	16	33
Healthcare	Compliance	What are the key factors influencing MediConn Solutions' decision to accept insurance for tele- health providers, considering HIPAA compliance and state-specific data privacy regulations?	easy	3	4	8
Healthcare	ITSM	How can we leverage ISTM and AI-driven ana- lytics to minimize IT service desk workload and improve response times in MediConn Solutions?	easy	3	6	12
Healthcare	Cybersecurity	What is the impact of third-party data breaches on MediConn's virtual healthcare platforms and patient data, and what new regulations can be implemented to defend against these breaches?	medium	4	7	14
Healthcare	CRM	How can MediConn Solutions leverage trendy new CRM solutions to improve patient engagement and retention, and what new CRM solutions are expected in 2025?	medium	4	12	24
Healthcare	Marketing	What are the most critical elements of a robust digital presence for a telehealth provider such as MediConn Solutions, and how can we optimize our website and content marketing strategy to attract digital-first patients?	hard	4	15	30
Electric Vehicle	Compliance	How can we balance the need for durability and warranty guarantees for EV batteries with evolving regulatory requirements, especially ACC regulations (ACC II), while staying on track with our production timelines through 2035?	easy	2	3	6
Electric Vehicle	Quality Assurance	How can Elexion Automotive's quality assurance processes be optimized to address the unique challenges of electric vehicle production, such as software and user experience issues, compared to gasoline cars?	easy	1	3	6
Electric Vehicle	Cybersecurity	How can Elexion Automotive effectively im- plement a cybersecurity strategy for its electric vehicles, considering the risks and challenges posed by connected and autonomous technologies?	medium	3	6	12
Electric Vehicle	Research	Can we leverage AI-enhanced battery management to improve EV battery lifespan by 15%?	medium	3	7	14
Electric Vehicle	CSM	How can Elexion Automotive increase customer trust through after-sales support while balancing the need for exceptional customer care with efficient and cost-effective service?	hard	4	15	30

D DRBENCH ENTERPRISE ENVIRONMENT

The *DRBench* Enterprise Environment provides a containerized simulation of realistic enterprise research settings where employees access confidential company information, personal files, and internal communications for comprehensive report generation. The environment simulates both a user's local machine filesystem and provides password-protected access to enterprise services.

To emulate realistic enterprise research settings, *DRBench* provides a self-contained Docker environment that integrates commonly used applications: Nextcloud for shared documents, Mattermost for internal chat, an IMAP server and Roundcube open-source client for emails, and Filebrowser to emulate local files. Each task is initialized by distributing its data across these services, enabling agents to retrieve, analyze, and cite information through enterprise-like interfaces rather than static files. This design ensures realistic interaction while maintaining reproducibility and controlled evaluation.

Table 8: Example Deep Research Question and Supporting Groundtruth Insights

Deep Research Question	Supporting Groundtruth Insight	Insight Category
	U.S. grocers are working to meet the FDA's FSMA 204 traceability rules	External
How can Lee's Market leverage FSMA	by January 2026, which require tracking lot codes and key data for high-risk	
204 regulations to enhance food safety and	foods to expedite recalls. This compliance is viewed as an "evolutionary	
customer trust?	step" to modernize grocery operations and enhance food safety.	
customer trust?	By capturing detailed traceability data, such as lot codes, at every step,	External
	retailers can meet regulations and gain inventory benefits. This allows grocers	
	to know exact expiration dates by lot, enabling them to discount items before	
	they expire, thus reducing food waste and keeping products fresher.	
	Regional grocers like Lunds & Byerlys and Raley's see FSMA 204 as	External
	a chance to enhance their systems and supply chain transparency. They	
	believe improved traceability will boost customer trust and could signal the	
	start of more extensive future food safety regulations.	
	Lee's Market tracks 250 high-risk food products as of Q3 2024, affecting	Internal
	30% of inventory.	
	Lee's Market reduced food waste by 8% in Q2 2024, saving \$1.2M.	Internal
	85% of high-risk food suppliers are FSMA 204 compliant, totaling 270	Internal
	vendors.	

D.1 ARCHITECTURE AND SERVICES

The environment implements a **multi-service architecture** within a single Docker container. This design prioritizes deployment simplicity and cross-platform compatibility while maintaining service isolation. The container orchestrates the following enterprise services:

- Nextcloud: Open-source file sharing and collaboration platform analogous to Microsoft SharePoint or Google Drive, providing secure document storage with user authentication.
- Mattermost: Open-source team communication platform simulating internal company communications similar to Microsoft Teams or Slack, with teams, channels, and persistent chat history.
- FileBrowser: Web-based file manager providing access to the container's local filesystem, simulating employee desktop environments and local document access.
- Email System: Roundcube webmail interface with integrated SMTP (postfix) and IMAP (dovecot) services for enterprise email communication simulation.
- VNC/NoVNC Desktop: Protocol and browser-based VNC access providing full desktop environment interaction within the container for comprehensive enterprise workflow simulation.

TASK LOADING AND DATA DISTRIBUTION

At initialization, the environment processes task configuration files (env. json) and distributes data across services through automated Python scripts and it makes sure that this source data is only accessible through the intended applications:

- File Distribution: Documents are placed in appropriate Nextcloud user folders and FileBrowser directories based on task specifications
- Communication Import: Chat histories and team conversations are imported into Mattermost channels with proper user attribution
- Email Integration: Email conversations are loaded into the mail system with realistic threading and metadata
- User Provisioning: Enterprise users are automatically created across all services with consistent authentication credentials

```
918
919
       from drbench import drbench_enterprise_space, task_loader
920
        # Load task configuration
921
        task = task_loader.get_task_from_id(task_id)
922
923
       # Initialize environment with automatic port allocation
924
       env = drbench_enterprise_space.DrBenchEnterpriseSearchSpace(
           task=task.get_path(),
925
           start_container=True,
926
           auto_ports=True # Prevents port conflicts in parallel execution
927
928
929
       # Environment provides service discovery
       available_apps = env.get_available_apps()
930
       # Returns: {'nextcloud': {'port': 8081, 'credentials': {...}}, ...}
931
932
        # Pass relevant information to the agent
933
934
       # Cleanup when research complete
       env.delete()
935
```

Listing 1: DrBench Environment Usage

D.3 PYTHON INTEGRATION

The DrBenchEnterpriseSearchSpace class provides programmatic container management with the following capabilities: container lifecycle management, service access information, task-specific data loading, and automatic cleanup. The typical usage pattern shown in Listing 1 demonstrates these integrated capabilities.

D.4 ENTERPRISE SERVICE APIS

Each service exposes both **web interfaces** for human and web-agent interaction, and **programmatic APIs** for agent access:

- Nextcloud: WebDAV API for file operations, sharing, and metadata retrieval
- Mattermost: REST API for message history, channel management, and user interactions
- Email: IMAP/SMTP protocols for message retrieval and sending
- FileBrowser: HTTP API for filesystem operations and file management

This dual-access model enables both agent-driven research and human verification of enterprise scenarios, supporting comprehensive evaluation of research capabilities across realistic enterprise information architectures.

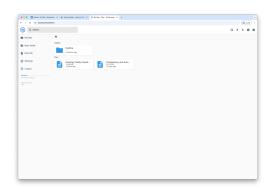
E DRBench EXAMPLES OF APPLICATION SCREENSHOTS

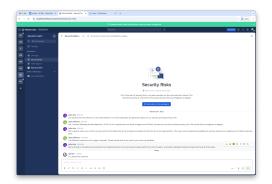
Figures 6 and 7 show the applications provided in *DRBench* environment: File Browser, Mattermost, Roundcube, and Nextcloud.

F DRBench AGENT IMPLEMENTATION DETAILS

F.1 DETAILED WORKFLOW

As depicted in Figure 3, the workflow begins with a Company Employee submitting an enterprise Deep Research Question along with Company Context. The *DRBench* agent processes this input through several key stages:

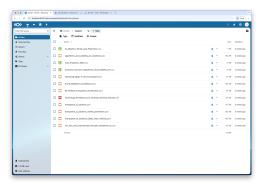


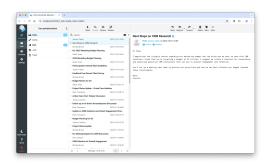


(a) Screenshot of the File Browser interface, displaying organized files and folders within the system.

(b) Screenshot of the Mattermost communication platform, showing a discussion channel and user interface elements.

Figure 6: Screenshots of Applications in *DRBench* environment (Part 1).





(a) Screenshot of the Nextcloud file management system, illustrating the file list view with various document types.

(b) Screenshot of Roundcube, an email client, it shows an open email in the user's inbox.

Figure 7: Screenshots of Applications in *DRBench* environment (Part 2).

Stage 1: Research Planning. The agent decomposes the research question into structured research investigation areas to guide subsequent action generation. This initial decomposition lays out a strategy to systematically cover the research space while maintaining focus on the initial deep research question.

Stage 2: Action Planning. The Action Planning stage translates the research objectives into executable actions through a planning subsystem. This component uses an LLM to create a prioritized sequence of actions. Each action is parameterized with specific tool selection and execution parameters, its dependencies to other actions in the plan, and a priority score.

Stage 3: Research Loop with Adaptive Execution. The Research Loop iterates over the following sub-stages until completion: (1) Tool Selection and Execution: The tool selection and execution subsystem implements a sophisticated priority-based selection of actions from the plan at each research iteration step and proceeds to execute it with the current research context. (2) Content Processing: If necessary, the action will make use of the content processing subsystem to extract, synthesize, and store retrieved documents and websites into a vector store to form a task-specific knowledge base that will grow on each iteration. (3) Adaptive Action Planning: After each execution round, the agent analyzes the most recent findings; if coverage gaps are detected, new actions are created and added to the plan at this point. This ensures that newly discovered knowledge is taken into account to answer the research question. (4) Iteration Findings Storage: Results from each iteration are stored in the vector store with rich metadata for traceability and later assessment.

Stage 4: Convergence and Completion. The research loop continues until all actions in the plan have been completed or the maximum iteration limit is reached.

1027

1028

1029

1030

1031

1032

1033 1034

1035 1036

1038

1040

1041

1042 1043

1045

1046

1047

1048

1049 1050

1051

1052

1053

1054

1055

1057

1058

1059

1060

1063

1064

1065

1068

1069

1070

1072

1074 1075 1076

1077

1078

1079

Stage 5: Report Writing. The report writing subsystem queries the vector store to synthesize the research findings and relevant retrieved content. This component generates a comprehensive report and uses its own citation tracking system to ensure proper attribution of claims within the report.

The vector store serves as the main knowledge integration component, maintaining embeddings of all processed content and enabling semantic retrieval at the report generation stage. This component is crucial in a deep research setting to prevent information loss from early research stages in arbitrarily long research sessions.

F.2 RESEARCH PLANNING IMPLEMENTATION

The Research Planning subsystem offers three operational modes to evaluate the impact of structured planning on research effectiveness:

- Complex Mode: Generates a comprehensive research plan with detailed investigation areas. These areas contain details about the specific research focus, expected information sources, and success criteria, among others. Each area includes an importance level and specific business intelligence objectives (see Listing 2).
- **Simple Mode:** Creates focused question decompositions with 4-10 self-contained subqueries derived directly from the main research question. Uses straightforward decomposition without the complex enterprise research structure of complex mode. See examples in Listing 4 and Listing 3 for comparison of different planning modes.
- None: Bypasses structured planning entirely, proceeding directly to action generation based on the original question. This mode serves as a baseline to measure the added-value of explicit planning stages.

The planning process begins with the enriched query (Prompt 1) and uses the research planning prompt (Prompt 2) to generate structured outputs. In Complex Mode, the system creates detailed investigation areas with enterprise-focused metadata, while Simple Mode produces straightforward question decompositions similar to existing multi-step reasoning approaches. The resulting plan structure directly feeds into the Action Planning System (Appendix F.3) for executable action generation.

```
"area_id": 1,
"research_focus": "Core strategic
  domain, market segment, or business hypothesis to investigate",
"information_needs": [
  "What specific intelligence is required for strategic decisions"
"knowledge_sources": ["internal", "external", "both"],
"research_approach
 ": "competitive_analysis | market_research | strategic_assessment
  | trend_analysis | risk_analysis | performance_benchmarking",
"key_concepts": ["concept1", "concept2"],
"business_rationale": "Why this investigation
  area is critical for enterprise strategy and decision-making",
"expected_insights": "What strategic
  understanding or competitive intelligence this area should provide",
"stakeholder_impact": "Which
  business units or decision-makers will benefit from these insights",
"importance_level": "critical | important | supplementary"
```

Listing 2: Investigation Area Structure for Full Planning Mode

```
{
  "query": "What are the new trends in the grocery retail market
  and what strategies can Lee's Market adopt to remain competitive?",
  "plan": {
    "research_investigation_areas": [
    {
```

```
1080
             "area_id": 1,
1081
             "research_focus": "Current trends in grocery retail market",
1082
             "information_needs": [
1083
             "Latest consumer preferences",
             "Emerging technologies influencing grocery shopping",
1084
             "Sustainability practices in grocery retail",
1085
             "Changes in supply chain dynamics"
1086
1087
             "knowledge_sources": ["external"],
1088
             "research_approach": "trend_analysis",
             "key_concepts
1089
           ": ["e-commerce growth", "sustainability in supply chains"],
1090
             "business_rationale": "Understanding consumer
1091
           trends and technological advancements that shape shopper behavior is
1092
           critical for adapting offerings and enhancing customer engagement.",
             "expected_insights": "Identify specific trends affecting customer
1093
           buying decisions, including the rise of online grocery shopping
1094
           and preferences for sustainable, local, or organic products.",
1095
             "stakeholder_impact
1096
           ": "Marketing, Product Development, Supply Chain Management",
             "importance_level": "critical"
1097
1098
1099
             "area_id": 2,
1100
             "research_focus": "Competitive analysis of grocery retailers",
1101
             "information needs": [
1102
             "Market share analysis",
             "Competitor strengths and weaknesses",
1103
             "Innovative strategies adopted by competitors"
1104
1105
             "knowledge_sources": ["external"],
1106
             "research_approach": "competitive_analysis",
1107
          key_concepts": ["market positioning", "competitive differentiation"],
1108
             "business_rationale": "A comprehensive understanding of
1109
           competitors allows for strategic positioning and the identification
1110
           of innovative practices that can be adopted or improved upon.",
1111
             "expected_insights": "Detailed profiles of key
1112
           competitors, including strategic moves they are making to capture
          market share, which can inform Lee's Market's competitive strategy.",
1113
             "stakeholder_impact
1114
           ": "Executive Leadership, Strategic Planning, Marketing",
1115
             "importance_level": "critical"
1116
           },
1117
           . . .
1118
           ]
1119
1120
```

Listing 3: Complex Mode Research Plan Example

1122 1123

```
1124
1125
         "query": "How can we leverage
1126
           data-driven loyalty programs to enhance customer engagement?",
         "plan": {
1127
           "mode": "simple",
1128
           "subqueries": [
1129
           "What are the key features
1130
           of successful data-driven loyalty programs in the retail industry?",
1131
           "How can data analytics be used to personalize rewards
1132
          and incentives in loyalty programs to increase customer engagement?",
1133
           "What types of customer data should be collected and analyzed
           to optimize loyalty programs for a company like Lee's Market?",
```

```
1134
           "Which
1135
           technology platforms and tools are most effective for implementing
1136
           and managing data-driven loyalty programs in the retail sector?",
           "How can Lee's Market measure the success of their
1137
           loyalty program in terms of customer engagement and sales growth?",
1138
           "What are the best practices for integrating a loyalty program
1139
           with existing marketing strategies to enhance customer experience?",
1140
           "How can Lee's Market ensure customer data
1141
          privacy and security while leveraging data-driven loyalty programs?",
           "What are the potential
1142
           challenges and limitations of implementing a data-driven loyalty
1143
          program for a retail company with Lee's Market's size and resources?"
1144
1145
           "research_methodology": {
1146
           "overall_approach": "Query decomposition into focused subqueries"
1147
1148
1149
```

Listing 4: Simple Mode Research Plan Example

F.3 ACTION PLANNING SYSTEM

The Action Planning System translates research objectives into executable actions through an intelligent planning subsystem that manages tools, prioritization, and dependencies.

