



# FINCH: Benchmarking Finance & Accounting across Spreadsheet-Centric Enterprise Workflows

Anonymous ACL submission

## Abstract

We introduce a finance & accounting benchmark (FINCH) for evaluating AI agents on real-world, enterprise-grade professional workflows that interleave data entry, structuring, formatting, web search, cross-file retrieval, calculation, modeling, validation, translation, visualization, and reporting. FINCH is sourced from authentic enterprise workspaces at Enron (15,000 Excel files and 500,000 emails from 150 employees) and other financial institutions, preserving in-the-wild messiness across multi-modal artifacts (text, tables, formulas, charts, and images) and spanning diverse domains such as budgeting, trading, asset management, and operational management.

We propose a workflow construction process that combines LLM-assisted discovery with expert annotation: (1) LLM-assisted, expert-verified derivation of workflows from real-world email threads and spreadsheet version histories, and (2) meticulous expert annotation requiring over 700 hours of domain-expert effort. This yields 172 composite workflows with 384 tasks, involving 1,710 spreadsheets with 27 million cells, along with PDFs and other artifacts, capturing the intrinsically messy, long-horizon, knowledge-intensive, and collaborative nature of real-world enterprise work.

We conduct both human and automated evaluations of frontier AI systems, including GPT 5.1, Claude Sonnet/Opus 4.5, Gemini 3 Pro, Grok 4, and Qwen 3 Max. Under human evaluation, GPT 5.1 Pro spends an average of 16.8 minutes per workflow yet passes only 38.4% of workflows. Comprehensive case studies further surface the challenges that real-world enterprise workflows pose for AI agents.

## 1 Introduction

Frontier AI systems are increasingly transforming professional workspaces. AI-assisted tools like ChatGPT (OpenAI, 2025), Claude (Anthropic, 2025), Gemini (Google, 2024), and Copilot (Microsoft, 2024) are now embedded in daily enterprise workflows—helping professionals draft documents, explore data, manipulate spreadsheets, and

generate reports. These tools are particularly impactful in finance and accounting (F&A), a high-stakes, knowledge- and labor-intensive domain critical to every organization.

However, real-world F&A work is inherently **messy** with substantial contextual complexity: artifacts are interconnected across heterogeneous spreadsheets, PDFs, and more, evolving through multiple versions with collaborative edits (Klimt and Yang, 2004); spreadsheets contain large, complex structures (Dong et al., 2024) with cross-sheet references, intricate layouts, inconsistent formatting, cryptic terms, erroneous formulas, and multimodal artifacts such as charts, images, and code. It is also **long-horizon** (Patwardhan et al., 2025): workflows demand multi-step reasoning spanning data entry, editing, retrieval, calculation, modeling, validation, reporting, and more.

This raises a key question: *Can today’s frontier AI agents actually handle the messy, long-horizon, and knowledge-intensive workflows that professionals face daily?*

To answer this, we introduce FINCH, a F&A benchmark sourced from authentic enterprise environments. FINCH captures the intrinsic complexity of professional work through:

- **In-the-wild enterprise sourcing:** FINCH is built around authentic enterprise spreadsheets, emails, and PDFs from real-world enterprise workspaces—primarily Enron (EnronData.org) (about 15,000 spreadsheet files and 500,000 emails from 150 employees) and EUSES (Fisher and Rothermel, 2005) (about 450 financial spreadsheet files from various sources), along with securities and asset management firms, global organizations such as World Bank (Bank, 2024), and Canadian and British governments (Department of Finance Canada, 2025; HM Treasury, 2023). Documents are large, cross-referenced, and messy—containing rich multimodal artifacts such as tables, formulas, charts, pivots, and images.
- **Rigorous construction process:** We pro-

090	pose a novel workflow construction pipeline	FINCH poses for AI agents. Comprehensive case	142
091	grounded in the real collaborative context	analyses further surface concrete challenges that	143
092	of emails and versioned artifacts. We in-	real-world enterprise workflows pose for AI agents.	144
093	duce workflows from enterprise email threads		
094	and attachments, where collaborators natu-	<b>2 FINCH: A Real-world Finance &amp;</b>	145
095	rally describe, discuss, and track workflows	<b>Accounting Workflow Benchmark</b>	146
096	as part of their daily work. We further use		
097	an LLM-assisted, expert-verified method to	<b>2.1 Dataset Construction</b>	147
098	derive workflows by analyzing changes across	We propose a novel workflow construction pipeline	148
099	versioned spreadsheets, surfacing the underly-	grounded in the real collaborative context of	149
100	ing goals that drive professionals’ work. An-	emails, versioned spreadsheets, and final deliver-	150
101	notators must reason over large multi-sheet	able spreadsheets and reports.	151
102	workbooks and subtle version deltas to infer	All workflows derived from these sources are	152
103	underlying workflows, making the annotation	consolidated into a unified schema with consis-	153
104	process substantially more difficult than curat-	tent fields (NL instruction, input files, reference	154
105	ing QA pairs over isolated tables.	outputs), and each workflow is tagged with task	155
106	We compile 172 meticulously annotated,	types (e.g., data entry/import, structuring, valida-	156
107	enterprise-grade workflows built on 1,710 spread-	tion) and business types (e.g., planning and bud-	157
108	sheets, along with PDFs and other artifacts,	geting, pricing and valuation, operations, asset	158
109	collectively capturing the compositional, messy,	management). Note that the reference outputs	159
110	knowledge-intensive, and collaborative nature of	may include both file-based reference answers (for	160
111	real work. Each workflow spans one or more in-	most generation/editing cases) and textual refer-	161
112	terdependent tasks—data entry, editing, retrieval,	ence answers (for a small number of QA and sum-	162
113	calculation, modeling, validation, translation, vi-	mary/visualization cases).	163
114	sualization, and reporting—mirroring how profes-		
115	sionals actually work on artifact manipulation and	<b>2.1.1 Workflow Derivation from Enterprise</b>	164
116	creation. These workflows cover a broad set of	<b>Email Threads</b>	165
117	enterprise domains such as trading and risk man-	We first mine real-world enterprise email threads to	166
118	agement, planning and budgeting, pricing and va-	surface workflows. Starting from the Enron Email	167
119	uation, and asset management.	Corpus, we prompt GPT-5 to identify collaborative	168
120	Evaluating such workflows poses nontrivial chal-	messages that (i) explicitly state a business goal	169
121	lenges, because FINCH tasks usually involve com-	(e.g., “update the RAC rankings” or “revise the	170
122	plex and large spreadsheets and require assessing	2002 allocations”) and (ii) reference one or more	171
123	both structural correctness and semantic fidelity be-	attached spreadsheets. For each selected thread,	172
124	yond exact matching. To address this, we provide	the model summarizes the communicative intent	173
125	both expert human evaluation and a scalable LLM-	and articulates a workflow description.	174
126	as-judge pipeline. Human evaluation serves as the	As illustrated in Figure 2, in the <i>strongly</i>	175
127	gold standard, judging whether a workflow has	<i>grounded</i> case, both the input and the reference	176
128	been satisfactorily completed end-to-end. In par-	artifacts for a workflow are already present as	177
129	allel, we introduce an efficient multimodal LLM-	attachments in the email thread (e.g., an initial	178
130	as-judge evaluation method that compares inputs,	ranking file and its corrected version). While	179
131	model outputs, and reference artifacts using struc-	strongly grounded cases are highly motivating for	180
132	tured diffs, compact snapshots, and multimodal	this benchmark, they are relatively rare. In the <i>par-</i>	181
133	renderings, enabling reliable assessment at scale.	<i>cially grounded</i> case, the email specifies a clear	182
134	We evaluate a spectrum of frontier AI systems—	goal, but only the input artifacts or the reference	183
135	including Claude Sonnet/Opus 4.5, GPT 5.1, Gem-	outputs are attached.	184
136	ini 3, Grok 4, and Qwen 3—using both expert eval-	Across both cases, human experts normalize con-	185
137	uation and a novel automated evaluation pipeline	versational email text and LLM-drafted descrip-	186
138	that closely aligns with expert judgments. Our	tions into workflow instructions and abstract away	187
139	experiments reveal that even frontier agents pass	idiosyncratic details while preserving the business	188
140	fewer than 50% of workflows under human eval-	intent. Annotators need to carefully verify and	189
141	uation, highlighting substantial challenges that	revise attachments to align with the instructions.	190

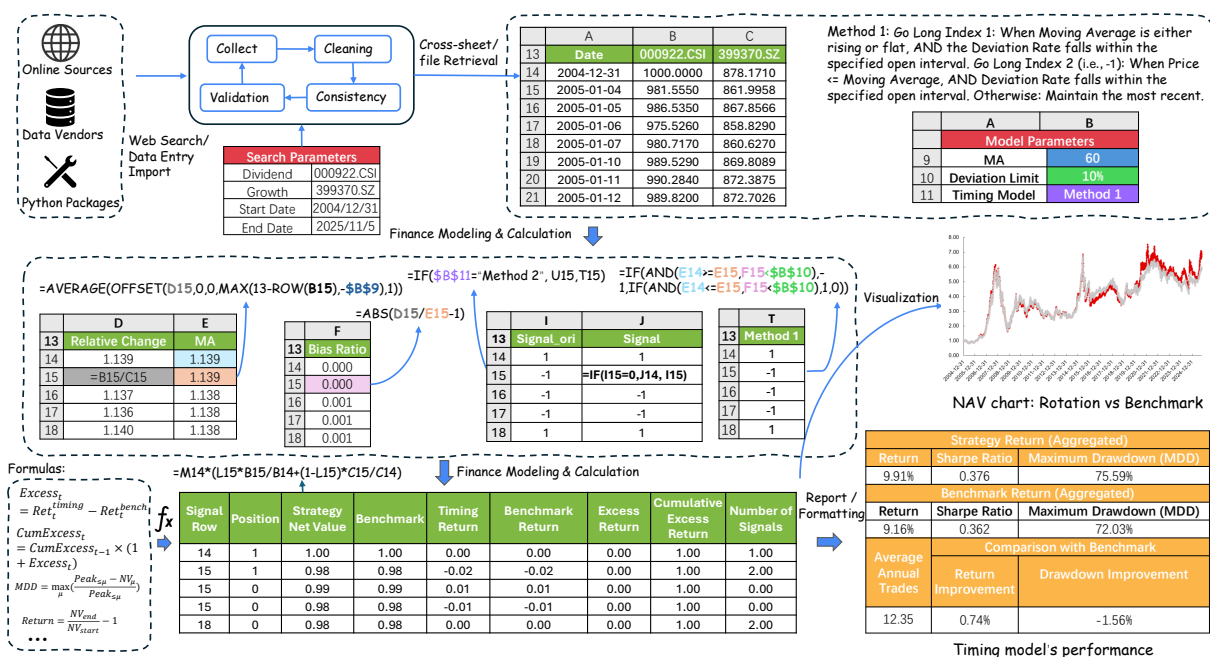


Figure 1: Illustration of an end-to-end predictive modeling workflow involving web search, data import, cross-sheet retrieval, calculation, and visualization. More illustrative examples in FINCH are presented in Appendix D.

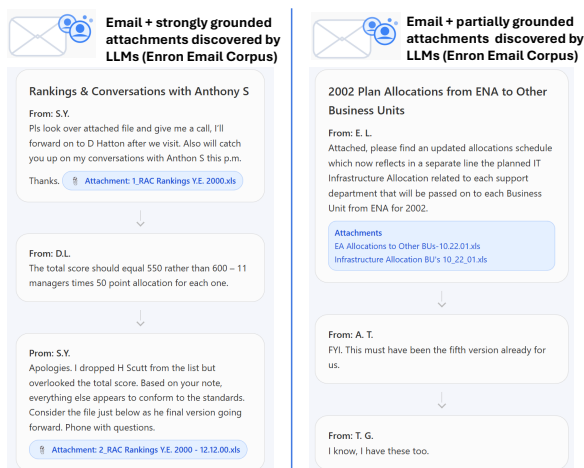


Figure 2: Illustration of real-world enterprise emails with attached artifacts.

## 2.1.2 Workflow Derivation from Versioned Spreadsheets

Beyond explicit messages in email threads, we propose to discover workflows that are implicitly captured in spreadsheet version histories, as illustrated in Figure 2(c). We collect families of versioned workbooks from the Enron and EUSES repositories and apply an LLM-based differencing procedure that recognizes consecutive versions and infers the underlying workflow.

For each aligned pair of versions, we prompt GPT-5 to propose (i) one or more workflow types (e.g., “date-stamped versioning, assumption updates, and error correction”, “data entry, structur-

ing, and visualization”) and (ii) a detailed NL description of all changes. Human experts then validate and refine these LLM-induced workflow candidates. They first determine whether the proposed diffs constitute a coherent and meaningful workflow rather than incidental churn. For accepted cases, they (i) rewrite the draft description into a precise task instruction that describes the transformation, and (ii) edit the corresponding workbook versions so that the input and reference files cleanly realize the described workflow without introducing out-of-scope changes beyond those specified in the instruction.

