# Towards Human-Guided, Data-Centric LLM Co-Pilots

**Evgeny Saveliev** *♣, **Jiashuo Liu** *♦, **Nabeel Seedat** *♣, **Anders Boyd** ‡, **Mihaela van der Schaar**♣
♣ University of Cambridge ♦ Tsinghua University ‡ Amsterdam University Medical Centers

## Abstract

Machine learning (ML) has the potential to revolutionize various domains and industries, but its adoption is often hindered by the disconnect between the needs of domain experts and translating these needs into robust and valid ML tools. Despite recent advances in LLM-based co-pilots to democratize ML for non-technical domain experts, these systems remain predominantly focused on model-centric aspects while overlooking critical data-centric challenges. This limitation is problematic in complex real-world settings where raw data often contains complex issues, such as missing values, label noise, and domain-specific nuances requiring tailored handling. To address this we introduce CliMB-DC, a human-guided, data-centric framework for LLM co-pilots that combines advanced data-centric tools with LLM-driven reasoning to enable robust, context-aware data processing. At its core, CliMB-DC introduces a novel, multi-agent reasoning system that combines a strategic coordinator for dynamic planning and adaptation with a specialized worker agent for precise execution. Domain expertise is then systematically incorporated to guide the reasoning process using a human-in-the-loop approach. To guide development, we formalize a taxonomy of key data-centric challenges that co-pilots must address. Thereafter, to address the dimensions of the taxonomy, we integrate state-of-the-art data-centric tools into an extensible, open-source architecture, facilitating the addition of new tools from the research community. Empirically, using real-world healthcare datasets we demonstrate CliMB-DC's ability to transform uncurated datasets into ML-ready formats, significantly outperforming existing co-pilot baselines for handling data-centric challenges. CliMB-DC promises to empower domain experts from diverse domains — healthcare, finance, social sciences and more — to actively participate in driving real-world impact using ML.

## 1 Introduction

Over the past decade, machine learning (ML) has evolved at a breathtaking pace, raising hopes that advanced ML methods can transform a wide range of domains and industries. However, for many domain experts — including medical researchers, social scientists, business analysts, environmental scientists, education researchers and more — conceiving a problem through which ML can provide a solution remains challenging. Despite having a deep understanding of their data and domain-specific challenges, these individuals often lack the programming or technical background needed to implement sophisticated ML pipelines (Pfisterer et al., 2019), and thus are considered *non-technical domain experts*. This gap in expertise creates a significant barrier to realizing the potential of ML across these domains.

Recent advancements in large language models (LLMs) have paved the way for AI co-pilots that promise to automate various aspects of ML development through natural language interaction (Hassan et al., 2023; Tu et al., 2024). However, current co-pilots remain predominantly focused on model-centric aspects—such as architecture selection and hyperparameter tuning—while overlooking the fundamental role of the data-centric side to ML. Since data-centric aspects largely determine the performance, fairness, robustness and safety of ML systems, ignoring the processes of constructing and handling data can negatively affect performance or worse lead to incorrect conclusions. Unfortunately, real-world data often contains missing values, inconsistencies, mislabeled records, and

---

*ES, JL and NS contributed equally.

domain-specific nuances (see Table 7) and thus the data is usually not ML-ready (Sambasivan et al., 2021; Balagopalan et al., 2024a). Furthermore, applying a "one-size-fits-all" data cleaning script from an LLM co-pilot that cannot be tailored to the varying structures of data risks erasing critical signals or introducing biases, and leaves the domain experts powerless to intervene.

There is indeed a growing interest in data-centric AI within the ML community — emphasizing the importance of ML to improve data quality, curation, and characterization (Seedat et al., 2022d; Zha et al., 2023; Liang et al., 2022). In particular, numerous data-centric ML tools and methods have been developed for handling common data issues, such as missing values, noisy labels, and data drift (Northcutt et al., 2021; Jarrett et al., 2022b; Seedat et al., 2022b; 2023b; Liu et al., 2023). However, for non-technical domain experts, these tools are often abstract to implement and remain out of reach to use. Integrating these tools into LLM-based co-pilots would not only allow tailored handling of data and thus empower domain experts, but also would broaden the use of data-centric AI research across various disciplines and application settings, including healthcare, finance, education etc.

Despite their value, data-centric tools are not a panacea in and of themselves and cannot be applied by co-pilots in isolation. Actions like imputing data or rectifying noisy labels require contextual understanding to avoid distorting critical domain-specific information. This underscores the need for expert oversight—guidance from individuals deeply familiar with the nuances of the data—to ensure that actions align with domain-specific goals and constraints. Such guidance is crucial in high-stakes fields like healthcare and finance, where improper data handling can lead to misleading conclusions or harmful decisions. This interplay between human expertise and data-centric automation presents a unique challenge for LLM-based co-pilots. Designing systems capable of nuanced reasoning and iterative planning, while effectively incorporating expert feedback, remains a significant hurdle. A co-pilot must not only execute tasks but also intelligently sequence and adapt data processing pipelines with a human-in-the-loop approach.

To address these challenges, we introduce **Cli**nical predictive **M**odel **B**uilder with **D**ata-**C**entric AI (**CliMB-DC**), a human-guided data-centric framework for LLM co-pilots. We advance upon and address limitations of other co-pilots by integrating advanced data-centric tools, along with a novel LLM-driven reasoning process to enable robust, context-aware data processing for real-world ML challenges. This enables robust, context-aware data processing for challenges faced when using ML approaches in real-world settings. Specifically, CliMB-DC introduces a novel, multi-agent reasoning system that combines a strategic coordinator for dynamic planning and adaptation with a specialized worker agent for precise execution. Domain expertise is then systematically incorporated to guide reasoning using a human-in-the-loop approach. Where CliMB established the foundation, CliMB-DC advances this vision by enabling sophisticated reasoning about data quality, integrity, and domain-specific constraints—essential capabilities for developing trustworthy ML systems when analyzing complex, real-world data.

**Contributions**: ① **Taxonomy of Challenges**: We formalize a taxonomy of data-centric challenges that co-pilots need to address. ② **Data-Centric Tools**: We integrate state-of-the-art, data-centric tools into an extensible and open-source framework. The broader accessibility for non-technical domain experts to these data-centric tools allows them more options when tailoring their data management accordingly. It additionally provides an opportunity for the data-centric ML research community to incorporate new tools or validate their tools more easily. ③ **Human-in-the-Loop Alignment**: We implement a human-in-the-loop system to ensure contextual alignment of data processing actions with domain-specific requirements. Moreover, we are able to incorporate domain expertise through natural language interaction, allowing experts to guide and assess data transformations without requiring coding experience. ④ **Multi-Agent Planning and Reasoning** : We introduce a novel multi-agent reasoning approach that combines a strategic coordinator agent with a specialized worker agent, enabling sophisticated planning and adaptation of data-centric workflows. ⑤ **Empirical Case Studies**: We conduct emprical case studies on real-world healthcare data, demonstrating where existing co-pilots fall short in handling the complexities of real-world data and illustrate the advantages of our approach.

CliMB-DC represents a significant step toward democratizing ML for non-technical domain experts, while ensuring the responsible and effective use of data-centric AI tools. By combining automation with expert oversight, it enables robust ML development that respects domain-specific knowledge. The target audience of CliMB-DC encompasses a wide range of users, including:

Table 1: Comparison of different co-pilots along different dimensions.

| Perspective | Capability | DS-Agent | AutoGen | Data-Interpreter | OpenHands | CliMB-DC (Ours) |
|---|---|---|---|---|---|---|
| Components | Data-centric tools | × | × | × | × | ✓ |
| | Clinical models | × | × | × | × | ✓ |
| Expert Input | Static integration | ✓ | × | ✓ | ✓ | ✓ |
| | Dynamic integration | × | × | × | ✓ | ✓ |
| Data-centric Reasoning | Setup refinement | × | × | × | × | ✓ |
| | Performance refinement | × | × | × | × | ✓ |
| Pipeline | End-to-end automation | × | ✓ | ✓ | ✓ | ✓ |
| | Replanning | × | × | × | × | ✓ |
| | Backtracking | × | × | × | × | ✓ |
| | Code refinement | ✓ | ✓ | ✓ | ✓ | ✓ |

- **Non-Technical Domain Experts**: While we contextualize and instantiate CliMB-DC as a tool for healthcare, we envision that such a system could be relevant to non-technical domain experts in other data-driven domains such as finance, environmental management, education etc.
- **Data-Centric ML Researchers**: CliMB-DC provides a platform for data-centric researchers to effortlessly compare tools, integrate and validate new ones, accelerating the advancement of data-centric AI.

## 2 RELATED WORK

This work engages not only with LLM Co-pilots tools, but also with the broader area of data-centric AI research. We provide a summary of existing LLM co-pilots in contrast to our proposal CliMB-DC in Table 1. An extended related work is provided in Appendix A.

*Remark:* We also formalize a taxonomy of data-centric challenges that future co-pilots need to address in Appendix C.1, which is then summarized in Table 7.

## 3 CLIMB-DC: AN LLM CO-PILOT FROM A DATA-CENTRIC PERSPECTIVE

**Problem Setting.** Denote $D = \{(\mathbf{x}^i, y^i)\}_{i=1}^n$, where $\mathbf{x} = (x_1, \ldots, x_p)$ with $x_d \in \mathcal{X}_d$ and $y \in \mathcal{Y}$, be a well-curated, "ML-ready" dataset suitable for training a given ML model (including an AutoML model) and achieving optimal target performance. Here, we consider a general scenario including both the supervised setting depending on the label types – such as classification $\mathcal{Y} = \{1, \ldots, C\}$, regression $\mathcal{Y} = \mathbb{R}$, and time-to-event analysis $\mathcal{Y} = (\{0, 1\}, \mathbb{R}_{\geqslant 0})$.

*Data corruption faced in Practice.* In real-world healthcare scenarios datasets have numerous challenges as discussed. Additionally, since non-technical domain experts have limited expertise in data science, it often results in uncurated datasets, denoted as $\tilde{D} = (\tilde{\mathbf{x}}^i, \tilde{y}^i)_{i=1}^{\tilde{n}}$, where $\tilde{\mathbf{x}} = (\tilde{x}_1, \cdots, \tilde{x}_{\tilde{p}})$ with $\tilde{x}_d \in \tilde{\mathcal{X}}_d$. These datasets are subject to various data-centric issues, as highlighted in Table 7. If left unprocessed, such issues can lead to undesired failures or suboptimal performance in downstream ML models. To clarify this concept, we formalize how a well-curated dataset, $D$, can be (unknowingly and unintentionally) transformed into an uncurated dataset, $\tilde{D}$, through a series of $L$ data corruption processes during real-world data collection.

$$\tilde{D} = g(D) = g_L \circ \cdots \circ g_1(D), \tag{1}$$

where $g_\ell$ represents the corruption applied at the $\ell$-th step. A well-curated dataset, $D$, can be corrupted in numerous ways, impeding the optimal performance and clinical impact of ML models. Based on our taxonomy in Appendix C.1, Table 7, we categorize the prominent data-centric issues commonly encountered in healthcare datasets, each representing a specific type of corruption function.

*Ideal data-centric curation.* Suppose it is feasible to revert the data corruption process applied to the well-curated dataset, $D$, from the given uncurated dataset, $\tilde{D}$. Ideally, the goal is to construct a series of $L$ data curation functions, $f_1, \ldots, f_L$, where each curation function is specifically designed to revert the corresponding data corruption function applied to $D$, i.e., $f_\ell = g_\ell^{-1}$.

*Domain-specific model learning.* Once the dataset is curated, the objective is to select an appropriate, context-dependent model class and train an ML model that achieves strong generalization performance for the user-defined task descriptions.

Then we move on to introduce in detail our proposed LLM co-pilot designed through a data-centric lens, named CliMB-DC (see Figure 1 for the overall architecture).
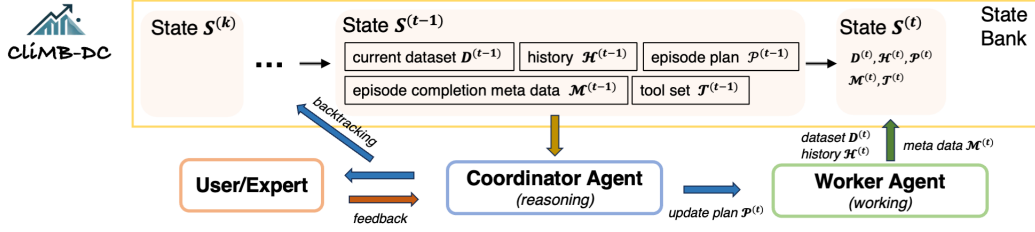
Figure 1: The overall architecture of CliMB-DC, which primarily consists of three components that interact with the state bank: (i) a coordinator agent responsible for reasoning and planning, (ii) a worker agent for code writing and execution, and (iii) the user or human experts.

## 3.1 OVERALL ARCHITECTURE

Recall that given an initial dataset $\mathcal{D}_0$, our goal is to find an optimal sequence of transformations $\mathbf{f} = (f_1, \ldots, f_L)$ that yields a curated dataset $\mathcal{D}^*$ suitable for downstream ML tasks. Each transformation $f_i \in \mathcal{F}$ is selected from a space of possible operations, guided by both LLM reasoning and expert feedback. The curated dataset is then used to select and train domain-specific ML models for prediction.

Our framework, CliMB-DC (see Figure 1), addresses challenges faced by existing co-pilots through a dual-agent architecture that combines the strengths of LLM-based reasoning with human domain expertise.

The system consists of a high-level *coordinator agent*, responsible for managing the overall data processing strategy, and a specialized *worker agent*, tasked with executing specific data transformations. This separation of responsibilities enables CliMB-DC to maintain strategic oversight while ensuring operational efficiency. Additionally, the framework supports continuous integration of expert feedback, facilitating iterative refinement of the strategy. As shown in Figure 1, the CliMB-DC primarily consists of three components that interact with the state bank:

**State bank.** The state bank stores the system states at each time step $t$. Formally, the system state $\mathcal{S}$ at time $t$ is defined by $\mathcal{S}^{(t)} = \{D^{(t)}, \mathcal{H}^{(t)}, \mathcal{P}^{(t)}, \mathcal{M}^{(t)}, \mathcal{T}^{(t)}\}$, where $D^{(t)}$ represents the current dataset state, $\mathcal{H}^{(t)}$ captures the interaction history including expert feedback, $\mathcal{P}^{(t)}$ maintains the dynamic episode plan, $\mathcal{M}^{(t)}$ tracks episode completion metadata, and $\mathcal{T}^{(t)}$ contains the available data-centric and modeling tools that can be used by the worker agent.

**Tool registry.** Recent advances in data-centric AI have led to the development of a variety of methods and tools from across the community. CliMB-DC integrates a large variety of diverse data-centric (and model-centric tools) from across the literature (see Table 2) which are available to the worker agent to utilize. Although not exhaustive, the current set of tools covers a diverse set of scenarios linked to the data-centic challenges taxonomy in Section C.1. Moreover, as outlined in Section C.2, we illustrate the extendable nature of the framework to easily integrate new tools from the ML community.

