A APPENDIX

649 650 651

648

A.1 DETAILS OF THE BREWALGORITHM

652 653

654

655

656

657

658

659 660

661

662 663

664 665

666

667

668 669

670

671

672

673

674

675

676

677 678

679 680

681

682

683

684

685

686

687

688

689

690

691 692

693

694

695

696

697

698

699

700

We provide pseudocode for the core components of BREW, aligning with the stages introduced in Section 3. Each algorithm plays a distinct role in constructing, organizing, or refining the knowledge base over iterative interactions. GenerateInsights (Alg. 2) produces concept-aligned insights from annotated rollouts using Reflagent. DeduplicateConcepts (Alg. 3) clusters semantically overlapping concepts into a compact meta-concept set. Integagent incrementally builds and updates per-concept documents using newly generated insights. Finally, Expanding (Alg. 4) performs MCTS-guided expansions to explore improved document variants, while Evaluate (Alg. 5) scores candidate KB states using correctness and retrieval-based rewards.

We specify the IntegAgent prompt below:

BREW Integrator Prompt

```
# Enhanced Documentation Editor Prompt
You are a meticulous documentation-level editor specializing in
   comprehensive technical reference materials. You will be given a
   list of topic nodes, each containing structured information that
   must be preserved and enhanced with maximum detail retention.
## Input Structure Analysis
Each node contains:
 **Title**: The primary topic identifier
 **Context**: Background information and conceptual foundation
- **How to Use**: Step-by-step instructions, commands, flags,
   parameters, and implementation details
 **When to Use**: Specific scenarios, conditions, and decision
   criteria
- **Best Practices**: Expert recommendations, optimization techniques,
    and common pitfalls to avoid
## Detailed Processing Requirements
### 1. Information Preservation (Zero Loss Policy)
 **Preserve every technical detail**: All command-line flags,
   parameter values, configuration options, file paths, URLs, version
    numbers, and exact syntax
- **Maintain all examples**: Keep every code snippet, sample input/
   output, file names, directory structures, and command sequences
   exactly as provided
- **Retain contextual nuances**: Preserve qualifying language like "
   typically, " "usually, " "in most cases, " "when available, " and
   conditional statements
 **Keep quantitative data**: Preserve all numbers, measurements,
   timeframes, limits, thresholds, and statistical information
  **Maintain cross-references**: Keep all mentions of related tools,
   dependencies, prerequisites, and interconnected concepts
### 2. Enhanced Detail Extraction
  **Expand abbreviations**: When encountering shortened forms, expand
   them naturally while preserving the original
- **Surface implicit knowledge**: Make obvious assumptions explicit (e
    .q., "this requires root permissions," "assumes default
   configuration")
 **Clarify relationships**: Explicitly describe how different
   components, options, or steps relate to each other
 **Highlight edge cases**: Emphasize special conditions, exceptions,
   or unusual scenarios mentioned in the source
  **Elaborate on consequences**: When the source mentions outcomes,
   expand on both success and failure scenarios
```

```
702
703
          ### 3. Prose Transformation Guidelines
704
          - **Bullet integration**: Transform each bullet point into 1-3
705
              complete sentences that naturally flow together
          - **Technical precision**: Use precise technical vocabulary while
706
              maintaining readability
707
          - **Logical flow**: Organize information within each section to follow
708
               a logical sequence (setup \rightarrowexecution \rightarrowverification)
709
            **Contextual embedding**: Weave code snippets and technical terms
710
              seamlessly into narrative sentences
          - **Comprehensive coverage**: Ensure every sub-bullet, nested item,
711
              and parenthetical note becomes part of the prose
712
713
          ### 4. Structural Requirements
714
          - **Heading hierarchy**: Use `# Title` for each node's main heading
          - **Section order**: Maintain Context \rightarrowHow to Use \rightarrowWhen to Use \rightarrowBest
715
               Practices sequence
716
          - **Paragraph organization**: Create substantial paragraphs (3-6
717
              sentences) rather than brief statements
718
          - **Transition quality**: Craft smooth bridges between sections and
719
              between different nodes
720
          - **Code formatting**: Preserve all inline code with backticks and
              maintain proper formatting for code blocks
721
722
          ### 5. Quality Assurance Checklist
723
          Before finalizing, verify:
724
          - [ ] Every piece of source information appears in the output
            [ ] All technical specifications, parameters, and examples are
725
              int.act.
726
            [ ] Code snippets maintain their exact syntax and formatting
727
          - [ ] Prose flows naturally without choppy or fragmented sentences
728
          - [ ] Each section provides comprehensive coverage of its topic area
729
          - [ ] Cross-references and dependencies are clearly explained
          - [ ] No section labels or formatting artifacts remain in the prose
730
731
          ## Output Specifications
732
          Generate a single, cohesive markdown document that reads as
733
              authoritative technical documentation. The result should be
734
              comprehensive enough that a reader could successfully implement
              the described tools or techniques using only the information
735
              provided, without referring back to the original nodes.
736
737
738
739
          **Input Nodes: **
          <NODES>
740
          {node list}
741
          </NODES>
742
743
744
          Now, produce the aggregated markdown reference sheet with maximum
745
              detail preservation and enhanced clarity.
746
```

Algorithm 2 GenerateInsights: Extract behavioral insights from trajectories

756

757

758

759

760

761

762

763 764

765

773

774

775

776

777

778

779

780

781

782

790

791

792

793

794

796

797

798

799

800

801

802

803

804

805

806

807

808

809

19: **end for**

```
Require: Queries Q, KB \mathcal{D}, rubrics
Ensure: Concept-insight pairs \mathcal{B}
 1: \mathcal{B} \leftarrow \emptyset
 2: for each query q \in \mathcal{Q} do
 3:
           \tau \leftarrow \text{LLM}(q, \mathcal{D})

⊳ Generate trajectory

 4:
           label \leftarrow Grade(\tau)
                                                                                                                          (c, i) \leftarrow \text{ReflAGENT}(\tau, \text{rubrics}, \text{label})
 5:
 6:
           \mathcal{B} \leftarrow \mathcal{B} \cup \{(c,i,q)\}

    Store with source query

 7: end for
 8: return \mathcal{B}
```