## 2.1.3 Workflow Sourced from Final Deliverable Spreadsheets and Reports

Third, we curate workflows from high-quality artifacts drawn from the Enron and EUSES corpora, as well as from various investment and securities companies (e.g., Figure 1), international organizations, and national governments. Domain experts author realistic workflow instructions and construct input and reference files based on final deliverable artifacts. For example, a valuation model from an investment firm can be turned into a financial modeling task; a World Bank report can be used to define a summarization and visualization task; and bilingual reports from the Canadian government can be used to construct translation tasks. To broaden coverage of web search and multi-source QA, we adapt 10 financial cases from WideSearch (Wong et al., 2025) into web-search-centric workflows and

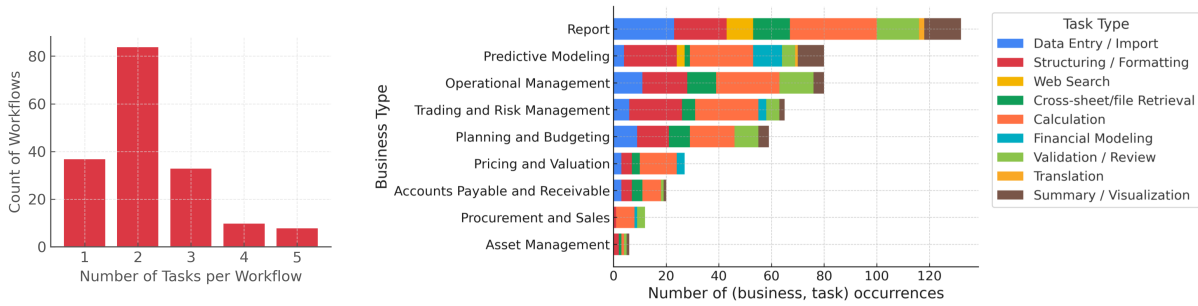


Figure 3: Distribution of number of tasks per workflow and task types across business types.

extend them into multi-step calculation and visualization pipelines. We further leverage 3 examples from DABStep (Egg et al., 2025) to construct multi-source question answering workflows.

### 2.1.4 Quality Control

Given the high complexity of each workflow, we adopt a rigorous multi-stage quality control process. All workflows are annotated and reviewed by a team of five experts: two annotators with interdisciplinary backgrounds in finance and computer science, and three annotators with computer science backgrounds. Among them, three annotators (including one female) have over nine years of industry experience, and two are outstanding Ph.D. or master’s students.

Annotators are instructed to skip workflows that are similar to existing ones to maximize diversity. Every workflow is validated by at least two independent annotators: approximately 40% of workflows undergo at least one round of revision, and particularly complex cases require multiple iterations. In addition, ChatGPT and Claude are used as secondary checkers to flag potential issues, which are always verified by human experts. Together, these procedures required over 700 hours of expert annotation and support the overall quality.

## 2.2 Dataset Characteristics

FINCH comprises 172 meticulously annotated, enterprise-grade workflows that collectively capture the compositional, messy, multimodal, and collaborative nature of real F&A work. All workflows are labeled exclusively for evaluation and are not intended for training. Across these workflows, the corpus contains 1,710 spreadsheets (956 distinct sheets) together with 17 PDFs, 12 images, 3 Word documents, and additional files such as JSON, CSV, and Markdown. This mixture reflects how real analysts coordinate over heterogeneous artifacts rather than clean, single-table inputs.

Figure 3 provides an overview of the task and

business coverage in FINCH. This distribution highlights that FINCH targets realistic, core enterprise workflows rather than curated toy tasks. More details can be found in Appendix A.

### 2.2.1 Task Compositionality

FINCH is explicitly designed around composite workflows rather than isolated tasks. As shown in Figure 3, only 37 workflows (21.5%) are single-task; the remaining 135 (78.5%) involve multiple tasks. Importantly, each task itself typically requires substantial multi-turn reasoning: for example, web search often entails many rounds of LLM calls to discover, filter, and verify evidence; cross-sheet retrieval requires iterative calls to read and locate key information across multiple sheets; and calculation usually spans many formulas distributed over different rows and columns.

### 2.2.2 Messiness

The source files in FINCH are large and interconnected. At the file level, 86.6% of workflows involve more than one file when counting both input and reference artifacts, and a workflow may touch up to 14 distinct files. At the spreadsheet level, 92.4% of workflows involve multiple input and reference sheets, with an average of 8 sheets and a long tail reaching up to 91 sheets. The median workflow covers 15K cells (157K on average), with the largest one scaling to 3.7 million cells.

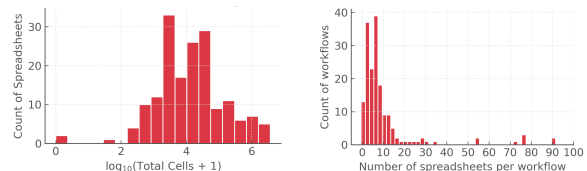


Figure 4: Distribution of the number of cells and sheets.

Moreover, most spreadsheets exhibit multimodal artifacts and intricate structures such as hierarchical headers, merged cells, and other irregularities. On average, a workflow involves 21.5K formulas (with a median of 212), reflecting deeply nested cal-

310 calculations and long dependency chains. In addition,  
 311 20.3% of workflows include charts, while 7.6%  
 312 explicitly require reasoning over PDFs or images.

## 313 2.3 Evaluation Method

### 314 2.3.1 Human Evaluation

315 We conduct human evaluation on all workflows  
 316 to directly assess model performance. For each  
 317 workflow, annotators read the NL instruction and  
 318 inspect the input, reference, and model output files  
 319 side by side (typically by aligning spreadsheets or  
 320 documents in adjacent tabs) to determine whether  
 321 the model has faithfully completed the requested  
 322 task. A workflow is marked as successful only if  
 323 the model generates or revises content and struc-  
 324 ture in accordance with the instruction, without in-  
 325 troducing critical errors, omissions, or unintended  
 326 changes; otherwise, it is labeled as a failure.

327 Importantly, evaluation is based on whether the  
 328 instruction has been satisfactorily fulfilled rather  
 329 than on a purely mechanical comparison between  
 330 the model and reference outputs, since multiple ac-  
 331 ceptable solutions may exist for summarization, vi-  
 332 sualization, formatting, structuring, formulas, and  
 333 related aspects. To reduce subjectivity and ambigu-  
 334 ity, annotators ultimately assign a binary pass/fail  
 335 label to each workflow.

### 336 2.3.2 LLM-as-Judge Evaluation

337 To scale evaluation, we employ an LLM-as-judge  
 338 framework that supports three file processing types:  
 339 *modification* (editing input artifacts), *generation*  
 340 (creating new workbooks or documents), and *QA*  
 341 (answering questions based on one or more arti-  
 342 facts). The framework accepts heterogeneous and  
 343 large inputs—including .xlsx, .txt, .docx, .md,  
 344 .pdf, and images—and normalizes them efficiently  
 345 into multimodal context for the judge model.

346 For spreadsheet modification tasks, instead of en-  
 347 coding full inputs and outputs separately—which  
 348 would consume a large number of tokens when  
 349 working with sizable artifacts—the framework  
 350 computes *structured diffs* between the input and the  
 351 reference output (`diff_ref`), as well as between  
 352 the input and the model output (`diff_model`). In  
 353 addition, it constructs a *compact input snapshot*  
 354 (snapshot) that, for each modified sheet, retains  
 355 only the first and last ten rows and the first five  
 356 columns—typically capturing table headers and  
 357 overall layout—along with all rows and columns  
 358 containing modified cells. This design preserves  
 359 the critical context required to interpret `diff_ref`

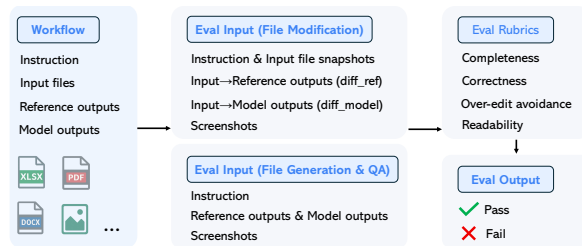


Figure 5: Illustration of automated evaluation pipeline.

360 and `diff_model` while substantially reducing to-  
 361 ken usage. In parallel, the framework renders  
 362 *screenshots* of sheets with changes from the input,  
 363 reference, and model output, enabling the judge  
 364 to perceive merged cells, conditional formatting,  
 365 charts, and other layout-sensitive properties.

366 For generation tasks, the method extracts all text,  
 367 cell values, and formulas from both the reference  
 368 and the model output, and captures screenshots of  
 369 every sheet/doc, since all contents should be ver-  
 370 ified rather than just local edits. For QA tasks,  
 371 it feeds the reference answer and the model’s re-  
 372 sponse, augmented with relevant input artifacts  
 373 when the question requires grounding.

374 Although different file processing strategies are  
 375 used, they share common evaluation rubrics on  
 376 (i) *completeness* with respect to the natural lan-  
 377 guage instruction, (ii) *numerical and logical cor-  
 378 rectness* of derived values and formulas, (iii) *over-  
 379 edit avoidance*, penalizing unnecessary or unex-  
 380 pected changes beyond the instruction, and (iv)  
 381 readability of the formatting and structure. Ex-  
 382 act cell-by-cell equality with the reference is not  
 383 required when multiple solutions are acceptable  
 384 (e.g., alternative layouts, equivalent formulas, or  
 385 different but semantically equivalent summaries);  
 386 instead, the judge determines whether the model  
 387 has satisfactorily fulfilled the instruction. To re-  
 388 duce subjectivity, the judge outputs a binary score  
 389 (pass/fail) along with a short natural language ra-  
 390 tionale for expert review and validation. For web  
 391 search tasks, the rubric permits a 1% tolerance  
 392 band, allowing for precision discrepancies across  
 393 different data sources.

394 This LLM-as-judge framework not only auto-  
 395 mates large-scale evaluation but also surfaces sub-  
 396 tle errors (such as formulas silently replaced with  
 397 static values) that are difficult to catch with GUI-  
 398 based human inspection alone. In Section 3.2, we  
 399 report the consistency between human and auto-  
 400 mated evaluations and show that the LLM-as-judge  
 401 scores closely align with human judgments.

### 3 Experiments

#### 3.1 Baselines

##### 3.1.1 Product-side Agents

We evaluate two frontier product-side agents: (i) *ChatGPT* using GPT 5.1 in Pro mode, and (ii) *Claude* using Opus 4.5 or Sonnet 4.5 in Thinking mode. We focus on these two systems (rather than Gemini or Grok) because they natively return downloadable artifacts (e.g., spreadsheets) for human inspection, instead of emitting executable code or intermediate markdown tables. For both agents, we enable external web browsing but disable access to historical chats, so that each workflow is evaluated independently without cross-run leakage.

##### 3.1.2 API-based Models

We evaluate five frontier LLMs via APIs (Table 5). We adopt SpreadsheetBench (Ma et al., 2024) as our baseline framework because it provides a principled code-generation paradigm for spreadsheet-centric tasks, enabling complex spreadsheet operations that cannot be reliably captured by text-only outputs.

Since SpreadsheetBench was originally designed for relatively small and clean spreadsheets, we extend it with richer encodings and multimodal input support to scale to the large and messy enterprise workflows in FINCH. We preserve its single-call setting to enable fair comparison across different LLMs without substantial framework changes. This setting may underestimate model performance, however, because it precludes iterative interaction, execution feedback, and self-correction mechanisms that product-side agents can leverage. Additional details are deferred to Appendix B.

**Spreadsheet Encoding.** SpreadsheetBench produces text tables without preserving cell addresses, data types, or formulas. We leverage SpreadsheetLLM (Dong et al., 2024)’s address-based encoding and introduce a *semantic-rich markdown encoding*. Each sheet begins with its name and the corresponding data range (e.g. ## Sheet: [name] (A1:Z100)), using a Markdown-based format for serialization. Each cell is encoded as a tuple (Address, Value, Type, Formula), where Address denotes the cell reference (e.g. A3), Type indicates the data type (T = Text, I = Integer, F = Float, D = Date, B = Boolean), and Formula records the cell formula (e.g., =SUM(A1:A10)->100).

**Multimodal Input Handling.** We extend the framework to support multimodal inputs involving images and PDFs, detailed in Appendix B.

Model	Pass Rate (%)
GPT-5.1 Pro (Product)	41.9
Claude Sonnet 4.5 (Product)	29.1
Claude Opus 4.5 (Product)	43.0
GPT-5.1 (API)	32.0
Claude Sonnet 4.5 (API)	20.3
Gemini 3 Pro Preview (API)	27.3
Grok 4 (API)	23.8
Qwen 3 Max (API)	14.5

Table 1: Automated LLM-as-judge evaluation results on FINCH workflows.