Available Tools Table 9 summarizes the available tools organized by category and their primary purposes.

Category	Tool	Purpose
Information Retrieval	Internet Search	External market research, competitive intelligence, and public data analysis. Ideal for market trends, competitor analysis, industry reports, news articles.
	Enterprise API	Access to proprietary internal data through extensible adapters (Nextcloud, Mattermost, email, FileBrowser). Ideal for internal metrics, communications, confidential documents.
	URL Fetch	Direct content extraction from specific URLs. Ideal for deep analysis of reports, whitepapers, case studies, competitor websites.
Analysis	Analyzer	AI-powered synthesis and analysis using vector search. Ideal for cross-referencing findings, identifying patterns, generating insights.
Local Processing	Local Document Search	Semantic search within locally ingested documents. Ideal for targeted retrieval from local files with source references.

Table 9: DrBench Agent Tool Categories and Purposes

Priority Scoring and Dependencies Actions are assigned priority scores (0.0-1.0 scale) based on strategic importance and expected information value. The priority assignment follows enterprise research principles:

- Source Type Prioritization: Enterprise and local sources receive higher priority than external sources, reflecting the strategic value of proprietary information in competitive analysis.
- Query Specificity: Targeted queries addressing specific investigation areas score higher than broad exploratory searches, ensuring focused research execution.

1188 • Dependency Management: Actions can specify prerequisite relationships where certain 1189 information gathering must precede analysis or synthesis tasks. The scheduler respects these 1190 dependencies while maximizing parallel execution within each iteration. 1191 1192 F.4 ENTERPRISE INTEGRATION ARCHITECTURE 1193 1194 **Service Adapters** The system implements extensible adapters for enterprise services including Nextcloud file server, Mattermost chat, IMAP email systems, and FileBrowser interfaces. Each adapter handles 1195 service-specific authentication, data retrieval, and metadata preservation for proper citation attribution. 1196 1197 **Source Prioritization Strategy** Enterprise and local sources receive priority multipliers in action scoring, 1198 reflecting the strategic value of proprietary information. The system maintains source type classification 1199 throughout the pipeline to ensure internal intelligence drives analytical conclusions while external sources 1200 provide context and validation. 1201 1202 F.5 TOOL SELECTION AND EXECUTION 1203 1204 The Tool Selection stage implements a priority-based action selection and tool invocation within each 1205 research iteration to execute the action plan. 1206 1207 **Action Selection Process** At each iteration, the system selects executable actions based on priority 1208 scores and dependency satisfaction: 1209 1210 • **Priority-Based Scheduling:** Actions are ranked by their priority scores, with enterprise and local sources prioritized over external sources to maximize the value of specific private information. 1211 1212 • Dependency Validation: The scheduler checks that all prerequisite actions have completed 1213 before making an action available for execution. 1214 Sequential Execution: Actions execute one at a time in priority order, maintaining research 1215 coherence and enabling each action to build upon previous findings. 1216 1217 Once selected, actions execute through a standardized interface and results are integrated into the research 1218 context, informing the following stage of adaptive action planning. 1219 1220 F.6 ADAPTIVE PLANNING 1221 1222 The Adaptive Planning system enables dynamic evolution of the action plan by analyzing results after each iteration to generate extra actions addressing information gaps. 1223 1224 It starts by analyzing the most recently completed actions and performs these two substages: 1225 1226 **Source analysis and gap classification.** The system evaluates the possible imbalances in the action 1227 completion if information came from internal or external sources and identifies possible scenarios to cover. 1228 1229 **Dynamic action generation.** After analyzing the sources and results from the previous actions, the 1230 system makes an LLM call to generate 1-5 extra candidate actions with a specific prioritization. After 1231 candidate actions are generated, they go through a deduplication process to make sure the plan didn't cover them already and incorporates the final subset into the priority action plan so they can be considered 1232 by the scheduler in the following iteration. 1233 1234 F.7 CONTENT PROCESSING AND VECTOR STORE 1235 1236 The Content Processing system implements a pipeline for unified ingestion of documents in multiple 1237 formats (PDF, docx, HTML, JSON, plain text formats, etc.) that normalizes and cleans text inside the documents and websites retrieved during the research. Content is then deduplicated and chunkized, and 1239 embeddings are computed for each of these chunks. 1240 The Vector Store implements the storage and retrieval of the content via JSON metadata and NumPy

embedding metrices, enabling semantic similarity and keyword based searches

1241

F.8 REPORT GENERATION

Multi-Stage Synthesis Pipeline The Report Generation stage implements a four-stage synthesis approach: (1) thematic content clustering via vector searches, (2) source prioritization and deduplication, (3) LLM Synthesis with source tracking, and (4) Final report writing and citation resolution.

The system generates targeted search queries based on the research plan, including specific to ensure that a predefined set of themes or sections (background, analysis, implementation, trends) are retrieved and written and prevents redundant analyses.

Unified Citation System The citation system implements deferred resolution to keep consistency of citations from section to section by referencing document IDs in the Vector Store and carrying them over for each piece of synthesized text. A final citation resolution stage will assign the correct numbering to each document in the final report.

F.9 DRBA PROMPTS

The *DRBench* agent relies on carefully designed prompts to orchestrate enterprise research workflows. These prompts implement the core architectural principles: enterprise context enrichment, structured planning decomposition, priority-based action generation, quantitative synthesis requirements, and adaptive research capabilities. The following five prompts represent the critical LLM interactions that enable systematic enterprise research with proper source prioritization and citation tracking:

- Enriched Query Generation (Prompt 1): Transforms basic research questions into enterprisecontextualized queries by incorporating company information, stakeholder personas, and business context to guide subsequent research activities.
- **Research Planning** (Prompt 2): Decomposes complex research questions into structured investigation areas with defined information needs, knowledge sources, and business rationales, enabling systematic coverage of the research space.
- Action Generation (Prompt 3): Converts research objectives into prioritized executable
 actions with tool specifications and dependency relationships, emphasizing enterprise source
 prioritization over external sources.
- Adaptive Action Generation (Prompt 4): Analyzes research progress to identify coverage gaps and source imbalances, generating complementary actions that enhance research depth and cross-validate critical findings.
- Report Synthesis (Prompt 5): Orchestrates quantitative-first content synthesis with strict citation requirements, ensuring numerical data leads analytical paragraphs and all claims are properly attributed to source documents.

Query Generation

```
Research Question: {dr_question}

Company Context: {company_name} is {company_desc}.

Persona Context: This analysis is requested by {name}, {role} in {department}.

Their responsibilities include: {responsibilities}.
```

Prompt 1: Query Enrichment with Enterprise Context.

G QUALITATIVE RESULTS

Shown below are some illustrative examples of how metrics are computed for different scenarios on a given test task.

```
1296
         7 Research Planning (Complex Mode)
1297
1298
        Design a comprehensive enterprise research strategy for:
1299
        "{question}"
1300
         {tools_section}
1301
1302
        As a senior enterprise researcher with deep business
1303
        intelligence expertise, create a thorough investigation
1304
        plan that combines rigorous research methodology with
        strategic business analysis. Your goal is to provide insights
1305
        that drive informed decision-making in complex enterprise
        environments.
        Enterprise Research Design Principles:
1309
        - Leverage proprietary internal data as competitive advantage while
1310
        ensuring external market context
1311
        - Design investigations that directly inform strategic decisions and
1312
        business outcomes
        - Prioritize research areas that maximize ROI and strategic value to the
        organization
        - Balance comprehensive analysis with focused insights relevant to
1315
        enterprise objectives.
1316
1317
        Generate a JSON object with strategic research investigation
1318
        areas:
1319
         {output_structure}
1320
1321
```

Prompt 2: Enterprise Research Planning with Investigation Areas. See example of the output structure in Listing 2

G.1 INSIGHTS RECALL

We showcase the insights recall metric using Task DR0002, whose DR Question is "How can personalization drive sales in the retail industry and what strategies can be used for Lee's Market in action?" (also shown in Table 7), which evaluates sales challenges and competitive landscape analysis.

Table 10 shows the overall performance comparison between using Llama 3.1 405B and GPT-5 in our agent. Using GPT-5 in our agent results in increasing the insights recall score from 0.14 to 0.43, successfully answering 3 out of 7 questions compared to Llama's 1 out of 7.

Table 10: Insights Recall Performance Comparison: Llama 3.1 405B vs GPT-5 (Task DR0002). We summarize the number of questions answered successfully and unsuccessfully as well as the overall insights recall score for the given task.

Metric	Llama 3.1 405B	GPT-5	Improvement
Insights Recall Score	0.14	0.43	+0.29
Questions Answered Successfully	1/7	3/7	+2
Questions Failed	6/7	4/7	-2

The question-by-question breakdown in Table 11 reveals the specific questions where each approach succeeded or failed. Both models successfully identified the insight related to online customer engagement, but only GPT-5 was able to identify the number of loyalty program members and the customer data collection rate. Neither model successfully answered the remaining 4 questions, indicating these insights may not have been readily available in the source materials or that agents struggled to find the right insights.

```
1350
         Action Generation
1351
1352
        Generate specific executable actions for this research investigation
        area with SOURCE PRIORITIZATION:
1353
1354
        Research Focus: {research_focus}
1355
1356
        Information Needs:
                             {information_needs}
1357
        Knowledge Sources:
                             {knowledge_sources}
1358
1359
        Research Approach:
                             {research_approach}
1360
        Available Tools: {available_tool_names}
1362
1363
        Tool Selection Guidelines: {tool_guidelines}
1364
        Return JSON array of actions with:
1365
        - "type": Action type (web_search, enterprise_api, url_fetch, analyzer,
        local_search)
1367
        - "description": Clear description of what this action will accomplish
        - "parameters": Tool-specific parameters including query, search_type,
1369
        - "priority": Float 0.0-1.0 (enterprise sources: 0.7-1.0, external:
1370
        0.4 - 0.7)
1371
        - "expected_output": What information this action should provide
1372
        - "preferred_tool": Specific tool class name to use
1373
```

Prompt 3: Priority-Based Action Generation with Source Awareness

Table 11: Question-by-Question Insights Recall Analysis (Task DR0002). We breakdown the results question by question for the given task, highlighting specifically which question is answered correctly or incorrectly for each model.

Question	Llama 3.1 405B	GPT-5	Δ
Online Customer Engagement with Personalized Recommendations	1.0	1.0	0.0
Number of Loyalty Program Members	0.0	1.0	+1.0
Customer Data Collection Rate	0.0	1.0	+1.0
Online Sales Growth	0.0	0.0	0.0
Effectiveness of Personalized Marketing	0.0	0.0	0.0
Personalized Promotions vs Mass Promotions	0.0	0.0	0.0
Retail Media Growth	0.0	0.0	0.0
Total Insights Recall Score	1.0/7.0	3.0/7.0	+2.0

In Table 12 we extend Table 4 and show all the groundtruth insights as well as each of the predicted insights from using both Llama 3.1 405B and GPT-5. As before, we highlight in bold where the model was able to accurately find details relevant to the expected insight, and show all the corresponding scores as given in Table 11.

G.2 FACTUALITY

The factuality metric evaluation uses the same Task DR0002 to assess the accuracy and reliability of generated content. Table 13 presents the factuality performance comparison, showing that while using Llama 3.1 405B achieved 0.41 factuality (7 factual claims out of 17 total claims), where as using GPT-5 reached 0.65 factuality (13 factual claim out of 20 total claims). This represents a significant improvement

```
1404
        Adaptive Action Generation
1405
1406
        Based on the research progress so far, suggest new actions with
1407
        INTELLIGENT SOURCE COMPLEMENTARITY:
1408
        Original Research Query: {research_query}
1409
1410
        Completed Actions Summary: {completed_actions_summary}
1411
1412
        Recent Findings Analysis: {findings_json}
1413
        Source Composition Analysis: {internal_findings}
1414
1415
        Available Tools: {available_tool_names}
1416
1417
        Generate 1-5 new actions that:
        1. **Address gaps** in current research coverage
1418
            **Balance source types** - if findings are heavily external,
1419
        prioritize internal sources and vice versa
1420
        3. **Build on discoveries** - leverage new information to explore
1421
        related areas
1422
        4. **Enhance depth** - dive deeper into promising
        findings
        5. **Cross-validate** - verify critical findings through alternative
1424
        sources
1425
1426
        Each action should have:
1427
        - Strategic rationale for why this action is needed
1428
        - Clear connection to research gaps or promising
1429
        leads
1430
        - Appropriate priority based on strategic value and source
1431
        balance
1432
        - Specific parameters that build on existing knowledge
1433
1434
```

Prompt 4: Gap-Driven Adaptive Action Generation

in content reliability. This also highlights that GPT-5 is much better prepared to make accurate claims that are sustained by evidence.