Algorithm 3 DeduplicateConcepts: Cluster similar concepts and map queries

```
Require: Concept-insight-query triples \mathcal{B}
Ensure: Meta-concepts \mathcal{K} with mapped queries and insights

1: Extract all concepts from \mathcal{B}

2: Embed and cluster concepts by similarity

3: \mathcal{K} \leftarrow cluster representatives

4: for each k \in \mathcal{K} do

5: \mathcal{Q}_k^{\text{train}} \leftarrow {training queries that contributed insights to k}

6: \mathcal{Q}_k^{\text{eval}} \leftarrow {held-out queries relevant to k}

7: \mathcal{I}_k \leftarrow {all insights mapped to concept k}

8: end for

9: return \mathcal{K} with associated queries and insights
```

Algorithm 4 ExpandNode: Generate and evaluate new document variants

```
Require: Node s, concept k, candidates h, current KB \mathcal{D}_{current}, best docs \mathcal{D}_{best}, tree
Ensure: Updated tree with new evaluated nodes

    □ Generate new insights from concept-relevant queries

 1:
 2: \mathcal{B}_{\text{new}} \leftarrow \varnothing
 3: for query q \in \mathcal{Q}_k^{\text{train}} do
           \tau \leftarrow \text{LLM}(q, \mathcal{D}_{\text{current}})
 5:
           (c, i) \leftarrow \text{ANNOTATE}(\tau, \text{rubrics}, \cdot)
 6:
           if c maps to k then
                 \mathcal{B}_{\text{new}} \leftarrow \mathcal{B}_{\text{new}} \cup \{i\}
 7:
 8:
           end if
 9: end for
10:
                                                                              11: for j = 1 to h do
12:
           d_{k,i} \leftarrow \text{INTEGAGENT}(k, \mathcal{I}_k \cup \mathcal{B}_{\text{new}}, d_k^s)
                                                    ⊳ Evaluate using hybrid KB with best docs from other concepts
13:
14:
           \mathcal{D}_{\text{hybrid}} \leftarrow \{d_{k,j}\} \cup \{d_{k'} \in \mathcal{D}_{\text{best}} : k' \neq k\}
15:
           R_{k,j} \leftarrow \text{EVALUATE}(d_{k,j}, \mathcal{D}_{\text{hybrid}}, \mathcal{Q}_k^{\text{eval}})
16:
                                                                                                  ▶ Add to tree and backpropagate
17:
           Add (d_{k,j}, R_{k,j}) as child of s in tree
18:
           Backpropagate R_{k,j} from new node to root
```

Algorithm 5 Evaluate: Score document using held-out queries

```
811
                     Require: Document d_k, hybrid KB \mathcal{D}_{hybrid}, eval queries \mathcal{Q}_k^{eval}
812
                      Ensure: Reward score R
813
                        1: R^{\text{corr}} \leftarrow 0
814
                        2: R^{\text{ret}} \leftarrow 0
                        3: for each q \in \mathcal{Q}_k^{\text{eval}} do
4: R^{\text{corr}} \leftarrow R^{\text{corr}} + \text{EVAL}(q, \underset{-}{\text{agent}} \oplus \mathcal{D}_{\text{hybrid}})
815
816
                                         R^{\text{ret}} \leftarrow R^{\text{ret}} + \text{MRR}(d_k, q, \mathcal{D}_{\text{hybrid}})
817
818
                        6: end for
                        7: R^{\text{corr}} \leftarrow \frac{R^{\text{corr}}}{|\mathcal{Q}_k^{\text{eval}}|}
8: R^{\text{ret}} \leftarrow \frac{R^{\text{ret}}}{|\mathcal{Q}_k^{\text{eval}}|}
819
820
821
                        9: return \lambda_{\text{corr}} \cdot R^{\text{corr}} + \lambda_{\text{ret}} \cdot R^{\text{ret}}
822
```

A.2 BREW CONFIGURATIONS

Base LLM Configuration For all BREWalgorithm steps, we use the OpenAI GPT-4.1-2025-04-14 model as the underlying language model. To balance exploration and stability, we set the temperature to 0.7 for the IntegAgent component to encourage diversity in sampled completions, while all other calls use a temperature of 0.1 for deterministic behavior. The search process employs an expansion width of e=3, a maximum search depth of k=3, and a maximum of n=10 iterations. Reward signals are weighted equally across correctness and retrieval relevance, with $\lambda_{corr}=\lambda_{ret}=0.5$.

A.3 BASELINE METHODS

We compare BREWagainst two common reasoning baselines. Step-Back Prompting encourages backward reasoning by guiding the model to work from the final task objective back to the initial actions. In-Context Learning augments the input prompt with successful trajectories from related tasks, enabling the model to benefit from relevant prior examples without additional fine-tuning.

A.4 BENCHMARK SPECIFICATIONS

A.4.1 OSWORLD: COMPUTER-USE AUTOMATION

Dataset Overview OSWorld Xie et al. (2024) comprises 369 real-world computer-use tasks spanning 10 distinct applications. The benchmark is divided into train and test sets, with the distribution of tasks across domains shown in Table 2.

Agent Specifications The UI-Tars-7B variant is a 7B-parameter multimodal transformer fine-tuned for graphical user interface understanding. It operates over an action space of PyAutoGUI commands (e.g., click, type, and key presses). The agent integrates a retrieval module that queries a task-relevant knowledge base using the user-provided description, with the top three retrieved items added to the system prompt. Inputs to the model consist of a screenshot of the active UI paired with the natural language task description.

The GTA1-7B configuration adopts a two-agent architecture, consisting of a planner and a grounding module. The planner (GTA-1-7B) generates the high-level action sequence, while the grounding module (OpenAI O3) verifies and refines each action before execution. Knowledge retrieval is incorporated differently for each component: the planner performs a single retrieval at the start of execution, which is persisted in its prompt, whereas the grounding module performs dynamic retrievals at each verification step.