#### 3.2 Experimental Results

**Product-side agents.** As shown in Table 1 (automated evaluation using GPT-5-mini as the judge) and Table 4 (human evaluation), ChatGPT 5.1 Pro and Claude Opus 4.5 achieve the strongest overall pass rates on FINCH. Their advantage largely comes from interactive affordances: they can iteratively inspect spreadsheets, revise intermediate states, and recover from partial errors through multiple tool calls. However, even these frontier agents solve fewer than 50% of workflows under human evaluation, suggesting that real-world finance and accounting work remains far from “solved.”

The human evaluation results in Table 2 highlight long-horizon composition as a key bottleneck: when a workflow contains more than two tasks, the pass rate drops sharply—GPT 5.1 Pro falls from 44.3% for workflows with  $\leq 2$  tasks to 23.5% for workflows with  $> 2$  tasks. We also examine how GPT 5.1 Pro’s completion time scales with workflow length (Table 2). The longest individual run takes roughly 60 minutes yet still fails, underscoring the difficulty of highly compositional workflows for current agents.

Tasks per workflow	# Workflows	Pass (%)	Time (min)
1	37	48.6	13.1
2	84	42.4	17.4
3	33	33.3	18.7
$\geq 4$	18	5.6	17.4

Table 2: GPT-5.1 Pro’s pass rates and completion time by workflow composition (human evaluation).

Pass rates also vary substantially by task type (Table 3). Data Entry / Import and Structuring / Formatting are among the most challenging cat-

egories, consistent with the messy layouts, irregular tables, and multi-sheet structures in FINCH. Data Entry / Import is also frequently coupled with web search, a known difficult setting (Wong et al., 2025). Notably, Translation—where modern LLMs typically excel in standard NLP benchmarks—performs poorly in FINCH, as finance-heavy tables make it easy to distort or omit critical structural cues (e.g., header hierarchies and row/column alignment). Detailed error analysis is provided in Section 3.3.

Pass Rate by Task Type(%)	GPT-5.1 Pro	Sonnet 4.5
Data Entry / Import	25.0	13.6
Structuring / Formatting	29.1	25.6
Web Search	9.1	0.0
Cross-sheet/file Retrieval	35.1	13.5
Calculation	39.2	24.2
Validation / Review	37.8	21.6
Financial Modeling	33.3	20.0
Translation	0.0	0.0
Summary / Visualization	36.4	27.3

Table 3: Pass rates by task type for GPT-5.1 Pro and Claude Sonnet 4.5 under human evaluation. For composite workflows, we use workflow-level correctness: a task is counted as correct only if the entire workflow succeeds; if a workflow fails, all tasks it contains are counted as incorrect.

**API-based models.** Table 1 shows that API-based baselines are generally weaker than official product-side agents. Under automated evaluation, GPT 5.1 Pro reaches 41.9% pass rate, whereas GPT 5.1 with our API-based agent design reaches 32.0%. A key limitation of the API baselines is the single-call setting, which precludes iterative interaction, execution feedback, and self-correction—an important direction for future work on enterprise-grade agents. Despite this constraint, our design narrows the gap by using more efficient spreadsheet encodings and task-appropriate tool outputs within the single-call budget.

**Consistency Between Human and Automated Evaluation** As shown in Table 4, the automated evaluation largely aligns with human judgments: for GPT 5.1 Pro and Claude Sonnet 4.5, the GPT-5-mini judge agrees with human labels on 82.1% and 90.2% of workflows, respectively. The judge also achieves high recall (83.3% and 88.4%), meaning it recovers most human-labeled passes. Replacing GPT-5-mini with the stronger GPT-5.1 judge yields lower automated pass rates (37.2% and 23.2%, respectively), bringing them closer to human scores.

Product	Automated (%)		Human (%)
	GPT-5-mini	GPT-5.1	
GPT 5.1 Pro	41.9	37.2	38.4
Sonnet 4.5	29.1	23.2	25.0

Table 4: Comparison of automated LLM-as-judge pass rates using different judges and human evaluation on product-side agents. Human evaluation serves as the gold standard. Since Opus 4.5 has been newly released, we did not perform human evaluation for it.

We adopt GPT-5-mini as the judge to avoid using the same model for both evaluation and the API-based baselines.

On the model side, the LLM judge can occasionally miss nuances in the rubric—either failing to catch subtle visual or numerical errors in large spreadsheets or, conversely, being overly literal about certain instructions (e.g., penalizing benign formula-to-constant conversions). However, we also observe cases where the LLM-based judge is correct but human raters are wrong—for example, when formulas are silently replaced with static values, which are difficult to detect through GUI-based inspection alone. On the system side, limitations of our spreadsheet tooling and data pipeline (e.g., incomplete support for corrupted but human-readable workbooks or uncommon file formats) can cause valid outputs to be marked as failures. Taken together, these factors mean that our automated scores should be interpreted as approximate rather than exact, and that human review remains important for borderline or high-impact workflows.

### 3.3 Error Analysis

To understand the sources of failure on FINCH, we conducted a qualitative error analysis of GPT 5.1 and Claude Sonnet 4.5, considering both product-side agents and API-based models. For all failed workflows in our evaluation, we manually inspected the trajectories and annotated the primary cause of failure.

From a workflow-centric perspective, we identify five dominant categories of error. Take the Claude Sonnet 4.5 product-side agent as an example. Across all examined failures, 10% stem from *task misunderstanding*: enterprise tasks often rely on implicit context in artifacts (e.g., spreadsheets), which models frequently overlook, leading them to misinterpret what is being asked and the required deliverable. 25% are *data retrieval errors*,

554 including selecting the wrong cross-sheet, cross-  
555 table, or intra-table row/column ranges. 35% arise  
556 from *formula reasoning errors*, such as failing to  
557 reconstruct the latent business logic encoded in for-  
558 mulas or deriving incorrect new formulas. 25%  
559 are due to *code generation errors*, where generated  
560 scripts (e.g., Python with spreadsheet APIs) are  
561 syntactically invalid or misaligned with the spread-  
562 sheet layout. The remaining 5% correspond to *data*  
563 *rendering errors*, including incorrect formatting,  
564 misconfigured charts, or flawed final reports that  
565 deviate from the requested layout or narrative—for  
566 example, creating a brand-new spreadsheet instead  
567 of modifying the original one as requested. We  
568 also compare error patterns between web-based  
569 agents and API-based setups, with details provided  
570 in Appendix C.1.

571 Notably, these errors largely arise from *generic*  
572 *capabilities* that modern LLMs already appear to  
573 master in isolation on many existing benchmarks.  
574 The sharp degradation on FINCH therefore stems  
575 not from the absence of these abilities, but from  
576 their *composition* within realistic enterprise Fi-  
577 nance & Accounting workflows. In FINCH, in-  
578 dividual workflows simultaneously involve large  
579 and fragmented spreadsheet ecosystems, dense and  
580 semantically homogeneous financial content, irreg-  
581 ular table structures, latent business logic encoded  
582 in formulas, and multimodal artifacts spanning  
583 spreadsheets, PDFs, charts, and other documents.  
584 When these challenges co-occur, small local er-  
585 rors—such as minor retrieval mistakes or misinter-  
586 preted structures—are easily amplified and propa-  
587 gate across long-horizon execution, ultimately lead-  
588 ing to workflow-level failure.

589 Taken together, these observations suggest that  
590 FINCH does not require fundamentally new model  
591 abilities; rather, it probes existing capabilities un-  
592 der an enterprise “extreme” regime characterized  
593 by high complexity, noise, and long-horizon de-  
594 pendencies that closely mirror real F&A work. A  
595 detailed breakdown and concrete examples of these  
596 error factors are provided in Appendix C.2.

## 597 4 Related Work

598 The integration of LLMs into enterprise produc-  
599 tivity tools has accelerated dramatically in recent  
600 years. The recently launched ChatGPT Agent (Ope-  
601 nAI, 2025) extends these capabilities to financial  
602 and accounting spreadsheet automation. Major pro-  
603 ductivity suites such as Microsoft 365 and Google

604 Workspace have also integrated LLMs (e.g., Copi-  
605 lot and Gemini) to support document drafting,  
606 spreadsheet analysis, and workflow automation  
607 through natural language interaction (Microsoft,  
608 2024; Google, 2024). Anthropic’s Claude Excel  
609 has similarly entered the enterprise space with  
610 spreadsheet automation capabilities (Anthropic,  
611 2025). In parallel, a growing body of research  
612 has emerged on agentic spreadsheet processing (Li  
613 et al., 2023a; Chen et al., 2025; Zhu et al., 2025;  
614 Ma et al., 2024; Wang et al., 2025c).

615 In the past few years, there has also been signifi-  
616 cant progress in benchmarks for financial reason-  
617 ing (Chen et al., 2021; Islam et al., 2023; Bigeard  
618 et al., 2025; Xie et al., 2024a; Liu et al., 2025a;  
619 Choi et al., 2025; Wang et al., 2025b; Li et al.,  
620 2025a; Liu et al., 2025b; Hu et al., 2025; Zhang  
621 et al., 2025; Xie et al., 2024b), office workspace  
622 automation (Wang et al., 2024, 2025a), spreadsheet  
623 reasoning (Ma et al., 2024; Indika and Molybog,  
624 2025; Li et al., 2023b; Dong et al., 2024; Wu et al.,  
625 2025; Li et al., 2025b; Dong et al., 2025; Li et al.,  
626 2024; Zhao et al., 2024), and multimodal document  
627 understanding (Mathew et al., 2021), driving ad-  
628 vances in LLM-based agents for enterprise tasks.  
629 In particular, GDPval (Patwardhan et al., 2025)  
630 evaluates economically valuable and professional  
631 tasks by comparing the quality of full work de-  
632 liverables. FINCH complements these efforts by  
633 emphasizing collaborative and long-horizon work-  
634 flows grounded in pre-existing enterprise artifacts  
635 (emails, versioned workbooks), where agents must  
636 coordinate and revise inherited materials rather  
637 than produce artifacts entirely from scratch.

## 638 5 Conclusion

639 We introduced FINCH, a new benchmark for real-  
640 world F&A enterprise workflows. FINCH derives  
641 workflows induced from enterprise email threads,  
642 version histories of spreadsheets, and high-quality  
643 financial artifacts with rigorous expert annotation  
644 and a calibrated LLM-as-judge framework, en-  
645 abling systematic evaluation of agents on diverse  
646 workflows that operate over large, messy, and  
647 multimodal enterprise artifacts and require long-  
648 horizon reasoning. Experiments show that even  
649 the strongest frontier systems pass fewer than 50%  
650 of workflows, revealing a gap between current AI  
651 capabilities and the demands of real enterprise prac-  
652 tice. We hope FINCH will inspire agents to tackle  
653 real, messy, and long-horizon professional work.

## 6 Limitations

**Access to contemporary enterprise data.** A central challenge in building realistic workflow benchmarks is that true enterprise work records are extremely difficult to obtain due to privacy, compliance, and contractual restrictions. FINCH is therefore anchored in publicly available real-world sources, most prominently the Enron corpus, which contains uniquely rich email threads and attached artifacts that reflect authentic collaborative work. However, Enron primarily reflects workplace practices in the early 2000s, and enterprise tooling and collaboration patterns have evolved over the past two decades. Fortunately, the core productivity substrates for finance and accounting—email as the communication backbone and spreadsheets as the primary computational artifact—were relatively mature by the early 2000s and remain central in modern enterprise operations.

To mitigate this limitation, we complement Enron with a large collection of more recent, publicly available enterprise artifacts from financial institutions, global organizations, and government agencies, including investment and securities firms, the World Bank, and Canadian and British government sources. These materials introduce contemporary data formats, reporting practices, and workflow patterns, helping FINCH better reflect modern finance and accounting work while respecting privacy and compliance constraints.

**From measurement to closing the gap.** Our benchmark surfaces both promising capabilities and substantial shortcomings of frontier AI systems on realistic workflows. However, FINCH is designed primarily as an evaluation and diagnostic resource rather than as a proposal of new modeling or training techniques. How to develop methodologically novel approaches that close the observed gaps—e.g., robust schema understanding under messy layouts, reliable formula-centric reasoning, and long-horizon execution with error recovery—is an important direction for future work.

## 7 Ethical considerations

The FINCH benchmark is constructed entirely from existing, publicly available data sources. Concretely, our workflows are derived from (1) the Enron email corpora, including the parsed Enron email dataset on Kaggle (released under the CC0 Public Domain dedication) and the Enron Email Dataset from EnronData.org (licensed under CC

BY 3.0 US); (2) the EUSES spreadsheet corpus and its modified variants (CC BY 4.0); and (3) a diverse collection of enterprise-like artifacts, including documents from investment and securities companies, the World Bank (CC BY 3.0), Canadian and British government websites (Open Government License), and public corpora such as WideSearch (MIT license) and DABStep (CC BY 4.0). We respect the original licenses of all upstream resources and only redistribute content within the terms they allow.