Table 2: Overview of tools available in CliMB-DC. This ensures data/model-centric tools are accessible to non-technical domain experts, while also providing data-centric ML researchers a platform for tool impact.

| Tool class | Available tools |
|---|---|
| Data understanding | Descriptive statistics, Exploratory data analysis (EDA), Feature selection (Remeseiro & Bolon-Canedo, 2019) |
| Feature extraction (from text) | spaCy Matcher |
| Data characterization | Data-IQ (Seedat et al., 2022a),TRIAGE (Seedat et al., 2023a) |
| Missing data | HyperImpute (Jarrett et al., 2022a) |
| Data valuation | KNN-Shapley (Jia et al., 2019) |
| Data auditing (outliers) | Confident Learning (Cleanlab) (Northcutt et al.; 2021) |
| Data imbalance | SMOTE (Chawla et al., 2002) |
| Model building | AutoPrognosis 2.0 (Imrie et al., 2023) (supports regression, classification, survival analysis) |
| Post-hoc interpretability | Permutation explainer (Breiman, 2001), SHAP explainer (Lundberg & Lee, 2017), AutoPrognosis 2.0 subgroup analysis (Imrie et al., 2023) |
| Test time risk or failure analysis | Data-SUITE (Seedat et al., 2022c), SMART Testing (Rauba et al., 2024) |

**Coordinator agent.** The Coordinator agent is the strategic planner of the system, responsible for maintaining a high-level view of the data processing pipeline and making decisions about task sequencing (i.e. the plan). It implements a three-stage reasoning process that continuously evaluates progress, identifies potential issues, and adapts the processing strategy based on both automated metrics and expert feedback. Operating through reasoning process $\pi_{\mathcal{C}}$, the Coordinator maps the current system state to the next-step strategic decisions/processing decision (i.e. plan $\mathcal{P}^{(t+1)}$):

$$\pi_{\mathcal{C}} : \mathcal{S}^{(t)} \to \mathcal{P}^{(t+1)}. \tag{2}$$

The detailed reasoning approach is demonstrated in Section 3.2.

**User/Expert integration.** Domain expertise is systematically integrated throughout the process through feedback that ensures all transformations align with domain-specific requirements and constraints. This integration occurs through a natural language feedback mechanism that evaluates proposed transformations and enriches the system's understanding of the domain. Additionally, it captures domain knowledge that enhances the future reasoning and decision making by the coordinator agent, creating a continuous learning loop that improves the system's performance over time. The detailed interaction mechanism is introduced in Section 3.2 (see Equation 5 in *State Observation*).

**Worker agent.** The Worker agent acts as the system's execution engine, translating high-level plans from the coordinator agent into concrete data transformations instantiated in code. It combines LLM capabilities with specialized data-centric tools to implement precise, context-aware transformations while maintaining interactions to integrate information from domain experts. The Worker's execution process is formalized as:

$$\pi_{\mathcal{W}} : (\mathcal{S}^{(t)}, \mathcal{P}^{(t)}, \mathcal{T}^{(t)}) \to (\mathcal{H}^{(t+1)}, \mathcal{D}^{(t+1)}, \mathcal{M}^{(t+1)}) \tag{3}$$

where $\mathcal{S}^{(t)}$ represents the current state, $\mathcal{T}^{(t)}$ indicates the current selected/available tool, and $\mathcal{D}^{(t+1)}$ denotes the resultant (new) dataset state, $\mathcal{H}^{(t+1)}$ and $\mathcal{M}^{(t+1)}$ denote the updated history records. The Worker operates at a granular level, focusing on individual data processing episodes and ensuring each transformation aligns with both technical requirements and domain expertise.

## 3.2 DETAILS OF CLIMB-DC'S REASONING PROCESS

Before introducing CliMB-DC's reasoning approach, we first highlight the challenges in our specific scenarios faced by an alternative approach — Monte Carlo Tree Search (MCTS).
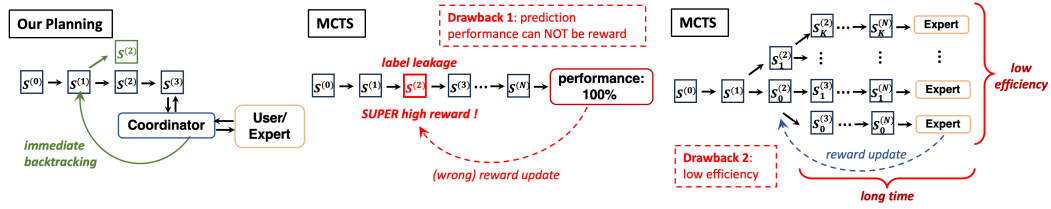


Figure 2: Challenges of Monte Carlo Tree Search (MCTS). We highlight two key drawbacks of MCTS. First, prediction performance cannot serve as a reliable reward, as it may favor data issues such as label leakage or meaningless problem setups (middle). Second, MCTS suffers from low efficiency, requiring experts to endure long waiting times and evaluate a large number of trials (right). In contrast, CliMB-DC's proposed reasoning approach enables immediate backtracking and replanning, significantly enhancing efficiency.

**Challenges of MCTS.** Monte Carlo Tree Search (MCTS) is a commonly used reasoning and planning mechanism that generates random paths to explore and evaluate potential plans based on a reward function. However, the complex nature of real-world data issues introduces several critical challenges, significantly limiting the practicality of MCTS in these contexts. We highlight the key challenges as follows:

- **Lack of intermediate reward model**: MCTS depends on a well-defined reward model. However, in our setting, there is no clear reward model, particularly for all the intermediate states that may arise. Even human experts are unable to provide such detailed rewards. For instance, given a dataset, it is challenging for experts to evaluate all the data issues listed in Table 7. As a result, MCTS would require complete model training and evaluation to obtain reward signals, making iterative data curation computationally prohibitive.
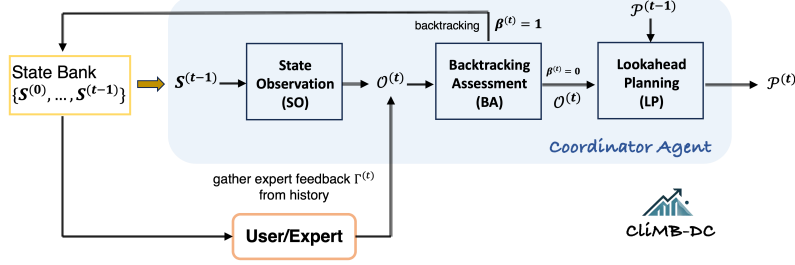
Figure 3: The framework of the coordinator agent in CliMB-DC, encompassing three parts named State Observation (SO), Backtracking Assessment (BA), and Lookahead Planning (LP).

- **Prediction performance is unsuitable as the final reward**: Another significant challenge for MCTS in our scenarios is that prediction performance cannot be directly used as the final reward. It is critical to first ensure that the entire data processing pipeline is valid and free from issues such as label leakage, which could render performance metrics unreliable. For example, as shown in Figure 2 (middle), if the state $\mathcal{S}^{(2)}$ introduces label leakage, relying solely on performance to determine the reward would assign an artificially high reward to this node. However, such a state should be avoided in the final path.
- **Low efficiency**: A further challenge is the low efficiency of the process. Since prediction performance is inadequate as a final reward, expert evaluation may be necessary. However, MCTS involves random exploration and requires numerous steps to transform a raw dataset into a well-trained prediction model. As shown in Figure 2 (right), this results in lengthy trials, and the large number of trials exacerbates the inefficiency. Additionally, in many scenarios, users may be unable to examine the details of all trials due to time constraints or limited expertise in data science and clinical domain knowledge. These limitations significantly restrict the applicability of MCTS in our scenarios.

The challenges associated with MCTS largely arise from its approach of treating all data processing steps as unknown and unexplored, attempting to navigate the entire sequence of actions, as is common in gaming scenarios. In contrast, data processing typically follows a well-established order of operations, making such exhaustive "searching" unnecessary. For example, addressing data missingness generally precedes other transformations or feature engineering steps, and exploring these well-known rules through extensive random trials is both inefficient and redundant.

A more practical alternative to MCTS is to focus on refining "local processing" within the general sequence of operations. Our framework incorporates automated planning combined with expert validation to ensure both technical quality and domain-specific appropriateness. The key innovation lies in enabling the method to backtrack after errors and consult experts when necessary, such as for decisions involving the meanings of features, handling label leakage, or determining whether to drop specific features, etc. As shown in Figure 2 (left), when combined with immediate expert feedback, the coordinator enables prompt backtracking, significantly improving efficiency.

**CliMB-DC's proposed multi-stage reasoning.** The reasoning mechanism of CliMB-DC is demonstrated in Figure 3. At time $t$, the coordinator agent takes the current state bank $\{\mathcal{S}^{(0)}, \ldots, S^{(t-1)}\}$ as input, and outputs the plan via:

$$\underbrace{P(\mathcal{P}^{(t)}|\{\mathcal{S}^{(0)}, \ldots, \mathcal{S}^{(t-1)}\})}_{\text{coordinator reasoning}} = \sum_{\mathcal{O}^{(t)}} \underbrace{P(\mathcal{O}^{(t)}|\mathcal{S}^{(t-1)})}_{\text{state observation}} \cdot \left( \sum_{\beta^{(t)}} \underbrace{P(\beta^{(t)}|\mathcal{O}^{(t)})}_{\text{backtracking}} \underbrace{P(\mathcal{P}^{(t)}|\beta^{(t)}, \mathcal{O}^{(t)})}_{\text{lookahead planning}} \right), \quad (4)$$

1. **State Observation (SO)**: $P(\mathcal{O}^{(t)}|\mathcal{S}^{(t-1)})$ denotes the state observation stage, where the coordinator analyzes the project state, focusing primarily on the last state $\mathcal{S}^{(t-1)}$. And it will gather expert human feedback $\Gamma^{(t)}$ from history interactions. Therefore, $P(\mathcal{O}^{(t)}|\mathcal{S}^{(t)})$ can be formulated as:

$$P(\mathcal{O}^{(t)}|\mathcal{S}^{(t-1)}) = \sum_{\Gamma^{(t)}} P(\Gamma^{(t)}|\mathcal{S}^{(t-1)})P(\mathcal{O}^{(t)}|\Gamma^{(t)}, \mathcal{S}^{(t-1)}), \quad (5)$$

where $\Gamma^{(t)}$ denotes the expert human feedback for current state. To inform its decision, the coordinator extracts information about the project state such as evaluating experiment outcomes, data quality metrics and expert feedback received, together denoted as $\mathcal{O}^{(t)}$.

2. **Backtracking Assessment (BA)**: Based on the analysis, the coordinator determines if previous decisions need revision or updating. If the project is not progressing satisfactorily (e.g., due to data
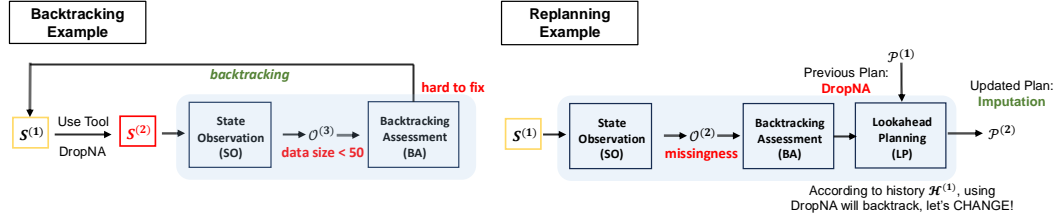
Figure 4: Example of backtracking and replanning in handling missing data, showcasing how the Backtracking Assessment (BA) and Lookahead Planning (LP) modules in the proposed reasoning approach collaborate to efficiently resolve data issues.

quality issues or expert feedback indicating problems), the coordinator identifies a backtracking point $k < t - 1$ and restores the project state to $\mathcal{S}^{(k)}$ (i.e. backtrack step), denoted by $P(\beta^{(t)}|\mathcal{O}^{(t)})$.

3. $M$**-Step Lookahead Planning (LP):** The coordinator evaluates the current plan, focusing specifically on the next $M$ episodes $e_t, \ldots, e_{t+M-1}$ within the plan $\mathcal{P}^{(t)}$. For each episode $e_i$, the coordinator assesses two key aspects: (i) Necessity and (ii) Appropriateness. Leveraging the history and user interactions (expert feedback) stored in $\mathcal{O}^{(t)}$, the coordinator refines the plan by excluding episodes considered unnecessary or inappropriate and incorporating new ones as needed. This process ensures the updated plan $\mathcal{P}^{(t)}$ remains aligned with the user's objectives. Specifically, this can involve the following types of updates: (i) Reordering episodes to better handle dependencies; (ii) Removing unnecessary episodes; (iii) Adding new episodes to address identified gaps; (iv) Modifying episode parameters based on expert input.

The coordinator then issues this updated plan to the worker agent to execute at the next iteration. When the Worker agent completes an episode, control returns to the Coordinator for the next iteration of plan analysis and refinement. Pseudo-code detailiing the interaction between components and the progression of transformations is provided in Algorithm 1.

**Demonstration of the reasoning process.** Let us consider the case of a dataset with missingness and how the reasoning process works. As shown in Figure 4, after loading the dataset, the co-pilot detects the missingness issue and initially plans to address it using the DropNA function, which removes all rows with missing values, resulting in the state $\mathcal{S}^{(2)}$. However, the State Observation (SO) module identifies a new problem: the reduced dataset size falls below 50, which is insufficient for subsequent processing and difficult to remedy. In response, the Backtracking Assessment (BA) module is triggered, rolling the state back to $\mathcal{S}^{(1)}$.

In the next step, the SO module detects the missingness issue, and the BA module is not triggered. Drawing on the history record, which indicates that using "DropNA" previously led to backtracking, the Lookahead Planning (LP) module revises the plan and selects an alternative approach—imputation—to address the missingness issue.

### 3.3 WORKER AGENT

The worker agent takes the updated plan and integrates it with the available tool set (see Table 2). It then generates and executes the necessary code to complete the plan. If execution fails, the agent autonomously updates the code to ensure successful execution. The worker agent also verifies the availability of required Python packages, installing them if needed before proceeding with execution. For the case studies (Sec 4), we perform an ablation, removing the coordinator & our tool set from CliMB-DC to emphasize the reliability of our worker agent in code generation and execution.

## 4 CASE STUDIES

We empirically investigate CliMB-DC's capability to handle real-world healthcare data challenges.
① **Does it work?**: We highlight for multiple data-centric challenges from our taxonomy, where vanilla co-pilots fail and the data-centric lens with human feedback helps.
② **Why does it work?**: We provide an in-depth analysis via various case studies to better understand why CliMB-DC succeeds and other co-pilots fail.

**Datasets.** We employ real-world tabular healthcare datasets with varying characteristics. i.e. different sample sizes, dimensionality, task types (classification, survival analysis) and task difficulty. These datasets reflect the following data challenges (as defined in our taxonomy): (i) Lung cancer: Data leakage (ii) PBC Dataset: Unaggregated data (based on identifiers) and (iii) Prostate cancer prediction: Ambiguous and Hard examples (mislabeled and outliers).

**Baselines.** We compare CliMB-DC as discussed in the related work to Data Interpreter and Open-Hands as representative co-pilots. Additionally, we perform an ablation of CliMB-DC. We remove the coordinator, highlighting its value, while assessing the worker agent's reliability in code generation and execution. This ablation is instantiated both with and without tools. We refer to these configurations as CliMB-DC (No coordinator, With tools) and CliMB-DC (No coordinator, No tools), respectively. Unless otherwise stated all results are averaged over 5 runs. In addition to the results in the main paper, Appendix E and F provides a detailed analysis of the logs of interactions for CliMB-DC and the baselines.