Evaluation Protocol Evaluation uses 134 task-specific scripts designed for automated verification. Success criteria include file state checks (e.g., validating .xlsx or .docx outputs), UI element validation to confirm correct interaction, and process completion checks to ensure that the intended automation sequence was executed successfully.

A.4.2 au^2 -Bench: Interactive Tool Usage

Dataset Overview τ^2 -Bench Barres et al. (2025b) extends τ -Bench by introducing bidirectional tool-calling capabilities. The dataset covers multiple service-oriented domains, with domain-level task distributions summarized in Table 3.

Domain Characteristics The benchmark spans several domains with distinct task characteristics. The Telecom domain focuses on connectivity troubleshooting, plan modifications, and service activation workflows. The Retail domain includes order processing, return handling, and inventory queries. The Airline domain emphasizes booking modifications and policy-compliant rescheduling scenarios.

Interaction Settings Two interaction modes are defined. In Easy mode, a human proxy (implemented via GPT-4.1) provides detailed guidance to the agent. The knowledge base is built exclusively from Easy mode trajectories, ensuring high-quality demonstrations for learning. In Hard mode, human intervention is minimized. The knowledge base combines both Easy and Hard trajectories, testing the agent's robustness to underspecified or noisy instructions.

Evaluation Criteria Task success is measured using domain-specific verification procedures. These include database state checks to validate final outcomes, status checks for confirming service or connection state, natural language verification to ensure correct confirmation statements appear in dialogue, and action matching to confirm that all required steps are completed. Each domain uses a tailored subset of these checks (e.g., Telecom relies primarily on status checks).

Domain	Test	Train
Calc	45	2
Chrome	44	2
Writer	21	2
Gimp	24	2
Impress	45	2
Os	22	2
Thunderbird	13	2
Multi-apps	99	2
VLC	15	2
VSCode	21	2
Total	349	20

Table 2: Test and Train samples across different domains in OSWorld.

Domain	Test	Train
Telecom	105	7
Retail	105	7
Airline	44	6
Total	254	20

Table 3: Task-wise breakdown for τ^2 -Bench with assumed 2-shot training samples per domain.

Domain Characteristics

- **Telecom:** Connectivity issues, plan management, service activation
- Retail: Order processing, returns, inventory queries
- Airline: Booking modifications, policy-compliant rescheduling

Evaluation Criteria Task success determined by:

• Database Checks: Final state verification

• Status Checks: Service/connection state validation

• NL Checks: Confirmation statements in dialogue

Action Matching: Required action sequence completion

921 922

923

Note: Each domain uses specific check combinations (e.g., Telecom uses only status checks).

924 925

926

SPREADSHEETBENCH: REAL-WORLD SPREADSHEET MANIPULATION

Dataset Overview SpreadsheetBench Ma et al. (2024) consists of 912 instructions collected from four major Excel forums and blogs. Each instruction is paired with spreadsheets reflecting authentic, complex user scenarios, often containing multiple tables and non-standard relational structures. The dataset totals 2,729 test cases, averaging three per instruction. A breakdown of cell-level and sheet-level manipulations is shown in Table 4.

931 932

Task Settings The benchmark defines two dimensions of evaluation:

933 934 935

• Granularity: Instructions involve either *cell-level* manipulations (specific ranges such as D2:D6) or *sheet-level* manipulations (entire tables or multi-sheet updates).

936 937

• Evaluation: Performance is measured using an Online Judge (OJ)-style protocol. The soft setting (IOI-style) awards partial credit when only some test cases are solved, while the hard setting (ICPC-style) requires solutions to succeed on all test cases.

938 939

940

941

942

943

944

Agent Configuration We evaluate

texttto4-mini using a function-calling agent connected to a single Python execution tool. The agent translates natural language instructions into Python code for spreadsheet manipulation (e.g., modifying cells, applying formulas, restructuring tables). After each tool call, all formulas in the spreadsheet are recalculated to ensure consistency before proceeding to the next step. This setup provides a controlled environment to assess reasoning, code generation, and execution robustness across diverse spreadsheet tasks.

945 946 947

948 949

Granularity	Instructions	Test Cases
Cell-Level	329	986
Sheet-Level	583	1,743
Total	912	2,729

950 951 952

Table 4: Cell-level vs. sheet-level distribution in SpreadsheetBench.

953 954 955

A.5 KB CONSTRUCTION AND RETRIEVAL DETAILS

956 957

Training Data Collection

958 959

• **OSWorld:** 20 successful trajectories (2 per application domain) and 10 for evals.

960 961

• τ^2 -Bench: 20 trajectories balanced across domains and difficulty settings and 10 for evals.

• SpreadsheetBench: Uniformly sample 30 trajectories for training and 10 for evaluation.

962 963

All numbers are reported on the remaining train set.

964 965

Retrieval Strategy

966 967

• Query Formation: For each task we take in the seed Natural Language query as the retrieval query.

968

• **Retrieval Count:** We take top-3 documents for all the retrieval steps

969 970

• Integration Point: For SPREADSHEET ENCH and OSWorld we insert retrievals in the system prompt augmentation. For τ^2 -bench we add perfrom retrieval after each user interaction.

	Baseline	max_width=3, max_depth=3	max_width=3, max_depth=10	max_width=10, max_depth=3
OSworld	44.20	47.56	43.83	49.32

Table 5: OSworld difference in MCTS parameters

B QUALITATIVE ANALYSIS

972973974975

976 977

978 979 980

981

982

983

984

985

986

987

988 989

990 991

992 993

994 995

996

998

Exploration on MCTS parameters WE evaluate OSworld on two different MCTS parameters.

- Increased Depth: To increase the depth we keep maximum width of the tree as 3 and depth as 10 with max number of iterations as 25. We observe that the Knowledge base over optimizes on the train set leading to a poorer performance on test set.
- Increased Width: For increased width we reverse the parameters where depth is 3 and maximum width is 10 with max iterations 25. We observe many different styles of KBs are generated storing very similar information, these different styles lead to a varied performance on both eval and test set notifying the importance of state search.