On top of these sources, we apply additional filtering, normalization, and expert annotation to organize spreadsheets and related documents into coherent workflows with task instructions, input files, and reference outputs. We do not introduce any new personally identifiable information. During curation, we remove obviously sensitive fields when they are not necessary for the task (e.g., personal contact information or signatures) and avoid annotating workflows whose successful completion would depend on sensitive personal attributes rather than business logic. The resulting FINCH dataset is released under the Creative Commons Attribution 3.0 United States license (CC BY 3.0 US), which permits broad reuse while requiring appropriate attribution.

Annotations were performed by volunteer domain experts, rather than paid crowdworkers, all of whom were informed about the research purpose and annotation requirements in advance. Detailed annotation guidelines, including task definitions, decision criteria, and representative examples, were provided to annotators and will be released as part of the dataset documentation.

AI-assisted tools were used to summarize workflows from emails and versioned spreadsheets, identify labeling issues, and improve language clarity during manuscript preparation. All technical content, experimental design, and conclusions were carefully determined by the authors.

The language in FINCH is primarily English, reflecting the dominant language of the underlying Enron and EUSES corpora and many of the public institutional sources. Because some artifacts originate from funds and securities institutions and from Canadian government materials, a small fraction of workflows include Chinese or French content.

Potential risks include over-reliance on agentic systems evaluated on enterprise-like benchmarks, where automation without sufficient human-in-the-loop validation could miss unexpected errors in professional workflows.

756  
757  
758  
  
759  
760  
761  
  
762  
763  
764  
  
765  
766  
767  
768  
  
769  
770  
771  
772  
773  
774  
  
775  
776  
777  
778  
779  
780  
781  
  
782  
783  
784  
785  
786  
787  
788  
  
789  
790  
791  
792  
793  
  
794  
795  
796  
797  
798  
799  
  
800  
801  
802  
803  
804  
  
805  
806  
807  
808

## References

Anthropic. 2025. Claude for excel. <https://claude.com/claude-for-excel>.

Anthropic. 2025. Introducing claude sonnet 4.5. <https://www.anthropic.com/news/claude-sonnet-4-5>. Accessed: 2025-12-14.

World Bank. 2024. *International Debt Report 2024*. World Bank, Washington, DC. World Bank’s annual publication on external debt statistics.

Antoine Bigeard, Langston Nashold, Rayan Krishnan, and Shirley Wu. 2025. Finance agent benchmark: Benchmarking llms on real-world financial research tasks. *arXiv preprint arXiv:2508.00828*.

Yibin Chen, Yifu Yuan, Zeyu Zhang, Yan Zheng, Jinyi Liu, Fei Ni, Jianye Hao, Hangyu Mao, and Fuzheng Zhang. 2025. Sheetagent: towards a generalist agent for spreadsheet reasoning and manipulation via large language models. In *Proceedings of the ACM on Web Conference 2025*, pages 158–177.

Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan R Routledge, and 1 others. 2021. Finqa: A dataset of numerical reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711.

Chanyeol Choi, Jihoon Kwon, Alejandro Lopez-Lira, Chaewoon Kim, Minjae Kim, Juneha Hwang, Jaeseon Ha, Hojun Choi, Suyeol Yun, Yongjin Kim, and 1 others. 2025. Finagentbench: A benchmark dataset for agentic retrieval in financial question answering. In *Proceedings of the 6th ACM International Conference on AI in Finance*, pages 632–637.

Department of Finance Canada. 2025. *Fiscal reference tables, november 2025*. Technical report, Government of Canada, Ottawa, Canada. Provides annual data on the financial position of the federal, provincial-territorial and local governments.

Haoyu Dong, Yue Hu, and Yanan Cao. 2025. Reasoning and retrieval for complex semi-structured tables via reinforced relational data transformation. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1382–1391.

Haoyu Dong, Jianbo Zhao, Yuzhang Tian, Junyu Xiong, Shiyu Xia, Mengyu Zhou, Yun Lin, José Cambronero, Yeye He, Shi Han, and 1 others. 2024. Spreadsheetllm: encoding spreadsheets for large language models. *arXiv preprint arXiv:2407.09025*.

Alex Egg, Martin Iglesias Goyanes, Friso Kingma, Andreu Mora, Leandro von Werra, and Thomas Wolf. 2025. Dabstep: Data agent benchmark for multi-step reasoning. *arXiv preprint arXiv:2506.23719*.

EnronData.org. Edo enron email pst dataset. <https://enrondata.readthedocs.io/en/latest/data/edo-enron-email-pst-dataset/>. Creative Commons Attribution 3.0 United States License. To provide attribution, please cite to “EnronData.org.”. 809  
810  
811  
812  
813

Marc Fisher and Gregg Rothermel. 2005. The euses spreadsheet corpus: a shared resource for supporting experimentation with spreadsheet dependability mechanisms. In *Proceedings of the first workshop on End-user software engineering*, pages 1–5. 814  
815  
816  
817  
818

Google. 2024. Gemini for google workspace. <https://workspace.google.com/solutions/ai/>. 819  
820

Google DeepMind. 2025. Gemini 3 pro. <https://deepmind.google/models/gemini/pro/>. Accessed: 2025-12-14. 821  
822  
823

HM Treasury. 2023. *Public expenditure statistical analyses 2023*. Technical report, HM Treasury, London, United Kingdom. UK public expenditure statistical release (PESA). 824  
825  
826  
827

Liang Hu, Jianpeng Jiao, Jiashuo Liu, Yanle Ren, Zhoufutu Wen, Kaiyuan Zhang, Xuanliang Zhang, Xiang Gao, Tianci He, Fei Hu, and 1 others. 2025. Finsearchcomp: Towards a realistic, expert-level evaluation of financial search and reasoning. *arXiv preprint arXiv:2509.13160*. 828  
829  
830  
831  
832  
833

Amila Indika and Igor Molybog. 2025. Sodbench: A large language model approach to documenting spreadsheet operations. *arXiv preprint arXiv:2510.19864*. 834  
835  
836  
837

Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. 2023. Financebench: A new benchmark for financial question answering. *arXiv preprint arXiv:2311.11944*. 838  
839  
840  
841

Bryan Klimt and Yiming Yang. 2004. The enron corpus: A new dataset for email classification research. In *European conference on machine learning*, pages 217–226. Springer. 842  
843  
844  
845

Haohang Li, Yupeng Cao, Yangyang Yu, Shashidhar Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen Jiang, Zining Zhu, Kp Subbalakshmi, Jimin Huang, and 1 others. 2025a. Investorbench: A benchmark for financial decision-making tasks with llm-based agent. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2509–2525. 846  
847  
848  
849  
850  
851  
852  
853

Hongxin Li, Jingran Su, Yuntao Chen, Qing Li, and Zhao-Xiang Zhang. 2023a. Sheetcopilot: Bringing software productivity to the next level through large language models. *Advances in Neural Information Processing Systems*, 36:4952–4984. 854  
855  
856  
857  
858

Jinyang Li, Nan Huo, Yan Gao, Jiayi Shi, Yingxiu Zhao, Ge Qu, Yurong Wu, Chenhao Ma, Jian-Guang Lou, and Reynold Cheng. 2024. Tapilot-crossing: Benchmarking and evolving llms towards interactive data analysis agents. *arXiv preprint arXiv:2403.05307*. 859  
860  
861  
862  
863

864	Peng Li, Yeye He, Cong Yan, Yue Wang, and Surajit Chaudhuri. 2023b. Auto-tables: Synthesizing multi-step transformations to relationalize tables without using examples. <i>Proceedings of the VLDB Endowment</i> , 16(11):3391–3403.	918
865		919
866		920
867		921
868		922
869	Zheng Li, Yang Du, Mao Zheng, and Mingyang Song. 2025b. Mimotable: A multi-scale spreadsheet benchmark with meta operations for table reasoning. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 2548–2560.	923
870		924
871		925
872		926
873		927
874	Shu Liu, Shangqing Zhao, Chenghao Jia, Xinlin Zhuang, Zhaoguang Long, Jie Zhou, Aimin Zhou, Man Lan, and Yang Chong. 2025a. Findabench: Benchmarking financial data analysis ability of large language models. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 710–725.	928
875		929
876		930
877		931
878		932
879		933
880	Zhaowei Liu, Xin Guo, Haotian Xia, Lingfeng Zeng, Fangqi Lou, Jinyi Niu, Mengping Li, Qi Qi, Jiahuan Li, Wei Zhang, and 1 others. 2025b. Visfineval: A scenario-driven chinese multimodal benchmark for holistic financial understanding. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 24099–24157.	934
881		935
882		936
883		937
884		938
885		939
886		940
887	Zeyao Ma, Bohan Zhang, Jing Zhang, Jifan Yu, Xiaokang Zhang, Xiaohan Zhang, Sijia Luo, Xi Wang, and Jie Tang. 2024. Spreadsheetbench: Towards challenging real world spreadsheet manipulation. <i>Advances in Neural Information Processing Systems</i> , 37:94871–94908.	941
888		942
889		943
890		944
891		945
892		946
893	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	947
894		948
895		949
896		950
897		951
898	Microsoft. 2024. Microsoft 365 copilot. <a href="https://www.microsoft.com/en-us/microsoft-365/copilot">https://www.microsoft.com/en-us/microsoft-365/copilot</a> .	952
899		953
900		954
901	OpenAI. 2025. Gpt-5. <a href="https://openai.com/index/introducing-gpt-5/">https://openai.com/index/introducing-gpt-5/</a> . Accessed: 2025-12-14.	955
902		956
903	OpenAI. 2025. Introducing chatgpt agent. <a href="https://openai.com/index/introducing-chatgpt-agent/">https://openai.com/index/introducing-chatgpt-agent/</a> .	957
904		958
905		959
906	Tejal Patwardhan, Rachel Dias, Elizabeth Proehl, Grace Kim, Michele Wang, Olivia Watkins, Simón Posada Fishman, Marwan Aljubeih, Phoebe Thacker, Laurence Fauconnet, and 1 others. 2025. Gdpval: Evaluating ai model performance on real-world economically valuable tasks. <i>arXiv preprint arXiv:2510.04374</i> .	960
907		961
908		962
909		963
910		964
911		965
912		966
913	Weixuan Wang, Dongge Han, Daniel Madrigal Diaz, Jin Xu, Victor Rühle, and Saravan Rajmohan. 2025a. Odysseybench: Evaluating llm agents on long-horizon complex office application workflows. <i>arXiv preprint arXiv:2508.09124</i> .	967
914		968
915		969
916		970
917		971
	Yan Wang, Keyi Wang, Shanshan Yang, Jaisal Patel, Jeff Zhao, Fengran Mo, Xueqing Peng, Lingfei Qian, Jimin Huang, Guojun Xiong, and 1 others. 2025b. Finauditing: A financial taxonomy-structured multi-document benchmark for evaluating llms. <i>arXiv preprint arXiv:2510.08886</i> .	972
		973
	Zilong Wang, Yuedong Cui, Li Zhong, Zimin Zhang, Da Yin, Bill Yuchen Lin, and Jingbo Shang. 2024. Officebench: Benchmarking language agents across multiple applications for office automation. <i>arXiv preprint arXiv:2407.19056</i> .	974
		975
		976
		977
		978
	Ziwei Wang, Jiayuan Su, Mengyu Zhou, Huaxing Zeng, Mengni Jia, Xiao Lv, Haoyu Dong, Xiaojun Ma, Shi Han, and Dongmei Zhang. 2025c. Sheetbrain: A neuro-symbolic agent for accurate reasoning over complex and large spreadsheets. <i>arXiv preprint arXiv:2510.19247</i> .	979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

974 spreadsheet formulas from natural language queries.  
975 In *Findings of the Association for Computational*  
976 *Linguistics: EACL 2024*, pages 2377–2388.

977 Ruiyan Zhu, Xi Cheng, Ke Liu, Brian Zhu, Daniel  
978 Jin, Neeraj Parihar, Zhoutian Xu, and Oliver Gao.  
979 2025. Sheetmind: An end-to-end llm-powered multi-  
980 agent framework for spreadsheet automation. *arXiv*  
981 *preprint arXiv:2506.12339*.