## 4.1 DOES IT WORK?

To demonstrate the effectiveness of CliMB-DC, we conduct the following case studies on healthcare datasets using different data challenges. Below, we summarize the results and the context and importance of these datasets:

- **Case study 1:** In the dataset on predictors of lung cancer, the primary challenge involves addressing data leakage. This scenario is particularly complex as it combines survival analysis with the need to identify and handle potentially leaked information from outcome-related variables. Table 4 demonstrates based on C-index that both Data-Interpreter and OpenHands were unable to provide valid results in almost all scenarios and produced several different reasons for failures, most of which related to data leakage. CliMB-DC was able to produce valid results without run failures.

- **Case study 2:** In the dataset on predictors of PBC, unaggregated data and potential data leakage are presented as simultaneous challenges. Table 3 demonstrates based on C-index that both Data-Interpreter and OpenHands were unable to provide vaild results, mainly due to failure to aggregate data per patient and identify data leakage. This scenario is especially relevant to healthcare settings where multiple measurements per patient are common. CliMB-DC was able to handle these issues and produce valid results, while maintaining temporal consistency and avoiding information leakage.

- **Case study 3:** In the datasets comparing predictors of prostate cancer mortality from the SEER (USA) and CUTRACT (UK) datasets, the challenge lies in handling data quality issues and data drifts, across different healthcare systems. Table 5 demonstrates based on AU-ROC that all three tools were able to produce valid results, while the accuracy and AU-ROC was slightly higher when using CliMB-DC. The results demonstrate our framework's robustness in managing dataset shifts while maintaining consistent performance across different healthcare contexts.

Table 3: Results on the PBC dataset, where the primary data challenges are addressing *unaggregated data* and *data leakage*. The prediction task in this case is *survival analysis*, a specialized and less common task compared to those typically encountered in general machine learning fields. The whole processing procedure of the proposed CliMB-DC is shown in Figure 5.

| Method | Human Assistance | Results Valid | C-Index | Failure Modes | % runs tested |
|---|---|---|---|---|---|
| Data-Interpreter | - | × | 0.789 | Failed to aggregate data per patient<br>Failed to identify data leakage<br>Failed to produce results<br>Failed to set up survival problem | 100%<br>100%<br>40%<br>20% |
| OpenHands | - | × | 0.468 | Failed to aggregate data per patient<br>Failed to identify data leakage<br>Failed to set up survival problem<br>Convergence issues causing task failure | 100%<br>100%<br>60%<br>20% |
| **CliMB-DC**<br>(No Coordinator<br>& No Tools) | - | × | 0.663 | Failed to aggregate data per patient<br>Failed to identify data leakage<br>Failed to produce results<br>Fail to solve convergence error | 100%<br>80%<br>60%<br>20% |
| **CliMB-DC** (No<br>Coordinator & With Tools) | - | × | 0.914 | Failed to aggregate data per patient<br>Failed to identify data leakage | 100%<br>100% |
| **CliMB-DC** | - | ✓ | **0.953** | (Successful) | 100% |

Table 4: Results on the Lung Cancer dataset, where the primary data challenge is addressing *data leakage*. The prediction task in this case is *survival analysis*, a specialized and less common task compared to those typically encountered in general machine learning fields. The whole processing procedure of the proposed CliMB-DC is shown in Appendix D, Figure 9.

| Method | Human Assistance | Results Valid | C-Index | Failure Modes | % runs tested |
|---|---|---|---|---|---|
| Data-Interpreter | - | ✗ | 0.625 | Failed to identify data leakage<br>Incorrect metric used<br>Failed to set up survival problem<br>PCA use degraded performance | 100%<br>20%<br>20%<br>60% |
| | Leakage features excluded | ✓<br>✗ | 0.738<br>0.995 | (Successful)<br>Label leakage reintroduced | 80%<br>20% |
| OpenHands | - | ✗ | N/A | Failed to identify data leakage<br>Incorrect metric reported | 100%<br>100% |
| | Specify cox model | ✗ | 0.496 | Failed to identify data leakage | 100% |
| | Leakage features excluded | ✓<br>✗ | 0.500<br>N/A | (Successful)<br>Stuck in a loop | 80%<br>20% |
| **CliMB-DC** (No Coordinator & No Tools) | - | ✗ | 0.765 | Failed to identify data leakage<br>Failed to solve convergence error | 100%<br>40% |
| | Leakage features excluded | ✓<br>✗ | 0.809<br>N/A | (Successful)<br>Failed to test on the test file | 80%<br>20% |
| **CliMB-DC** (No Coordinator & With Tools) | - | ✗ | 0.871 | Failed to identify data leakage | 100% |
| **CliMB-DC** | - | ✓ | **0.848** | (Successful) | 100% |

Table 5: Results on cross cancer mortality prediction (SEER from the USA to CUTRACT from the UK), where the primary data challenges are addressing *data quality/hardness* and *data drifts*. The prediction task in this case is *classification*, a common task in general machine learning fields. The whole processing procedure of the proposed CliMB-DC is shown in Appendix D, Figure 10.

| Method | Human Assistance | Results Valid | Accuracy | AUC-ROC | Failure Modes | % runs tested |
|---|---|---|---|---|---|---|
| Data-Interpreter | - | ✓<br>✗ | 66.5<br>N/A | 0.727<br>N/A | (Successful)<br>Failed in preprocessing | 80%<br>20% |
| OpenHands | - | ✓ | 67.5 | 0.729 | (Successful) | 100% |
| **CliMB-DC** (No Coordinator & Tool Set) | - | ✓<br>✗ | 67.8<br>68.3 | 0.749<br>0.683 | (Successful)<br>Failed to compute AUC-ROC | 80%<br>20% |
| **CliMB-DC** (No Coordinator & With Tools) | - | ✓ | 69.4 | 0.765 | (Successful) | 100% |
| **CliMB-DC** | - | ✓ | **69.9** | **0.771** | (Successful) | 100% |

## 4.2 Why does it work?

To provide a deeper understanding of how CliMB-DC can excel in specific, data-centric challenges, we describe various facets of using CliMB-DC in comparison to other co-pilots. These case studies illustrate how our multi-agent architecture, reasoning processes and human-in-the-loop feedback can provide specific advantages, and where problems when using other co-pilots arise. Figure 5 illustrates the specific reasoning and planning mechanisms through which CliMB-DC reasons, adapts the plan, engages with the domain expert and then resolves these challenges. Further deepdives can be found in Appendix D (Figures 9-10)

**Case study 1:** In many datasets from healthcare settings, multiple measurements are recorded for each patient over time, leading to unaggregated data. Baseline co-pilots treat each row as an independent patient observation, creating two issues: (1) data leakage between training and test sets when measurements from the same patient appear in both, and (2) an ill-posed modeling setup that violates the independence assumptions of survival analysis. In contrast, as shown in Figure 5, CliMB-DC identifies the structure of the dataset, leveraging its state observation and human-guided feedback mechanisms to aggregate measurements correctly. This ensures that the modeling process aligns with the true data generation process, avoiding leakage.

**Case study 2:** Survival analysis tasks are particularly vulnerable to label leakage from other features or covariates that can compromise model validity. Specifically, time-dependent variables like "time_to_lung_cancer" inherently leak information about the outcome. Other co-pilots fail to account for such features, resulting in inflated performance metrics and compromising the model's real-world applicability and generalization. In contrast, as shown in Appendix D, Figure 9, CliMB-DC's dynamic
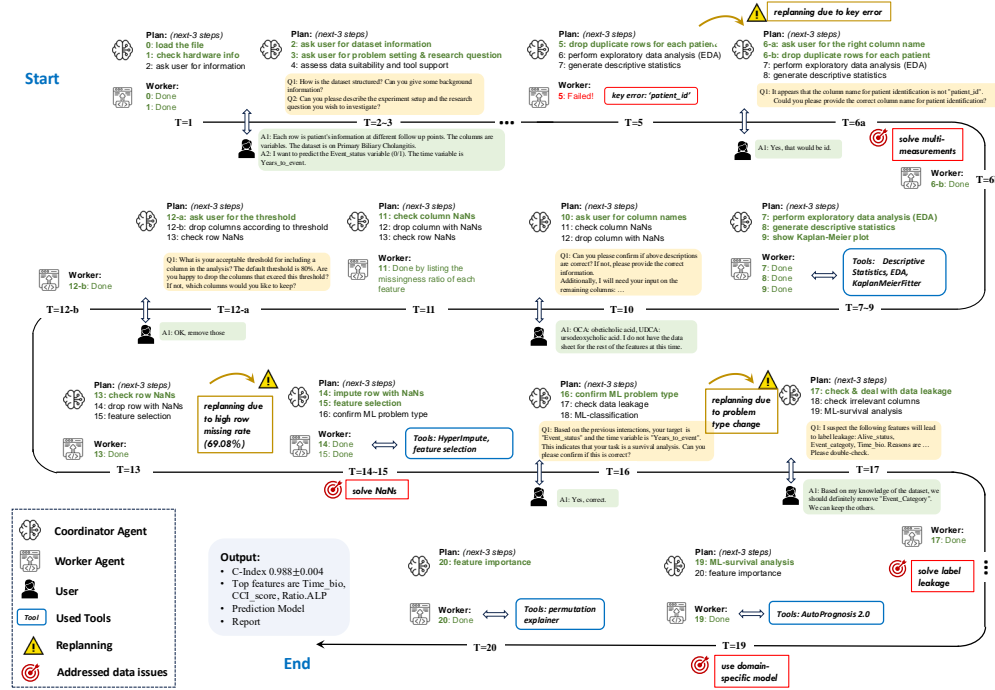
Figure 5: The processing workflow of CliMB-DC on the PBC dataset, illustrating how the coordinator agent, worker agent, and user interact at each processing stage $T$. The data issues include multiple measurements, missingness, and label leakage. The prediction task is survival analysis, requiring domain-specific model classes.

reasoning and iterative planning allows detection and mitigation of label leakage. Through domain expert feedback, the system removes problematic features like *"time_to_lung_cancer"*. This ensures that the resulting models are valid and generalizable.

**Case study 3:** Data from healthcare settings can often have data quality challenges such as hard examples (mislabeled, heterogenous outcomes etc). These observations can affect model training (Seedat et al., 2022b) and can be sculpting or filtered from the dataset to improve generalization. In addition, when models are used across countries, as is the case in the two prostate cancer datasets from SEER and CUTRACT, distribution shift could occur. In contrast, as shown in Appendix D, Figure 10, CliMB-DC's dynamic reasoning and data-centric tool usage allow it to understand that data quality is a challenge, run a method for data characterization (e.g. Data-IQ (Seedat et al., 2022b)) and based on the output, engage with the human expert to remove ambiguous observations that may cause downstream problems during modeling. We show that this improves model generalization cross-domain (i.e. in different countries).

## 5 CONCLUSION

We introduced CliMB-DC, a human-guided, data-centric framework for LLM-based co-pilots. Importantly, this framework addresses a critical gap in current LLM co-pilots: their inability to effectively handle real-world data challenges while incorporating domain expertise. Our contributions span multiple dimensions, from a taxonomy of data-centric challenges to developing a novel multi-agent architecture enabling sophisticated reasoning about data quality and processing workflows.

Our empirical evaluations highlight several key advantages of CliMB-DC when handling key data challenges, allowing it to achieve robust ML outcomes where existing co-pilots may come across problems. Beyond these technical contributions, the open-source nature of CliMB-DC encourages the broader research community to extend its capabilities, ensuring its relevance across diverse data structures and modeling applications. By highlighting the importance of data-centric aspects to AI co-pilots, CliMB-DC represents a critical step towards democratizing ML for non-technical domain experts (in a variety of fields), ensuring that data quality and contextual understanding are central to ML workflows. We envision this framework as a foundational tool for democratizing the adoption of ML across a variety of problem settings and domains.

REFERENCES

Andrea Apicella, Francesco Isgrò, and Roberto Prevete. Don't push the button! exploring data leakage risks in machine learning and transfer learning. *arXiv preprint arXiv:2401.13796*, 2024.

Anand Avati, Martin Seneviratne, Emily Xue, Zhen Xu, Balaji Lakshminarayanan, and Andrew M Dai. Beds-bench: Behavior of ehr-models under distributional shift–a benchmark. *arXiv preprint arXiv:2107.08189*, 2021.

Aparna Balagopalan, Ioana Baldini, Leo Anthony Celi, Judy Gichoya, Liam G McCoy, Tristan Naumann, Uri Shalit, Mihaela van der Schaar, and Kiri L Wagstaff. Machine learning for healthcare that matters: Reorienting from technical novelty to equitable impact. *PLOS Digit. Health*, 3(4): e0000474, April 2024a.

Aparna Balagopalan, Ioana Baldini, Leo Anthony Celi, Judy Gichoya, Liam G McCoy, Tristan Naumann, Uri Shalit, Mihaela van der Schaar, and Kiri L Wagstaff. Machine learning for healthcare that matters: Reorienting from technical novelty to equitable impact. *PLOS Digital Health*, 3(4):e0000474, 2024b.

Xiao-Yuan Bao, Wan-Jing Huang, Kai Zhang, Meng Jin, Yan Li, and Cheng-Zhi Niu. A customized method for information extraction from unstructured text data in the electronic medical records. *Beijing da xue xue bao. Yi xue ban= Journal of Peking University. Health Sciences*, 50(2):256–263, 2018.

Brett K Beaulieu-Jones, Jason H Moore, and POOLED RESOURCE OPEN-ACCESS ALS CLINICAL TRIALS CONSORTIUM. Missing data imputation in the electronic health record using deeply learned autoencoders. In *Pacific symposium on biocomputing 2017*, pp. 207–218. World Scientific, 2017.

Louise Bloch, Christoph M Friedrich, and Alzheimer's Disease Neuroimaging Initiative. Data analysis with shapley values for automatic subject selection in alzheimer's disease data sets using interpretable machine learning. *Alzheimer's Research & Therapy*, 13:1–30, 2021.

Sabri Boughorbel, Fethi Jarray, Neethu Venugopal, and Haithum Elhadi. Alternating loss correction for preterm-birth prediction from ehr data with noisy labels. *arXiv preprint arXiv:1811.09782*, 2018.

Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.

Angel Alexander Cabrera, Minsuk Kahng, Fred Hohman, Jamie Morgenstern, and Duen Horng Chau. Discovery of intersectional bias in machine learning using automatic subgroup generation. In *ICLR Debugging Machine Learning Models Workshop*, 2019.

Tiffany Tianhui Cai, Hongseok Namkoong, and Steve Yadlowsky. Diagnosing model performance under distribution shift. *arXiv preprint arXiv:2303.02011*, 2023.

Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.

Davide Chicco, Luca Oneto, and Erica Tavazzi. Eleven quick tips for data cleaning and feature engineering. *PLOS Computational Biology*, 18(12):e1010718, 2022.

Nastaran Enshaei, Moezedin Javad Rafiee, Arash Mohammadi, and Farnoosh Naderkhani. Data shapley value for handling noisy labels: An application in screening covid-19 pneumonia from chest ct scans. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1381–1385. IEEE, 2022.

Hossein Estiri and Shawn N Murphy. Semi-supervised encoding for outlier detection in clinical observation data. *Computer methods and programs in biomedicine*, 181:104830, 2019.

Pablo Ferri, Nekane Romero-Garcia, Rafael Badenes, David Lora-Pablos, Teresa García Morales, Agustín Gómez de la Cámara, Juan M García-Gómez, and Carlos Sáez. Extremely missing numerical data in electronic health records for machine learning can be managed through simple imputation methods considering informative missingness: A comparative of solutions in a covid-19 mortality case study. *Computer Methods and Programs in Biomedicine*, 242:107803, 2023.