Table 14 provides a detailed breakdown of factual versus unfactual claims. The agent using GPT-5 generated 6 additional factual claims while producing 3 fewer unfactual claims, resulting in a net improvement in accuracy percentage. This demonstrates that GPT-5 may generate a higher proportion of factual information than Llama 3.1 405B.

The impact of these factuality improvements of using GPT-5 over Llama 3.1 405B is summarized in Table 15. The 24.0 percentage point improvement in factuality represents enhanced content quality and research reliability. The increase in 3 total claims shows that GPT-5 can generate more content overall.

H EXPERIMENTAL SETTINGS

All experiments were conducted on a cluster of 8 NVIDIA A100 GPUs (80GB each). For file generation and task construction, we primarily used the Llama-3.1-405B model, with decoding performed using nucleus sampling at a temperature of 0.7 unless otherwise specified. For larger-scale evaluations of the DRBench Agent, we also used closed-source models such as GPT-40 and GPT-5, alongside DeepSeek models, to enable comparison across open- and closed-source backbones.

To ensure reproducibility, the DRBench environment was deployed as a self-contained Docker container with all supporting applications (Nextcloud, Mattermost, and Filebrowser) pre-configured. Each task was

Table 12: Insights Recall Improvement Areas (Task DR0002). We highlight in bold where each model was able to accurately find details relevant to the groundtruth insight. We also show the corresponding score where 1.0 is considered a successful recall and 0.0 an unsuccessful recall.

Groundtruth Insight	Insight Predicted by Llama 3.1 405B	Insight Predicted by GPT-5	
45% of our online customers have interacted with personalized product recommendations,	45% of Lee's Market online customers engage with personalized product recom-	45% of online customers engaged with personalized product recommendations	
resulting in a 25% increase in average order	mendations, resulting in a 25% increase in	and among those engagers average order	
value.	average order value. (Score = 1.0)	value increased by 25%. (Score = 1.0)	
As of Q3 2024, Lee's Market has 1.2 million	85% of transactions are linked to loyalty	From Q2 2024 to Q3 2024, loyalty members	
loyalty program members.	accounts at Lee's Market, providing a solid foundation for personalized marketing and improving customer engagement. (Score = 0.0)	increased from 1,050,000 to 1,200,000 (+150,000; +14.29%), average spend per member rose from 24 to 25 (+4.17%), and total member spend increased from 25,200,000 to 30,000,000 (+19.05%). (Score = 1.0)	
85% of Lee's Market transactions are linked	85% of transactions are linked to loyalty	As of Q2 2024, 85% of transactions were	
to customer loyalty accounts as of Q2 2024.	accounts at Lee's Market, providing a solid foundation for personalized marketing and	linked to loyalty accounts, leaving a 15 unlinked identity gap. (Score = 1.0)	
	improving customer engagement. (Score = 0.0)		
Lee's Market online sales grew 12% in Q2 2024 compared to Q2 2023.	45% of Lee's Market online customers engage with personalized product recommendations,	A naive blended online AOV upside of approximately 11.25% is derived from 45%	
2024 compared to Q2 2023.	resulting in a 25% increase in average order	engagement multiplied by a 25% AOV lift	
	value. (Score = 0.0)	among engagers. (Score = 0.0)	
Retailers excelling in personalized marketing	Retailers have seen a consistent 25% increase	Grocers running data-driven loyalty cam	
are growing revenues about 10 percentage points faster than their peers, according to	in revenue due to advanced personalization capabilities.	paigns have realized an average 3.8 like-for-like sales uplift. (Score = 0.0)	
BCG. By effectively using first-party customer data, these leaders could unlock an estimated	(Score = 0.0)		
\$570 billion in additional sales, highlighting			
the importance of data-driven sales strategies for growth.			
Personalized promotions can deliver returns	Retailers have seen a consistent 25% increase	External sources indicate POS-enabled	
three times higher than mass promotions, yet many retailers allocate under 5% of their	in revenue due to advanced personalization capabilities.	personalization can lift revenue 5%-15% and advocate personalized e-receipts with relevant	
promo budgets to personalization. One major chain increased its personalized promo spend	(Score = 0.0)	offers and coupons to extend post-purchase engagement.	
from 1% to 10% by establishing a "customer		(Score = 0.0)	
investment council," resulting in \$250 million in incremental sales.			
Retail media networks are expanding rapidly,	85% of transactions are linked to loyalty	As of Q2 2024, 85% of transactions were	
with retail media growing at approximately 25% annually, offering retailers a profitable	accounts at Lee's Market, providing a solid foundation for personalized marketing and	linked to loyalty accounts , leaving a 15% unlinked identity gap.	
revenue stream to reinvest in technology, data,	improving customer engagement.	(Score = 0.0)	
and personnel. By integrating loyalty data, retailers like Sephora, which links 95% of	(Score = 0.0)		
transactions to loyalty accounts, enhance preci-			
sion in product recommendations and provide a seamless omnichannel experience, boosting			
conversion rates and customer lifetime value.			

Table 13: Factuality Performance Comparison: Llama 3.1 405B vs GPT-5 (Task DR0002). We show the number of factual and unfactual claims made by each model, as well as the overall factuality score for the given task.

Metric	Llama 3.1 405B	GPT-5	Improvement
Factuality Score	0.41	0.65	+0.24
Factual Claims	7	13	+6 claims
Unfactual Claims	10	7	-3 claims

```
1512
        Report Synthesis
1513
1514
        As an expert research analyst, synthesize the following content into
1515
        a coherent, insightful, and well-supported analysis for the theme:
        "\{\text{theme}\}" directly related to the overarching research question:
1516
        "{original_question}"
1517
1518
        Source Priority Guidelines:
1519
1520
        1. **"internal"**: Highest priority (internal company
        documents, proprietary files, confidential reports,
1521
        enterprise chat messages, local documents, CRM data,
1522
        internal APIs, project management tools).
                                                     Insights from
1523
        these sources should form the primary foundation of the
1524
        analysis.
1525
        2. **"external"**: Medium priority (public web sources,
        academic papers, industry reports, news articles). Use these
1526
        to provide broader context, external validation, or contrasting
        perspectives.
1528
1529
        **Synthesis Requirements:**
1530
        * **QUANTITATIVE PRIORITY: ** Lead with numerical data, calculations, and
1531
        aggregations
        * Extract ALL percentages, costs, metrics, and performance
1532
1533
        * Perform mathematical operations: aggregate percentages, calculate
1534
        increases, sum totals
1535
        * Example: "Finance customers (35%) combined with healthcare
        (40%) represent 75% of regulated industry concerns
1536
        [DOC:doc_1] [DOC:doc_2] "
1537
1538
        * **FACT VERIFICATION (CRITICAL): **
1539
        * ONLY state what documents explicitly contain - no inference or
1540
        extrapolation
        * Use exact quotes for key numerical claims: "As stated in the
1541
        document: '[exact quote]' [DOC:doc_id]"
1542
1543
        * **Citation Usage (Critical):**
        * **Format:** Reference sources by their document ID: "Internal review
1545
        shows 15% increase [DOC:doc_079c2e0f_1752503636]"
        * **NEVER HALLUCINATE CITATIONS:** Only use provided doc_id
        values
1547
        * **Cite every numerical claim and calculation with source
1549
1550
        Generate 2-4 paragraphs of synthesized analysis with proper inline
        citations.
1551
```

Prompt 5: Report Section Synthesis with Citation Requirements

1559

1560

1561

1562

1563

1564

1565

executed by instantiating a fresh container to avoid state leakage across runs. We capped the number of agent iterations according to the settings described in Section 5, with each iteration limited by a fixed computational budget.

For model outputs, we standardized all prompts and evaluation pipelines across backbones, using identical research questions, company contexts, and injected insight sets. To avoid stochastic variability, we repeated generation three times per task and reported averaged scores.

Finally, all supporting scripts, environment configurations, and evaluation code are fully containerized, enabling consistent replication of our reported results across hardware setups.

Table 14: Factuality Content Analysis (Task DR0002). We show the number of factual and unfactual claims made by each model, highlighting the factuality accuracy of each model for the given task.

Agent	Factual	Unfactual	Accuracy
Llama 3.1 405B	7 claims	10 claims	41.0%
GPT-5	13 claim	7 claims	65.0%
Change	+6	-3	+24.0%

Table 15: Factuality Summary (Task DR0002). We summarize the factuality result improvements made by using GPT-5 over Llama 3.1 405B.

Impact Category	Value
Factuality Improvement	+24.0 percentage points
Claims Added	3 total claims
Accuracy Enhancement	From 41.0% to 65.0%
Content Quality	More grounded information
Task Domain	Sales
Task Industry	Retail

I DATA SYNTHESIS AND DRBA COST DETAILS

To generate the DRBench tasks, we combined external insight extraction with internal file synthesis. Each task included an average of 10 supporting files spanning heterogeneous formats (PDF, DOCX, PPTX, XLSX, and JSONL), with each file containing roughly 5 paragraphs of content. Files were designed to embed injected insights while mixing in distractor material, ensuring realistic enterprise complexity without exceeding practical runtime or storage budgets.

Data synthesis was primarily powered by GPT-4o. During task construction, GPT-4o was responsible for (1) extracting structured insights from public web sources, (2) adapting these insights into enterprise-grounded interpretations, (3) generating persona-specific deep research questions, and (4) producing file-level content with a balanced mix of insights and distractors. For evaluation, the DRBench Agent (DRBA) used GPT-4o as its backbone, with each task typically requiring 15 iterations and approximately 120–150 model calls.

In terms of cost, GPT-4o-based synthesis of a single task (10 files, 5 paragraphs each, plus metadata) consumed about 30k–40k tokens, while DRBA execution required an additional 50k–70k tokens per task. At current GPT-4o API pricing (\$5 per million input tokens and \$15 per million output tokens), this corresponds to a per-task cost of approximately \$1.5–\$3.5 depending on the mix of input/output tokens and the iteration budget. This makes large-scale benchmarking feasible at moderate cost, while still being significantly cheaper than manual authoring or annotation.

We also note that smaller open-source models such as Llama-3.1-8B-Instruct perform well for file generation. Unlike GPT-4o, which requires API usage, Llama-3.1-8B can be hosted locally and runs efficiently on a single NVIDIA A100 40GB GPU. This provides a cost-effective alternative for generating large numbers of supporting documents, especially when full closed-source quality is not required.

J WEB AGENTS FOR DEEP RESEARCH TASKS

Since each of the environments can be access directly through a web user interface (UI), we also experimented with an agent that can directly interact with the webpages through common browser actions like click, input and scroll, which are executed through *playwright*⁴. We implement our web agent using the AgentLab and BrowserGym frameworks (Chezelles et al., 2025) with a GPT-4.1⁵ backbone. Our agent is implemented from AgentLab's *GenericAgent*, which achieves respectable performance when used

⁴https://playwright.dev

⁵https://openai.com/index/gpt-4-1/

 with GPT-40⁶ as a backbone; it completes 45.5% of the tasks in WorkArena (Drouin et al., 2024), 31.4% in WebArena (Zhou et al., 2024) and achieves a step-level reward of 13.7% on WebLINX (Lù et al., 2024).

Hyperparameters and Prompt We present the hyperparameters for the agent in tables 16 and 17, which are in majority set to the default hyperparameters, except for the maximum number of input tokens (bound to a reasonable maximum length) and a higher maximum number of steps (to allow the agent to perform more actions required to write the report). We further update the agent's action space on the last step to only allow it to reply to the user with a report, ensuring that each trajectory terminates with a report. To ensure that the agent is aware of the tools it can use, we modify the default system prompt (see prompt 6). Additionally, each task intent is provided alongside information about the user and company (see prompt 7).

Results We find that the GPT-4.1-powered web agent achieves an insights recall and factuality of 1.11% and 6.67% respectively and a report quality score of 33.07%. Although the high report quality indicates that the agent can properly formulate a report, the insights quality is severely limited, with none of the claims being backed by useful sources. For example, a DRBench Agent powered by GPT-5 may answer the question *What is Lee's Market's current food waste reduction rate as of Q2 2024?* with *An 8% reduction in food waste in Q2 2024 saved Lee's Market \$1.2 million, indicating that better inventory control can yield both safety and financial benefits.*, which achieves a score of 1.0 for the question. On the other hand, a GPT-4.1-powered agent will provide an unsatisfactory answer, thus achieving an insights recall of 0.0. The most likely cause of this poor performance is the model's limited capability to properly interact with web interfaces when encountering unfamiliar tools. For instance, the agent may be unfamiliar with the VNC and file browser applications, making it harder for it to correctly select the file it needs to use. Moreover, whenever the agent ends up performing an ineffective action (e.g. click on an element that does not trigger any change to the page), it tends to persist by reiterating the same action (see Table 18), or the same sequence of ineffective actions, despite not achieving anything in the previous steps. As a result, despite a large number of steps, most of the agent's actions are not helpful towards solving the task.

Table 16: Web Agents Boolean Hyperparameters

Value	Flags
True	vision_support, use_ax_tree, use_tabs, use_focused_element, use_error_logs, use_history, use_action_history, use_screenshot, use_som, extract_visible_tag, extract_clickable_tag, use_thinking, use_concrete_example, use_abstract_example, use_hints, be_cautious, add_missparsed_messages
False	use_html, use_past_error_logs, use_think_history, use_diff, filter_visible_elements_only, filter_with_bid_only, filter_som_only, multiaction, strict, long_description, individual_examples, use_plan, use_criticise, use_memory, enable_chat

K STANDARD ERROR

Restricting to the MinEval subset, we average the results on each task across 3 different runs in Table 19. We give both the means and standard errors for the insight recall, factuality, distractor avoidance, and report quality.