We report the numbers on table ??

C EXEMPLAR KNOWLEDGE BASES

C.1 KNOWLEDGE BASE LEARNED FOR OSWORLD

We showcase a small part of knowledge base learned thought BREW. This demonstrate 3 major parts on which each document is aggregated. These parts discuss when to use a piece of information, why to use the information, how to use the information/tool.

```
999
          ## Search and Open Files
1000
          **When to use**: Locating documents, spreadsheets, images, or
1001
              downloads for editing, conversion, or attachment.
1003
          ### How to Perform
1004
          - Open **File Manager (Nautilus) ** from launcher or system dock
          - Press 'Ctrl + F' or click the search icon
1005
          - Enter part of filename, full name, or wildcard ('*.pdf', 'report*')
1006
          - Use right-click →**Open With** to choose the desired application
          - Use the sidebar to navigate to **Downloads**, **Documents**, or
              custom folders
1009
1010
          ### Additional Actions

    Right-click →**Properties** to check modification date or file type

1011
          - Sort results by Date, Type, or Name from the top-right dropdown
1012
          - Use 'F2' to rename files inline
1013
1014
          ### Example
          - Task: "Edit the file titled `sales_report_march.ods`"
1015
            - Search for 'sales' in File Manager
1016
           - Confirm '.ods' type and open with LibreOffice Calc
1017
1018
1019
1020
          ## Insert Images
1021
          **When to use**: Adding visual elements to documents, presentations,
1022
              emails, or templates.
1023
1024
          ### How to Perform
1025
          - Navigate to **Insert →Image →From File** (in Writer, Impress,
              Thunderbird)
```

```
1026
           - Select an image file ('.png', '.jpg', '.svg') from the file dialog
1027
           - Use drag handles to resize; right-click \rightarrow **Wrap** or **Alignment**
1028
               for layout
1029
           ### Additional Actions
1030
           - In GIMP: **File →Open as Layers** to insert image as a new layer
1031
           - Use drag-and-drop from file manager into open document windows
1032
           - Use **Format \rightarrowImage** to apply borders, shadows, or color
1033
               corrections (in Writer/Impress)
           ### Example
1035
           - Task: "Insert the logo.png image into the title slide"
1036
             - Open '.odp' file in Impress \rightarrowGo to Slide 1 \rightarrowInsert \rightarrowImage \rightarrow
1037
                Select 'logo.png'
1038
1039
           . . .
1040
           ## Export as PDF
1041
1042
           **When to use**: Required submission format
1043
1044
           ### How to Perform
           - Go to **File →Export As PDF**
1045
           - Choose output folder (usually **Documents** or **Downloads**)
1046
           - Click **Save**, then confirm the exported file opens correctly
1047
1048
           ### Additional Actions
1049

    In GIMP or Impress: choose **File →Export As**, then select `.pdf`

               from format list
1050
           - Use **Save As** to preserve both editable and exported versions
               separately
1052
1053
           ### Example
           - Task: "Export the flyer.xcf as a PDF"
1054
             - Open in GIMP \rightarrowFile \rightarrowExport As \rightarrowRename to 'flyer.pdf' \rightarrowClick
1055
                 Export
1056
```