## A Detailed Task and Business Type Distribution

Figure 3 summarizes the coverage of task and business types in FINCH. On the task side, the benchmark includes calculation tasks (119 workflows) such as filling in formulas or computing figures; structuring and formatting tasks (86) involving table reorganization, layout adjustment, and row/column insertion or deletion; data entry and import tasks (44) that transcribe or ingest data from spreadsheets, PDFs, images, or external sources; validation and review tasks (37) that check consistency and reconcile calculations within or across sheets and files; cross-sheet or cross-file retrieval tasks (36) that aggregate values from multiple sources into a target workbook; summary and visualization tasks (33) that produce financial summaries or charts; financial modeling tasks (15) that extend or calibrate valuation and timing models; web search tasks (11) that collect and integrate financial data from online sources; and a small number of translation tasks (3) that translate spreadsheets or reports while preserving structure, formatting, and layout.

On the business side, workflows span reporting (48 workflows), trading and risk management (35), predictive modeling (33), operational management (36), planning and budgeting (26), pricing and valuation (15), accounts payable and receivable (10), as well as procurement and sales (7) and asset management (3). Some workflows are tagged with multiple business types, reflecting the cross-functional nature of real-world enterprise finance and accounting work.

## B Experiment Details

### B.1 API-based Models

**Execution Paradigm.** We frame the API evaluation as a **code generation task**. Models are instructed to solve spreadsheet manipulation and generation workflows by generating executable Python scripts, which are then executed in a sandboxed environment to produce output artifacts. This paradigm aligns with SpreadsheetBench’s philosophy of treating model-written code as the primary action space, but is adapted here to accommodate long-horizon, multimodal enterprise tasks.

- **Action Space:** Models generate Python code using standard libraries including `openpyxl` (for Excel manipulation), `pandas` (for data pro-

cessing), `matplotlib` (for visualization), and `scikit-learn` (for statistical analysis).

- **Output Format:** Models must produce complete, self-contained Python scripts wrapped in markdown code blocks (“python . . .”).
- **Sandboxed Execution:** Generated code is extracted via regex parsing and executed in isolated Docker containers running Jupyter Kernel Gateway. Each container mounts the dataset volume at `/mnt/data/` with a 10-minute session timeout.
- **Single-shot Protocol:** We employ a strict one-shot generation protocol without iterative refinement—each model produces exactly one solution per workflow. If the generated code fails to execute (e.g., due to syntax errors or runtime exceptions), the workflow is marked as failed without retry. This strict setting is designed to evaluate the model’s raw code generation capability under realistic deployment constraints.

This unified code-as-action setting ensures that the measured performance reflects the model’s inherent competence on complex workflows rather than benefits derived from interactive debugging.

**Prompting Strategy.** We employ a **zero-shot** setting with a structured system prompt comprising:

1. A role definition: “You are an expert who can manipulate spreadsheets through Python code.”
2. A detailed description of the compact spreadsheet encoding format with illustrative examples.
3. The task instruction and explicit input/output file paths.
4. Library-specific best practices (e.g., `openpyxl` chart creation patterns) to mitigate common code errors.
5. An explicit directive to generate Python code as the final output.

This structured design explicitly guides models toward generating valid, context-aligned Python code, minimizing ambiguity in task interpretation. However, for models that support reasoning traces (GPT 5.1, Gemini 3 Pro), we request explicit reasoning via the `include_reasoning` API parameter, enabling us to capture the model’s internal deliberation process for subsequent qualitative error analysis. Temperature is set to 0.7 across all models.

Model	Provider	Context	Max Output	Vision	Native PDF
GPT 5.1 (OpenAI, 2025)	OpenAI	400K	128K	✓	✓
Claude Sonnet 4.5 (Anthropic, 2025)	Anthropic	1M <sup>†</sup>	64K	✓	✓
Grok 4 (xAI, 2025)	xAI	256K	256K	✓	—
Qwen 3 Max (Yang et al., 2025)	Alibaba	256K	32.8K	—	—
Gemini 3 Pro Preview (Google DeepMind, 2025)	Google	1.05M	65.5K	✓	✓

Table 5: API-based model configurations. Context and output limits are measured in tokens. Vision indicates native image input support, while Native PDF refers to direct PDF file ingestion via the provider’s API without explicit text extraction. <sup>†</sup>Available via long-context beta API mode.

**Model License** All models evaluated in this work are accessed through their official APIs or product interfaces and are subject to the corresponding providers’ terms of service and licensing agreements. We do not redistribute any model weights, training data, or proprietary outputs.

GPT 5.1 is provided by OpenAI and used in accordance with OpenAI’s API and product usage policies. Claude Sonnet 4.5 is provided by Anthropic and accessed via its official API under Anthropic’s commercial license. Grok 4 is provided by xAI and evaluated through its public API under xAI’s terms of service. Qwen 3 Max is provided by Alibaba Cloud and accessed subject to Alibaba’s licensing terms. Gemini 3 Pro Preview is provided by Google and evaluated via Google’s official API under the applicable Google AI service terms.

All evaluations are conducted for research purposes. We do not claim ownership over any model outputs, and the results reported in this paper reflect observed model behavior under controlled experimental settings rather than any guarantees of model performance. Any use of the evaluated models beyond this study must comply with the respective providers’ licenses and usage policies.

**Multimodal Input Handling.** We extend the framework to support multimodal inputs involving images and PDFs. For vision-capable models (GPT 5.1, Claude Sonnet 4.5, Grok 4, and Gemini 3 Pro), we use each provider’s official multimodal API to transmit visual inputs alongside text prompts. For PDF documents, we adopt a tiered strategy. Models with native PDF support—GPT 5.1, Claude Sonnet 4.5, and Gemini 3 Pro Preview—directly ingest PDF files via their file upload interfaces, enabling analysis of both textual and visual elements without pre-extraction. For Grok 4, which lacks native PDF support, we extract text using PyMuPDF and include it in the pdf\_content field. For Qwen 3 Max, which lacks multimodal support entirely, both image and PDF content are converted

to textual descriptions. While this fallback retains semantic cues, it loses layout and visual context.

**Context Management.** To handle large spreadsheets that may exceed model context limits, we implement automatic truncation. We reserve 32K tokens for model output—sufficient for comprehensive code generation and analysis while remaining within the output limits of all evaluated models. Truncation is triggered when input exceeds the remaining capacity, removing content from the end of the spreadsheet data with an explicit notice appended to inform the model of data loss.

## C Detailed Analysis

### C.1 Agent Behavior Analysis

**Product-side agents.** ChatGPT 5.1 Pro tends to decompose workflows into more, smaller steps, with explicit reasoning, tool calls, execution, and self-checking at each step. This leads to longer traces and noticeably higher latency, but also more opportunities for intermediate validation (e.g., sanity-checking partial results). However, the code it generates is often hidden behind tool abstractions, so our error attribution is limited to observed behavior and natural language reasoning rather than the exact implementation details. Claude Sonnet 4.5 typically uses fewer steps and produces more direct solutions. In visualization-heavy workflows, its generated charts are often both more accurate and more aesthetically polished than those produced by ChatGPT 5.1 Pro, leading to relatively fewer failures in the data visualization sub-tasks.

ChatGPT 5.1 Pro and Claude Sonnet 4.5 agents can explore Excel files through many API calls within a single workflow, but their encoding methods are not well-suited to spreadsheets with complex layouts and structures. Thanks to efficient encoding and appropriate tool use, the following single-call API-based method achieves a pass rate that is much closer to that of product-side agents.



Figure 6: Real-world F&A work is messy, spanning heterogeneous and large-scale artifacts such as spreadsheets and PDFs. It’s also long-horizon and knowledge-intensive: workflows interleave multiple tasks and span diverse domains such as budgeting, trading, asset management, and operations.

**API-based models** Our API-based runs are single-call: they leverage the models’ underlying reasoning capabilities but lack two crucial affordances that web agents exploit. (i) interleaved code execution with feedback, and (ii) explicit reflection based on intermediate tool outputs. As a result, the API agents must generate the entire plan, code, and outputs within a single LLM call. When their initial structural assumptions about a spreadsheet are slightly off, they have no mechanism to detect or correct the mistake, leading to a significantly higher error rate, particularly in categories related to schema understanding and table manipulation. It’s desirable for future work to explore agentic methods with multiple rounds of API calls.

## C.2 Detailed Error Factors in Enterprise F&A Workflows

First, FINCH workflows routinely involve *large, fragmented spreadsheet ecosystems*: dozens of interlinked workbooks and thousands of rows distributed across many sheets. Executing these workflows accurately requires long-range cross-sheet navigation and precise referencing, which substantially increases the likelihood of small retrieval errors. Second, the *content is dense and semantically homogeneous*: many cells contain domain-specific financial concepts that are subtly different yet lexically similar (e.g., variants of revenue/expense items, adjusted vs. unadjusted metrics), making entity disambiguation and cell grounding unusually difficult. Third, the *table layouts and structures are complex and often irregular*, including multi-level headers, merged cells, nested subtotals, and

bespoke layouts that force the model to infer structure from noisy contents and ad hoc formatting. For example, at the code level, even tiny misinterpretations of these layouts (e.g., off-by-one errors when specifying ranges) can then propagate into globally incorrect outputs, especially when logic is applied in batch across many such sheets. Fourth, *formulas encode latent structure and logic*. In the FINCH dataset, each sheet contains a large number of formulas that encode latent business logic, temporal assumptions, and fine-grained dependencies that are not visible from displayed values alone; yet models typically prioritize cell values and under-use formulas, leading to systematic misinterpretations. For example, in a pricing sheet with the column header IF NGPL MidContinent index (@ Baker), the apparent semantics from the header alone suggest a daily exposure metric. However, inspecting the associated formula ( $25 * V21 + C41 * C22$ ) reveals that the column in fact encodes a 55-day payment timing. Models that ignore or under-utilize formulas systematically misinterpret such columns’ roles in downstream calculations, and this misinterpretation then propagates through subsequent steps. Finally, many workflows involve *multimodal artifacts and chat-centric tasks* such as combining spreadsheets with PDFs, charts, and screenshots, requiring the agent to jointly reason over heterogeneous formats. For example, tables embedded in PDFs are often only partially referenced, with key entries missing or truncated.

## D Examples

This section presents representative examples from FINCH to illustrate the diversity and complexity of spreadsheet-centric enterprise workflows covered by the benchmark. The examples span a wide range of realistic professional tasks, including cross-sheet verification, formula auditing, document-grounded extraction, schema transformation, financial modeling, and reporting.

Note that the figures shown in this section are illustrative excerpts rather than complete workflow inputs or full task specifications. Many FINCH workflows involve large, multi-sheet or multi-file inputs and extensive supporting documents, which cannot be fully visualized in static screenshots. Accordingly, the figures are intended to highlight representative data characteristics, structural complexity, and reasoning challenges that arise in the underlying workflows, rather than to fully specify the tasks themselves.

Across these workflows, AI agents are required to jointly reason over heterogeneous artifacts, while preserving structural consistency and semantic correctness. Many tasks demand multi-step reasoning, precise interpretation of formulas, and careful cross-referencing across sheets or external documents. Collectively, these examples highlight the core challenges targeted by FINCH: messy and multimodal document processing, long-horizon reasoning, robust code generation under schema variation, and faithful execution of complex financial and accounting logic.

	B	C	D	E	F	K
1	Houston					
2	Gas Trading (including NGLs)				43	
3	Gas Origination (includes Mexico)				34	
4	Admins				12	
5						
6	Power Trading				43	
7	Power Origination				14	
8	Development Systems				3	
9	Admins				4	
10						
11	Portland					
12	Power Trading & Origination				36	
13						
14	Admins				3	
15	Canada					
16	Trading				11	
17	Origination				15	
18	Admins				5	
19						
20						
21	Fundamentals					
22	Houston (includes Competitor Analysis = 3)				23	
23						
24	Structuring				3	
25						
26	Weather				5	
27						
28	Credit				15	
29						
30	Regulatory Affairs				9	
31						
32						
33	Energy Operations					
34	Gas Logistics				33	
35	Gas Risk				41	
36	Gas Settlements				7	
37	Gas Volume Management				21	
38	Documentation				30	
39	Power Logistics				13	
40	Power Book Running				11	
41	Power Settlements & Confirms				9	
42	Power Volume Management				9	
43	Financial Settlements				2	
44	Management					
45						
46	Technology					
47	Infrastructure				60	
48	IT				141	
49						
50	OOO				4	
51	EnronOnline				57	
52						
53	HR				12	
54						
55	Legal				20	
56	Accounting				44	
57	Transacion Support				2	
58	Cash Operations				6	
59						
60	SAP				-3	
61	Canada Support (including Legal but not HR)				33	
62						
63						
64						
65						
66	Total					