L COMPLEX MODEL ABLATION RESULTS

Extending our discussion in Section 5.3 to more GPT and Llama models in Table 20, we see that the smaller GPT-5-mini model lags behind but still outperforms earlier closed-source backbones such as GPT-40 and GPT-40-mini, particularly in terms of harmonic mean. In addition, smaller variants of Llama

⁶https://openai.com/index/hello-gpt-4o/

1723 1724

172517261727

```
1675
1676
1677
         Web Agents System Prompt
1678
1679
        You are an agent trying to solve a web task based on the content of the
1680
        page and user instructions. You can interact with the page and explore,
1681
        and send messages to the user. Each time you submit an action it will
1682
        be sent to the browser and you will receive a new page.
        You will be solving tasks that involve Deep Research, by
1683
        navigating websites and web apps with information useful
1684
        solving the task. You will need to gather insights from
1685
        data contained in the services provided and the internet
1686
        to complete your report. You have a maximum of 50 steps
1687
        to complete the task. Before the end, you must use the
        action {send_msg_to_user}, which should contain a final Deep
1688
        Research report detailing everything the user needs to
1689
1690
1691
        You must sustain your claims in files, chats, or emails
1692
        from the enterprise environment or in websites you
1693
        searched. You must provide a citation (or an inline
        citation, that works too) with the source of those claims
1694
         (e.g. << add citation examples from the documentation
1695
        >>). Do not make up citations if you haven't retrieved its
1696
        content.
1697
        Here are some expected agent behavior:
1698
1699
        **Enterprise Environment Interaction:**
1700
        - Access Nextcloud files, Mattermost chats, emails, VNC desktop,
1701
        etc.
1702
        - Extract relevant information from multiple sources
        - Navigate complex enterprise data landscapes
1703
1704
        **Report Requirements:**
1705
        - Synthesize findings into comprehensive research
1706
        report
1707
        - Include proper citations for all claims (flexible format -
1708
        auto-normalized)
        - Draw meaningful insights and conclusions
1709
        - Ground all statements in available evidence
1710
1711
1712
        **Citation Format (Flexible - Auto-Normalized):**
        - Files: 'quarterly_report.pdf', 'shared/budget-analysis.xlsx',
1713
        'Analysis document (reports/analysis.docx)'
1714
        - URLs: Direct links or '[Article Title] (https://example.com)'
1715
        - Emails: 'Email from alice@company.com on Jan 20,
1716
        2025 \
1717
        - Chat:
                 'Mattermost message from john.doe in Compliance team, General
        channel'
1718
1719
1720
        If you need navigate the internet (outside of the designated websites),
1721
        you can use the browser inside noVNC Desktop.
1722
```

Prompt 6: Extended instructions given to the Deep Research web agent.

Table 17: Web Agents Non-Boolean Hyperparameters

Parameter	Value
chat_model.model_name chat_model.max_total_tokens chat_model.max_input_tokens	gpt-4.1 32768 28672
chat_model.max_new_tokens chat_model.temperature action.action_set.subsets	4096 0 webarena
action.action_set.retry_with_force flags.max_prompt_tokens flags.max_trunc_itr	true 28672 20
env.max_steps	50

```
How can Lee's Market leverage FSMA 204 regulations to enhance food safety and customer trust?

Here is some information about the user:
ID: MNG0003
First name: John
Last name: Doe
...
Justification: As a regulatory affairs manager, John
...

Here is some information about the company:
Name: Lee's Market
Annual revenue: $500M - $600M
...
Target markets: Asian communities in the U.S. and ...
```

Prompt 7: Web Agents Task Intent Prompt

Table 18: Web Agents tends to get stuck on cycles of actions, and are unable to backtrack or to restart with a different application.



degrade metric results further. This further substantiates the claim that larger and more advanced models tend to offer a better balance between recall, factuality, and overall report quality.

M QUANTITATIVE RESULTS PER TASK

We show a detailed breakdown of the insights recall in Table 21, factuality in Table 22, distractor avoidance in Table 23 and report quality in Table 24 on the MinEval subset for a variety of models.

Table 19: DRBA performance with different planning configurations on MinEval. We compare the base agent with variants using Simple Research Planning (SRP), Complex Research Planning (CRP), Adaptive Action Planning (AAP), and their combinations. Scores are reported for insight recall, factuality, distractor avoidance, report quality, and the overall harmonic mean.

Configuration	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
Base DRBA	.188 ± .038	$.665 \pm .09$.981 ± .01	.912 ± .004	.448
+ SRP	.137 ± .04	.622 ± .11	$1.00 \pm .00$.901 ± .006	.365
+ CRP	$.154 \pm .03$	$.705 \pm .05$	$1.00 \pm .00$	$.917 \pm .004$.400
+ AAP	$.197 \pm .03$	$.604 \pm .08$	$.995 \pm .39$	$.928 \pm .005$.454
+ SRP + AAP	.142 ± .04	.504 ± .10	.990 ± .39	.906 ± .007	.359
+ CRP + AAP	$.188 \pm .05$	$.691 \pm .03$	$1.00 \pm .00$	$.923 \pm .006$.453

Table 20: Performance of DRBA on the MinEval subset using different backbone language models and planning strategies. Scores are reported for insight recall, factuality, distractor avoidance, report quality, and harmonic mean. Note that higher numbers corresponds to better scores, and the best result on each metric is bolded.

Model	Planning	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
GPT-5	None	38.33	74.52	95.14	94.56	79.80
GPT-5	Simple	37.86	72.09	97.14	95.34	78.92
GPT-5	Complex	39.63	65.17	92.86	93.42	77.74
GPT-5-mini	None	25.68	58.76	96.51	84.48	60.21
GPT-5-mini	Simple	28.37	58.96	95.81	85.74	63.73
GPT-5-mini	Complex	26.57	51.07	94.48	86.28	58.89
GPT-40	None	17.53	65.43	99.05	92.58	48.47
GPT-40	Simple	20.37	62.35	98.57	93.12	53.06
GPT-40	Complex	17.31	60.84	98.33	91.62	47.35
GPT-4o-mini	None	13.75	46.68	99.05	84.72	38.33
GPT-4o-mini	Simple	13.67	55.59	97.14	85.86	39.39
GPT-4o-mini	Complex	13.08	48.95	97.14	86.04	37.28
GPT-OSS-120B	None	22.40	29.24	97.14	84.24	44.82
GPT-OSS-120B	Simple	17.42	27.48	97.14	84.36	38.38
GPT-OSS-120B	Complex	18.31	38.92	98.1	85.44	44.14
Llama-3.1-405B-Instruct	None	17.37	78.91	100.00	90.48	49.78
Llama-3.1-405B-Instruct	Simple	16.97	79.27	98.10	92.34	48.87
Llama-3.1-405B-Instruct	Complex	20.16	69.75	97.90	91.26	53.86
Llama-3.1-70b-Instruct	None	18.54	64.64	96.70	85.32	50.08
Llama-3.1-70b-Instruct	Simple	16.28	69.43	97.62	84.96	46.41
Llama-3.1-70b-Instruct	Complex	17.02	52.82	98.10	86.16	45.46
DeepSeek-V3.1	None	25.15	72.66	97.43	86.52	62.59
DeepSeek-V3.1	Simple	25.56	73.45	96.67	87.36	63.29
DeepSeek-V3.1	Complex	30.26	70.27	96.67	86.88	69.28
Qwen-2.5-72B-Instruct	Complex	26.82	58.35	97.65	89.64	61.75
Qwen-2.5-72B-Instruct	None	25.55	69.39	98.10	90.24	62.64
Qwen-2.5-72B-Instruct	Simple	23.20	67.23	98.10	88.14	58.58

N EFFECT OF NUMBER OF ITERATIONS

We next analyze the effect of varying the number of research loop iterations when using DRBA with GPT-5 as the backbone language model. Results for both the baseline configuration without explicit planning and the complex planning setup are shown in Table 25. Overall, increasing the iteration budget does not guarantee consistent improvements. With no planning, performance initially drops when the agent executes more iterations, as additional exploration often introduces noise and distracts from key insights. However, with a larger budget the agent partially recovers, suggesting that a small number of additional iterations can help refine factual grounding, while excessive exploration reduces focus.

For the complex planning setting, higher iterations improve certain metrics such as factuality, but this comes at the cost of lower insight recall and reduced overall balance. This indicates that while more steps allow the agent to verify citations more carefully, they can also lead to fragmented reasoning and overfitting

gpt-5 enterprise_fact 0.597 external_fact 0.0 qwen-2.5-72b-instruct enterprise_fact 0.417 external_fact 0.0			
enterprise_fact 0.597 external_fact 0.00 qwen-2.5-72b-instruct enterprise_fact 0.417	336		
enterprise_fact 0.597 external_fact 0.00 qwen-2.5-72b-instruct enterprise_fact 0.417	337		
enterprise_fact 0.597 external_fact 0.00 qwen-2.5-72b-instruct enterprise_fact 0.417	338		
external_fact 0.0 qwen-2.5-72b-instruct enterprise_fact 0.417	339		gpt-5
qwen-2.5-72b-instruct	340	enterprise_fact	0.597
enterprise_fact 0.417	341	external_fact	0.0
enterprise_fact 0.417	342		
enterprise_fact 0.417	3		
enterprise_fact 0.417	14		
enterprise_fact 0.417	5		
enterprise_fact 0.417	6		
	7		qwen-2.5-72b-instruct
external_fact 0.0		enterprise_fact	0.417
		external_fact	0.0

eepseek-chat-v3.1		gpt-5-mini
0.472	enterprise_fact	0.444
0.0	external_fact	0.0

	llama-3.1-405b-instruct
enterprise_fact	0.347
external_fact	0.0

	gpt-oss-120b
enterprise_fact	0.333
external_fact	0.0

	gpt-4o-mini
enterprise_fact	0.194
external_fact	0.0

	llama-3.1-70b-instruct
enterprise_fact	0.194
external_fact	0.0

	gpt-4o		
enterprise_fact	0.182		
external_fact	0.0		

Figure 8: Total average Insight Recall scores per model and insight source type computed on all the results available for each model running in Complex Research Plan mode. Insights embedded in enterprise sources are more easily retrieved by DRBA in all the models.

to peripheral evidence. The best overall performance emerges at moderate iteration counts, highlighting the importance of carefully tuning the iteration budget rather than simply scaling up the number of steps.

O DATA GENERATION PROMPTS

O.1 COMPANY AND PERSONA DATA GENERATION

In this section we give the prompts used for company generation 8 and persona generation 9.

```
1873
1874
         7 Company Generation Prompt
1875
1876
         Generate a realistic company structure for {company_name} in the
         {industry} industry.
1877
1878
         Company size: {size} ({employee_range} employees)
1879
1880
         The company should focus on this domain {domain}
1881
         EXTERNAL INSIGHTS: {external_insights}
1882
1883
         Return ONLY a valid JSON object with this structure:
1884
         {output_structure}
1885
1886
        Make it realistic for the {industry} industry.
1887
```

Prompt 8: Company Generation Prompt Template.

Table 21: Mean and standard error of the insight recall metric for the first five tasks, obtained from three runs of on *DRBench* our agent (DRBA) using 15 iterations across different backbone models.

Configuration	Plan	DR0001	DR0002	DR0003	DR0004	DR0005
GPT-5	None	$.222 \pm .056$.429 ± .000	$.467 \pm .033$.519 ± .098	.281 ± .035
GPT-5	Simple	$.222 \pm .056$	$.381 \pm .048$	$.533 \pm .033$	$.370 \pm .098$	$.386 \pm .046$
GPT-5	Complex	$.278 \pm .056$	$.381 \pm .126$	$.567 \pm .067$	$.370 \pm .037$	$.386 \pm .046$
GPT-5-mini	None	$.111 \pm .056$	$.286 \pm .000$	$.367 \pm .033$	$.222 \pm .064$	$.298 \pm .018$
GPT-5-mini	Simple	$.222 \pm .056$	$.286 \pm .082$	$.333 \pm .033$	$.296 \pm .074$	$.281 \pm .063$
GPT-5-mini	Complex	$.000 \pm .000$	$.381 \pm .126$	$.367 \pm .067$	$.370 \pm .098$	$.211 \pm .053$
GPT-4o	None	.111 ± .056	$.238 \pm .048$	$.167 \pm .033$	$.185 \pm .074$	$.175 \pm .018$
GPT-4o	Simple	$.167 \pm .000$	$.333 \pm .048$	$.300 \pm .000$	$.148 \pm .037$	$.070 \pm .018$
GPT-4o	Complex	$.111 \pm .056$	$.238 \pm .095$	$.300 \pm .100$	$.111 \pm .000$	$.105 \pm .000$
GPT-4o-mini	None	$.000 \pm .000$	$.095 \pm .048$	$.267 \pm .033$	$.185 \pm .037$	$.140 \pm .035$
GPT-4o-mini	Simple	$.056 \pm .056$	$.190 \pm .048$	$.167 \pm .067$	$.148 \pm .074$	$.123 \pm .018$
GPT-4o-mini	Complex	$.000 \pm .000$	$.190 \pm .048$	$.267 \pm .033$	$.074 \pm .037$	$.123 \pm .018$
GPT-OSS-120B	None	.111 ± .111	$.143 \pm .000$	$.433 \pm .033$	$.222 \pm .000$	$.211 \pm .030$
GPT-OSS-120B	Simple	$.056 \pm .056$	$.143 \pm .000$	$.367 \pm .033$	$.148 \pm .098$	$.158 \pm .061$
GPT-OSS-120B	Complex	$.000 \pm .000$	$.190 \pm .048$	$.333 \pm .067$	$.111 \pm .064$	$.281 \pm .063$
Llama-3.1-405B-Instruct	None	$.167 \pm .000$	$.143 \pm .000$	$.233 \pm .088$.185 ± .074	$.140 \pm .035$
Llama-3.1-405B-Instruct	Simple	$.222 \pm .056$	$.190 \pm .048$	$.167 \pm .033$	$.111 \pm .064$	$.158 \pm .053$
Llama-3.1-405B-Instruct	Complex	$.111 \pm .056$	$.238 \pm .048$	$.300 \pm .058$	$.148 \pm .098$	$.211 \pm .030$
Llama-3.1-70B-Instruct	None	$.167 \pm .000$	$.381 \pm .048$	$.200 \pm .000$	$.074 \pm .037$	$.105 \pm .030$
Llama-3.1-70B-Instruct	Simple	$.056 \pm .056$	$.286 \pm .082$	$.200 \pm .058$	$.185 \pm .098$	$.088 \pm .046$
Llama-3.1-70B-Instruct	Complex	$.167 \pm .000$	$.238 \pm .048$	$.267 \pm .033$	$.074 \pm .074$	$.105 \pm .000$
DeepSeek-V3.1	None	$.167 \pm .000$	$.286 \pm .082$	$.300 \pm .058$	$.259 \pm .037$	$.246 \pm .046$
DeepSeek-V3.1	Simple	$.278 \pm .056$	$.238 \pm .048$	$.333 \pm .033$	$.148 \pm .074$	$.281 \pm .046$
DeepSeek-V3.1	Complex	$.167 \pm .000$	$.095 \pm .048$	$.600 \pm .058$	$.370 \pm .037$	$.281 \pm .035$
Qwen-2.5-72B-Instruct	None	.222 ± .056	$.238 \pm .048$	$.400 \pm .058$	$.259 \pm .037$.158 ± .061
Qwen-2.5-72B-Instruct	Simple	.111 ± .056	$.333 \pm .048$	$.333 \pm .088$	$.259 \pm .037$	$.123 \pm .018$
Qwen-2.5-72B-Instruct	Complex	$.167 \pm .096$	$.143 \pm .082$	$.400 \pm .058$	$.333 \pm .000$	$.298 \pm .076$

Persona Generation Prompt

Generate $\{persona_count\}$ diverse employee personas for $\{company_name\}$ in the $\{industry\}$ industry.