C.2 BREW KNOWLEDGE BASE FOR τ^2 -BENCH

1057 1058

10591060

1061

1062

BREW enable use to learn relevant information for tau bench for across the domains in a single knowledge base. This knowledge base is helpful to use relevant actions from the action pool.

```
1063
1064
          ### Additional Actions
1065
          * Inform the user:
1066
            - Refunds via gift card = immediate.
1067
            - Refunds via other methods = -57 business days.
          ### Example
1069
1070
          * Task: "Cancel a T-shirt order placed yesterday"
1071
            * Validate: Status is 'pending'
1072
            * Reason: "no longer needed"
1073
            * Confirm
1074
            * Execute tool call
1075
1076
          # Exchange Delivered Order
1077
1078
          **When to use**:
1079
          User wants to swap delivered items for a different variant (e.g., size
               or color).
```

```
1080
1081
          **Why to use it**:
1082
          To fix sizing or option errors without needing a new purchase.
1083
          ### How to Perform
1084
          - Authenticate user
1085
          - Confirm order status is 'delivered'
1086
          - Get full list of exchange items
1087
          > "Please ensure all items for exchange are listed. This step 'cant be
1088
               repeated."
          - Ask for refund/payment method
1089
          - Confirm:
1090
           > "'Youre exchanging item X for same product, different option.
1091
               Proceed?"
1092
          - On confirmation:
            '''python
1093
           request_exchange(order_id="45678", item_exchanges=[...],
1094
            payment_method="paypal")
1095
1096
1097
          ### Additional Actions
1098
          * Mention: An email will be sent with return instructions
1099
          * Validate that the new variant is from the same product
1100
1101
          ### Example
1102
          * Task: "Exchange red shirt for blue in Order #45678"
1103
           * Confirm all exchange items
1104
            * Confirm payment method for difference
1105
            * Execute tool call
1106
1107
          ### Example
1108
          * Task: "Show me my last 2 orders"
1109
            * Authenticate
1110
            * Retrieve and present info
1111
1112
          # Deny Unsupported Request
1113
          **When to use**:
1114
          User asks for an unsupported action (e.g., cancel processed order,
1115
              exchange to different product type, help another user).
1116
1117
          **Why to use it**:
          To stay compliant with platform policy.
1118
1119
          ### How to Perform
1120
          - Politely reject:
1121
           > "'Im sorry, but I 'cant process that request. 'Its outside the
1122
               allowed scope."
1123
          ### Example
1124
1125
          * Task: "Cancel a processed order"
1126
           * Respond with denial message
1127
          # Transfer to Human Agent
1128
          **When to use**:
1129
          User needs help outside the 'assistants permitted capabilities.
1130
1131
          **Why to use it**:
1132
          To ensure user gets the right help from trained staff.
1133
          ### How to Perform
```

```
1134
          - Make tool call:
1135
            '''python
1136
            transfer_to_human_agents()
1137
          - Then inform user:
1138
            > "YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON."
1139
1140
          ### Example
1141
1142
           * Task: "Delete a task"
            * Deny deletion
1143
            * Transfer to human
1144
```

```
1146
      C.3 BREW Knowledge Base for SpreadsheetBench
1147
             Header Extraction
1148
          1. Detecting Header Rows
1149
          Overview:
1150
          To accurately identify header rows, scan the initial region of your
1151
             dataset. This process is crucial for mapping column information
1152
             for further processing.
1153
          Approaches:
1154
          - Heuristic Checks:
1155
          - Look for rows where all cells are strings (e.g., "Name", "Date", "
1156
             Region", "Amount").
1157
          - Identify rows with distinctive formatting such as bold text or
             background color.
1158
          - Example:
1159
          | Name | Date | Region | Amount | |--
1160
             ----| John | 2024-01-01 | North
1161
              | 100 |
1162
          - Pattern Recognition:
          - Use regex to match typical header patterns, such as column names
1163
             starting with uppercase letters.
1164
          - Score candidate rows based on the likelihood of being headers.
1165
          - Multi-Table Sheets:
1166
          - Detect gaps, empty rows, or separators indicating a new table.
1167
          - Assign a Table ID to each detected table for later reference.
1168
          Edge Cases:
1169
          - Merge multi-row headers (e.g., "Sales" over "2024", "2025" becomes "
1170
             Sales 2024", "Sales 2025").
1171
          - Fill in missing headers by inferring from context.
1172
          2. Assigning and Validating Headers
1173
          Overview:
1174
          Once headers are detected, assign them programmatically and ensure
1175
             they match expected schema and data types.
1176
          Implementation:
1177
          - Column Naming:
1178
          - Set names in code, e.g., df.columns = ["Name", "Date", "Region", "
1179
             Amount"].
1180
          - Schema Mapping:
1181
          - Map headers to a standardized schema, using external files or user
1182
             prompts.
          - Example:
1183
          - Raw header: "Amt"; Mapped header: "Amount"
1184
          - Quality Checks:
1185
          - Detect duplicate or empty headers ("Date", "Date" becomes "Date_1",
1186
             "Date_2").
1187
          - Validate each column's expected data type.
```

```
1188
          3. Automation and Usability Enhancements
1189
          Overview:
1190
          Enhance usability and automation to streamline header extraction and
1191
              user interaction.
1192
          Features:
1193
          - Freeze Panes:
1194
          - Automatically freeze header rows in Excel for easier navigation.
1195
          - Highlighting:
1196
          - Use colored formatting to visually distinguish headers.
          - Example:
1197
          - Yellow fill for header row.
1198
          - Documentation:
1199
          - Log extraction logic and confidence scores for each detected header.
1200
          - Integration:
          - Build header extraction into ETL pipelines and record process
1201
              metadata.
1202
1203
          Block Detection
1204
          1. Identifying Block Boundaries
1205
          Overview:
1206
          Block detection segments data into logical units or tables.
1207
          Methods:
1208
          - Boundary Detection:
1209
          - Find empty rows, repeated labels, or formatting changes.
1210
          - Example:
          | Name | Amount | |-----| | John | 100 | | | | <-- Empty row
1211
               indicates new block | Name | Amount | | Alice| 200 |
1212
          - Machine Learning:
1213
          - Train classifiers to detect block boundaries based on cell patterns.
1214
1215
          Advanced:
          - Detect nested blocks or hierarchies using indentation or merged
1216
              cells.
1217
          - Identify summary blocks with keywords like "Total" or "Summary".