Validate

Cross-sheet retrieval

	B	C	D	E	F	K
1	Houston					
2	Gas Trading (including NGLs)				43	
3	Gas Origination (includes Mexico)				35	
4	Admins				12	
5						
6	Power Trading				45	
7	Power Origination				14	
8	Development Systems				3	
9	Admins				4	
10						
11	Portland					
12	Power Trading & Origination				37	
13						
14	Admins				6	
15	Canada					
16	Canada				11	
17	Trading				11	
18	Origination				11	
19	Toronto				4	
20						
21						
22	Fundamentals					
23	Houston (includes Competitor Analysis = 3)				24	
24						
25	Structuring				3	
26						
27	Weather				5	
28						
29	Credit				15	
30						
31	Regulatory Affairs				9	
32						
33						
34	Energy Operations					
35	Gas Logistics				32	
36	Gas Book Running				41	
37	Gas Settlements				21	
38	Gas Volume Management				7	
39	Documentation				30	
40	Power Logistics				14	
41	Power Risk				11	
42	Power Settlements & Confirms				9	
43	Power Volume Management				13	
44	Financial Settlements				9	
45	Management				2	
46						
47	Technology					
48	Infrastructure				61	
49	IT				143	
50						
51	OOO				4	
52	EnronOnline				4	
53						
54	HR				12	
55						
56	Legal				23	
57						
58	Accounting				44	
59	Cash Operations				6	
60	Canada Support (including Legal but not HR)				33	
61						
62	Tax				4	
63						
64						
65						
66	Total					170

Example Sheets

A	B	C	D	E	F
1	Lawyers	ok		Bonus	
2	Houston				
3	Hansen, Leslie	Sr. Counsel		75000	
4	Cook, Mary	Sr. Counsel			R
5	Headrick, Mark	MD			
6	Hodge, Jeff	VP	125000		
7	Koehler, Anne	Sr. Counsel		75000	
8	McCullough, Travis	VP			R
9	Hermes, Gerald	Sr. Counsel		25000	
10	Natanson, Marcus	Sr. Counsel		75000	
11	Prins-Lewis, Francisco	Sr. Counsel		75000	
12	Sager, Elizabeth	VP			R
13	St. Clair, Carol	VP			R
14	Taylor, Mark	VP			
15	Van Hoser, Steve	Asst. General Counsel		75000	
16	Portland				
17	East, Steve	Sr. Counsel		100000	
18	Rasmussen, Dale	Sr. Counsel		75000	
19	Yoder, Christian	Sr. Counsel		75000	
20	Legal Specialists				
21	Fitzgerald, Genia	Sr. Counsel		15000	
22	Hearn, Mone	Sr. Counsel		15000	
23	Jones, Tana	Sr. Counsel		15000	
24	Admins				
25	Elberston, Janette	Admin		5000	
26	Kaiser, Holly	Admin		4000	
27	Adams, Suzanne	Admin		4000	
28	Simmons, Linda	Admin		3000	
29	Total		21	836,000	
30					
31					
32	Removed				
33	Cash				
34					
35	Add				
36					
37					
38	Adams, Suzanne	Admin		4000	
39	Simmons, Linda	Admin		3000	

Sheet: Legal

Sheet: Gas Trading

Sheet: Portland Trading & Origination

Sheet: Canda

A	B	C	D	E	F	G
1	Regulatory	ok				
2						
3	Commes	Alan	West Power	Director		25000
4	Young	Charles	ERCOT	Director		40000
5	Lindberg	Susan		Director	3	40000
6	Calciagno	Suzanne		Manager	1	
7	Dasovich	Jeff	West			
8	Nicolay	Christi	Sr. Director			R
9	Walton	Steve	Sr. Director			R
10	Novosell	Sarah	Sr. Director	4		50000
11	Montbano	Steve	East Gas	Vice Pres	1	75000
12						
13						
14						
15						230,000
16						
17	Staffies	James	VP			

Sheet: Regulatory

Sheet: Market Risk & Research

Sheet: Portland Fundies & Structuring

Figure 7: For this task, the model must verify the department headcount summary by cross-checking each of the 39 departments against its detailed roster sheet. It should correct discrepancies such as miscounts and missing or duplicate entries. The summary must be updated by fixing incorrect totals, removing departments that no longer exist, and adding any omitted departments. Furthermore, the underlying schema varies slightly across departments, which challenges reliable code generation.

Table 5.1 Public sector expenditure on services by departmental group and function, 2024-25

Departmental Grouping	Function	Accredited Official Statistics										Public sector expenditure on services for each department									
		1. General public services	2. Defence	3. Public order and safety	4. Economic affairs	5. Environment	6. Housing and community amenities	7. Health	8. Recreation, culture and religion	9. Education	10. Social protection										
Health and Social Care		-	-	-	-	-	-	-	-	-	-	-	193,350								
Education		-	-	-	-	-	-	-	-	-	-	-	52,215								
Home Office		-	-	-	-	-	-	-	-	-	-	-	6,679								
Justice		-	-	-	-	-	-	-	-	-	-	-	13,721								
Law Officers' Departments		-	-	-	-	-	-	-	-	-	-	-	954								
Defence		-	-	-	-	-	-	-	-	-	-	-	61,138								
Single Intelligence Account		-	-	-	-	-	-	-	-	-	-	-	4,442								
Foreign, Commonwealth and Development Office		-	-	-	-	-	-	-	-	-	-	-	10,852								
Housing, Communities and Local Government		-	-	-	-	-	-	-	-	-	-	-	10,889								
Culture, Media and Sport		-	-	-	-	-	-	-	-	-	-	-	10,808								
Science, Innovation and Technology		-	-	-	-	-	-	-	-	-	-	-	13,797								
Transport		-	-	-	-	-	-	-	-	-	-	-	31,097								
Energy Security and Net Zero		-	-	-	-	-	-	-	-	-	-	-	9,379								
Environment, Food and Rural Affairs		-	-	-	-	-	-	-	-	-	-	-	6,432								
Business and Trade		-	-	-	-	-	-	-	-	-	-	-	2,985								
Work and Pensions		-	-	-	-	-	-	-	-	-	-	-	282,927								
HM Revenue and Customs		-	-	-	-	-	-	-	-	-	-	-	33,993								
HM Treasury		-	-	-	-	-	-	-	-	-	-	-	124,967								
Cabinet Office		-	-	-	-	-	-	-	-	-	-	-	4,022								
Scottish Government		-	-	-	-	-	-	-	-	-	-	-	40,583								
Welsh Government		-	-	-	-	-	-	-	-	-	-	-	16,387								
Northern Ireland Executive		-	-	-	-	-	-	-	-	-	-	-	28,307								
Small and Independent Bodies		-	-	-	-	-	-	-	-	-	-	-	3,109								
Local Government		-	-	-	-	-	-	-	-	-	-	-	193,846								
<b>Public sector expenditure on services for each</b>		<b>157,591</b>	<b>22,342</b>	<b>10,533</b>	<b>124,715</b>	<b>63,648</b>	<b>51,370</b>	<b>96,793</b>	<b>20,422</b>	<b>9,677</b>	<b>3,983</b>	<b>6,296</b>	<b>46,415</b>	<b>17,141</b>	<b>22,319</b>	<b>241,835</b>	<b>14,522</b>	<b>118,674</b>	<b>383,934</b>	<b>-1,433</b>	<b>1,156,393</b>

Figure 1 Public and publicly guaranteed debt, by creditor and creditor type in 2023, including IMF credit

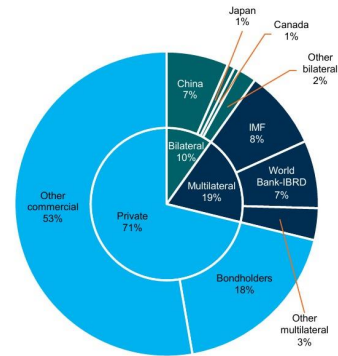


Figure 8: An example of extracting data from tables and charts in PDFs and saving it to a spreadsheet. AI agents must understand the layout, parse hierarchical structures, interpret how values map to specific cells, and reconstruct formulas from aggregated values.

	Enron Amounts	AEP Amounts	Difference	Adjusted Difference
<b>Working Gas</b>				
Invoiced by Enron	\$ 100,067,365.68			
Actual	\$ 82,294,269.66			
<b>Working Gas Adjustment Due Buyer</b>	<b>\$ 17,773,096.02</b>	\$ 17,773,096.02	\$ -	
<b>Add(Deduct):</b>				
Texas General Land Office	\$ 264,392.52	\$ 264,392.52	\$ -	
Cannon	\$ 393,750.00	\$ 393,750.00	\$ -	
Lyondell Cigo Adjustment	\$ (1,762,925.00)	\$ (1,762,925.00)	\$ -	
Centana Gas Payment	\$ (2,936,179.25)	\$ (2,936,179.25)	\$ -	
Centana Ad Valorem Tax Proration	\$ 40,059.11	\$ 64,094.57	\$ (24,035.47)	\$ (8,011.82)
Centana July Reimbursement	\$ -	\$ 95,130.00	\$ (95,130.00)	\$ -
Gas Lift Deposits	\$ 25,000.00	\$ 25,225.00	\$ (225.00)	\$ (22.00)
Specialty Sands	\$ 91,000.00	\$ 91,000.00	\$ -	
SAP to PeopleSoft Conversion	\$ (43,001.25)	\$ (43,001.25)	\$ -	
<b>Adjusted Working Gas Payment</b>	<b>\$ 13,845,192.15</b>	<b>\$ 13,964,582.61</b>	<b>\$ (119,390.47)</b>	<b>\$ (8,033.82)</b>
Interest Payment (4)	\$ 373,725.36	\$ 486,415.08	\$ (112,689.72)	\$ (112,689.72)
<b>Adjusted Payment</b>	<b>\$ 14,218,917.51</b>	<b>\$ 14,450,997.69</b>	<b>\$ (232,080.19)</b>	<b>\$ (120,723.54)</b>
(1) Estimated Bammel Storage volume as of 05/31/01				
(2) Actual Bammel Storage volume as of 05/31/01				
(3) Published Price in Inside FERC (HSC) for the month of July 2001				
(4) Interest Payment calculated based on an assumed pay date of 10/31/01 for Daily Prime for 153 days (Enron calculate on "Additional Payment" and AEP calculated on "Working Gas Adjustment")				

Figure 9: Cross-sheet reference validation. This example is relatively easy for frontier AI agents.

	A	B	F
7	<b>Base Case Variables</b>		
8	Plant Capacity (PC)		263,000
9	Discount Rate		5.00%
10	<b>Question 1</b>		
12	Include Interest Rate Adjustment?	No	
13	Include Apache Savings?	No	
14	<b>Question 2</b>		
16	Include Interest Rate Adjustment?	Yes	
17	Include Apache Savings?	Yes	
18	<b>Question 3</b>		
20	Include Interest Rate Modification?	Yes	
21	Include Apache Savings?	Yes	
22	<b>Question 4</b>		
24	Include Interest Rate Adjustment?	No	
25	Include Apache Savings?	Yes	

**Cleburne Plant**  
Tenaska's Interest Rate Adjustment

THIS INFORMATION IS CONFIDENTIAL PER SECTION 9 OF THE FACILITATION AGREEMENT

	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018
<b>Base Case</b>	0.00%	1.25%	1.75%	2.25%	2.75%	3.25%	3.75%	4.25%	4.75%	5.25%	5.75%	6.25%	6.75%	7.25%	7.75%	8.25%	8.75%	9.25%	9.75%	10.25%	10.75%	11.25%
<b>Scenario 1</b>	0.00%	1.25%	1.75%	2.25%	2.75%	3.25%	3.75%	4.25%	4.75%	5.25%	5.75%	6.25%	6.75%	7.25%	7.75%	8.25%	8.75%	9.25%	9.75%	10.25%	10.75%	11.25%
<b>Scenario 2</b>	0.00%	1.25%	1.75%	2.25%	2.75%	3.25%	3.75%	4.25%	4.75%	5.25%	5.75%	6.25%	6.75%	7.25%	7.75%	8.25%	8.75%	9.25%	9.75%	10.25%	10.75%	11.25%
<b>Scenario 3</b>	0.00%	1.25%	1.75%	2.25%	2.75%	3.25%	3.75%	4.25%	4.75%	5.25%	5.75%	6.25%	6.75%	7.25%	7.75%	8.25%	8.75%	9.25%	9.75%	10.25%	10.75%	11.25%
<b>Scenario 4</b>	0.00%	1.25%	1.75%	2.25%	2.75%	3.25%	3.75%	4.25%	4.75%	5.25%	5.75%	6.25%	6.75%	7.25%	7.75%	8.25%	8.75%	9.25%	9.75%	10.25%	10.75%	11.25%