Marzyeh Ghassemi, Tristan Naumann, Peter Schulam, Andrew L Beam, Irene Y Chen, and Rajesh Ranganath. A review of challenges and opportunities in machine learning for health. *AMIA Summits on Translational Science Proceedings*, 2020:191, 2020.

Karan Goel, Albert Gu, Yixuan Li, and Christopher Re. Model patching: Closing the subgroup performance gap with data augmentation. In *International Conference on Learning Representations*, 2020.

Lea Goetz, Nabeel Seedat, Robert Vandersluis, and Mihaela van der Schaar. Generalization—a key challenge for responsible ai in patient-facing clinical applications. *npj Digital Medicine*, 7(1):126, 2024.

Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Ds-agent: Automated data science by empowering large language models with case-based reasoning. *arXiv preprint arXiv:2402.17453*, 2024.

Udo Hahn and Michel Oleynik. Medical information extraction in the age of deep learning. *Yearbook of medical informatics*, 29(01):208–220, 2020.

Sebastien Haneuse, David Arterburn, and Michael J Daniels. Assessing missing data assumptions in ehr-based studies: a complex and underappreciated task. *JAMA Network Open*, 4(2):e210184–e210184, 2021.

Md Mahadi Hassan, Alex Knipper, and Shubhra Kanti Karmaker Santu. Chatgpt as your personal data scientist. *arXiv preprint arXiv:2305.13657*, 2023.

Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. *Advances in Neural Information Processing Systems*, 36, 2024.

Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei, Danyang Li, Jiaqi Chen, Jiayi Zhang, et al. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*, 2024.

Fergus Imrie, Bogdan Cebere, Eoin F McKinney, and Mihaela van der Schaar. AutoPrognosis 2.0: Democratizing diagnostic and prognostic modeling in healthcare with automated machine learning. *PLOS Digit. Health*, 2(6):e0000276, June 2023.

Daniel Jarrett, Bogdan C Cebere, Tennison Liu, Alicia Curth, and Mihaela van der Schaar. HyperImpute: Generalized iterative imputation with automatic model selection. In *International Conference on Machine Learning*, pp. 9916–9937. PMLR, June 2022a.

Daniel Jarrett, Bogdan C Cebere, Tennison Liu, Alicia Curth, and Mihaela van der Schaar. Hyperimpute: Generalized iterative imputation with automatic model selection. In *International Conference on Machine Learning*, pp. 9916–9937. PMLR, 2022b.

Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gurel, Bo Li, Ce Zhang, Costas Spanos, and Dawn Song. Efficient task-specific data valuation for nearest neighbor algorithms. *Proceedings of the VLDB Endowment*, 12(11):1610–1623, 2019.

Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pp. 5637–5664. PMLR, 2021.

Moritz Lehne, Julian Sass, Andrea Essenwanger, Josef Schepers, and Sylvia Thun. Why digital medicine depends on interoperability. *NPJ digital medicine*, 2(1):79, 2019.

Weixin Liang, Girmaw Abebe Tadesse, Daniel Ho, Li Fei-Fei, Matei Zaharia, Ce Zhang, and James Zou. Advances, challenges and opportunities in creating data for trustworthy ai. *Nature Machine Intelligence*, 4(8):669–677, 2022.

Jiashuo Liu, Jiayun Wu, Renjie Pi, Renzhe Xu, Xingxuan Zhang, Bo Li, and Peng Cui. Measure the predictive heterogeneity. In *The Eleventh International Conference on Learning Representations*, 2022.

Jiashuo Liu, Tianyu Wang, Peng Cui, and Hongseok Namkoong. On the need for a language describing distribution shifts: Illustrations on tabular datasets. *Advances in Neural Information Processing Systems*, 36, 2023.

Tianyang Liu, Fei Wang, and Muhao Chen. Rethinking tabular data understanding with large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 450–482, 2024.

Andrew Low and Z Yasemin Kalender. Data dialogue with chatgpt: Using code interpreter to simulate and analyse experimental data. *arXiv preprint arXiv:2311.12415*, 2023.

Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 2017.

Chuizheng Meng, Loc Trinh, Nan Xu, James Enouen, and Yan Liu. Interpretability and fairness evaluation of deep learning models on mimic-iv dataset. *Scientific Reports*, 12(1):7166, 2022.

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pp. 220–229, 2019.

Aída Muñoz Monjas, David Rubio Ruiz, David Pérez Del Rey, and Matvey B Palchuk. Enhancing real world data interoperability in healthcare: A methodological approach to laboratory unit harmonization. *International Journal of Medical Informatics*, 193:105665, 2025.

Yang Nan, Javier Del Ser, Simon Walsh, Carola Schönlieb, Michael Roberts, Ian Selby, Kit Howard, John Owen, Jon Neville, Julien Guiot, et al. Data harmonisation for information fusion in digital healthcare: A state-of-the-art systematic review, meta-analysis and future research directions. *Information Fusion*, 82:99–122, 2022.

Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021.

Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.

Luke Oakden-Rayner, Jared Dunnmon, Gustavo Carneiro, and Christopher Ré. Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, pp. 151–159, 2020.

Nassim Oufattole, Teya Bergamaschi, Aleksia Kolo, Hyewon Jeong, Hanna Gaggin, Collin M Stultz, and Matthew McDermott. Meds-tab: Automated tabularization and baseline methods for meds datasets. *arXiv preprint arXiv:2411.00200*, 2024.

Konstantin D Pandl, Fabian Feiland, Scott Thiebes, and Ali Sunyaev. Trustworthy machine learning for health care: scalable data valuation with the shapley value. In *Proceedings of the Conference on Health, Inference, and Learning*, pp. 47–57, 2021.

Kayur Patel, James Fogarty, James A Landay, and Beverly Harrison. Investigating statistical machine learning as a tool for software development. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 667–676, 2008.

Florian Pfisterer, Janek Thomas, and B Bischl. Towards human centered AutoML. *ArXiv*, abs/1911.02391, November 2019.

Oleg S Pianykh, Georg Langs, Marc Dewey, Dieter R Enzmann, Christian J Herold, Stefan O Schoenberg, and James A Brink. Continuous learning ai in radiology: implementation principles and early applications. *Radiology*, 297(1):6–14, 2020.

Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33: 17044–17056, 2020.

Paulius Rauba, Nabeel Seedat, Max Ruiz Luyten, and Mihaela van der Schaar. Context-aware testing: A new paradigm for model testing with large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

Beatriz Remeseiro and Veronica Bolon-Canedo. A review of feature selection methods in medical applications. *Computers in Biology and Medicine*, 112:103375, 2019.

Jenna Rychert. In support of interoperability: A laboratory perspective. *International Journal of Laboratory Hematology*, 45(4):436–441, 2023.

Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. "everyone wants to do the model work, not the data work": Data cascades in high-stakes AI. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, May 2021. ACM.

Leonard Sasse, Eliana Nicolaisen-Sobesky, Juergen Dukart, Simon B Eickhoff, Michael Götz, Sami Hamdan, Vera Komeyer, Abhijit Kulkarni, Juha Lahnakoski, Bradley C Love, et al. On leakage in machine learning pipelines. *arXiv preprint arXiv:2311.04179*, 2023.

Bey-Marrié Schmidt, Christopher J Colvin, Ameer Hohlfeld, and Natalie Leon. Definitions, components and processes of data harmonisation in healthcare: a scoping review. *BMC Medical Informatics and Decision Making*, 20:1–19, 2020.

Nabeel Seedat, Jonathan Crabbe, Ioana Bica, and M Schaar. Data-IQ: Characterizing subgroups with heterogeneous outcomes in tabular data. *Neural Inf Process Syst*, abs/2210.13043:23660–23674, October 2022a.

Nabeel Seedat, Jonathan Crabbé, Ioana Bica, and Mihaela van der Schaar. Data-iq: Characterizing subgroups with heterogeneous outcomes in tabular data. *Advances in Neural Information Processing Systems*, 35:23660–23674, 2022b.

Nabeel Seedat, Jonathan Crabbé, and Mihaela van der Schaar. Data-suite: Data-centric identification of in-distribution incongruous examples. In *International Conference on Machine Learning*, pp. 19467–19496. PMLR, 2022c.

Nabeel Seedat, Fergus Imrie, and Mihaela van der Schaar. Dc-check: A data-centric ai checklist to guide the development of reliable machine learning systems. *arXiv preprint arXiv:2211.05764*, 2022d.

Nabeel Seedat, Jonathan Crabbé, Zhaozhi Qian, and Mihaela van der Schaar. Triage: Characterizing and auditing training data for improved regression. *Advances in Neural Information Processing Systems*, 36:74995–75008, 2023a.

Nabeel Seedat, Fergus Imrie, and Mihaela van der Schaar. Dissecting sample hardness: Fine-grained analysis of hardness characterization methods. In *The Twelfth International Conference on Learning Representations*, 2023b.

Nabeel Seedat, Fergus Imrie, and Mihaela van der Schaar. Navigating data-centric artificial intelligence with dc-check: Advances, challenges, and opportunities. *IEEE Transactions on Artificial Intelligence*, 2023c.

Arindam Sett, Somaye Hashemifar, Mrunal Yadav, Yogesh Pandit, and Mohsen Hejrati. Speaking the same language: Leveraging llms in standardizing clinical data for ai. *arXiv preprint arXiv:2408.11861*, 2024.

Janmajay Singh, Masahiro Sato, Tomoko Ohkuma, et al. On missingness features in machine learning models for critical care: observational study. *JMIR Medical Informatics*, 9(12):e25022, 2021.

Daniel J Stekhoven and Peter Bühlmann. Missforest—non-parametric missing value imputation for mixed-type data. *Bioinformatics*, 28(1):112–118, 2012.

Harini Suresh, Jen J Gong, and John V Guttag. Learning tasks for multitask learning: Heterogenous patient populations in the ICU. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 802–810, 2018.

Ana Szarfman, Jonathan G Levine, Joseph M Tonning, Frank Weichold, John C Bloom, Janice M Soreth, Mark Geanacopoulos, Lawrence Callahan, Matthew Spotnitz, Qin Ryan, et al. Recommendations for achieving interoperable and shareable medical data in the usa. *Communications Medicine*, 2(1):86, 2022.

Siyi Tang, Amirata Ghorbani, Rikiya Yamashita, Sameer Rehman, Jared A Dunnmon, James Zou, and Daniel L Rubin. Data valuation for medical imaging using shapley value and application to a large-scale chest x-ray dataset. *Scientific reports*, 11(1):8366, 2021.

Nenad Tomašev, Xavier Glorot, Jack W Rae, Michal Zielinski, Harry Askham, Andre Saraiva, Anne Mottram, Clemens Meyer, Suman Ravuri, Ivan Protsyuk, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature*, 572(7767):116–119, 2019.

Andrej Tschalzev, Sascha Marton, Stefan Lüdtke, Christian Bartelt, and Heiner Stuckenschmidt. A data-centric perspective on evaluating machine learning models for tabular data. *arXiv preprint arXiv:2407.02112*, 2024.

Xinming Tu, James Zou, Weijie Su, and Linjun Zhang. What should data science education do with large language models? 2024.

Boris van Breugel, Nabeel Seedat, Fergus Imrie, and Mihaela van der Schaar. Can you rely on your model evaluation? improving model evaluation with synthetic test data. *Advances in Neural Information Processing Systems*, 36, 2024.

Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.

Yishu Wei, Yu Deng, Cong Sun, Mingquan Lin, Hongmei Jiang, and Yifan Peng. Deep learning with noisy labels in medical prediction problems: a scoping review. *Journal of the American Medical Informatics Association*, pp. ocae108, 2024.

Steven Euijong Whang, Yuji Roh, Hwanjun Song, and Jae-Gil Lee. Data collection and quality challenges in deep learning: A data-centric ai perspective. *The VLDB Journal*, 32(4):791–813, 2023.

Jenny Yang, Hagen Triendl, Andrew AS Soltan, Mangal Prakash, and David A Clifton. Addressing label noise for electronic health records: Insights from computer vision for tabular data. *medRxiv*, pp. 2023–10, 2023.

Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyou Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. *Advances in Neural Information Processing Systems*, 35:32598–32611, 2022.

Karina Zadorozhny, Patrick Thoral, Paul Elbers, and Giovanni Cinà. Out-of-distribution detection for medical applications: Guidelines for practical evaluation. In *Multimodal AI in healthcare: A paradigm shift in health intelligence*, pp. 137–153. Springer, 2022.

Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, and Xia Hu. Data-centric ai: Perspectives and challenges. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pp. 945–948. SIAM, 2023.

Boyang Zhao. Clinical data extraction and normalization of cyrillic electronic health records via deep-learning natural language processing. *JCO Clinical Cancer Informatics*, 3:1–9, 2019.

Yue Zhao, Zain Nasrullah, and Zheng Li. Pyod: A python toolbox for scalable outlier detection. *Journal of machine learning research*, 20(96):1–7, 2019.

# Appendix: Towards Human-Guided, Data-Centric LLM Co-Pilots

## Table of Contents

**Code will be released upon acceptance.**

## A  EXTENDED RELATED WORK

Below we provide a further assessment of LLM-based code interpreters.

### A.1  RELATED WORK OVERVIEW

**LLM-Based Co-pilots.**    The rapid advancements in LLMs have paved the way for various stages of ML and data science workflows to be automated by co-pilots and code interpreters that leverage the reasoning and code generation capabilities of LLMs (Tu et al., 2024; Hollmann et al., 2024; Low & Kalender, 2023). These tools allow users to specify their requirements for data science pipelines via natural language and thus offer greater flexibility compared to traditional AutoML systems.

Prominent examples include systems that chain task execution (e.g., AutoGPT, DS-Agent (Guo et al., 2024)), modular frameworks for multi-step reasoning (e.g., OpenHands (Wang et al., 2024)), and graph-based workflow decomposition (e.g., Data Interpreter (Hong et al., 2024)).

Despite recent progress, significant challenges remain in addressing the complex, data-centric aspects of real-world ML workflows. Many co-pilots operate within predefined pipelines or task hierarchies, making them ill-suited for dynamic, data-centric workflows. Furthermore, while these systems excel in automating code generation, they often lack mechanisms for robust data reasoning, such as diagnosing data issues, incorporating domain-specific knowledge, or addressing contextual nuances (e.g., data leakage or feature importance validation). These gaps are especially pertinent to real-world datasets that commonly exhibit variability and noise-common to data from healthcare settings (as summarized in Table 7).

Additionally, the effectiveness of these co-pilots to empower non-technical domain experts, particularly in healthcare, remains a significant challenge. Healthcare data is often characterized by heterogeneity, complexity, and susceptibility to biases and data quality issues. Hence, a co-pilot blindly applying generic data processing techniques to raw clinical data can lead to the introduction of errors and the loss of important clinical information. For example, detecting and removing outliers based on the percentile of a variable distribution might remove extreme lab values that are clinically meaningful, as they could represent a critical underlying condition rather than noise. In another example, correcting suspected label errors without domain-specific knowledge risks obscuring meaningful patterns or rare cases that are needed in downstream decision-making. These challenges underscore that when using co-pilots with non-technical domain experts, there is a need for co-pilots to reason and update via expert human guidance along with incorporating data-centric tools.

Among existing frameworks, *OpenHands* (Wang et al., 2024) and *Data Interpreter* (Hong et al., 2024) are the closest to incorporating data-centric aspects and are particularly relevant due to their emphasis on multi-step reasoning and dynamic task execution.