1218
1219
          2. Processing and Tracking Blocks
1220
          Overview:
          Once blocks are detected, assign IDs and enable block-level analysis.
1221
1222
          Actions:
1223
          - Block ID:
1224
          - Assign unique IDs (e.g., Block_001, Block_002).
1225
          - Analysis:
          - Perform group-by or aggregation within each block.
1226
          - Example:
1227
          - Sum "Amount" for Block_001: 100 + 150 = 250
1228
1229
          3. Additional Block Actions
1230
          Overview:
          Enable modular analysis and reporting at the block level.
1231
1232
          Features:
1233
          - Summary Rows:
1234
          - Add computed totals/averages for each block.
1235
          - Export/Save:
          - Save blocks as separate files or sheets.
1236
          - Example:
1237
          - Export Block_001 to "block1.csv"
1238
1239
          Search for Values or Patterns
1240
          1. Search Execution Methods
1241
          Overview:
          Efficiently locate specific values or patterns in your data.
```

```
1242
1243
          Techniques:
1244
          - Manual Tools:
1245
          - Use Ctrl + F in Excel for quick lookups.
          - Programmatic Search:
1246
          - Scan all cells using loops or vectorized code.
1247
          - Example:
1248
          - Find all instances of "North" in the "Region" column.
1249
          - Pattern Matching:
1250
          - Support exact, wildcard (*Total*), and regex (\d{4}-\d{2}-\d{2} for
              dates).
1251
1252
          2. Recording and Highlighting Results
1253
          Overview:
1254
          Log and visualize search matches for user review.
1255
          Actions:
1256
          - Logging:
1257
          - Record coordinates (e.g., Sheet1, Row 3, Col "Region").
1258
          - Highlighting:
1259
          - Apply conditional formatting to search hits.
1260
          3. Advanced Search Scenarios
1261
          Overview:
1262
          Handle complex or large-scale search requirements.
1263
1264
          Scenarios:
1265
          - Merged Cells:
          - Search within merged cells or across multiple sheets.
1266
          - Export:
1267
          - Export found results for further analysis.
1268
          - Example:
1269
          - Export all rows containing "John" to "john_results.csv"
1270
          Writeback Results
1271
          1. Output Placement
1272
          Overview:
1273
          Choose where and how to insert results.
1274
          Options:
1275
          - Target Columns:
1276
          - Select existing or blank columns for output.
1277
          - Appending:
1278
          - Add new columns for flags, counts, or statuses.
1279
          - Example:
          - Add "Approved_Flag" column next to "Status".
1280
1281
          2. Writing and Styling Results
1282
1283
          Automate and style the output for visibility.
1284
          Methods:
1285
          - Formulas/Code:
1286
          - Use code (e.g., ws.cell(row, col).value = result) to insert results.
1287
          - Styling:
1288
          - Bold, borders, or colors for output cells.
1289
          - Example:
          - Green fill for "Success", red for "Error".
1290
1291
          3. Audit and Protection
1292
          Overview:
1293
          Maintain the integrity and traceability of results.
1294
1295
          Measures:
          - Lock Columns:
```

```
1296
          - Prevent edits to output columns.
1297
          - Timestamps/User Info:
1298
          - Add audit trail for writebacks.
1299
          - Example:
          - "2024-06-01, User: admin"
1300
1301
          Difference in State
1302
          1. Sheet Comparison
1303
          Overview:
1304
          Identify changes between input and output sheets.
1305
          Process:
1306
          - Load Sheets:
1307
          - Read both sheets into memory.
1308
          - Compare Cells:
          - Detect differences by position and value.
1309
1310
          2. Recording and Reporting Differences
1311
          Overview:
1312
          Log and report all detected changes.
1313
          Actions:
1314
          - Log Mismatches:
1315
          - Record cell coordinates and values.
1316
          - Example:
1317
          - Cell B3: "North" \rightarrow "South"
1318
          - Export Diff Report:
          - List all detected differences for review.
1319
1320
          3. Visualization and Automation
1321
          Overview:
1322
          Make changes visible and automate validation.
1323
1324
          Features:
          - Highlight Changes:
1325
          - Color code changed cells.
1326
          - Automate Checks:
1327
          - Integrate diff comparisons into test scripts.
1328
          Column Selection
1329
          1. Selection Criteria
1330
          Overview:
1331
          Choose relevant columns for analysis.
1332
1333
          Methods:
          - Labels/Indices:
1334
          - Select by name or position.
1335
          - Dynamic Rules:
1336
          - E.g., all numeric columns.
1337
          - Assign Roles:
1338
          - Example: "ID", "Date", "Metric"
1339
          2. Preparation and Validation
1340
          Overview:
1341
          Prepare columns for consistent use.
1342
1343
          Actions:
          - Rename/Relabel:
1344
          - Standardize column names.
1345
          - Validate Types:
1346
          - Ensure columns are of expected type.
1347
          - Example:
1348
          - "Date" column as datetime.
1349
          3. Reusability
```

```
1350
          Overview:
1351
          Save and reuse column selections.
1352
1353
          Features:
          - Presets:
1354
          - Save selection profiles.
1355
          - Downstream Use:
1356
          - Use validated columns in subsequent processes.
1357
1358
          Filter Rows
          1. Filtering Methods
1359
          Overview:
1360
          Refine your dataset with filters.
1361
1362
          Techniques:
          - Spreadsheet Tools:
1363
          - Use built-in filters.
1364
          - Code Logic:
1365
          - Filter with code (e.g., df[df['Status'] == 'Approved']).
1366
          - Multiple Criteria:
1367
          - Combine conditions (AND/OR).
1368
          - Example:
          - Status = "Approved" AND Amount > 100
1369
1370
          2. Helper Columns and Complex Filters
1371
          Overview:
1372
          Simplify filtering using helper columns.
1373
          Actions:
1374
          - Helper Columns:
1375
          - Compute intermediate flags.
1376
          - Document Logic:
1377
          - Record filtering rules for audit.
1378
          3. Post-Filter Actions
1379
          Overview:
1380
          Visualize and export filtered data.
1381
1382
          Features:
          - Highlighting:
1383
          - Grey-out filtered-out rows.
1384
          - Export:
1385
          - Save the filtered dataset.
1386
1387
          Merge Tables
          1. Key-Based Merging
1388
          Overview:
1389
          Combine tables using shared keys.
1390
1391
          Techniques:
1392
          - Join Operations:
          - Use VLOOKUP, JOIN, or code merges.
1393
          - Example:
1394
          - Merge "Customer_ID" from two tables.
1395
          - Align Data:
1396
          - Match on columns like "ID", "Name".
1397
          2. Stack-Based Merging
1398
          Overview:
1399
          Append tables when keys 'arent needed.
1400
1401