Tenaska's i Adjustment Sheet

=+'Tenaska's i Adjustment'!D\$44\*G\$46

=+'Damage Calculations'!G\$85\*Damage Calculations!\$C\$33

	A	B	C	G	H	I	J	K		
				1997	1998	1999	2000	2001	2002	2003
44	<b>Question 2 - At Inception</b>									
45	Contract Capacity Rates	=+IF(\$B\$1 6="yes",1, 0) \$/kW-mo.		\$ 13.60	\$ 14.6	\$ 14.57	\$ 15.09	\$ 15.65	\$ 16.21	\$ 16.87
46	Interest Rate Adjustment	=+IF(\$B\$1 7="yes",1, 0) \$/kW-mo.		\$ 1.20	\$ 1.19	\$ 1.18	\$ 1.17	\$ 1.15	\$ 1.13	\$ 1.11
47	Apache Savings	\$/kW-mo.		\$ -	\$ -	\$ (1.08)	\$ (1.13)	\$ (1.19)	\$ (1.25)	
48	O & M Adjustment	\$/kW-mo.		\$ (1.96)	\$ (2.12)	\$ (3.94)	\$ (2.24)	\$ (3.42)	\$ (5.95)	\$ (2.65)
49	Adjusted Capacity Rate	\$/kW-mo.		\$ 13.11	\$ 11.81	\$ 12.94	\$ 12.25	\$ 10.20	\$ 14.08	
50	Months in the Year			12	12	12	12	12	12	12
51	Plant Capacity	kW		263,000	263,000	263,000	263,000	263,000	263,000	263,000
52	Yearly Capacity Payments			\$40,510,657	\$41,370,159	\$37,268,061	\$40,824,319	\$38,657,491	\$32,191,827	\$44,437,196
53	Monthly Capacity Payments			\$3,375,888	\$3,447,513	\$3,105,672	\$3,402,027	\$3,221,458	\$2,682,652	\$3,703,100

Figure 10: This task requires deriving the XNPV5 of the contract under different combinations of assumptions. The analysis uses contract capacity rates, plant capacity, and the specified discount rate provided in the table. While key adjustment components—namely the Interest Rate Adjustment, Apache Savings, and O&M Adjustment—must be retrieved from supporting documents and applied according to each scenario (included, excluded). For each assumption set, the analyst must then construct the annual capacity payment cash flows by deriving the adjusted capacity rate, converting it into monthly and annual capacity payments, and assembling the full month-by-month cash-flow schedule. Only after these intermediate steps are completed can the cash flows be discounted to the valuation date (e.g., December 31, 2000) to compute XNPV5.

Financial Data Schedule - Revenue & Expenses					Financial Indicator Components						
Account Description	Low Rent Public Housing	Housing Choice Voucher	CFP	TOTAL	1	2	3	4	5	6	
	14,850a	14,871	14,872		Current Ratio	Number of Months Expendable Fund Balance	Tenant Receivable Outstanding	Occupancy Loss	Expense mgmt/ Utility Consumption,.... [for Low Rent Program Only - CFDA# 14.850a]	Net Income (Loss)	
<b>REVENUE:</b>											
703 Net tenant rental revenue		70,020		70,020							
704 Tenant revenue - other											
705 Total tenant revenue		70,020		70,020			70,020	B			
706 HUD PHA operating grants	106,447	211,473	62,422	380,342							
706.1 Capital Grants			53,752	53,752						53,752 A	
<b>Equity Roll Forward Test:</b>											
128 Calculation from R/E Statement	139,662	7,905	53,752	201,319							
129 B'S Line 513	139,662	7,905	53,752	201,319							
130 Difference											
					SUM OF A: (numerator)	145,026	83,606	541	1,178	37,692	95,079
					SUM OF B: (denominator)	62,020	235,221	70,020	1,224	1,178	83,606

Figure 11: The sum of A&B and the equity roll-forward test requires cross-sheet retrieval and calculation.

**Exposition liée aux comptes commerciaux et souverains**

(en millions de dollars canadiens)

	31 mars 2024			31 mars 2023		
	Prêts commerciaux	Prêts souverains	Total	Prêts commerciaux	Prêts souverains	Total
<b>Prêts</b>						
Concessionnels – CUEC	8 508	-	8 508	40 153	-	40 153
Concessionnels	11	422	433	11	455	466
Non concessionnels	425	16 992	17 417	130	16 252	16 382
	8 944	17 414	26 358	40 294	16 707	57 001
<b>Engagements de financement et passifs éventuels</b>						
Engagements de prêts	941	2 273	3 214	1 007	3 039	4 046
Garanties de prêts	-	18 500	18 500	-	11 500	11 500
	941	20 773	21 714	1 007	14 539	15 546
<b>Total</b>	<b>9 885 \$</b>	<b>38 187 \$</b>	<b>48 072 \$</b>	<b>41 301 \$</b>	<b>31 246 \$</b>	<b>72 547 \$</b>
Pourcentage	21 %	79 %	100 %	57 %	43 %	100 %

**Exposition liée aux comptes commerciaux et souverains, par industrie et par pays**

(en millions de dollars canadiens)

	31 mars 2024		31 mars 2023	
	Total	%	Total	%
<b>Prêts commerciaux :</b>				
CUEC (diverses industries)	8 508	18	40 153	55
Services publics	1 131	3	1 000	2
Fabrication	169	-	-	-
Information	39	-	35	-
Autres	38	-	113	-
	9 885	21	41 301	57
<b>Prêts souverains :</b>				
Canada	37 670	78	30 670	42
Chine	255	1	279	1
Turquie	68	-	72	-
Moroc	51	-	54	-
Irak	46	-	58	-
Inde	31	-	32	-
Autres	66	-	81	-
	38 187	79	31 246	43
<b>Total</b>	<b>48 072 \$</b>	<b>100</b>	<b>72 547 \$</b>	<b>100</b>

La baisse de l'exposition relative aux comptes commerciaux s'explique surtout par la diminution des prêts du CUEC. Quant à l'exposition relative aux prêts souverains, elle a augmenté en raison surtout de l'augmentation des garanties de prêts pour l'oléoduc Trans Mountain.

**Commercial and Sovereign Exposure**

(in millions of Canadian dollars)

	Mar 2024		Mar 2023		Total
	Commercial	Sovereign	Commercial	Sovereign	
<b>Loans receivable</b>					
Concessional - CEBA	8,508	-	8,508	-	40,153
Concessional	11	422	433	11	466
Non-concessional	425	16,992	17,417	130	16,252
	8,944	17,414	26,358	40,294	57,001
<b>Financing commitments and contingent liabilities</b>					
Loan commitments	941	2,273	3,214	1,007	3,039
Loan guarantees	-	18,500	18,500	-	11,500
	941	20,773	21,714	1,007	14,539
<b>Total</b>	<b>9,885</b>	<b>38,187</b>	<b>48,072</b>	<b>41,301</b>	<b>72,547</b>
Percentage	21%	79%	100%	57%	43%

**Commercial and Sovereign Exposure by Industry and Country**

(in millions of Canadian dollars)

	Mar 2024		Mar 2023	
	Total	%	Total	%
<b>Commercial:</b>				
CEBA (various)	8,508	18	40,153	55
Utilities	1,131	3	1,000	2
Manufacturing	169	-	-	-
Information	39	-	35	-
Other	38	-	113	-
	9,885	21	41,301	57
<b>Sovereign:</b>				
Canada	37,670	78	30,670	42
China	255	1	279	1
Turkey	68	-	72	-
Morocco	51	-	54	-
Iraq	46	-	58	-
India	31	-	32	-
Other	66	-	81	-
	38,187	79	31,246	43
<b>Total</b>	<b>48,072</b>	<b>100</b>	<b>72,547</b>	<b>100</b>

The decrease in commercial exposure was primarily due to the decrease in loans receivable for the CEBA program. Sovereign exposure increased mainly as a result of the increase in the TMP loan guarantee.

Figure 12: A workflow that translates a French report into English while preserving its format and structure. The report contains many tables to translate, along with text, notes, and even charts.

	A	B	C	D	E	F	G
1							
2	TRANSWESTERN PIPELINE - SUMMARY OF OBA BALANCES					Index	
3						SJ	\$2.09
4	Positive=due Transwestern		Negative = due operator			AVG	\$2.08
5						NTXPH	\$2.08
6							
7	<b>Operator</b>	<b>Dollars</b>	<b>Volume</b>	<b>Date</b>	<b>Imbal Type</b>	<b>Mktg Rep</b>	<b>MS rep</b>
8	PNM	\$879,575	422,873	2/19	Dollar Valued		
9	Conoco	\$484,069	223,110	2/19	Dollar Valued		
10	Mojave Pipeline	\$360,739	173,432	2/19	Volumetric		
11	OneOk Westex-Ward	\$328,563	157,963	2/18	Dollar Valued		
12	Mewborne	\$326,518	156,980	2/18	Dollar Valued		
13	NGPL	\$196,353	89,593	2/19	Volumetric		
14	Dominion Gas Ventures	\$172,975	83,161	2/19	Dollar Valued		
15	Amoco Abo	\$167,604	80,579	2/18	Dollar Valued		
16	SoCal	\$166,015	79,815	2/19	Volumetric		
17	El Paso Field Services	\$143,001	68,750	2/19	Dollar Valued		
18	Red Cedar	\$129,691	62,053	2/19	Volumetric		
19	Agave	\$108,157	51,999	2/19	Dollar Valued		
20	Amarillo Nat Gas	\$102,694	49,372	2/17	Dollar Valued		
21	Citizens-Griffith	\$96,384	46,339	2/19	Dollar Valued		
22	Lonestar	\$87,495	42,065	2/18	Volumetric		
23	PG&E Topock	\$87,416	42,027	2/19	Volumetric		
24	Plains Gas Farmers Co-Op	\$63,242	30,405	1/31	Dollar Valued		
25	Calpine	\$50,583	24,319	2/18	Dollar Valued		
26	Panhandle Eastern	\$49,402	23,751	2/18	Volumetric		
27	Stalland Exploration	\$48,490	23,313	1/31	Dollar Valued		
28	El Paso	\$48,144	23,343	2/19	Volumetric		
29	Continental	\$46,769	22,485	2/18	Dollar Valued		
30	Receivable imbalances	\$4,257,409	2,046,730				
31							
32	<b>Operator</b>	<b>Dollars</b>	<b>Volume</b>	<b>Date</b>	<b>Imbal Type</b>	<b>MS rep</b>	<b>MS rep</b>
33	Citizens Communications	(\$563,447)	(270,888)	2/17	Dollar Valued		
34	North Star Steel	(\$269,783)	(129,703)	2/18	Dollar Valued		
35	MaVida/Richardson Gas Treating	(\$192,286)	(92,445)	1/31	Dollar Valued		
36	Crosstex Energy Serv	(\$134,414)	(64,622)	2/17	Dollar Valued		
37	Duke Energy Field Services	(\$128,990)	(62,014)	2/17	Dollar Valued		
38	Burlington	(\$56,909)	(27,229)	2/18	Dollar Valued		
39	SW Gas Transmission	(\$27,828)	(13,379)	2/18	Dollar Valued		
40	PG&E Topock	(\$27,828)	(13,379)	2/18	Dollar Valued		
41							
42	by type_area	area info	summary	williams	Lonestar	PG&E	SoCal

	A	B	C	D	E	F
1						
2	TRANSWESTERN PIPELINE - SUMMARY OF OBA BALANCES					Index
3						SJ
4	Positive=due Transwestern		Negative = due operator			AVG
5						NTXPH
6						
7	<b>DOLLAR VALUED IMBALANCES</b>					
8		<b>Prod Month</b>	<b>Volume</b>	<b>Accum Prod</b>	<b>As of</b>	
9	<b>Operator</b>	<b>\$ Value</b>	<b>Equivalent</b>	<b>Mo Volume</b>	<b>Date</b>	
10						
11	<b>West of Thoreau</b>					
12	Calpine	\$50,583	24,319	111,418	2/18	
13	North Star Steel	(\$269,783)	(129,703)	(3,264)	2/18	
14	Citizens Communications	(\$563,447)	(270,888)	(49,205)	2/17	
15	<b>Total WOT</b>	(\$714,091)	(343,313)	108,546		
16						
17	<b>San Juan</b>					
18	TransColorado	(\$374)	(180)	(56,126)	2/18	
19	Williams Field Services	(\$18,950)	(9,067)	(9,067)	2/19	
20	Burlington	(\$56,909)	(27,229)	(27,371)	2/18	
21	<b>Total SJ</b>	(\$76,234)	(36,476)	(92,564)		
22						
23	<b>Total \$ Value</b>	\$1,666,771	801,507	554,989		
24						
25						
26						
27	<b>VOLUMETRIC IMBALANCES</b>					
28		<b>Prod Mo</b>	<b>Value @curr</b>	<b>Accum Prod</b>	<b>As of</b>	
29	<b>Operator</b>	<b>Volume</b>	<b>Mo prices</b>	<b>Mo Value</b>	<b>Date</b>	
30						
31	<b>West of Thoreau</b>					
32	Mojave Pipeline	173,432	\$360,739	\$171,960	2/19	
33	SoCal	79,815	\$166,015	\$280,738	2/19	
34	El Paso - Window Rock	64,269	\$133,680	(\$1,582,961)	2/19	
35	PG&E Topock	42,027	\$87,416	(\$115,995)	2/19	
36						
37	by type_area	area info	summary	williams	Lonestar	PG&E
38						SoCal
39						PG&E
40						SoCal
41						PG&E
42						SoCal
43						PG&E
44						SoCal
45						PG&E
46						SoCal
47						PG&E
48						SoCal
49						PG&E
50						SoCal
51						PG&E
52						SoCal
53						PG&E
54						SoCal
55						PG&E
56						SoCal
57						PG&E
58						SoCal
59						PG&E
60						SoCal
61						PG&E
62						SoCal
63						PG&E
64						SoCal
65						PG&E
66						SoCal
67						PG&E
68						SoCal
69						PG&E
70						SoCal
71						PG&E
72						SoCal
73						PG&E
74						SoCal
75						PG&E
76						SoCal
77						PG&E
78						SoCal
79						PG&E
80						SoCal
81						PG&E
82						SoCal
83						PG&E
84						SoCal
85						PG&E
86						SoCal
87						PG&E
88						SoCal
89						PG&E
90						SoCal
91						PG&E