The personas should focus on this domain: {domain}

Create diverse roles across seniority levels: Junior, Mid, Senior, Executive $% \left(1\right) =\left(1\right) +\left(1\right$

Return ONLY a valid JSON array with this exact format:
{output_structure}

Make personas realistic with appropriate responsibilities for their roles.

Prompt 9: Persona Generation Prompt.

O.2 QUESTION GENERATION

In this section we give the prompt used to generate our deep research questions 10.

Table 22: Mean and standard error of the factuality metric for the first five tasks, obtained from our agent (DRBA) using 15 iterations across different backbone models.

Configuration	Plan	DR0001	DR0002	DR0003	DR0004	DR0005
GPT-5	None	.761 ± .072	$.504 \pm .075$	$.866 \pm .007$.833 ± .019	$.762 \pm .077$
GPT-5	Simple	$.714 \pm .050$	$.384 \pm .076$	$.848 \pm .034$	$.812 \pm .070$	$.846 \pm .029$
GPT-5	Complex	$.730 \pm .060$	$.291 \pm .094$	$.782 \pm .012$	$.782 \pm .064$	$.674 \pm .121$
GPT-5-mini	None	$.585 \pm .045$	$.297 \pm .119$	$.705 \pm .039$	$.704 \pm .067$	$.647 \pm .098$
GPT-5-mini	Simple	.647 ± .120	$.299 \pm .056$	$.624 \pm .041$	$.694 \pm .028$	$.683 \pm .020$
GPT-5-mini	Complex	$.309 \pm .126$	$.381 \pm .161$	$.699 \pm .047$	$.692 \pm .111$	$.472 \pm .114$
GPT-40	None	.792 ± .150	$.490 \pm .110$	$.827 \pm .056$	$.570 \pm .058$	$.593 \pm .20$
GPT-40	Simple	.485 ± .262	$.512 \pm .131$	$.813 \pm .041$	$.693 \pm .139$	$.614 \pm .12$
GPT-4o-mini	Simple	.475 ± .166	$.653 \pm .097$	$.704 \pm .037$	$.542 \pm .110$	$.406 \pm .02$
GPT-40	Complex	$.828 \pm .043$	$.265 \pm .133$	$.800 \pm .000$	$.690 \pm .128$	$.459 \pm .23$
GPT-4o-mini	None	.611 ± .056	$.429 \pm .092$	$.622 \pm .062$	$.481 \pm .209$.191 ± .04
GPT-4o-mini	Complex	$.557 \pm .030$	$.324 \pm .169$	$.580 \pm .075$	$.642 \pm .119$.344 ± .14
GPT-OSS-120B	None	.144 ± .099	$.150 \pm .035$	$.386 \pm .040$	$.337 \pm .117$	$.445 \pm .05$
GPT-OSS-120B	Simple	$.074 \pm .074$	$.128 \pm .072$	$.410 \pm .090$	$.311 \pm .155$	$.451 \pm .08$
GPT-OSS-120B	Complex	$.368 \pm .061$	$.178 \pm .078$	$.564 \pm .064$	$.400 \pm .076$	$.435 \pm .19$
Llama-3.1-405B-Instruct	None	.852 ± .087	.726 ± .158	.803 ± .028	$.820 \pm .066$.745 ± .02
Llama-3.1-405B-Instruct	Simple	$.802 \pm .125$	$.638 \pm .202$	$.800 \pm .074$	$.892 \pm .035$	$.832 \pm .08$
Llama-3.1-405B-Instruct	Complex	$.789 \pm .053$	$.392 \pm .154$	$.792 \pm .055$	$.771 \pm .073$	$.745 \pm .10$
Llama-3.1-70B-Instruct	None	.618 ± .109	$.431 \pm .160$	$.684 \pm .104$	$.812 \pm .021$	$.687 \pm .07$
Llama-3.1-70B-Instruct	Simple	$.608 \pm .173$	$.681 \pm .069$	$.800 \pm .069$	$.826 \pm .067$	$.557 \pm .14$
Llama-3.1-70B-Instruct	Complex	$.588 \pm .082$	$.286 \pm .108$	$.686 \pm .011$	$.522 \pm .270$	$.559 \pm .24$
DeepSeek-V3.1	None	$.860 \pm .014$	$.518 \pm .085$	$.818 \pm .041$	$.679 \pm .095$	$.757 \pm .05$
DeepSeek-V3.1	Simple	.696 ± .041	$.531 \pm .086$	$.922 \pm .056$	$.769 \pm .035$	$.754 \pm .08$
DeepSeek-V3.1	Complex	$.581 \pm .042$	$.657 \pm .024$	$.838 \pm .050$	$.774 \pm .053$	$.662 \pm .09$
Qwen-2.5-72B-Instruct	None	.674 ± .077	.493 ± .109	$.866 \pm .002$.741 ± .060	.696 ± .06
Qwen-2.5-72B-Instruct	Simple	$.806 \pm .049$	$.540 \pm .174$	$.741 \pm .074$	$.724 \pm .101$	$.550 \pm .14$
Qwen-2.5-72B-Instruct	Complex	.626 ± .114	$.396 \pm .056$	$.723 \pm .053$	$.587 \pm .139$	$.586 \pm .05$

O.3 PUBLIC SOURCE AND INSIGHT COLLECTION

In this section we give the prompt used to generate external insights 11. The URLs used for external insight extraction and deep research question creation can be found in Table 26.

O.4 INTERNAL INSIGHT GENERATION

In this section we give the prompts to generate both internal insights 12 and internal distractors 13.

O.5 FILE GENERATION

In this section we give the prompts used for generating each of the file types used in our tasks, which we list as follows:

• **PDF**: Prompts 18, 19, and 20

• Excel: Prompts 21, 22, and 23

• Powerpoint: Prompts 24, 25, and 26

• Email: Prompts 27, 28, and 29

• Chat: Prompts 30, 31, 29, and 33

Table 23: Mean and standard error of the distractor avoidance metric for the first five tasks, obtained from three runs of on *DRBench* our agent (DRBA) using 15 iterations across different backbone models.

Configuration	Plan	DR0001	DR0002	DR0003	DR0004	DR0005
GPT-5	None	$.857 \pm .000$	$.900 \pm .058$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-5	Simple	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-5	Complex	$.905 \pm .048$	$.900 \pm .058$	$.933 \pm .000$	$.905 \pm .024$	$1.00 \pm .000$
GPT-5-mini	None	$.905 \pm .048$	$.967 \pm .033$	$.978 \pm .022$	$.976 \pm .024$	$1.00 \pm .000$
GPT-5-mini	Simple	$.857 \pm .000$	$.933 \pm .033$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-5-mini	Complex	$.857 \pm .000$	$.867 \pm .133$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-40	None	$.952 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-40	Simple	$.952 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$.976 \pm .024$	$1.00 \pm .000$
GPT-4o	Complex	$.952 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$.964 \pm .036$	$1.00 \pm .000$
GPT-4o-mini	None	$.952 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-4o-mini	Simple	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-4o-mini	Complex	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-OSS-120B	None	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-OSS-120B	Simple	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
GPT-OSS-120B	Complex	$.905 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
Llama-3.1-405B-Instruct	None	$1.00 \pm .000$				
Llama-3.1-405B-Instruct	Simple	$.905 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
Llama-3.1-405B-Instruct	Complex	$.952 \pm .048$	$.967 \pm .033$	$1.00 \pm .000$	$.976 \pm .024$	$1.00 \pm .000$
Llama-3.1-70B-Instruct	None	$.905 \pm .048$	$1.00 \pm .000$	$.978 \pm .022$	$.952 \pm .048$	$1.00 \pm .000$
Llama-3.1-70B-Instruct	Simple	$.905 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$.976 \pm .024$	$1.00 \pm .000$
Llama-3.1-70B-Instruct	Complex	$.905 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
DeepSeek-V3.1	None	$.905 \pm .048$	$.967 \pm .033$	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$
DeepSeek-V3.1	Simple	$.857 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$.976 \pm .024$	$1.00 \pm .000$
DeepSeek-V3.1	Complex	$.905 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$.929 \pm .000$	$1.00 \pm .000$
Qwen-2.5-72B-Instruct	None	$1.00 \pm .000$	$1.00 \pm .000$	$1.00 \pm .000$	$.905 \pm .024$	$1.00 \pm .000$
Qwen-2.5-72B-Instruct	Simple	$.952 \pm .048$	$1.00 \pm .000$	$1.00 \pm .000$	$.952 \pm .024$	$1.00 \pm .000$
Qwen-2.5-72B-Instruct	Complex	$.952 \pm .048$	$1.00 \pm .000$	$.978 \pm .022$	$.952 \pm .048$	$1.00 \pm .000$

P EVALUATION PROMPTS

In this section we give the prompts for decomposing reports into atomic insights 14, computing insight recall 15, computing factuality 16, and computing report quality 17. These prompts are discussed in detail in Section 5.1.

Q HUMAN PREFERENCE EVALUATION

We calculate a human score for model a task t as:

$$S_{a,t} \!=\! \frac{1}{n} \!\sum_{i=1}^n \! s_{a,i} \text{, where } s_{a,i} \!=\! \begin{cases} 1, & \text{if human_choice is } a \text{ or "both good"} \\ 0, & \text{otherwise} \end{cases}$$

, where n is the number of gold insights in task t, $s_{a,i}$ is the human score of model a on insight i. Figure 9 shows that the insight recall metric is on par with human decision.

Table 24: Mean and standard error of the report quality metric for the first five tasks, obtained from three runs of *DRBench* with 15 iterations across different backbone models.

Configuration	Plan	DR0001	DR0002	DR0003	DR0004	DR0005
GPT-5	None	$.936 \pm .008$.918 ± .004	$.924 \pm .005$.909 ± .001	.927 ± .009
GPT-5	Simple	$.942 \pm .000$	$.927 \pm .008$	$.936 \pm .006$	$.915 \pm .002$	$.933 \pm .007$
GPT-5	Complex	.948 ± .007	$.921 \pm .001$	$.940 \pm .008$	$.922 \pm .009$	$.929 \pm .008$
GPT-5-mini	None	.892 ± .002	$.884 \pm .007$	$.879 \pm .004$.891 ± .000	$.886 \pm .005$
GPT-5-mini	Simple	.901 ± .009	$.889 \pm .002$	$.887 \pm .001$	$.895 \pm .008$	$.892 \pm .003$
GPT-5-mini	Complex	$.896 \pm .001$	$.882 \pm .006$	$.884 \pm .003$	$.889 \pm .009$	$.890 \pm .001$
GPT-4o	None	$.927 \pm .000$.911 ± .003	.903 ± .001	.918 ± .009	$.909 \pm .002$
GPT-4o	Simple	$.934 \pm .008$	$.919 \pm .000$	$.911 \pm .009$	$.923 \pm .001$	$.916 \pm .000$
GPT-4o	Complex	.929 ± .009	$.914 \pm .002$	$.905 \pm .000$	$.920 \pm .008$	$.913 \pm .009$
GPT-4o-mini	None	$.886 \pm .004$	$.874 \pm .008$.861 ± .007	$.872 \pm .002$	$.879 \pm .006$
GPT-4o-mini	Simple	$.893 \pm .002$	$.881 \pm .005$	$.867 \pm .004$	$.878 \pm .001$	$.884 \pm .003$
GPT-4o-mini	Complex	$.889 \pm .003$	$.877 \pm .006$	$.864 \pm .005$	$.875 \pm .000$	$.882 \pm .004$
GPT-OSS-120B	None	.872 ± .007	.861 ± .001	$.849 \pm .009$	$.858 \pm .006$	$.866 \pm .008$
GPT-OSS-120B	Simple	$.878 \pm .006$	$.867 \pm .009$	$.854 \pm .008$	$.863 \pm .004$	$.872 \pm .007$
GPT-OSS-120B	Complex	$.874 \pm .007$	$.863 \pm .000$	$.851 \pm .008$	$.860 \pm .005$	$.869 \pm .006$
Llama-3.1-405B-Instruct	None	$.914 \pm .000$	$.903 \pm .003$.897 ± .001	$.909 \pm .008$.902 ± .002
Llama-3.1-405B-Instruct	Simple	$.921 \pm .008$	$.910 \pm .001$	$.904 \pm .000$	$.915 \pm .009$	$.908 \pm .000$
Llama-3.1-405B-Instruct	Complex	$.917 \pm .009$	$.906 \pm .002$	$.899 \pm .001$	$.911 \pm .008$	$.905 \pm .001$
Llama-3.1-70B-Instruct	None	$.889 \pm .003$	$.877 \pm .007$	$.869 \pm .005$.881 ± .001	$.873 \pm .004$
Llama-3.1-70B-Instruct	Simple	$.895 \pm .002$	$.883 \pm .005$	$.874 \pm .004$	$.886 \pm .000$	$.878 \pm .003$
Llama-3.1-70B-Instruct	Complex	$.891 \pm .003$	$.879 \pm .006$	$.871 \pm .005$	$.883 \pm .001$	$.875 \pm .004$
DeepSeek-V3.1	None	.884 ± .004	$.872 \pm .008$	$.864 \pm .006$	$.876 \pm .002$	$.869 \pm .005$
DeepSeek-V3.1	Simple	$.890 \pm .003$	$.878 \pm .006$	$.870 \pm .005$	$.881 \pm .001$	$.874 \pm .004$
DeepSeek-V3.1	Complex	$.886 \pm .004$	$.874 \pm .007$	$.866 \pm .006$	$.878 \pm .002$	$.871 \pm .005$
Qwen-2.5-72B-Instruct	None	.901 ± .001	$.889 \pm .005$	$.881 \pm .002$.893 ± .009	$.885 \pm .003$
Qwen-2.5-72B-Instruct	Simple	$.908 \pm .000$	$.896 \pm .003$	$.888 \pm .001$	$.899 \pm .008$	$.891 \pm .002$
Qwen-2.5-72B-Instruct	Complex	$.904 \pm .001$	$.892 \pm .004$	$.884 \pm .002$	$.895 \pm .009$	$.888 \pm .003$

Table 25: Effect of the number of research loop iterations on the performance of DRBA with GPT-5 as the backbone model, on MinEval.