**Challenges with Existing Co-pilots.**    Despite the progress demonstrated by OpenHands, Data Interpreter, and similar systems, several key challenges remain (C1-C4):

- **(C1) Overlooking Data-Centric Challenges:** Existing co-pilots often overlook data quality issues such as multi-measurements, noise, outliers, and missingness. In particular, they don't integrate state-of-the-art data-centric tools. They also fail to incorporate domain-specific reasoning for tasks requiring contextual interpretation, such as deciding how to deal with multiple measurements or whether a statistical anomaly is meaningful or erroneous. The integration of human expertise is vital for this contextual reasoning.

- **(C2) Static Workflow Architectures:** Many systems operate with predefined task structures, making them ill-suited for workflows where the data science flow is dependent on the unique challenges in the data or can be influenced via human expertise in a dynamic manner.

- **(C3) Healthcare-Specific Challenges:** The inability of these systems to contextualize healthcare data poses risks to using currently available frameworks. Some examples already mentioned include erroneous exclusion of clinically meaningful extreme lab values or data redundancy when retrieving data from electronic medical records or bioinformatic pipelines. Again, the integration of human expertise along with data-centric tools is vital in this regard.

- **(C4) Shallow Reasoning:** While these systems excel at automating task execution, they lack mechanisms for higher-level reasoning about data, such as validating correlations, diagnosing feature leakage, or ensuring robustness after data transformations.

**Data-centric AI.** Data-centric AI represents a paradigm shift in ML in which assessing and improving the quality of the data are prioritized over model-specific tasks (Liu et al., 2022; Liang et al., 2022; Zha et al., 2023; Whang et al., 2023; Seedat et al., 2023c). This paradigm has gained increasing importance within the ML community and has led to advances in methods and tools to systematically address issues, such as mislabeled samples (Seedat et al., 2023b; Northcutt et al., 2021; Pleiss et al., 2020), missing data (Jarrett et al., 2022b; Stekhoven & Bühlmann, 2012), outliers (Zhao et al., 2019; Yang et al., 2022), data leakage (Mitchell et al., 2019; Seedat et al., 2022d), and data drifts (Cai et al., 2023; Liu et al., 2023). These approaches have demonstrated improvements in model generalization in the context of real-world scenarios characterized by noisy and heterogeneous data.

Despite as the benefits of data-centric AI being demonstrated, existing LLM-based co-pilots adopt a model-centric perspective-focusing on automation of the model building pipeline, while neglecting the underlying data challenges. These limitations make existing co-pilots less effective for real-world applications where data quality directly impacts further modelling. We posit that the inclusion of data-centric AI tools in LLM-based co-pilots could significantly enhance their utility by automating the process of identifying data issues, improving dataset quality, and ensuring robust ML workflows. However, the autonomous application of these tools without contextual oversight can have unintended consequences.

Consequently, we advocate that data-centric AI tools are integrated into LLM-based co-pilots, while emphasizing the importance of the human-in-the-loop to contextualize and guide their usage.

## A.2 LLM-BASED CODE INTERPRETERS

- **GPT-Code Interpreter:** Aimed at simplifying tasks such as data visualization, basic modeling, and statistical analysis, this tool allows users to query and interact with datasets dynamically. However, its design is limited to single-step tasks and lacks support for multi-stage workflows, iterative refinement, or complex reasoning across interdependent tasks. Its applicability to real-world datasets with evolving requirements is minimal due to its static nature.

- **AutoGPT:** AutoGPT generalizes task execution by chaining multiple steps through autonomous prompts. It explores iterative workflows but relies heavily on predefined task templates. This rigidity makes it ill-suited for dynamic, data-centric environments where task dependencies evolve unpredictably. Moreover, AutoGPT lacks mechanisms to diagnose or correct data quality issues during execution.

- **BLADE:** Designed primarily as a benchmarking framework, BLADE evaluates LLM agents on open-ended scientific analyses and decision-making tasks. It provides insights into flexibility and task accuracy but does not address robustness or the ability to adapt workflows based on intermediate results. Furthermore, its scope is confined to task execution, neglecting data-centric challenges such as feature leakage or outlier handling.

- **DS-Agent:** Integrating case-based reasoning (CBR) with LLMs, DS-Agent automates ML workflows by leveraging prior knowledge from human-curated cases (e.g., Kaggle). It iteratively refines workflows by incorporating execution feedback. However, its dependency on retrieved cases limits its adaptability to novel or unstructured problems. DS-Agent's reliance on historical cases also makes it less effective for workflows requiring real-time adaptability or dynamic reasoning about data.

- **OpenHands:** OpenHands introduces a modular architecture for multi-agent collaboration and secure task execution in sandboxed environments. Its strengths lie in its flexibility and support for multi-step workflows, including software engineering tasks and web interaction. However, it lacks built-in tools for diagnosing and resolving data-centric issues, such as missing data or noise, and offers limited support for domain-specific reasoning, which is critical for high-stakes domains like healthcare.

- **Data Interpreter:** Data Interpreter employs hierarchical graph-based reasoning to model workflows as interdependent tasks, allowing for iterative refinement and robust task decomposition. This makes it highly effective for structured ML pipelines. However, its reliance on predefined task graphs limits its generalization to exploratory workflows or tasks with undefined dependencies. Additionally, it lacks direct integration of domain-specific insights, such as clinical knowledge for healthcare datasets.

## A.3 Limitations Across Approaches

Despite their individual strengths, these systems share several overarching limitations that hinder their applicability to real-world, data-centric workflows:

- **Static Pipeline Architectures:** Most interpreters rely on predefined templates or fixed task hierarchies, which restrict their ability to adapt to evolving requirements. For example, AutoGPT and DS-Agent struggle with workflows where task dependencies are contingent on intermediate results.

- **Insufficient Data Reasoning:** While these tools excel at executing predefined workflows, they lack higher-level data-centric reasoning capabilities, such as identifying feature correlations, addressing data drift, or diagnosing feature leakage. For instance, GPT-Code Interpreter and OpenHands fail to contextualize data preprocessing steps to account for domain-specific nuances.

- **Healthcare-Specific Challenges:** Healthcare datasets present unique challenges, including heterogeneity, noise, and biases. Generic preprocessing approaches risk introducing errors or obscuring critical clinical information. For example, extreme lab values might appear as statistical outliers but could signify a critical medical condition. Current systems fail to incorporate the domain expertise required to navigate such complexities.

- **Limited Adaptability to Data Evolution:** Real-world datasets often exhibit evolving distributions, feature sets, or objectives. Most interpreters, including BLADE and Data Interpreter, are designed for static workflows and do not account for the dynamic nature of these datasets.

- **Lack of Control in Open-Ended Scenarios:** Systems with open-ended prompting, such as AutoGPT and DS-Agent, can generate uncontrolled outputs when used by non-experts. This is particularly problematic in sensitive domains like healthcare, where errors can lead to significant consequences.

### A.4 TERMINOLOGY DIFFERENCES BETWEEN MACHINE LEARNING (ML) AND BIOSTATISTICS/EPIDEMIOLOGY

Table 6 demonstrates the different terminologies between communities that it is useful for a co-pilot to handle.

Table 6: Comparison of terminology in Machine Learning (ML) and Biostatistics/Epidemiology, which is considered in the design of CliMB-DC.

| Machine Learning (ML) | Biostatistics/Epidemiology |
|---|---|
| Model/Algorithm | Statistical/Predictive Model |
| Features | Covariates/Covariables |
| Targets | Outcomes/Endpoints |
| Training | Model Fitting/Estimation |
| Test Set | Validation Data |
| Overfitting | Overparameterization |
| Hyperparameters | Tuning Parameters |
| Performance Metrics | Goodness-of-Fit Measures |
| Cross-Validation | Internal Validation |
| Bias-Variance Tradeoff | Model Complexity |
| Generalization | External Validity |
| Feature Selection | Variable Selection |

## B DATASET DESCRIPTIONS

### B.1 LUNG CANCER DATASET

The dataset consists of **216714** records, capturing baseline and follow-up data related to lung cancer risk factors, smoking history, and demographic attributes. It includes **31 features**, broadly categorized as follows:

- **Demographic and Administrative:** This category includes age at recruitment, sex, ethnicity, and highest qualifications attained.
- **Smoking History:** Features include the number of cigarettes smoked per day, age at which smoking started and stopped, smoking duration, years since quitting, and pack-years (a cumulative measure of smoking exposure).
- **Respiratory and Comorbid Conditions:** This includes self-reported history of respiratory diseases such as asbestosis, pneumonia, chronic obstructive pulmonary disease (COPD), emphysema, chronic bronchitis, asthma, and allergic conditions (eczema, allergic rhinitis, hay fever).
- **Cancer History:** The dataset captures both personal and family history of lung cancer, including lung cancer diagnoses in parents (mother and father) and siblings, as well as the number of self-reported cancers and prior personal history of cancer.
- **Occupational and Environmental Exposure:** Presence of asbestos exposure is recorded as a binary indicator.
- **Lung Cancer Outcomes and Time-to-Event Data:** The dataset includes indicators for lung cancer diagnosis and related outcomes, along with survival-related features such as time to lung cancer diagnosis, time to death from lung cancer, and time to event development.

***Task.*** Lung cancer risk prediction and survival analysis

### B.2 PRIMARY BILIARY CHOLANGITIS (PBC) DATASET

The dataset consists of **43,834** records across 2181 patients, capturing baseline and follow-up data for individuals diagnosed with **Primary Biliary Cholangitis (PBC)**. It includes **33 features**, broadly categorized as follows:

- **Demographic and Administrative:** This category includes patient ID, sex, age, visit type, and time-related variables, which provide essential context for each recorded observation.
- **Clinical Outcomes:** Features in this category capture event status, survival status, and liver transplantation (LT) status, allowing for disease progression analysis.
- **Clinical Complications:** These features focus on manifestations of liver dysfunction, including decompensation (Decomp), variceal hemorrhage (VH), ascites, and hepatic encephalopathy (HE).
- **Treatment Variables:** This category records the use of Ursodeoxycholic Acid (UDCA), Obeticholic Acid (OCA), and Bezafibrate (BZF), which are commonly used interventions in PBC management.
- **Laboratory Measurements:** Example biomarkers include Albumin, Bilirubin, ALP, ALT, Platelets, Hemoglobin, White Cell Count, Urea, Creatinine, Sodium, Potassium, IgM, IgG, IgA
- **Comorbidity Assessment:** The Charlson Comorbidity Index (CCI) score is included as a prognostic measure for patient risk stratification.

Note the PBC dataset has repeat measurements for each patient that need to be aggregated before modelling.

***Task.*** Time-to-event modeling and survival analysis.

### B.3 PROSTATE CANCER PREDICTION: SEER AND CUTRACT DATASETS

This task focuses on **10-year prostate cancer mortality prediction** using two datasets: **SEER** (Surveillance, Epidemiology, and End Results) from the **United States** and **CUTRACT** from the **United Kingdom**. The goal is to assess how models trained in one country generalize when deployed in another, particularly in handling distribution shifts across different healthcare systems.

Both datasets are balanced to 2000 patient records each, with **10 features** related to patient demographics, cancer severity, treatment, and outcomes. The features are categorized as follows:

- **Demographic and Clinical Characteristics:** Age at diagnosis (`age`), baseline Prostate-Specific Antigen (`psa`), and presence of comorbidities (`comorbidities`).
- **Treatment Variables:** Indicators of different treatment modalities, including conservative management (`treatment_CM`), primary hormone therapy (`treatment_Primary hormone therapy`), radical radiotherapy (`treatment_Radical Therapy-RDx`), and radical surgery (`treatment_Radical therapy-Sx`).
- **Tumor Characteristics:** Cancer grade (`grade`) and stage (`stage`), which define disease severity.
- **Outcome Variable:** The binary label `mortCancer` represents **10-year mortality due to prostate cancer**, serving as the prediction target.

*Task.* 10-year prostate cancer prediction (binary classification). Given the multi-country for train and test it serves as an example of cross-national deployment of prediction models and their robustness to distribution shift.

# C CLiMB-DC ADDITIONAL DETAILS

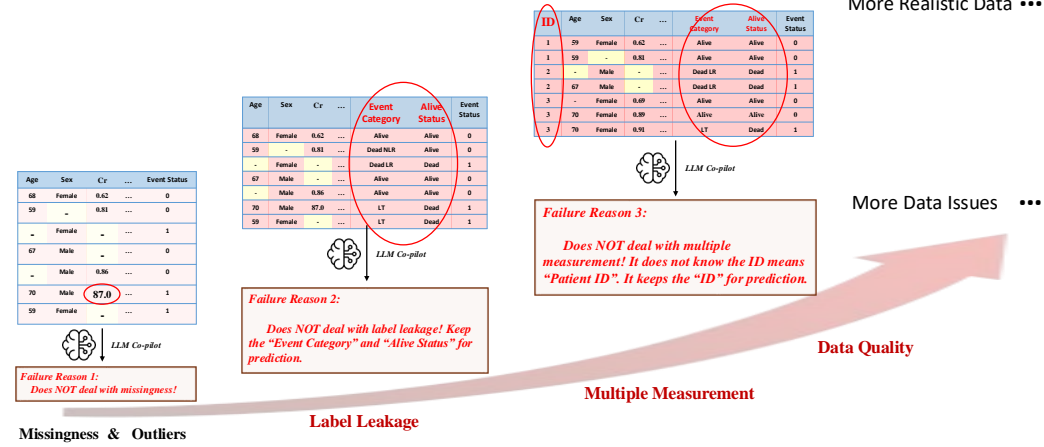## C.1 TAXONOMY OF DATA-CENTRIC ISSUES FACING CO-PILOTS



Figure 6: Illustrative examples of potential data issues in real-world healthcare scenarios, highlighting challenges at various levels and demonstrating how the current LLM co-pilot struggles to address these issues.

Recent advancements in LLM-based agents used in co-pilots have largely concentrated on code generation for model-centric issues, such as algorithm selection, hyperparameter tuning and performance evaluation. These processes take on datasets that have been assumedly clean (e.g., outliers removed, missingness assessed and handled, data errors removed, etc.) and problem setups that are well-defined for an ML task. However, transforming raw, sometimes disorganized, real-world datasets into clean, structured ones, while at the same time defining a clear problem setup is not necessarily trivial and can be complex, particularly for non-technical domain experts with limited experience in data science. Such data-centric challenges are precisely the area where co-pilots are expected to provide significant support, yet have been overlooked.

**Motivated examples from healthcare.** In healthcare scenarios, it is common for some variables collected during data acquisition to be highly correlated with the outcome or to have been measured only after the outcome occurred. Including such variables in predictive models can lead to label leakage, compromising the model's validity. Consequently, these variables must be carefully excluded during model construction. As illustrated in Figure 6 (middle), current LLM co-pilot fail to exclude variables such as "Event Category" and "Alive Status", which are highly correlated with the outcome "Event Status". Including these variables results in exceptionally high predictive performance, which is a misleading conclusion for users. Similarly, healthcare datasets often contain multiple records for a single patient, as one patient could come to the hospital multiple times for follow-ups or during a chronic condition. . However, current LLM co-pilots do not automatically perform aggregation to handle such cases. Being unable to appropriately account for these data structures can result in severe label leakage and render the problem setup meaningless, as demonstrated in Figure 6 (right).

Beyond label leakage, data-centric challenges in ML —including issues with data quality, preprocessing, and curation—are particularly pronounced in healthcare. These datasets are often collected by clinicians with limited data science expertise, rather than by experienced data scientists. Some datasets are retrieved from bioinformatic pipelines, which could have problems with certain reads or even produce invalid measurements. As a result, data are frequently incomplete and noisy, but usually in a context-dependent manner. The complexity when processing these data necessitates domain-specific expertise, assessment and handling. However, such challenges remain under-explored in the field of co-pilots.