          Methods:
1402
          - Vertical Append:
          - Combine rows from similar tables.
1403
          - Deduplicate:
```

```
1404
          - Remove duplicate records.
1405
1406
          3. Tracking and Audit
1407
          Overview:
          Track source and unmatched records.
1408
1409
          Actions:
1410
          - Source Column:
1411
          - Add "Source" to indicate origin.
1412
          - Highlight Unmatched:
          - Mark or export mismatched rows.
1413
1414
          Pivot or Unpivot
1415
          1. Pivoting Data
1416
          Overview:
1417
          Summarize data using pivots.
1418
          Methods:
1419
          - PivotTables:
1420
          - Group by \operatorname{row/column} dimensions.
1421
          - Example:
          - Sum "Amount" by "Region".
1422
          - Aggregation:
1423
          - Choose SUM, AVG, COUNT, etc.
1424
1425
          2. Unpivoting (Melting) Data
1426
          Overview:
1427
          Reshape data from wide to long format.
1428
          Techniques:
1429
          - Melt Operations:
1430
          - Convert columns into rows.
1431
          - Example:
1432
          | Year | Sales_2019 | Sales_2020 | |-----|
1433
1434
          | Year | Sales_Year | Value |
1435
          - Flexible Restructuring:
1436
          - Selectively unpivot non-ID columns.
1437
          3. Post-Pivot Actions
1438
          Overview:
1439
          Prepare pivoted data for export.
1440
1441
          Features:
          - Flatten Pivot Table:
1442
          - Convert back to flat for further analysis.
1443
          - Reorder/Rename:
1444
          - Clarify pivoted fields.
1445
1446
          Map with Lookup Tables
          1. Mapping Techniques
1447
          Overview:
1448
          Standardize data using lookups.
1449
1450
          Methods:
1451
          - Functions:
          - Use VLOOKUP, merge with dictionaries.
1452
          - Code-to-Label:
1453
          - Example:
1454
          - Code "N" →Label "North"
1455
1456
          2. Application and Fallbacks
1457
          Overview:
          Apply lookups and handle missing values.
```

```
1458
1459
          Actions:
1460
          - Apply Mappings:
          - Across selected columns.
1461
          - Handle Missings:
1462
          - Use defaults for missing codes.
1463
1464
          3. Audit and Display
1465
          Overview:
1466
          Ensure mapping transparency.
1467
          Features:
1468
          - Cache Mappings:
1469
          - Store for repeated use.
1470
          - Display Codes/Labels:
1471
          - Show both for clarity.
1472
          Fill Missing Data
1473
          1. Choosing Fill Methods
1474
          Overview:
1475
          Impute missing data appropriately.
1476
          Techniques:
1477
          - Forward/Backward Fill:
1478
          - Fill gaps with prior/next value.
1479
           - Default Values:
1480
          - Use fixed placeholder (e.g., 0, "Unknown").
          - Contextual Example:
1481
          - Dates: Fill missing month with last known month.
1482
1483
          2. Application and Auditing
1484
          Overview:
1485
          Apply fills and flag for review.
1486
          Actions:
1487
           - Targeted Filling:
1488
          - Apply to specific columns/rows.
1489
          - Flag Filled Cells:
1490
          - Highlight for later review.
1491
          3. Documentation
1492
          Overview:
1493
          Keep fill logic transparent.
1494
1495
          Features:
1496
          - Record Logic:
          - Document assumptions and methods.
1497
          - Audit Trail:
1498
          - Track all changes.
1499
1500
          Flag Rows or Cells
          1. Defining Flag Rules
1501
1502
          Establish criteria for flagging.
1503
1504
          Examples:
1505
          - Simple Rule:
          - Flag where Amount < 0
1506
          - Complex Rule:
1507
          - Flag where Status = "Pending" and Amount > 1000
1508
1509
          2. Applying Flags
1510
          Overview:
          Insert flags and summarize.
1511
```

```
1512
          Actions:
1513
          - Flag Column:
1514
          - Add "Flag" column with "Yes"/"No".
1515
          - Export Flagged Rows:
          - Save for further inspection.
1516
1517
          3. Advanced Flagging
1518
          Overview:
1519
          Use multiple criteria and document.
1520
          Features:
1521
          - Multi-Criteria:
1522
          - Combine several rules for granular checks.
1523
          - Notes:
1524
          - Document flagging rationale.
1525
          Sort Data
1526
          1. Setting Sort Criteria
1527
          Overview:
1528
          Organize data for analysis.
1529
1530
          Options:
          - Sort Columns:
1531
          - By value, ascending/descending.
1532
          - Multi-Level:
1533
          - E.g., sort by "Region", then by "Amount".
1534
1535
          2. Applying Sorts
          Overview:
1536
          Implement sorting programmatically or manually.
1537
1538
          Methods:
1539
          - Spreadsheet Tools:
1540
          - Built-in sort features.
          - Code:
1541
          - E.g., df.sort_values(['Region', 'Amount'])
1542
1543
          3. Post-Sort Actions
1544
          Overview:
          Finalize sorted data.
1545
1546
          Actions:
1547
          - Renumber Rows:
1548
          - Update indices.
1549
          - Highlight Extremes:
          - Mark top/bottom values.
1550
1551
          Validate Data
1552
          1. Validation Checks
1553
          Overview:
1554
          Ensure data meets required standards.
1555
          Checks:
1556
          - Type:
1557
          - Ensure numeric columns contain numbers.
          - Range:
1559
          - E.g., "Amount" > 0.
          - Pattern:
1560
          - Date columns match YYYY-MM-DD.
1561
          - Business Rule Example:
1562
          - "Start Date" < "End Date"
1563
1564
          2. Marking and Reporting
1565
          Overview:
          Visualize and report errors.
```

```
1566
1567
          Actions:
1568
          - Highlight Invalids:
          - Color-code errors.
1569
          - Export Summary:
1570
          - Table of error counts and locations.
1571
1572
          3. Integration in Workflow
1573
          Overview:
          Make validation a routine part of processing.
1574
1575
          Features:
1576
          - Pre-Processing Step:
1577
          - Validate before analysis.
1578
          - Automation:
          - Integrate into data pipelines.
1579
1580
          Split Sheets or Data
1581
          1. Defining Split Rules
1582
          Overview:
1583
          Segment data for modular analysis.
1584
          Methods:
          - By Category:
1586
          - E.g., split by "Region".
1587
          - By Date Range:
1588
          - E.g., split by year.
1589
          2. Exporting Segments
1590
          Overview:
1591
          Save segments for separate use.
1592
1593
          Actions:
          - Export Files:
1594
            "North_Region.csv", "South_Region.csv"
1595
          - Consistent Formatting:
1596
          - Ensure identical columns and styling.
1597
1598
          3. Automation and Documentation
1599
          Automate splitting and track provenance.
1601
          Features:
1602
          - Automation:
1603
          - Use scripts/macros for repeated splits.
          - Documentation:
1604
          - Record rules and export logs.
1605
```