Planning Method	# Iterations	Insight Recall	Factuality	Distractor Avoidance	Report Quality	Harmonic Mean
None	15	39.45	72.65	93.14	94.56	66.20
None	30	28.80	69.03	98.57	96.12	57.34
None	50	37.10	78.84	100.00	93.48	66.30
Complex	15	44.44	62.51	90.95	93.12	66.41
Complex	30	31.61	73.94	94.38	94.76	60.32
Complex	50	38.16	66.05	94.38	92.64	63.76

```
2107
2108
         Deep Research Question Generation Prompt
2109
2110
        Generate 3 Deep Research (DR) questions for the following business
        context:
2111
2112
        Persona: {persona_name} - {persona_role}
2113
        Department: {persona_department}
2114
        Responsibilities: {persona_responsibilities}
2115
        Company: {company_name} ({company_industry})
        Domain: {domain}
2116
2117
        External Insights: {external_insights}
2118
2119
        Generate 3 Deep Research questions that:
2120
        1. Are appropriate for the persona's role and department
        2. Require analysis of the provided internal insights
2121
        3. Consider the external market context ...
2122
2123
        Each question should be 1 sentence of 15 words max, in
2124
        plain english, and end with a question mark. Do not
2125
        include any preamble or explanation - return only the JSON
        array.
2126
2127
        Return ONLY a valid JSON array with this structure:
2128
         {output_structure}
2129
2130
2131
```

Prompt 10: Deep Research Question Generation Prompt Template. Subquestions are generated to help human annotators select good DR questions.

```
You will be given a report (with url, and date). Based on the report, generate 3 external insights that summarize important findings, trends, or takeaways from the report.

Output Format {output_structure}

Url: {url}
Industry: {industry}
Domain: {domain}
Company Information: {company_information}

Important notes
Focus only on insights that are external and grounded in the report.
Insights should be concise, factual, and directly tied to the retail industry context.
```

Prompt 11: External Insight Extraction Prompt.

Table 26: Public URLs For Deep Research Task Creation.

Industry	Domain	Reference
Retail	Compliance	Grocers on FSMA-204 Compliance (GroceryDive)
Retail	CRM	Grocery Loyalty & Inflation (EagleEye)
Retail	Market Analysis	Grocery Trends Outlook 2025 (GroceryDive)
Retail	ITSM	Retail IT Optimization (Thirdera)
Retail	CSM	Chatbots & Grocery Interactions (GroceryDoppio)
Retail	Knowledge Mgmt	Retail Knowledge Management (Knowmax)
Retail	Sales	Personalization in Action (BCG)
Retail	Cybersecurity	Retail Cybersecurity Threats (VikingCloud)
Retail	Public Relations	Walmart CSR Strategy (SunriseGeek)
Healthcare	Compliance	Telehealth Regulations (HealthcareDive)
Healthcare	CRM	Future of Healthcare CRM (WTT Solutions)
Healthcare	Market Analysis	Future of Telehealth (CHG Healthcare)
Healthcare	ITSM	Healthcare ITSM (Topdesk)
Healthcare	CSM	Patient Engagement Tech (TechTarget)
Healthcare	Knowledge Mgmt	Knowledge Mgmt in Healthcare (C8Health)
Healthcare	Sales	Sales for Digital Health (Medium)
Healthcare	Marketing	Marketing Telehealth Services (MarketingInsider)
Healthcare	Cybersecurity	Healthcare Cybersecurity 2024 (AHA)
Automobiles	Compliance	Evolving EV Regulations (WardsAuto)
Automobiles	CRM	Salesforce Automotive Cloud (TechTarget)
Automobiles	CSM	EV Aftersales Support (EVReport)
Automobiles	Sales	Tesla vs Dealerships (TheWeek)
Automobiles	Research	AI for EV Optimization (Here.com)
Automobiles	Cybersecurity	Cybersecurity Risks in Cars (HelpNetSecurity)
Automobiles	Quality Assurance	EV Quality Issues (GreenCars)
Automobiles	Asset Mgmt	Digital Twins in Autos (RTInsights)
Automobiles	Market Analysis	Global EV Outlook 2024 (IEA)

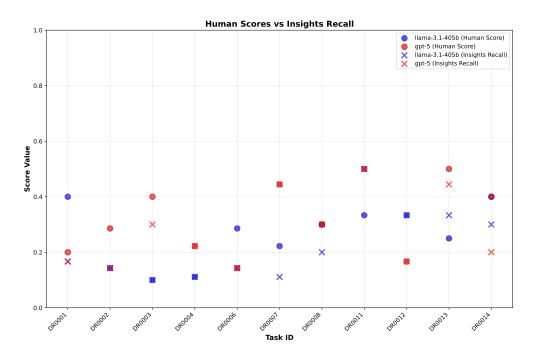


Figure 9: Comparison of Human Scores and Insight Recall Scores. As can be seen the human evaluation results are aligned with our automated evaluation.

2261 2262

```
2215
2216
2217
2218
2219
        Internal Supporting Insight Generation Prompt
2220
        Based on the Deep Research (DR) Question, external
2221
        market insights, company context, and previous internal
2222
        insights, generate 3 specific QA pair that an expert data
2223
        scientist would need to get in order to address the DR
2224
        Ouestion.
2225
        Company: {company_name} - {company_description}
2226
        Industry: {industry}
2227
        Company Size: {company_size} ({employee_count})
2228
        Annual Revenue: {annual_revenue}
2229
        Persona Context: {persona_context}
        DR Question: {dr_question}
2230
        External Market Context (for inspiration): {external_context}
2231
        QA History: {qa_list}
2232
2233
        Insight Requirements:
2234
        - Use the DR Question as the central theme for the
2235
        - Draw inspiration and supporting details from the other internal
2236
        insights and external insights provided.
2237
        - Include QUANTITATIVE DATA in the answer: metrics, percentages, dollar
2238
        amounts, timeframes, KPIs.
2239
        Specific Question Instructions
2240
        - the specific question should be a question that would be a step
2241
        towards resolving the DR Question
2242
        {additional_question_instructions}
2243
2244
        Answer Instructions
2245
        - the answer should be 12 words max and minimum 5
        words
2246
        {additional_answer_instructions}
2247
2248
        Justification Instructions
2249
        - the justification should directly explain how the specific question,
2250
        answer pair help address the DR Question in 15 words
        max
2251
2252
        Misc Instructions
2253
        - the filename should be 3 words max with dashes
2254
        in between, do not mention the file type in the
        filename
2255
        - use the example below as inspiration but do not use it
2256
        directly
2257
2258
        Return ONLY a valid JSON object with this exact structure:
2259
        {output_structure}
2260
```

Prompt 12: Internal Supporting Insight Generation Prompt Template. *specific_questions* and *justification* help human annotators to select good insights.

```
2268
2269
         Internal Distractor Insight Generation Prompt
2270
2271
        Based on the Deep Research (DR) Question, external
2272
        market insights, company context, and previous internal
2273
        insights, generate 3 specific QA pairs that are DISTRACTOR
2274
        questions
        - these should be questions that an expert data scientist might
2275
        ask about the company but are NOT relevant to addressing the DR
2276
        Question.
2277
2278
        Company: {company_name} - {company_description}
2279
        Industry: {industry}
        Company Size: {company_size} ({employee_count})
2280
        Annual Revenue: {annual_revenue}
2281
        Persona Context: {persona_context}
2282
        DR Question: {dr_question}
2283
        External Market Context (for inspiration): {external_context}
        QA History: {qa_list}
2284
2285
        DISTRACTOR Requirements:
2286
        - Generate questions that are plausible for this
2287
        company and industry but DO NOT help address the DR
2288
        Ouestion.
        - The questions should be about the company's operations, metrics, or
2289
        business but tangential to the DR Question.
2290
        - Include QUANTITATIVE DATA in the answer: metrics, percentages, dollar
2291
        amounts, timeframes, KPIs.
        - Focus on different business areas that are NOT central to the
2293
        DR Question (e.g. if DR Question is about pricing, ask about HR
2294
        metrics).
        Specific Question Instructions
2295
        - the specific question should be a plausible business
2296
        question for this company but NOT related to the DR
2297
        Question
2298
        - the specific_question should lead to a quantitative answer and should
2299
        be dated such as Q3 of 2025, the question should contain a date like Q2
        of 2024
        - the specific_question should be 10 words max
2301
        - make sure to be different from any question in the
2302
        ga_list
2303
        - choose business areas like: HR metrics, facility costs, IT
2304
        infrastructure, compliance, training, etc. that are UNRELATED to the
        DR Question
2305
        - make sure to be different from any question in the QA
2306
        History
2307
2308
        Answer Instructions {answer_instructions}
2309
        Justification Instructions
2310
        - the justification should explain why this specific_question
2311
        is NOT relevant to addressing the DR Question in 15 words
2312
        max
2313
        Misc Instructions {misc_instructions}
2314
2315
        Return ONLY a valid JSON object with this exact structure:
2316
         {output_structure}
2317
```

Prompt 13: Internal Distractor Insight Generation Prompt. *specific_questions* and *justification* help human annotators to select good insights.

2320

```
2322
2323
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
        Breaking Report into Insights Prompt
2336
        Please break down the following report text into insight claims. Each
2337
        insight claim should be:
2338
2339
        1. A single insight, that might include multiple statements and
2340
        claims
2341
        2. Independent and self-contained
        3. Each claim can have more than one sentence, but should be focused
2342
        on a single insight
2343
        4. Support each insight with citations from the report text following
2344
        these specific rules: {rules}
2345
        5. Citations should be in one of these formats (various formats will
        be automatically normalized):
2346
        {citation_formats} 6. Do not include general summaries, opinions,
2347
        or claims that lack citation, just the sentences that are
2348
2349
        7. Each claim should be a concise but complete
2350
        sentence.
2351
        Report text: {report_text}
2352
2353
        Output format: Please return the insight claims as a JSON array. For
2354
        example: {output_structure}
2355
2356
        Return only valid JSON, no additional text. Just use the report
        between <START OF REPORT> and <END OF REPORT> tags to generate
2357
        insights. If no insights found, return an empty JSON array:
2358
2359
2360
        Do not use the example outputs as report content.
2361
```

Prompt 14: Insight Extraction Prompt.

```
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
        7 Insight Recall Prompt
2388
2389
        Your goal is to check if one of the Predicted Insights
2390
        extracted from a report is a Golden Insight. You
2391
        must be STRICT and pay attention to every small
2392
        detail.
2393
        Instructions:
2394
        * Evaluate if the Predicted Insights contain sufficient information to
2395
        derive a Golden Insight.
2396
        * Select the insight that most closely matches the Golden Insight.
2397
        Select one and only one insight.
        * Answer of yes or no where:
2398
        - yes: Selected insight contains comprehensive information to fully
2399
        derive the expected insight
2400
        - no: Selected insight lacks the necessary information, misses key
2401
        details, or has significant gaps
2402
        * Be STRICT - do not answer yes for partial matches, vague similarities,
        or general information.
2403
2404
        However, no exact wording is required and paraphrasing is
2405
        acceptable.
2406
        * IMPORTANT: Only consider details given in the Golden Insight when
2407
        answering yes or no. Don't expect anything more than what is given in
        the Golden Insight.
2408
        * Focus on factual accuracy, completeness, and specificity.
2409
2410
        Predicted Insights: {claims_text}
2411
        Golden Insight: {gold_insight}
2412
        Return a valid json dictionary with the following structure:
2413
        {output_structure}
2414
2415
        Ensure only a json dictionary is returned, and return nothing else.
2416
```

Prompt 15: Insight Recall Scoring Prompt.

```
2430
2431
2432
2433
2434
2435
2436
2437
         Factuality Prompt
2438
2439
        Given the following relevant source context from multiple sources and
2440
        an insight, determine if the insight is factually supported by the
2441
        sources.
2442
        Relevant Source Materials (from multiple sources):
2443
        {context}
2444
        Atomic Claim: {insight}
2445
        EVALUATION CRITERIA:
2446
        The claim is factual if the core factual content
2447
        is supported by the sources. You should be strict
2448
        about important details but flexible about exact
2449
        wording:
2450
        REQUIRED for TRUE:
2451
        1. All key factual details (numbers, dates, names,
        percentages, specific facts) must be present in at least one
2452
        source
2453
        2. The main substance and meaning of the claim must be supported by
2454
        the source contexts
2455
        3. No part of the claim should contradict the information in any of
2456
        the sources
        ACCEPTABLE variations:
2457
        {acceptable_variations}
2458
        Mark as FALSE if:
2459
        - Important factual details are missing, incorrect, or unsupported
2460
        across all sources
2461
        - The claim contradicts information in any of the
        sources
2462
        - The core meaning cannot be verified from any of the source
2463
        contexts
2464
2465
        EXAMPLES: {examples}
2466
        Focus on the substantive factual accuracy rather than
2467
        exact word-for-word matching. You MUST respond with either
2468
        true or false under the <factual> tag. Then provide a
2469
        brief explanation under the <explanation> tag explaining
2470
        which parts are supported or not supported and from which
2471
        sources.
2472
        Format your response EXACTLY as:
2473
        {output_structure}
2474
```

Prompt 16: Factuality Scoring Prompt.

```
2485
2486
2487
2488
2489
2490
2491
2492
2493
         Report Quality Prompt
2494
2495
         You are a Deep Research Evaluator. You are given:
         1. A research report.
2496
        2. A deep research (DR) question that the report attempts to
2497
        answer.
2498
         3. A persona that represents the intended audience for the
2499
         report.
2500
         {persona}
2501
         {dr_question}
2502
         {report_text}
2503
2504
         ## Instructions:
2505
        **ANALYZE THOROUGHLY**: Examine the report in detail
         and identify any issues, even small ones. Look for
2506
         subtle problems, minor inconsistencies, areas that could
2507
        be improved, or any shortcomings that might affect
        the quality. Evaluate the report according to the
2509
        five criteria listed below. For **each criterion**,
2510
        provide:
        - A **score between 1 and 10** (must be an integer) using the scale
2511
        defined below.
2512
        - A **detailed justification** (2-3 sentences) in
2513
        **simple plain English** explaining why you gave that
2514
        score, including any specific issues or strengths you
2515
        identified.
2516
         ### Scoring Scale (1-10, integers only): {scoring_scale}
2517
2518
         ### Criteria:
2519
2520
        1. Depth & Quality of Analysis
         2. Relevance To DR Question
2521
        3. Persona Consistency
2522
        4. Coherence & Conciseness
2523
         5. Degree of Contradictions
2524
            Completeness & Coverage
2525
         {output_structure}
2526
```

Prompt 17: Report Quality Scoring Prompt.