	A	B	C	D	E	F	G	H	I
			IF NGPL	IF NGPL	IF NGPL	IF CIG Rocky	PG&E Topock	Gas Daily EI	Gas Daily
			MidContinent	MidContinent	MidContinent	Mtns. index	index minus	Paso- San	NWPL
			index (@	index (@	Index minus	minus \$0.03	\$0.02	Juan index	Wyoming Pool
			Forgan)	Baker)	\$0.01	(Proposed)	(Proposed)	minus \$0.10	index minus
									\$0.10
2									(Proposed)
3	31	Dec-01	155,147	155,147	-	1,387,125	898,812	480,067	431,317
4	31	Jan-02	370,351	370,351	-	3,536,910	2,117,920	1,149,369	1,058,919
5	28	Feb-02	401,692	401,692	-	4,112,550	2,257,575	1,258,026	1,180,826
6	31	Mar-02	380,019	380,019	-	3,880,530	2,119,363	1,199,536	1,109,736
7	30	Apr-02	392,883	392,883	-	3,843,585	2,167,842	1,235,554	1,128,504
8	31	May-02	386,271	386,271	-	3,596,625	2,121,738	1,200,934	1,088,184
9	30	Jun-02	406,844	406,844	1,702,339	2,057,625	2,295,482	684,037	620,487
10	31	Jul-02	415,605	415,605	3,807,843	-	2,441,730	-	-
11	31	Aug-02	437,663	437,663	3,995,353	-	2,606,812	-	-
12	30	Sep-02	445,239	445,239	4,059,253	-	2,615,090	-	-
13	31	Oct-02	444,413	444,413	4,047,111	-	2,548,440	-	-

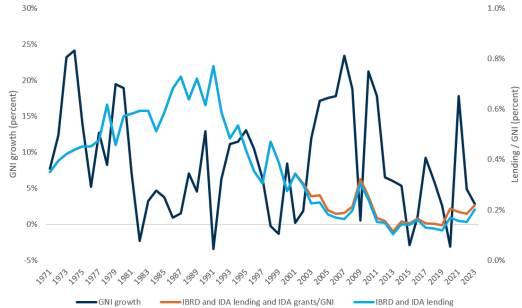
Figure 14: The apparent semantics from the headers suggest a monthly/daily exposure metric. However, inspecting the underlying formula (e.g., C5=\$A4\*Volumes!B6\*Curves!G7+25\*Volumes!B7\*Curves!G8) reveals that it actually encodes a 55-day payment timing schedule. Models that ignore or underutilize formula information, therefore systematically misattribute the column's role in downstream computations, and this misinterpretation then propagates through subsequent steps.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	<b>SUMMARY OF CANADA'S TRADING INCOME BY TRADER - IN US\$</b>														
2															
3	<b>FX - AVG OF MNTH</b>	1.4511	1.4493	1.4607	1.4677	1.4943	1.4761	1.4778	1.4822	1.4833	1.5108	1.5429	1.5396		
4													To Dec 6	1.48	
5															
6															
7	<b>BY BOOK:</b>														
8	<b>TERM:</b>	<i>(Note: Q1 Origination has been manually backed out)</i>													
9	Alberta Term	3,458,525	5,544,777	2,686,484	469,449	4,631,088	(1,097,125)	(5,566,310)	5,913,482	777,220	(3,599,231)	(7,518,121)	2,503,740	8,203,978	
10	BC Term	(879,015)	8,011	902,041	710,889						126,089	(\$5,525,759)	\$7,648,004	2,990,260	
11	EOL - Term				7,208	936,819	(1,582,742)	(512,592)	471,676	2,538,389	1,319,205	(298,040)	2,011,563	4,891,486	
12	Options	(224,086)	685,015	(711,171)	=SummaryE10/SUM-USD!E3			(1,140,535)	(2,145,933)	(5,898,385)	(\$1,050,467)	2,645,718	(1,111,159)	(2,863,509)	
13	<b>CASH</b>														
14	Alberta Cash	(1,649,853)	692,072	785,943	703,120	2,076,946	2,115,658	1,174,214	1,353,510	503,658	743,525	766,044	(162,227)	9,102,609	
15	BC Cash	(143,902)	176,401	(\$19,854)	50,999	73,039	317,584	295,608	327,887	673,449	522,897	341,336	238,176	2,853,619	
16	BC Pipe Cash													3,278,295	
17	Alberta Term - GD	(242,131)	300,898	30,625	88,895	398	(142)	150	(365,403)	(962,508)	124,583	(560,897)	(4,123)	(1,589,854)	
18	BC Term - GD	(615,175)	21,265	(4)										(593,915)	
19	Options - GD	286,232	517,629	348,138	424,851	354,022	1,007,495	853,328	555,567	1,935,097	(175,224)	(1,554,150)	2,667,285	7,220,271	
20	Power	246,906	183,840	1,261,789	1,453,448	1,254,449	(998,132)	(416,028)	(1,002,687)	1,471,034	(136,978)	47,101		3,364,742	
21	PMA						(812,960)							(812,960)	
22	<b>TOTAL CANADA</b>	237,501	8,129,908	5,283,991	4,248,065	13,745,994	278,691	(5,312,165)	5,108,098	1,037,955	(2,511,333)	(8,133,204)	13,931,724	36,045,223	
23															
24															
25	<b>BY RISK TYPE:</b>														
26	Total Term	2,355,424	6,237,804	2,877,354	1,526,752	9,987,140	(1,350,812)	(7,219,438)	4,239,225	(2,582,776)	(3,204,404)	(10,696,202)	11,052,148	13,222,216	
27	Check														
28	Total Cash	(2,117,923)	1,892,104	2,406,637	2,721,313	3,758,853	1,629,502	1,907,272	868,873	3,620,730	693,071	2,562,998	2,879,576	22,823,007	
29	Check														
30	<b>TOTAL CANADA</b>	237,501	8,129,908	5,283,991	4,248,065	13,745,994	278,691	(5,312,165)	5,108,098	1,037,955	(2,511,333)	(8,133,204)	13,931,724	36,045,223	
31															
32															
33	<b>BY AREA/TRADER:</b>														
34	West Term - Lavorato	3,216,394	5,845,675	2,717,109										11,779,178	
35	West Term - Mckay	(1,638,092)	205,677	902,037	1,269,232	4,631,487	(1,097,266)	(5,566,160)	471,676	2,538,389	1,059,563	(2,300,237)	9,800,031	10,276,335	
36	West Term - Lambie				7,208	936,819	(1,582,742)	(512,592)	5,548,079	(185,288)	(3,474,648)	(8,079,018)	2,499,617	(4,842,564)	
37	Options - Disturnal	62,146	1,202,645	(363,033)	764,058	4,773,255	2,336,550	(287,207)	(1,590,366)	(3,963,287)	(1,225,692)	1,091,569	1,556,126	4,356,762	
38	Alberta Cash - Cowan	(1,649,853)	692,072	785,943	703,120	2,076,946	1,302,698	1,174,214	1,353,510	503,658	743,525	766,044	(162,227)	8,289,648	
39	BC Cash - Clark			(19,854)	50,999	73,039	317,584	295,608	327,887	673,449	522,897	341,336	238,176	2,821,121	
40	Power - Greenizan	246,906	183,840	1,261,789	1,453,448	1,254,449	(998,132)	(416,028)	(1,002,687)	1,471,034	(136,978)	47,101		3,364,742	
41	<b>TOTAL CANADA</b>	(2,978,893)	2,284,233	2,566,882	4,248,065	13,745,994	278,691	(5,312,165)	5,108,098	1,037,955	(2,511,333)	(8,133,204)	13,931,724	36,045,223	

Figure 15: This workflow requires creating a new spreadsheet with all values converted to USD. It also requires correct in-sheet and cross-sheet formula references while preserving the original spreadsheet layout.

Figure 1.4 GNI Growth Versus Ratio of New World Bank Lending to Gross National Income, Low- and Middle-Income Countries, 1971-2023

Percent	1971	1972	1973	1974	1975	1976	1977	1978	1979	1980	1981	1982	1983	1984	1985	1986	1987
GNI growth	8%	12%	23%	24%	14%	5%	13%	8%	19%	19%	7%	-2%	3%	5%	4%	1%	2%
IBRD and IDA lending and IDA grants/GNI	0.35%	0.40%	0.42%	0.44%	0.45%	0.45%	0.48%	0.62%	0.46%	0.57%	0.58%	0.59%	0.59%	0.51%	0.59%	0.68%	0.73%
IBRD and IDA lending	0.35%	0.40%	0.42%	0.44%	0.45%	0.45%	0.48%	0.62%	0.46%	0.57%	0.58%	0.59%	0.59%	0.51%	0.59%	0.68%	0.73%



but also implicit ex ante debt relief and financial support. Most IDA credits carry a zero or very low interest rate, and repayments typically extend over 30–50 years; however, more than one-third of IDA-eligible countries receive all or part of their IDA resources in the form of grants that carry no repayments in the future. Whereas IDA focuses on the most impoverished nations, the World Bank's other lending arm, IBRD, has played a crucial role in coordinating responses to regional and global challenges by providing loans and financial services to middle-income and creditworthy low-income countries (figure 1.4). IBRD was created to support countries rebuilding after World War II and has continued its crisis and emergency support through increased lending to countries affected by other crises since then, including the 2008–09 financial crisis, the 2014 Ebola outbreak, and the COVID-19 pandemic. Since inception, IBRD and IDA lending has responded positively to adverse external shocks affecting the economies of countries eligible for such financing, and this countercyclical lending has been a recurring and stabilizing response to dramatic drops in economic growth in these economies over the years.

Figure 16: Generating reports from tabular data requires financial knowledge of data analysis, financial events, and visualization. For example, one may plot two series with different units on a single chart (e.g., using a secondary y-axis) to reveal their correlation.

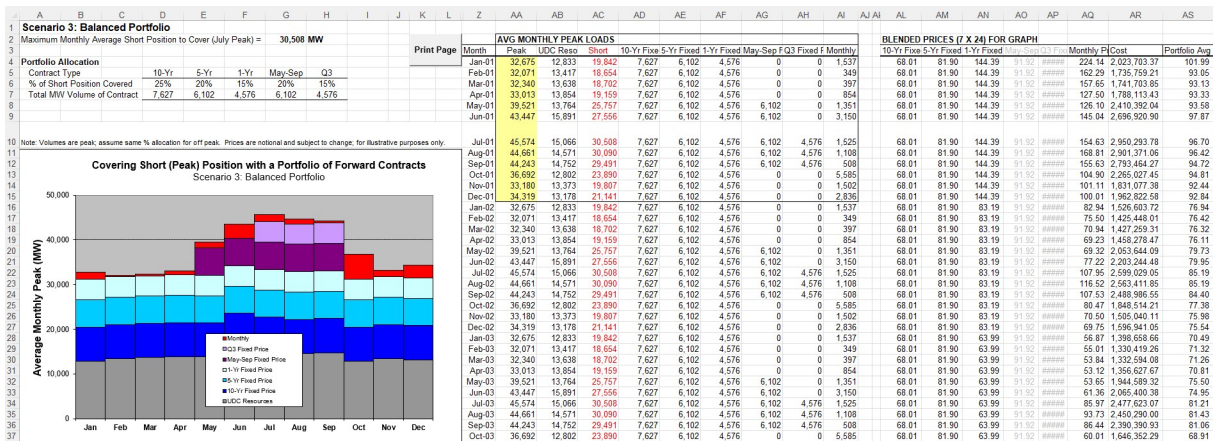


Figure 17: This Excel sheet shows an assumption-update workflow, where a mix of forward contracts is used to cover monthly peak-load short positions. It lists the contract allocations and MW volumes, along with monthly peak loads and the resulting short MW. A table on the right computes blended prices and portfolio costs, and the stacked chart visualizes coverage by contract type over the year.