**Key perspectives for ensuring reliable LLM co-pilots.** In this work, we present a formalized taxonomy of key issues that LLM co-pilots must address to enable reliable deployment in healthcare scenarios. Our taxonomy follows a bottom-up approach, drawing on a broad survey of literature where these challenges have been extensively documented and analyzed (Zadorozhny et al., 2022;

Table 7: Taxonomy of key data-centric challenges frequently encountered in healthcare machine learning pipelines. While not exhaustive, these categories represent a significant fraction of issues that co-pilots must address to ensure strong predictive performance, robustness, fairness, and clinical feasibility.

| Category | Issues | Description | References highlighting issue | Resultant Issues |
|---|---|---|---|---|
| Data-Centric (Formatting) | Multiple measurements | Challenges from datasets including multiple observations for a single individual, requiring aggregation and standardization. | (Tschalzev et al., 2024; Liu et al., 2024; Oufattole et al., 2024; Sett et al., 2024) | Ill-posed problem setup, Temporal misalignment, Potential data inflation |
| | Multiple files | Datasets from different sources/periods of time need to be correctly aggregated or harmonized across files. | (Balagopalan et al., 2024b; Lehne et al., 2019; Nan et al., 2022; Schmidt et al., 2020) | Ill-posed problem setup, Inconsistent representation, Duplication risk |
| | Inconsistent data | Data might be inconsistent based on units or how data might be represented. | (Rychert, 2023; Monjas et al., 2025; Szarfman et al., 2022) | Ill-posed problem setup, Label leakage, Reduced reproducibility |
| | Data extraction | Data might be stored in heterogenous text fields and needs to be extracted as features. | (Bao et al., 2018; Zhao, 2019; Hahn & Oleynik, 2020) | Ill-posed problem setup, Inconsistent representation |
| | Feature redundancy | Multiple features conveying similar information in a dataset. | (Chicco et al., 2022; Apicella et al., 2024; Meng et al., 2022; Sasse et al., 2023) | Poor generalization, Poor interpretability, Label leakage |
| Data-Centric (Statistical - Train) | Outliers | Extraordinary values (leading to soft outliers) or mistakes in the data creation process (possibly leading to hard outliers). | (Zadorozhny et al., 2022; Avati et al., 2021; Estiri & Murphy, 2019) | Overfitting, Misleading performance metrics, Potential data misinterpretation |
| | Label leakage | Features can include future information or tests dependent on the outcome, or datasets can have multiple correlated outcome variables. | (Tomašev et al., 2019; Ghassemi et al., 2020) | Ill-posed problem setup, Over-optimistic performance, Failed clinical deployment |
| | Missingness | Missing values caused by not being recorded (MCAR), later feature aggregation (MAR), or differing clinical practices (MNAR). | (Beaulieu-Jones et al., 2017; Ferri et al., 2023; Singh et al., 2021; Haneuse et al., 2021) | Imputation risk, Model bias, Reduced external validity |
| | Noisy labels | Incorrect labels caused by erroneous annotation, recording mistakes, or difficulty in labeling. | (Yang et al., 2023; Wei et al., 2024; Boughorbel et al., 2018) | Poor generalization, Compromised interpretability, Unstable model calibration |
| | Data valuation | General data quality issues impacting model performance. | (Bloch et al., 2021; Enshaei et al., 2022; Tang et al., 2021; Pandl et al., 2021) | Poor generalization, Suboptimal performance, High curation overhead |
| Data-Centric (Statistical - Test) | Subgroup challenges | Poor performance or generalization on certain subgroups (in-distribution heterogeneity). | (Oakden-Rayner et al., 2020; Suresh et al., 2018; Goel et al., 2020; Cabrera et al., 2019; van Breugel et al., 2024) | Poor generalization, Fairness concerns |
| | Data shift | Changes due to novel equipment, different measurement units, or clinical practice evolution over time. | (Pianykh et al., 2020; Koh et al., 2021; Patel et al., 2008; Goetz et al., 2024) | Poor generalization, Model bias, Need for continuous monitoring |

Avati et al., 2021; Estiri & Murphy, 2019; Tomašev et al., 2019; Ghassemi et al., 2020; Beaulieu-Jones et al., 2017; Ferri et al., 2023; Singh et al., 2021; Haneuse et al., 2021). After synthesizing insights from these diverse studies and their practical applications, we present a structured taxonomy, highlighting the most pressing data-centric challenges affecting ML workflows. As shown in Table 7, these perspectives address both data-centric and model-centric aspects.

On the *data-centric* side, we highlight elements related to data formatting, as well as statistical (both training and test). When an LLM co-pilot fails to address these data issues effectively, it can lead to a range of problems. These include issues with the final ML model (e.g., overfitting, model bias, poor generalization, and limited interpretability) along with flaws in experimental setups (e.g., improper formulation of the problem and label leakage (see case study 1 and 2 in Section 4.2)).

While not included in the table, we also note there do remain model-centric challenges. While algorithm selection, hyperparameter tuning, and performance evaluation, have been frequently discussed and relatively well-covered in recent LLM co-pilots, there should also be a focus on *domain-specific model classes* and *model interpretability*. Different from typical data science tasks that mainly focus on classification and regression: domain-specific model classes account for temporal dependencies, hierarchical structures, and clinical context, ensuring that models are both accurate and practically applicable. These issues arise frequently in data from healthcare settings. For instance,

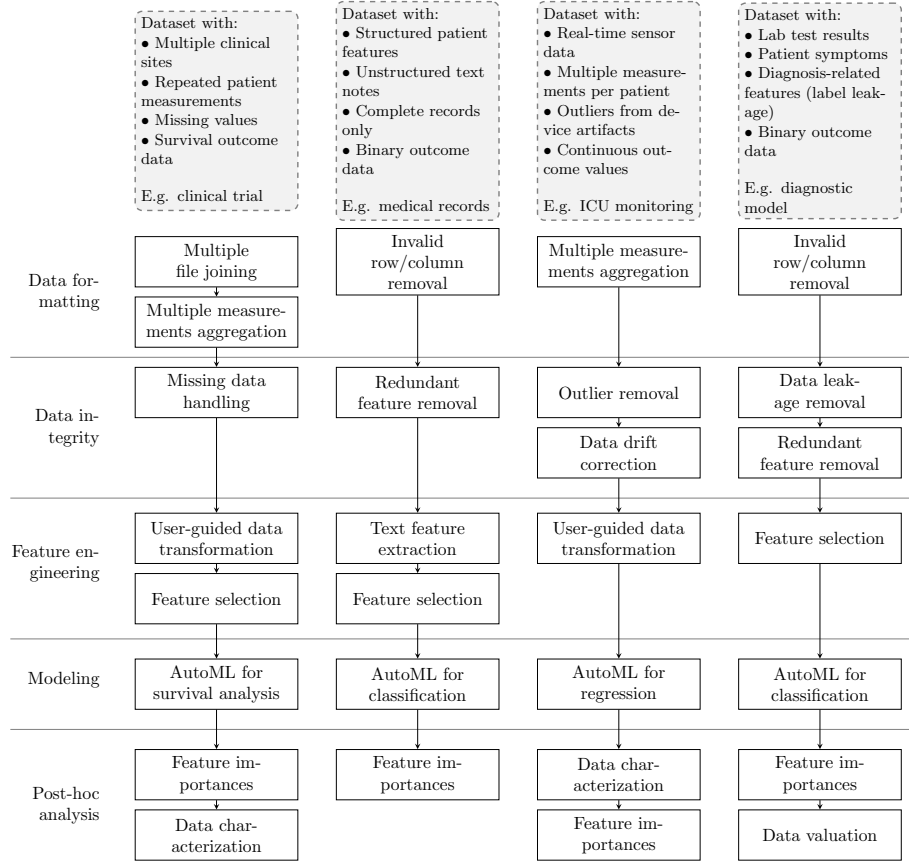| | Dataset with:<br>• Multiple clinical sites<br>• Repeated patient measurements<br>• Missing values<br>• Survival outcome data<br><br>E.g. clinical trial | Dataset with:<br>• Structured patient features<br>• Unstructured text notes<br>• Complete records only<br>• Binary outcome data<br><br>E.g. medical records | Dataset with:<br>• Real-time sensor data<br>• Multiple measurements per patient<br>• Outliers from device artifacts<br>• Continuous outcome values<br><br>E.g. ICU monitoring | Dataset with:<br>• Lab test results<br>• Patient symptoms<br>• Diagnosis-related features (label leakage)<br>• Binary outcome data<br><br>E.g. diagnostic model |
|---|---|---|---|---|
| Data formatting | Multiple file joining<br>Multiple measurements aggregation | Invalid row/column removal | Multiple measurements aggregation | Invalid row/column removal |
| Data integrity | Missing data handling | Redundant feature removal | Outlier removal<br>Data drift correction | Data leakage removal<br>Redundant feature removal |
| Feature engineering | User-guided data transformation<br>Feature selection | Text feature extraction<br>Feature selection | User-guided data transformation | Feature selection |
| Modeling | AutoML for survival analysis | AutoML for classification | AutoML for regression | AutoML for classification |
| Post-hoc analysis | Feature importances<br>Data characterization | Feature importances | Data characterization<br>Feature importances | Feature importances<br>Data valuation |

Figure 7: Addressing real data challenges is complex and requires multi-step reasoning.

specialized models are designed for survival analysis, a critical and widely applied task in healthcare. Moreover, the role of interpretability is to ensure that predictive models can provide transparent and actionable insights, which is crucial for enabling clinicians to trust and validate their decision-making.

Our taxonomy consequently offers a systematic foundation for challenges that co-pilots should address and hence should impact the design and evaluation of LLM co-pilots. Specifically, we posit that the structured taxonomy will enable the development of co-pilots that are better equipped to handle real-world data issues, ultimately fostering more reliable, interpretable, and impactful ML systems in high-stakes domains like healthcare.

**Challenges vary by problem and context.** While the taxonomy of challenges describes each issue in an isolated manner, real-world scenarios often require a more integrated approach. When a co-pilot addresses a user's task, the challenges are inherently problem- and context-dependent, requiring end-to-end consideration. As illustrated in Figure 7, there can be multiple data issues, which must be handled in a nuanced manner, thereby making real-world applications complex. Consequently, systems must be capable of reasoning about these challenges autonomously, while gathering and considering expert human feedback.

## C.2 CLIMB-DC: OPEN-SOURCE SOFTWARE TOOLKIT

Beyond usage by diverse users and improved performance, a an important aspect for CliMB-DC for impact in healthcare is its role as a software toolkit to empower domain experts. Consequently, a key aspect is the open-source nature of the framework, which enables the community to contribute and integrate new tools to extend its capabilities.

To achieve this goal of empowerment for diverse users, three software challenges are vital to address: *extensibility to new tools*, *human integration*, and *support for diverse predictive tasks* in medicine, specifically *classification*, *survival analysis*, and *regression*. This enables a more robust and user-friendly system for clinical predictive modeling.

### C.2.1 EXTENSIBILITY TO NEW TOOLS

The diversity and rapid development of data-centric tools means the framework must be capable of integrating new tooling from the community with minimal effort.

**Data-Centric Tool Support:** CliMB-DC emphasizes a data-centric AI approach by integrating specialized tools (see Table 2) that enhance dataset quality including:

- *Imputation Tools*: Frameworks like HyperImpute handle missing data with advanced iterative imputation techniques.

- *Exploratory Analysis and data quality evaluation*: Tools such as Data-IQ enable detailed subgroup analysis, data heterogeneity and noisy labels.

- *Interpretability*: Built-in post-hoc interpretability methods like SHAP and permutation explainers ensure models remain transparent and actionable.

**Extensibility:** (i) We have a *tool registry* which catalogs the available tools, their supported predictive tasks, and their data requirements, enabling users to easily incorporate new methods without modifying core system logic. (ii) Modular APIs allow developers to register new tools, ensuring that CliMB-DC evolves alongside data-centric advances. (iii) The Open-source architecture encourages community contributions to expand the ecosystem of available data-centric tools to the co-pilot. This enables broader accessibility to data-centric tools for non-technical domain experts thereby empowering them. Additionally, it provides an opportunity for the data-centric ML research community to easily incorporate new tools and/or validate their tools, facilitating research impact via usage on diverse applications.

### C.2.2 HUMAN INTEGRATION THROUGH UI AND FEEDBACK

More complex ML frameworks generally require a wider range of skill sets that are often lacking by non-technical domain experts, whereas the setup of more complex biological research questions risk being misunderstood by technical domain experts. One way of minimizing the impact of these limitations is creating a user interface (UI) that allows both mutual understanding of the tasks between users and integrates specific feedback from the type of user.

**User Interface** The UI for CliMB-DC combines output from natural, conversational language, along with updates on the progress of the task pipeline accompanied with visualizations (see Figure 8). This type of interface provides non-technical domain experts the opportunity to perform tasks that they might not be able to do directly with an ML tool and technical domain experts the opportunity to examine more closely which ML procedures were effectuated.

Oftentimes, users with technical and non-technical domain expertise employ different terminology for the same task or problem at hand, obfuscating the processes needed to complete the task or solve the given problem. For example, users from the ML community might adapt a feature selection process prior to or during predictive modeling, while those from the epidemiology community would refer to this process as model building (see Table 6 for more examples). The fact that users can communicate desired processes in their natural language makes it possible to carry out the intended task, supporting its nuances and bringing about a more fluid user experience.

**Dynamic Plan Refinement via expert guidance** CliMB-DC importantly incorporates *human feedback* into its reasoning. The user can refine the data science pipeline in an iterative manner by weighing in their expertise on a variety of processes (e.g., on data transformations, feature selection, or model evaluation). The iterative feedback dynamically adjusts the plan, ensuring alignment with domain-specific goals and that outputs are clinically relevant.
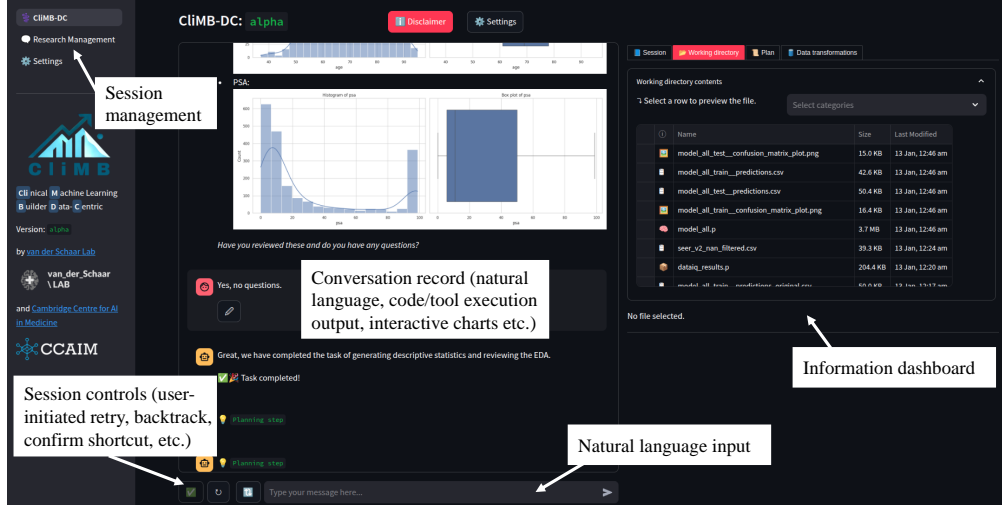
Figure 8: The user interface of CliMB-DC, which supports natural language input, multi-modal conversation record and dashboard, session controls (including user-initiated retry and backtracking), and session management across multiple conversations.