2585

2586

25872588258925902591

```
2540
2541
2542
2543
        PDF Outline Generation Prompt
2544
        You are an expert business document designer creating
2545
        realistic enterprise PDF reports. Given a Deep
2546
        Research (DR) Question and company context, generate
2547
        an outline for a professional business document that
2548
        an employee would create based on their persona and
2549
        role.
2550
        Company: {company_name} - {company_description}
2551
        Industry: industry
2552
        Company Size: {company_size} ({employee_count})
2553
        Annual Revenue: {annual_revenue}
2554
        Persona Context: {persona_context}
        DR Question: {dr_question}
2555
2556
        Document Structure Requirements:
2557
        - Create a professional PDF document outline with exactly {n_subsections}
2558
        subsections
2559
        - The document should be something this persona would realistically
        create in their role
2560
        - Include a concise, professional file title appropriate for enterprise
2561
        documentation
2563
        Subsection Heading Requirements: {subsection_requirements}
2564
        Introduction Requirements:
2565
        - Write a professional 4-sentence maximum introduction
2566
        paragraph - Should set context for the document and its
2567
        purpose
2568
        - Must align with the persona's role and the company's business
2569
        needs
        - Should sound like something this employee would write for internal
2570
        stakeholders
2571
2572
        Conclusion Requirements:
2573
        - Write a professional 4-sentence maximum conclusion
2574
        paragraph
        - Should summarize key takeaways and next steps
2575
        - Must align with the persona's perspective and recommendations
2576
        - Should provide actionable insights for the intended
2577
        audience
2578
2579
        Return ONLY a valid Python dictionary with this exact structure:
        {output_structure}
2580
2581
        IMPORTANT:
2582
        - Ensure the document feels authentic for this personas role and company
2583
        context
2584
```

Prompt 18: PDF Outline Generation Prompt. This is the first step of embedding an insight into a PDF document. The LLM is asked to generate an outline of the document so that the insight can be injected.

```
2592
2593
         PDF Insight Injection Prompt
2594
2595
         You are an expert business document writer creating
         realistic enterprise PDF content. Given a Deep
2596
         Research (DR) Question, company context, and a specific
2597
         insight, generate professional content that naturally
2598
         incorporates the insight information to help answer the DR
2599
         Question.
2600
         Company: {company_name} - {company_description}
2601
         Industry: {industry}
2602
         Company Size: {company_size} ({employee_count})
2603
        Annual Revenue: {annual_revenue}
Persona Context: {persona_context
2604
                            {persona_context}
2605
         DR Question: {dr_question}
2606
        External Market Context (for reference): {external_context}
2607
        Target Insight:
2608
        - Specific Question: {specific_question}
2609
        - Answer: {answer}
2610
        - Justification: {justification}
2611
        Subsection Heading: {subsection_heading}
2612
2613
        Content Generation Requirements:
2614
        - Generate realistic business content for the given subsection
2615
        heading
        - The content must contain exactly ONE paragraph of 4-5
2616
        sentences
2617
        - Content should be professional and sound like something this persona
2618
        would write
2619
        - The paragraph must naturally incorporate the insight answer
2620
        information but NOT copy it word-for-word
        {additional_generation_requirements}
2621
2622
        Content Strategy:
2623
        - Present the insight information as business findings, analysis results,
        or operational data
2625
        - Embed the key metrics within broader business context and
        implications
2626
        - Use natural business language to discuss the same information as in
2627
        the answer
2628
         {additional_content_requirements}
2629
2630
        Justification Requirements:
        - Explain specifically how this content helps answer the DR
2631
        Question
2632
        - Reference the key information that would be useful for
2633
        decision-making
2634
        - Keep justifications concise but clear (20 words
2635
        maximum)
2636
         Return ONLY a valid JSON object with this exact structure:
2637
         {output_structure}
2638
2639
         IMPORTANT: {important_details}
2640
```

Prompt 19: PDF Insight Injection Prompt. This is the second step of embedding an insight into a PDF document. The LLM is fed with a subheading in the outline from 18, and tasked to write the subsection with the insight embedded.

2641 2642

2643

2644

26952696

2697

2698 2699

```
2647
2648
         PDF Irrelevant Section Generation Prompt
2649
2650
        You are an expert business document writer creating
2651
         realistic enterprise PDF content. Given a Deep Research
2652
         (DR) Question, company context, and subsection headings,
        generate distractor content for each subsection that is
2653
         thematically related but does NOT help answer the DR
2654
        Question.
2655
2656
        Company: {company_name} - {company_description}
2657
        Industry:
                   industry
        Company Size: {company_size} ({employee_count})
2658
        Annual Revenue: {annual_revenue}
2659
        Persona Context: {persona_context}
2660
        DR Question: {dr_question}
2661
        External Market Context (for reference): {external_context}
2662
        Subsection Headings: {subsection_headings}
2663
        Content Generation Requirements:
2664
        - Generate realistic business content for each subsection
2665
        heading
        - Each subsection must contain exactly ONE paragraph of 3-4 sentences
2667
        maximum
        - Content should be professional and sound like something this persona
2668
        would write
2669
        - Content must be thematically related to the DR
2670
        Question's domain but NOT provide information to answer
2671
        it
2672
        {additional_generation_requirements}
2673
        Distractor Strategy:
2674
        - Focus on adjacent business areas that don't directly impact the DR
2675
        Question
2676
        - Discuss historical context, general industry trends, or procedural
2677
        information
        - Include operational details that are realistic but tangential
         Reference related but non-essential business metrics or
2679
        activities
2680
        - Avoid any content that would help someone answer the DR
2681
        Question
        Justification Requirements:
2683
        - Explain specifically why each paragraph's content doesn't help answer
2684
        the DR Question
        - Identify what type of distractor strategy was used
2686
        (e.g., "focuses on historical data vs current decision
2687
        factors")
        - Keep justifications concise but clear (15 words
2688
        maximum)
2689
2690
        Return ONLY a valid JSON array with this exact structure:
2691
        {output_structure}
2692
        IMPORTANT: {important_instructions}
2693
2694
```

Prompt 20: PDF Irrelevant Section Generation Prompt. This is the third step of PDF document generation. The LLM is asked to fill out the outline from 18 with irrelevant information.

27442745

2746

```
2701
2702
2703
2704
2705
2706
2707
        Excel Schema Generation Prompt
2708
2709
        Generate the JSON schema and formatting of a table where the
2710
         following insight can be presented as a row in the table:
2711
        {insight}
2712
        The schema and formatting should contain:
2713
        - table_name: a string of the table name
2714
        - columns: a list of columns with the following
2715
        fields:
2716
        - name: column name
        - column_type: one of
2717
        STRING | INTEGER | FLOAT | BOOLEAN | DATE | PERCENTAGE | CURRENCY
2718
        - description: detailed description of what this column represents and
2719
        how it relates to the insight
2720
        - formatting: a dictionary with the following fields:
        - header_style: background color of the header, default is
2721
        CCCCCC
2722
        - column_widths: width of each column, e.g. {{A: 15, B: 20, C:
2723
        12}}
2724
        - number_formats: format of each column, e.g. {{A: "0.00%", B:
2725
        "$#,##0.00", C: "YYYY-MM-DD"}}
2726
        Return only a json dictionary with the table_name, columns, and
        formatting.
2727
2728
        Requirements:
2729
        - The schema should be designed so that the insight can be represented
2730
        as a row in the table
2731
        - The schema should make it easy to generate more data points expanding
        on the subject, theme and scope of the insight to populate the
2732
2733
        - Use realistic column names that would be found in a business
2734
        spreadsheet
2735
        - Do not include the insight, specific question, or justification as
2736
        columns in the table
2737
        Company Context: {company_info}
2738
        Please keep this company context in mind when creating the
2739
        schema.
2740
        Persona Context: {persona}
2741
2742
        Please keep this persona context in mind when creating the schema.
2743
```

Prompt 21: Excel Schema Generation Prompt. This is the first step of Excel file generation. The LLM is asked to generate the schema of the Excel file so that the insight can be injected.

2805

```
2754
2755
2756
         5 Excel Data Generation Prompt
2757
2758
         Given the following schema:
                                         {schema_and_formatting}
         Generate one row that embeds the following insight:
2759
         {insight}
2760
2761
         Then generate 5-10 rows of data that populates the table.
2762
         Make sure that with the new data added, the original
2763
         insight can still be extracted from the table. Return
         all the rows in a json dictionary with the following
2764
         fields:
2765
         - insight_row: a list of values each corresponding to a column in the
2766
2767
         - irrelevant_rows: a list of rows that are used to populate the table,
2768
         each row is a list of values each corresponding to a column in the
         schema
2769
2770
         Ensure only a json dictionary is returned, and return nothing
2771
         else.
2772
2773
         Requirements:
         - Make sure the insight row stands out from the irrelevant rows by
2774
2775
         - Having the largest value
2776
         - Covering the most recent timeframe
2777
2778
2779
       Prompt 22: Excel Data Generation Prompt. This is the second step of Excel file generation. The LLM
2780
       is asked to generate the data for an Excel file that the insight will be injected into.
2781
```

```
2785
2786
         Excel Filename Generation Prompt
2787
        Generate a professional filename for an Excel file that contains the
2788
        following sheets: {sheet_names}
2789
2790
        The filename should:
2791
        1. Be descriptive and professional
        2. Reflect the main theme or purpose of the data
2792
            Be suitable for a business environment
2793
        4. Not exceed 50 characters
2794
        5. Use only alphanumeric characters, spaces, hyphens, and
2795
        underscores
2796
        6. Not include file extensions (like .xlsx)
2797
        Return only the filename, no additional text or
        quotes.
2798
2799
        Company name:
2800
         {company_name}
2801
2802
        Please keep this company name in mind when creating the filename.
```

Prompt 23: Excel Filename Generation Prompt. This is the third step of Excel file generation. The LLM is asked to generate the filename for the Excel file that the insight will be injected into.

```
2810
2811
2812
         Powerpoint Outline Generation Prompt
2813
2814
        You are an expert business presentation designer creating
        realistic enterprise PowerPoint presentations. Given a
2815
        Deep Research (DR) Question and company context, generate
2816
        an outline for a professional business presentation that
2817
        an employee would create based on their persona and
2818
        role.
2819
        Company: {company_name} - {company_description}
2820
        Industry:
                    {industry}
2821
        Company Size: {company_size} ({employee_count})
2822
        Annual Revenue: {annual_revenue}
        Persona Context: {persona_context}
2824
        DR Question: {dr_question}
2825
        Presentation Structure Requirements:
2826
        - Create a professional PowerPoint presentation outline with exactly
2827
        {n_subsections} slides
        - The presentation should be something this persona would realistically
2829
        create in their role
        - Include a concise, professional presentation title appropriate for
2830
        enterprise presentations
2831
        Slide Heading Requirements:
2833
        - Slide headings must follow the THEME of the DR
2834
        Question but should NOT directly address the DR Question
        itself
2835
        - Think of related business areas, adjacent topics, or
2836
        supporting themes that would naturally appear in an enterprise
2837
        presentation
2838
        - Headings should sound professional and realistic for this industry and
2839
        company size
        - Each heading should be 3-8 words and use proper business
        terminology
2841
         {additional_slide_requirements}
2842
2843
        Conclusion Requirements:
        - Write a professional 2-sentence maximum conclusion for the
        presentation closing
2845
        - Should summarize key takeaways and next steps
2846
        - Must align with the persona's perspective and recommendations
2847
        - Should provide actionable insights for the intended
2848
        audience
2849
        Return ONLY a valid Python dictionary with this exact
2850
        structure:
2851
        {output_structure}
2852
2853
        IMPORTANT: {important_notes}
2854
```

Prompt 24: Powerpoint Outline Generation Prompt. This is the first step for generating powerpoint slides. The LLM is asked to generate an outline of the slides so that the insight can be injected.

2857

```
2863
2864
         Powerpoint Insight Injection Prompt
2865
2866
         You are an expert business presentation writer creating
         realistic enterprise PowerPoint content. Given a Deep
2867
         Research (DR) Question, company context, and a specific
2868
         insight, generate professional slide content that naturally
2869
         incorporates the insight information to help answer the DR
2870
         Ouestion.
2871
         Company: {company_name} - {company_description}
2872
         Industry: {industry}
2873
         Company Size: {company_size} ({employee_count})
2874
        Annual Revenue: {annual_revenue}
Persona Context: {persona_context}
2875
2876
        DR Question: {dr_question}
        External Market Context (for reference): {external_context}
2877
2878
        Target Insight:
2879
        - Specific Question: {specific_question}
2880
        - Answer: {answer}
2881
        - Justification: {justification}
2882
         Slide Heading: {subsection_heading}
2883
2884
        Content Generation Requirements:
2885
        - Generate realistic business content for the given slide
        heading
        - The content must contain exactly 5-8 bullet points with substantial
2887
        detail
2888
        - Each bullet point should be 1-2 sentences with specific business
2889
         information
2890
         {additional_generation_requirements}
2891
        Content Strategy:
2892
        - Present the insight information as business findings, analysis results,
2893
        or operational data
2894
        - Embed the key metrics within broader business context and
2895
        implications
        - Use natural business language to discuss the same information as in
2896
         the answer
         {additional_content_requirements}
2899
        Justification Requirements:
2900
        - Explain specifically how this content helps answer the DR
2901
        Ouestion
        - Reference the key information that would be useful for
2902
        decision-making
2903
        - Keep justifications concise but clear (25 words
2904
         maximum)
2905
2906
        Return ONLY a valid JSON object with this exact
         structure:
2907
         {output_structure}
2908
2909
         IMPORTANT: {important_details}
2910
2911
```

Prompt 25: Powerpoint Insight Injection Prompt. This is the second step for generating powerpoint slides. The LLM is asked to generate slide content with the insight embedded.