D QUALITATIVE ANALYSIS OF BREW-GENERATED KNOWLEDGE BASES

1616

1617

1618

1619

This section presents a comprehensive qualitative analysis of knowledge bases generated through the BREW technique applied to two distinct agent training environments: OSWorld and τ^2 Bench described in the section before. The analysis examines knowledge representation patterns, procedural sophistication, and domain-specific learning characteristics extracted from CUA agent behaviors, providing insights into the effectiveness and scope of knowledge distillation techniques across diverse task environments.

D.1 CROSS-DOMAIN KNOWLEDGE BASE ANALYSIS

D.1.1 BASE STRUCTURE & ORGANIZATION

Schema Consistency and Evolution: Both knowledge bases demonstrate consistent structural schemas, though adapted to their respective domains. The OSWorld KB employs a four-part schema (contextual triggers, procedural steps, extended capabilities, concrete instantiation), while the τ^2 Bench KB extends this to a five-part structure, adding explicit purpose rationale ("Why to use it"). This evolution suggests that BREW adapts its extraction patterns to domain-specific requirements—conversational commerce demands explicit justification for actions due to customer interaction contexts.

Taxonomic Organization Principles: The OSWorld KB reveals a capability-based taxonomy organized around computational tasks: file operations, document processing, inter-application workflows, and data visualization. Each category represents a distinct computational domain with specific tool requirements and interaction patterns. In contrast, the τ^2 Bench KB employs a **lifecycle-based taxonomy** structured around transactional states: order creation, modification, fulfillment, and post-delivery operations. This organizational difference reflects fundamental domain characteristics—desktop automation focuses on tool orchestration, while conversational commerce centers on process management.

Hierarchical Task Decomposition: Both KBs demonstrate sophisticated hierarchical reasoning, but through different decomposition strategies. OSWorld exhibits **technical decomposition**, breaking complex operations like "Create Charts from Data" into constituent technical steps (data selection, chart insertion, customization, formatting). τ^2 Bench shows **process decomposition**, structuring operations like order modification into authentication, validation, confirmation, and execution phases. This suggests BREW successfully identifies domain-appropriate decomposition strategies rather than applying uniform patterns.

Knowledge Boundary Definition: Both KBs explicitly encode operational boundaries, but through contrasting mechanisms. OSWorld boundaries are **capability-constrained**—determined by available applications and system resources. τ^2 Bench boundaries are **policy-constrained**—explicitly defined through "Deny Unsupported Request" patterns and escalation protocols. This difference highlights how knowledge extraction adapts to domain-specific constraint types.

D.1.2 PROCEDURAL KNOWLEDGE GROUNDING

Context-Dependent Action Selection: Both domains demonstrate sophisticated context awareness, but grounded in different environmental factors. OSWorld exhibits application-context sensitivity, where identical operations (e.g., image insertion) require different procedures across LibreOffice Writer, Impress, GIMP, and Thunderbird. The agent learned application-specific affordances and interaction patterns rather than generic command sequences. τ^2 Bench demonstrates state-context sensitivity, where available actions depend on order status (pending vs. delivered), payment methods, and authentication levels. This reveals learned understanding of business process constraints and temporal operation windows.

Error Prevention and Validation Workflows: Both KBs incorporate sophisticated error prevention mechanisms, but grounded in domain-specific failure modes. OSWorld emphasizes technical validation: file integrity checks ("confirm the exported file opens correctly"), application state verification, and multi-step confirmation for irreversible operations. τ^2 Bench emphasizes transactional validation: authentication cascades, confirmation dialogues with standardized templates, and explicit user consent protocols. The emergence of defensive programming practices across both domains suggests these represent fundamental principles of reliable agent behavior.

State-Dependent Decision Logic: The procedural knowledge in both domains demonstrates sophisticated state machine reasoning. OSWorld exhibits application state awareness—understanding when applications are ready for input, when files are loaded, and when operations can be safely executed. Window management and application switching reveal learned understanding of desktop metaphors and resource constraints. τ^2 Bench demonstrates business process state awareness—finite state machine reasoning where order lifecycle states determine available operations. The agent learned that pending orders enable modification while delivered orders unlock return workflows, indicating internalized understanding of business logic constraints.

Security and Authentication Grounding: While OSWorld operates in a trusted desktop environment with minimal explicit security concerns, τ^2 Bench reveals pervasive authentication-first paradigms. Nearly every transactional operation begins with identity verification through email, name, and zip code combinations. The KB demonstrates graduated security reasoning: information retrieval requires basic authentication while financial transactions trigger rigorous verification protocols. This contrast highlights how procedural knowledge adapts to domain-specific security requirements.

 Cross-Application vs. Cross-Process Orchestration: OSWorld demonstrates technical orchestration—coordinating multiple applications (Chrome, LibreOffice suite, File Manager, GIMP) to accomplish complex workflows. The "Navigate Between Applications" section reveals learned behaviors for window management, application switching, and resource coordination. τ^2 Bench exhibits process orchestration—coordinating authentication, validation, confirmation, and execution phases across different operational contexts. Both forms of orchestration require sophisticated temporal reasoning and constraint management, but applied to different environmental complexity types.

Failure Mode Internalization: Both KBs reveal learned understanding of domain-specific failure modes. OSWorld incorporates file validation, application crash recovery suggestions, and verification steps for critical operations. τ^2 Bench includes explicit escalation protocols ("Transfer to Human Agent"), policy compliance mechanisms, and irreversibility warnings for financial operations. The consistent emergence of failure-aware procedures suggests that agents successfully internalize risk assessment and mitigation strategies during training.

Domain-Specific Communication Patterns: The procedural knowledge reveals distinct communication paradigms appropriate to each domain. OSWorld procedures are **task-oriented** with minimal user interaction—focusing on efficient command execution and verification. τ^2 Bench procedures are **dialogue-oriented** with standardized customer interaction templates, confirmation protocols, and expectation management communications. This adaptation demonstrates that BREW extracts not just procedural logic but domain-appropriate interaction modalities.

The cross-domain analysis reveals that BREW successfully extracts procedural knowledge that is both **structurally consistent** (following learnable organizational patterns) and **contextually grounded** (adapted to domain-specific constraints, failure modes, and interaction requirements). This dual capability suggests significant potential for knowledge transfer across related domains while maintaining appropriate domain-specific adaptations.