### C.2.3  SUPPORT FOR DIVERSE PREDICTIVE TASKS

Clinical datasets require predictive modeling across varied tasks, including *classification*, *survival analysis*, and *regression*, each with distinct data processing and modeling needs. Recall that in carrying out these analytics, non-technical domain experts like clinical researchers, biostatisticians, epidemiologists etc, do not need to do ML analytics. Rather CliMB-DC facilitates this with domain expertise used to guide and validate the process.

### C.3  CLIMB ALGORITHM

---

**Algorithm 1** CliMB-DC Optimization with Expert Integration

---

**Require:** Initial dataset $\mathcal{D}^{(0)}$, tools $\mathcal{T}$
**Ensure:** Curated dataset $\mathcal{D}^*$
1: Initialize $\mathcal{S}^{(0)} = \{\mathcal{D}^{(0)}, \varnothing, \mathcal{P}^{(0)}, \varnothing, \mathcal{T}\}$
2: $\mathcal{D}^* \leftarrow \mathcal{D}^{(0)}$
3: **while** not converged **do**            ▷ Coordinator reasoning phase
4:      $\mathcal{O}^{(t)} \leftarrow \text{STATEOBSERVE}(\mathcal{S}^{(t)})$            ▷ see Section 3.2
5:      $\beta^{(t)} \leftarrow \text{ASSESSBACKTRACK}(\mathcal{O}^{(t)})$
6:      **if** $\beta^{(t)} = 1$ **then**
7:          $(\mathcal{D}^*, \mathcal{S}^{(t)}) \leftarrow \text{RESTORECHECKPOINT}(\mathcal{H}^{(t)})$
8:          continue
9:      **end if**
10:      $\mathcal{P}^{(t)} \leftarrow \text{PLANNING}(\mathcal{O}^{(t)})$
11:      **while** not episode_complete **do**            ▷ Worker execution phase
12:          $f_t \leftarrow \mathcal{W}.\text{PROPOSETRANSFORM}(\mathcal{S}^{(t)}, \mathcal{P}^{(t)}, \mathcal{T}^{(t)})$
13:          $\mathcal{D}^* \leftarrow f_t(\mathcal{D}^*)$
14:      **end while**
15:      $\mathcal{S}^{(t+1)} \leftarrow \text{UPDATESTATE}(\mathcal{S}^{(t)}, \mathcal{D}^*, \mathcal{H}^{(t)})$            ▷ Update system state
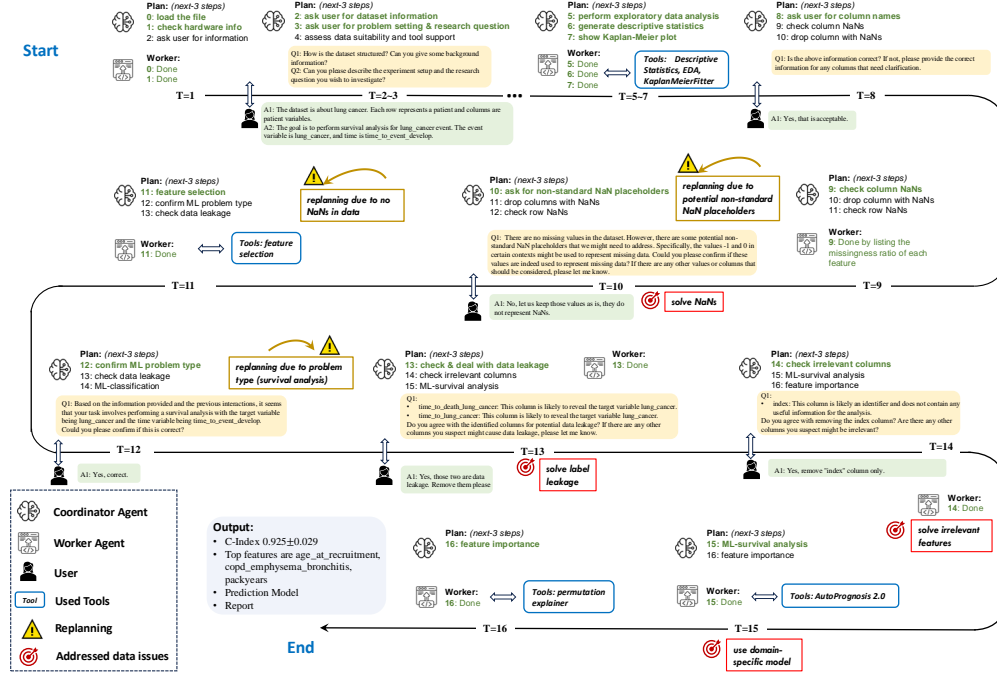16: **end while**

---

# D CLiMB-DC CASE STUDY DEEP DIVES



Figure 9: The processing workflow of CliMB-DC on the Lung Cancer dataset, illustrating how the coordinator agent, worker agent, and user interact at each processing stage $T$. The data issues include missingness, feature redundancy, and label leakage. The prediction task is survival analysis, requiring domain-specific model classes.
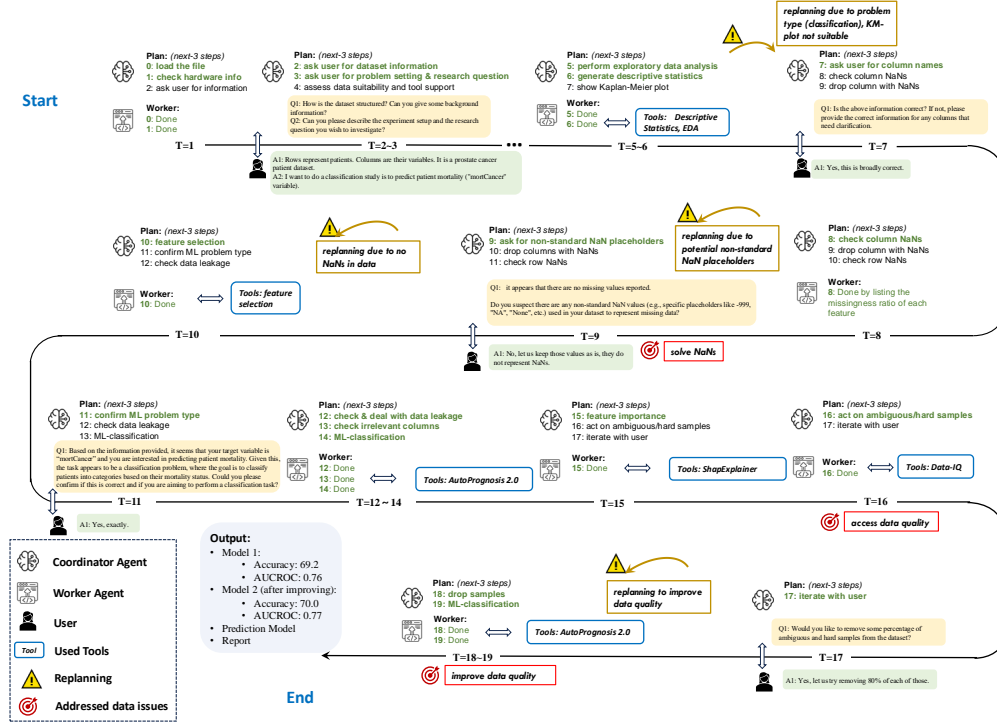


Figure 10: The processing workflow of CliMB-DC on the Prostate Cancer dataset, illustrating how the coordinator agent, worker agent, and user interact at each processing stage $T$. The data issues include missingness and data quality/drifts. The prediction task is classification.

# E  INTERACTION EXAMPLES

In this section, we present execution logs from both the baselines and CliMB-DC to illustrate two key points: (1) why the baselines fail to address data issues effectively, and (2) how CliMB-DC successfully tackles these challenges in healthcare. We first show examples of baselines, i.e. Data-Interpreter and OpenHands, and then provide the corresponding successful examples of CliMB-DC.

## E.1  DATA-INTERPRETER

### E.1.1  EXAMPLE: INCORRECT METRIC USED

In the following example (lung cancer dataset), we observe a failure case with Data-Interpreter, where it used the ROC-AUC score as a proxy for the C-index score. This issue arises due to the absence of Python packages for survival analysis. Since Data-Interpreter lacks user interaction capabilities, the system failed to detect and address this error.

---

**Summary of work up to this point:** The datasets have been successfully loaded and initially preprocessed. However, despite multiple attempts to utilize Python packages for survival analysis, such as `lifelines` and `sksurv`, the installation process has not been triggered.

...

**Generated code:**

```
1   import xgboost as xgb
2   from sklearn.model_selection import train_test_split
3   from sklearn.metrics import roc_auc_score
4   import numpy as np
5   # Prepare the data for survival analysis
6   X_train = train_df_copy.drop(columns=['time_to_event_develop', '
        lung_cancer'])
7   y_train = train_df_copy[['time_to_event_develop', 'lung_cancer']]
8   X_eval = eval_df_copy.drop(columns=['time_to_event_develop', '
        lung_cancer'])
9   y_eval = eval_df_copy[['time_to_event_develop', 'lung_cancer']]
10  # Convert the target to a structured array for xgboost
11  y_train_structured = np.array([(time, event) for time, event in zip
        (y_train['time_to_event_develop'], y_train['lung_cancer'])],
12                                  dtype=[('time', 'float'), ('event', '
                                      bool')])
13  y_eval_structured = np.array([(time, event) for time, event in zip(
        y_eval['time_to_event_develop'], y_eval['lung_cancer'])],
14                                  dtype=[('time', 'float'), ('event', '
                                      bool')])
15  # Create DMatrix for xgboost
16  dtrain = xgb.DMatrix(X_train, label=y_train['lung_cancer'])
17  deval = xgb.DMatrix(X_eval, label=y_eval['lung_cancer'])
18  # Set up parameters for xgboost
19  params = {
20      'objective': 'binary:logistic',
21      'eval_metric': 'auc',
22      'max_depth': 3,
23      'eta': 0.1,
24      'subsample': 0.8,
25      'colsample_bytree': 0.8,
26      'seed': 42
27  }
28  # Train the xgboost model
29  bst = xgb.train(params, dtrain, num_boost_round=100)
30  # Predict the risk scores for the evaluation set
31  risk_scores = bst.predict(deval)
32  # Calculate the concordance index
33  c_index = roc_auc_score(y_eval['lung_cancer'], risk_scores)
```

```
34    print(f"C-index: {c_index:.4f}")
```

**Code execution output:**
Code execution finished successfully ✓

```
1   C-index: 0.5337
```

**Reason behind this.** Data-Interpreter lacks the ability to interact with users, making it difficult to trigger the installation of required Python packages when they are missing. Additionally, since survival analysis is not as common as standard machine learning tasks like classification and regression, the system fails to recognize how to compute the C-index score. As a result, it incorrectly applies an inappropriate metric without detecting the issue.

### E.1.2 EXAMPLE: FAIL TO SETUP SURVIVAL ANALYSIS

In the following example (lung cancer dataset), we observe a failure case with Data-Interpreter, where it *failed to setup the survival analysis*. The system "auto-degraded" to XGBoost and failed to address this error.

---

**Summary of work up to this point:** The datasets have been successfully loaded and initially pre-processed. Python packages for survival analysis, such as `lifelines` and `sksurv`, are *manually* installed. However, it encountered "Convergence error" when fitting `CoxPHFitter` due to the high collinearity.

. . .

**Generated code:**

```
1    from xgboost import XGBRegressor
2    from lifelines.utils import concordance_index
3    import pandas as pd
4    # Prepare the data
5    X_train = train_data_var.drop(columns=['time_to_event_develop', '
         lung_cancer'])
6    y_train = train_data_var[['time_to_event_develop', 'lung_cancer']]
7    X_eval = eval_data_var.drop(columns=['time_to_event_develop', '
         lung_cancer'])
8    y_eval = eval_data_var[['time_to_event_develop', 'lung_cancer']]
9    # Train the XGBoost model
10   xgb_model = XGBRegressor(objective='survival:cox', n_estimators
         =100, learning_rate=0.1, max_depth=3, subsample=0.8,
         colsample_bytree=0.8)
11   xgb_model.fit(X_train, y_train['time_to_event_develop'])
12   # Predict risk scores
13   eval_data_combined['risk_score'] = xgb_model.predict(X_eval)
14
15   # Calculate the C-index
16   c_index = concordance_index(y_eval['time_to_event_develop'], -
         eval_data_combined['risk_score'], y_eval['lung_cancer'])
17   print(f"C-index: {c_index}")
```

**Code execution output:**
Code execution finished successfully ✓

```
1   C-index: 0.5000
```

**Reason behind this.** Data-Interpreter lacks both data-centric tools and sufficient knowledge of data preprocessing to address collinearity effectively. As a result, even when provided with the appropriate

survival analysis package (`lifelines`), it repeatedly encounters "Convergence errors." Eventually, it abandons the use of `lifelines` and reverts to using XGBoost.

### E.1.3 Example: Fail to address label leakage

In the following example (PBC dataset), we observe a failure case with Data-Interpreter, where it *failed to address the data aggregation as well as the label leakage*. The system did not recognize such issues.

---

**Summary of work up to this point:** The datasets have been successfully loaded and initially pre-processed. Python packages for survival analysis, such as `lifelines` and `sksurv`, are *manually* installed.