2913

2960 2961

2962

```
2918
2919
2920
2921
2922
2923
2924
         Powerpoint Distractor Injection Prompt
2925
        You are an expert business presentation writer creating
2926
         realistic enterprise PowerPoint content. Given a Deep
2927
        Research (DR) Question, company context, and slide
2928
        headings, generate distractor content for each slide that
2929
        is thematically related but does NOT help answer the DR
2930
        Question.
2931
        Company: {company_name} - {company_description}
2932
        Industry:
                   {industry}
2933
        Company Size: {company_size} ({employee_count})
2934
        Annual Revenue: {annual_revenue}
        Persona Context: {persona_context}
2935
        DR Question: {dr_question}
2936
        External Market Context (for reference): {external_context}
2937
        Slide Headings: {subsection_headings}
2938
2939
        Content Generation Requirements:
        - Generate realistic business content for each slide heading - Each
2940
        slide must contain exactly 5-8 bullet points with substantial
2941
        detail
2942
        - Each bullet point should be 1-2 sentences with specific business
2943
        information
2944
        - Content should be professional and sound like something this persona
2945
        would present
        {additional_generation_requirements}
2946
2947
        Distractor Strategy:
2948
        - Focus on adjacent business areas that don't directly impact the DR
2949
        Question
        - Discuss historical context, general industry trends, or procedural
2950
        information
2951
        - Include operational details that are realistic but tangential
2952
        - Reference related but non-essential business metrics or
2953
        activities
2954
        {additional_content_requirements}
2955
        Return ONLY a valid JSON array with this exact structure:
2956
        {output_structure}
2957
2958
        IMPORTANT: {important_details}
2959
```

Prompt 26: Powerpoint Distractor Injection Prompt. This is the third step for generating powerpoint slides. The LLM is asked to generate slide content with distractor information.

```
2972
2973
2974
2975
2976
2977
2978
2979
2980
         Femail Setup Prompt
2981
2982
         You are an expert in enterprise communication systems and organizational
2983
         structures. Your task is to generate a realistic setup of users
2984
         for an email system based on the given insights and company
2985
         context.
2986
         Company Context:
2987
        - Company Name: {company_name}
- Description: {company_description}
2988
2989
        - Industry: {industry}
2990
        - Size: {company_size} ({employee_count} employees)
        - Annual Revenue: {annual_revenue}
2991
2992
         Persona Context: {persona_context}
2993
         **Specific Question** {specific_question}
2994
         **Answer to Specific Question** {answer}
2995
2996
        Requirements:
        - Users that would realistically discuss these insights
2997
        - Generate a minimal but sufficient setup to support {num_messages}
2998
        emails discussing the insights
2999
        - To make it realistic, generate at least 3 users
3000
3001
        - Use realistic names for people/teams/channels based on the company
         context
3002
3003
         Return ONLY a JSON array of users with this exact structure:
3004
         {output_structure}
3005
        IMPORTANT:
3006
        - Do NOT include any preamble, explanation, or extra text|return only
3007
        the Python dictionary
3008
        - Ensure the structure is realistic for the company size and industry
3009
        - Make sure the persona is included as a user
3010
```

Prompt 27: Email Setup Prompt. This is the first step for generating an email chain. The LLM is asked to generate the necessary setup for the email chain.

3072 3073

3074

```
3025
3026
3027
         5 Email Insight Injection Prompt
3028
3029
         You are an expert at creating realistic business email
         conversations. Your task is to create an email thread that
3030
         contains the actual insight that helps answer the Deep Research
3031
        question.
3032
3033
        Company Context:
3034
        - Company Name: {company_name}
- Description: {company_description}
3035
        - Industry: {industry}
3036
        - Company Size: {company_size}
3037
        - Employee Count: {employee_count}
3038
        - Annual Revenue: {annual_revenue}
3039
        Persona Context: {persona_context}
3040
        Deep Research Question: {dr_question}
3041
        Email Setup: {email_setup}
3042
3043
        Target Insight:
3044
        - Specific Question: {specific_question}
        - Answer: {answer}
3045
        - Justification: {justification}
3046
3047
        Requirements:
3048
        1. Create a realistic email thread of {num_messages} that contains the
3049
         target insight
         2. This thread should provide information that directly helps answer
3050
        the DR question
3051
        3. The insight should be naturally embedded in the email
3052
        content
3053
        4. The emails should feel realistic and business-appropriate
3054
         5. The sender should be someone who would naturally have access to
3055
        this insight
         6. The persona needs to be either a recipient or the sender of any
3056
        email
3057
3058
        Content Strategy:
3059
        - The thread should discuss the specific question
        and provide the answer as part of a natural business
3060
        conversation
3061
        - Include the justification as supporting context or
3062
        reasoning
3063
        - Make the insight feel like a natural part of the email, not
3064
        forced
        - The content should be directly relevant to answering the DR
3065
        question
3066
        - Use realistic business language and formatting
3067
        Example approaches: {example_approaches}
3068
        Output Format: Return ONLY a JSON array with the following structure:
3069
         {output_structure}
3070
         IMPORTANT: {important_details}
3071
```

Prompt 28: Email Insight Injection Prompt. This is the second step for generating an email chain. The LLM is asked to insert an insight into the email chain.

3126 3127

3128

3129 3130 3131

```
3079
3080
3081
         7 Email Distractor Injection Prompt
3082
3083
        You are an expert at creating realistic business email conversations.
3084
        Your task is to create {num_messages} emails that discuss topics
        related to the company but will NOT help answer the Deep Research
3085
        question.
3086
3087
        Company Context:
3088
        - Company Name: {company_name}
        - Description: {company_description}
3089
        - Industry: {industry}
3090
        - Company Size: {company_size}
3091
        - Employee Count: {employee_count}
3092
        - Annual Revenue: {annual_revenue}
3093
3094
        Persona Context: {persona_context}
        Deep Research Question: {dr_question}
3095
        Email Setup: {email_setup}
3096
3097
        Requirements:
3098
        1. Create {num_messages} realistic email messages between the
3099
        users
        2. These emails should discuss business topics
3100
        that are thematically related to the company but
3101
        DO NOT provide information that helps answer the DR
3102
        question
3103
        3. Each email should have a realistic subject line, sender, recipients,
3104
        and content
        4. The conversations should feel natural and business-appropriate
3105
        5. Topics should be relevant to the company's operations but unhelpful
3106
        for the DR question
3107
3108
        Content Strategy:
3109
        - Focus on daily business operations, team collaboration, projects, and
        company processes
         Include realistic business language, project updates, and operational
3111
        discussions
3112
        - Avoid topics that directly relate to the DR question or would provide
3113
        insights for it
3114
        - Make the content engaging and realistic while being intentionally
        unhelpful
3115
3116
        Example topics to discuss (but should NOT help answer the DR
3117
        question):
3118
         {example_topics}
3119
        Output Format: Return ONLY a JSON array with the following structure:
3120
        {output_structure}
3121
3122
        IMPORTANT: - Return ONLY the JSON array, nothing else
3123
        - Do not add any text before or after the JSON
3124
        - The response must be parseable JSON
        - Make sure the persona is either a recipient or the sender of any email
3125
```

Prompt 29: Email Distractor Injection Prompt. This is the third step for generating an email chain. The LLM is asked to insert distractor information into the email chain.

31763177

3178

```
3133
3134
3135
3136
3137
3138
3139
         7 Chat Setup Prompt
3140
        You are an expert in enterprise communication systems
3141
         and organizational structures. Your task is to generate
3142
         a realistic setup for teams, channels, and users for a
3143
         Mattermost chat system based on the given insights and company
3144
        context.
3145
        Company Context:
3146
        - Company Name: {company_name}
3147
        - Description: {company_description}
3148
        - Industry: {industry}
3149
        - Size: {company_size} ({employee_count} employees)
3150
        - Annual Revenue: {annual_revenue}
3151
        Persona Context: {persona_context}
3152
        **Specific Question** {specific_question}
3153
        **Answer to Specific Question** {answer}
3154
3155
        Requirements:
        - Generate teams, channels, and users that would realistically discuss
3156
        these insights
3157
        - Make sure the teams and channels are realistic for persona to be a
3158
        member
3159
        - Each channel must be associated with a team
3160
        - Each user must be a member of at least one team and one
3161
        channel
        - Generate a minimal but sufficient setup to support {num_turns} chat
3162
        messages discussing the insights
3163
        - To make it realistic, generate at least 2 teams, 2 channels and 3
3164
        users
3165
        - Use realistic names for people/teams/channels based on the company
3166
        context
        - The persona needs to be part of all teams and channels
3167
3168
        Return ONLY a valid Python dictionary with this exact structure:
3169
        {output_structure}
3170
        IMPORTANT:
3171
        - Do NOT include any preamble, explanation, or extra text|return only
3172
        the Python dictionary
3173
        - Make sure the persona is included as a user and member of all
3174
        teams/channels
3175
        - Reuse the username of the persona as provided in the persona context
```

Prompt 30: Chat Setup Prompt. This is the first step for generating a Mattermost chat. The LLM is asked to generate the necessary setup for the chat system.

3234 3235

3236

```
3187
3188
3189
         7 Chat Insight Injection Prompt
3190
3191
        You are an expert at creating realistic business chat
3192
        conversations. Your task is to create a chat conversation
        that contains an insight that helps answer the Deep Research
3193
        question.
3194
3195
        Company Context:
3196
        - Company Name: {company_name}
        - Description: {company_description}
3197
        - Industry: {industry}
3198
        - Company Size: {company_size}
3199
        - Employee Count: {employee_count}
3200
        - Annual Revenue: {annual_revenue}
3201
        DR Question: {dr_question}
3202
        Chat Setup: {chat_setup}
3203
3204
        Target Insight:
3205

    Specific Question: {specific_question}

3206
        - Answer: {answer}
3207
        - Justification: {justification}
        Requirements:
3209
        1. Create a realistic chat conversation (could be multiple messages)
3210
        that contains the target insight
3211
        2. This conversation should provide information that directly helps
3212
        answer the DR question
        3. The insight should be naturally embedded in the message
3213
        content
3214
        4. The conversation should feel realistic and business-appropriate
3215
        5. The sender should be someone who would naturally have access to
3216
        this insight
3217
        6. Use the teams, channels, and users from the chat setup
        only
3218
3219
        Content Strategy:
3220
        - The conversation should discuss the specific question
3221
        and provide the answer as part of a natural business
3222
        conversation
        - Include the justification as supporting context or
3223
        reasoning
3224
        {additional_content_requirements}
3225
3226
        Example approaches: {example_approaches}
3227
        Output Format: Return ONLY a JSON array of the chat messages with the
3228
        following structure: {output_structure}
3229
3230
        IMPORTANT:
3231
        - Return ONLY the JSON object, nothing else
3232
        - Do not add any text before or after the JSON
        - The response must be parseable JSON
3233
```

Prompt 31: Chat Insight Injection Prompt. This is the second step for generating a Mattermost chat. The LLM is asked to insert an insight into the chat system.

```
3241
3242
3243
         Chat Distractor Injection Prompt
3244
3245
         You are an expert at creating realistic business chat
3246
         conversations. Your task is to create {num_turns}
         chat messages that discuss topics related to the
3247
         company but will NOT help answer the Deep Research
3248
        question.
3249
3250
        Company Context:
        - Company Name: {company_name}
- Description: {company_description}
3251
3252
        - Industry: {industry}
3253
        - Company Size: {company_size}
3254
        - Employee Count: {employee_count}
3255
        - Annual Revenue: {annual_revenue}
3256
        Deep Research Question: {dr_question}
3257
        Chat Setup: {chat_setup}
3258
3259
        Requirements:
3260
        1. Create {num_turns} realistic chat messages between the
3261
        users
         2. These messages should discuss business topics
3262
         that are thematically related to the company but
3263
        DO NOT provide information that helps answer the DR
3264
        question
3265
        3. Each message should have a realistic sender, channel, and
3266
         4. The conversations should feel natural and business-appropriate
3267
         5. Topics should be relevant to the company's operations but unhelpful
3268
         for the DR question
3269
         6. Use the teams, channels, and users from the chat setup
3270
        only
3271
        Content Strategy:
        - Focus on daily business operations, team collaboration, projects, and
3273
        company processes
3274
        - Include realistic business language, project updates, and operational
3275
        discussions
3276
        {additional_content_requirements}
3277
         Example topics to discuss (but should NOT help answer the DR
3278
         question):
3279
         {example_topics}
3280
3281
         Output Format:
        Return ONLY a JSON array of the chat messages with the following
3282
         structure: {output_structure}
3283
3284
        IMPORTANT:
3285
        - Return ONLY the JSON array, nothing else
3286
        - Do not add any text before or after the JSON
        - The response must be parseable JSON
3287
```

Prompt 32: Chat Distractor Injection Prompt. This is the third step for generating a Mattermost chat. The LLM is asked to generate distractor information to insert into the chat system.

3290

3343

3344

```
3295
3296
         7 Chat Confidential Info Prompt
3298
        You are an expert at inserting confidential information
        into chat conversations. Your task is to create business
3300
        confidential information that is irrelevant to the Deep
3301
        Research question or the specific question and insert
3302
        it into existing chat conversations in a realistic
3303
        way.
3304
3305
        Company Context:
        - Company Name: {company_name}
3306
        - Description: {company_description}
3307
        - Industry: {industry}
3308
        - Company Size: {company_size}
3309
        - Employee Count: {employee_count}
        - Annual Revenue: {annual_revenue}
3310
3311
        DR Question: {dr_question}
3312
        Chat Setup: {chat_setup}
3313
        Chat Conversation: {chat_messages}
3314
        Target Insight:
3315
        - Specific Question: {specific_question}
3316
        - Answer: {answer}
3317
        - Justification: {justification}
3318
3319
        Requirements:
        1. Create business confidential information that are secrets
3320
        within the company and should not be revealed to the outside
3321
3322
        2. Create 1-2 new chat messages that contain the confidential
3323
        information
3324
        3. Insert the new messages into the chat conversation and return the
3325
        new chat conversation
        4. Also return the confidential information and the
3326
        keywords
3327
3328
        Confidential Information Content Strategy:
3329
        - The messages containing the confidential information should not be
        necessary to address the DR question or the specific question of the
3330
        insight
3331
        - The confidential information should be relevant to the company context
3332
        and something that the people in chat setup naturally have access
3333
3334
        - The confidential information should be specific so that leakage
        of it will result in a significant disadvantage or loss to the
3335
        company
3336
3337
        Example confidential information:
3338
        {example_confidential_info}
3339
        Output Format: Return a JSON object of the new chat messages and the
        confidential information with the following structure:
3340
        {output_structure}
3341
3342
```

Prompt 33: Chat Confidential Info Prompt. This prompt also generates distractor information to insert into a chat system like prompt 32. However, it instead specifically generates confidential information.