**Generated code:**

```python
1   # Evaluate the model using 5-fold cross-validation and report the C
        -index
2   import pandas as pd
3   from metagpt.tools.libs.data_preprocess import FillMissingValue
4   from sklearn.model_selection import KFold
5   from sklearn.preprocessing import LabelEncoder, StandardScaler
6   from lifelines import CoxPHFitter
7   from lifelines.utils import concordance_index
8   import numpy as np
9   # Load and preprocess the data
10  data_path = '/shared/share_mala/jiashuo/MetaGPT/data/di_dataset/
        copilot/PBCData.csv'
11  df = pd.read_csv(data_path, encoding='ISO-8859-1')
12
13  # Fill missing values
14  df_copy = df.copy()
15  numeric_features = df_copy.select_dtypes(include=['float64', 'int64
        ']).columns.tolist()
16  fill_missing_numeric = FillMissingValue(features=numeric_features,
        strategy='mean')
17  df_copy = fill_missing_numeric.fit_transform(df_copy)
18  categorical_features = df_copy.select_dtypes(include=['object']).
        columns.tolist()
19  fill_missing_categorical = FillMissingValue(features=
        categorical_features, strategy='most_frequent')
20  df_copy = fill_missing_categorical.fit_transform(df_copy)
21  # Encode categorical features
22  label_encoders = {}
23  for col in ['Visit_Type', 'Sex', 'Event_Category', 'Alive_status',
        'LT', 'Decomp', 'VH', 'Ascites', 'HE']:
24      le = LabelEncoder()
25      df_copy[col] = le.fit_transform(df_copy[col].astype(str))
26      label_encoders[col] = le
27  # Fill remaining missing values
28  df_copy.fillna(df_copy.mean(numeric_only=True), inplace=True)
29  for col in ['Visit_Type', 'Sex', 'Event_Category', 'Alive_status',
        'LT', 'Decomp', 'VH', 'Ascites', 'HE']:
30      df_copy[col].fillna(df_copy[col].mode()[0], inplace=True)
31
32  # Drop low variance columns
33  low_variance_cols = [col for col in df_copy.columns if df_copy[col
        ].var() < 1e-5]
34  df_copy.drop(columns=low_variance_cols, inplace=True)
35
36  # Prepare features and target
37  X = df_copy.drop(columns=['Event_status', 'Years_to_event'])
38  y = df_copy[['Event_status', 'Years_to_event']]
39  # Standardize features
```

```
40    scaler = StandardScaler()
41    X_scaled = scaler.fit_transform(X)
42    X_scaled = pd.DataFrame(X_scaled, columns=X.columns)
43
44    # Perform 5-fold cross-validation
45    kf = KFold(n_splits=5, shuffle=True, random_state=42)
46    c_indices = []
47
48    for train_index, test_index in kf.split(X_scaled):
49        X_train, X_test = X_scaled.iloc[train_index], X_scaled.iloc[
              test_index]
50        y_train, y_test = y.iloc[train_index], y.iloc[test_index]
51
52        train_data = X_train.copy()
53        train_data['Event_status'] = y_train['Event_status']
54        train_data['Years_to_event'] = y_train['Years_to_event']
55        test_data = X_test.copy()
56        test_data['Event_status'] = y_test['Event_status']
57        test_data['Years_to_event'] = y_test['Years_to_event']
58
59        cph = CoxPHFitter(penalizer=0.1)
60        cph.fit(train_data, duration_col='Years_to_event', event_col='
              Event_status')
61
62        test_predictions = cph.predict_partial_hazard(test_data)
63        c_index = concordance_index(test_data['Years_to_event'], -
              test_predictions, test_data['Event_status'])
64        c_indices.append(c_index)
65    mean_c_index = np.mean(c_indices)
66    print(f'Mean C-index from 5-fold cross-validation: {mean_c_index}')
```

**Code execution output:**
```
Code execution finished successfully ✓
```
```
1  Mean C-index from 5-fold cross-validation: 0.7750170186750802
```

**Reason behind this.** Although Data-Interpreter successfully fitted a Cox regression model for survival analysis, it failed to address key data issues in the PBC dataset. Specifically, the PBC data contains multiple observations per patient ID, requiring aggregation prior to analysis. Additionally, besides the "Event_status" and "Years_to_event" features, there are "Event_Category" and "Alive_status" features that are highly correlated with the target variable, leading to label leakage. However, Data-Interpreter did not account for these issues, resulting in flawed analysis.

### E.2 OPENHANDS

In the following example (lung cancer dataset), we observe a failure case with OpenHands, where it (1) *failed to appropriately preprocess the data* and (2) *failed to deal with label leakage*. The system did not recognize or deal with such issues.

**Summary of work up to this point:** The datasets have been successfully loaded and initially preprocessed.

```
1  import pandas as pd
2  from lifelines import CoxPHFitter
3  from lifelines.utils import concordance_index
4
5  # Load the training and testing datasets
6  train_data = pd.read_csv('/workspace/copilot/lung_train.csv')
7  test_data = pd.read_csv('/workspace/copilot/lung_test.csv')
8
```

```
9  # Prepare the data for survival analysis
10 X_train = train_data.drop(columns=['time_to_event_develop', '
       lung_cancer'])
11 y_train = train_data[['time_to_event_develop', 'lung_cancer']]
12
13 # Fit the Cox Proportional Hazards model
14 cph = CoxPHFitter()
15 cph.fit(y_train, duration_col='time_to_event_develop', event_col='
       lung_cancer')
16
17 # Prepare the test data
18 X_test = test_data.drop(columns=['time_to_event_develop', 'lung_cancer
       '])
19 y_test = test_data[['time_to_event_develop', 'lung_cancer']]
20
21 # Predict the risk scores for the test data
22 risk_scores = cph.predict_partial_hazard(X_test)
23
24 # Calculate the C-index
25 c_index = concordance_index(y_test['time_to_event_develop'], -
       risk_scores, y_test['lung_cancer'])
26
27 c_index
```

**Code execution output:**
```
Code execution finished successfully ✓
```
```
1  np.float64(0.5)
```

**Assistant:** The survival prediction model for lung cancer has been built using the Cox Proportional Hazards model. The C-index reported on the test dataset is 0.5. However, the model summary is empty, which may indicate that the model did not find any significant predictors in the training data. If you need further analysis or adjustments, please let me know!

**Reason behind this.**   The prediction model developed by OpenHands demonstrates subpar performance. This is primarily due to the lack of any data preprocessing, resulting in an ineffective and meaningless survival model. Besides, OpenHands did not recognize the label leakage problem caused by feature "time_to_death_lung_cancer" and "time_to_lung_cancer". Moreover, OpenHands failed to identify the label leakage issue caused by the features "time_to_death_lung_cancer" and "time_to_lung_cancer". The failure scenarios closely resemble those observed with Data-Interpreter.

### E.3   SUCCESSFUL CASES WITH CLIMB-DC

Here, we demonstrate how CliMB-DC effectively addresses the previously discussed failure cases, using the lung cancer dataset as an example.

#### E.3.1   EXAMPLE: ADDRESS LABEL LEAKAGE

This example demonstrates how CliMB-DC effectively addresses label leakage issues via advanced reasoning and human interaction.

---

**Summary of Progress to Date:** The datasets have been successfully loaded and appropriately preprocessed, including steps such as imputation, normalization, and other necessary adjustments. Before initiating experiments, CliMB-DC ensures alignment with the user by confirming the problem setup, while also checking for potential label leakage.

. . .

**Assistant:** I will now consult the message history to check the meaning and details of each of these columns and identify any potential data leakage.

**Suspected Data Leakage Columns**:

- "time_to_death_lung_cancer": This column is likely to reveal the target variable "lung_cancer".
- "time_to_lung_cancer": This column is likely to reveal the target variable "lung_cancer".

Data leakage is a problem because it can lead to overly optimistic results and a model that cannot be used in practice. By including information that would not be available at the time of prediction, the model can learn patterns that are not generalizable to new data.

---

**Assistant:** Do you agree with the identified columns for potential data leakage? If there are any other columns you suspect might cause data leakage, please let me know.

---

**User:** Yes, those two are data leakage. Remove them please

---

**Assistant:** I will now generate the code to remove the identified columns ("time_to_death_lung_cancer" and "time_to_lung_cancer") from both the training and test datasets.

**Generated code:**

```python
import pandas as pd

# Load the training dataset
train_df = pd.read_csv("lung_train.csv")

# Drop the identified columns
train_df = train_df.drop(columns=["time_to_death_lung_cancer", "time_to_lung_cancer"])

# Save the modified training dataset
train_df.to_csv("lung_train_no_leakage.csv", index=False)

# Load the test dataset
test_df = pd.read_csv("lung_test.csv")

# Drop the identified columns
test_df = test_df.drop(columns=["time_to_death_lung_cancer", "time_to_lung_cancer"])

# Save the modified test dataset
test_df.to_csv("lung_test_no_leakage.csv", index=False)

print("Columns 'time_to_death_lung_cancer' and 'time_to_lung_cancer' have been removed from both datasets.")
```

**Code execution output:**
Code execution finished successfully ✓

```
Columns ``time_to_death_lung_cancer'' and ``time_to_lung_cancer'' have been removed from both datasets.
```

. . .

**Summary of Subsequent Work:** CliMB-DC then went on to identifying irrelevant features to remove. An illustrative example is provided below.

### E.3.2 EXAMPLE: IDENTIFY IRRELEVANT FEATURES

This example demonstrates how CliMB-DC identified irrelevant features and interacted with the user to remove them.

---

**Summary of Progress to Date:** The datasets have been appropriately preprocessed and the label leakage issue has been addressed based on expert human feedback. Moving forward, CliMB-DC focuses on identifying and addressing irrelevant features.

. . .

---

**Assistant:** I will now check if there are any meaningless/irrelevant columns still left in the dataset.

**Suspected Irrelevant Columns:**

- "index": This column is likely an identifier and does not contain any useful information for the analysis.

Including irrelevant columns can lead to overfitting and misleading feature importance. It is important to remove them before continuing with the machine learning study.

---

**Assistant:** Do you agree with removing the "index" column? Are there any other columns you suspect might be irrelevant?

---

**User:** Yes, remove "index" column only.

---

**Assistant:** I will now generate the code to remove the "index" column from both the training and test datasets.
**Generated code:**

```python
import pandas as pd

# Load the training dataset
train_df = pd.read_csv("lung_train_no_leakage.csv")

# Drop the identified irrelevant column
train_df = train_df.drop(columns=["index"])

# Save the modified training dataset
train_df.to_csv("lung_train_no_leakage_no_index.csv", index=False)

# Load the test dataset
test_df = pd.read_csv("lung_test_no_leakage.csv")

# Drop the identified irrelevant column
test_df = test_df.drop(columns=["index"])

# Save the modified test dataset
test_df.to_csv("lung_test_no_leakage_no_index.csv", index=False)

print("Column 'index' has been removed from both datasets.")
```

**Code execution output:**
Code execution finished successfully ✓

```
1  Column ``index'' has been removed from both datasets.
```

**Assistant:** The index column has been removed from both the training and test datasets. The modified datasets have been saved as lung_train_no_leakage_no_index.csv and lung_test_no_leakage_no_index.csv.

**Feature Columns Left**: Here are the feature columns that are left in the dataset:

- age_at_recruitment
- sex
- body_mass_index
- n_cig_per_day
- ...

Are you happy to use all of these features in the machine learning study?

**User:** Yes

. . .

**Summary of Subsequent Work:** CliMB-DC then went on to utilizing domain-specific models tailored for survival analysis. An illustrative example is provided below.

### E.3.3 EXAMPLE: SETUP DOMAIN-SPECIFIC MODEL & METRIC

Following the above examples, this example shows how CliMB-DC leveraged domain-specific tools to build survival models effectively.

**Summary of Progress to Date:** The datasets have been appropriately preprocessed and the label leakage issue has been addressed based on expert human feedback. Irrelevant features are removed.

. . .

**Assistant:** Tool `autoprognosis_survival_train_test` completed successfully ✓
**Tool logs:**

```
1  Trying models:
2  [
3    "cox_ph",
4    "deephit"
5  ]
6  Running the survival analysis study, this may take several minutes
       ...
7  AutoPrognosis 2.0: Duration of this iteration: 15.6s
8  AutoPrognosis 2.0: Investigating model: data_cleanup->cox_ph
9  AutoPrognosis 2.0: C-index achieved: 0.917
10 AutoPrognosis 2.0: Brier score achieved: 0.003
11 AutoPrognosis 2.0: Duration of this iteration: 4.1s
12 AutoPrognosis 2.0: Investigating model: nop->feature_normalizer->
       data_cleanup->cox_ph
13 AutoPrognosis 2.0: C-index achieved: 0.896
14 AutoPrognosis 2.0: Brier score achieved: 0.003
15 AutoPrognosis 2.0: Duration of this iteration: 3.7s
16 AutoPrognosis 2.0: Investigating model: nop->normal_transform->
       data_cleanup->cox_ph
```

```
17  AutoPrognosis 2.0: C-index achieved: 0.892
18  AutoPrognosis 2.0: Brier score achieved: 0.003
19  AutoPrognosis 2.0: Duration of this iteration: 3.4s
20  AutoPrognosis 2.0: Investigating model: data_cleanup->cox_ph
21  AutoPrognosis 2.0: C-index achieved: 0.898
22  AutoPrognosis 2.0: Brier score achieved: 0.007
23  AutoPrognosis 2.0: Duration of this iteration: 3.7s
24  AutoPrognosis 2.0: Investigating model: nop->feature_normalizer->
        data_cleanup->cox_ph
25  AutoPrognosis 2.0: C-index achieved: 0.862
26  AutoPrognosis 2.0: Brier score achieved: 0.007
27  AutoPrognosis 2.0: Duration of this iteration: 3.4s
28  AutoPrognosis 2.0: Investigating model: nop->normal_transform->
        data_cleanup->cox_ph
29  AutoPrognosis 2.0: C-index achieved: 0.876
30  AutoPrognosis 2.0: Brier score achieved: 0.007
31  AutoPrognosis 2.0: Duration of this iteration: 2.9s
32  AutoPrognosis 2.0: Investigating model: data_cleanup->cox_ph
33  AutoPrognosis 2.0: C-index achieved: 0.843
34  AutoPrognosis 2.0: Brier score achieved: 0.012
35  AutoPrognosis 2.0: Duration of this iteration: 3.2s
36  AutoPrognosis 2.0: Investigating model: nop->feature_normalizer->
        data_cleanup->cox_ph
37  AutoPrognosis 2.0: C-index achieved: 0.803
38  AutoPrognosis 2.0: Brier score achieved: 0.012
39  AutoPrognosis 2.0: Duration of this iteration: 3.0s
40  AutoPrognosis 2.0: Investigating model: nop->normal_transform->
        data_cleanup->cox_ph
41  AutoPrognosis 2.0: C-index achieved: 0.830
42  AutoPrognosis 2.0: Brier score achieved: 0.012
43  AutoPrognosis 2.0: Duration of this iteration: 7.6s
44  AutoPrognosis 2.0: Duration of this iteration: 7.6s
45  AutoPrognosis 2.0: Investigating model: ['1.0 * data_cleanup->
        cox_ph']
46  AutoPrognosis 2.0: C-index achieved: 0.952 +/- 0.0
47  AutoPrognosis 2.0: Brier score achieved: 0.003 +/- 0.0
48  AutoPrognosis 2.0: Investigating model:
49  ['1.0 * data_cleanup->cox_ph']
50  AutoPrognosis 2.0: C-index achieved: 0.952 +/- 0.0
51  AutoPrognosis 2.0: Brier score achieved: 0.003 +/- 0.0
52  AutoPrognosis 2.0: Duration of this iteration: 7.6s
53  AutoPrognosis 2.0: Duration of this iteration: 7.6s
54  AutoPrognosis 2.0: Investigating model: ['1.0 * data_cleanup->
        cox_ph']
55  AutoPrognosis 2.0: C-index achieved: 0.936 +/- 0.0
56  AutoPrognosis 2.0: Brier score achieved: 0.007 +/- 0.0
57  AutoPrognosis 2.0: Investigating model:
58  ['1.0 * data_cleanup->cox_ph']
59  AutoPrognosis 2.0: C-index achieved: 0.936 +/- 0.0
60  AutoPrognosis 2.0: Brier score achieved: 0.007 +/- 0.0
61  AutoPrognosis 2.0: Duration of this iteration: 7.1s
62  AutoPrognosis 2.0: Duration of this iteration: 7.1s
63  AutoPrognosis 2.0: Investigating model: ['1.0 * data_cleanup->
        cox_ph']
64  AutoPrognosis 2.0: C-index achieved: 0.87 +/- 0.0
65  AutoPrognosis 2.0: Brier score achieved: 0.012 +/- 0.0
66  AutoPrognosis 2.0: Investigating model:
67  ['1.0 * data_cleanup->cox_ph']
68  AutoPrognosis 2.0: C-index achieved: 0.87 +/- 0.0
69  AutoPrognosis 2.0: Brier score achieved: 0.012 +/- 0.0
70  AutoPrognosis 2.0: Duration of this iteration: 609.3s
71  AutoPrognosis 2.0: Duration of this iteration: 609.3s
```

```
72  AutoPrognosis 2.0: Investigating model: ["1.0 * data_cleanup->
        cox_ph({'cox_ph': {}})", "1.0 * data_cleanup->cox_ph({'cox_ph':
        {}})", "1.0 * data_cleanup->cox_ph({'cox_ph': {}})"]
73  AutoPrognosis 2.0: C-index achieved: 0.903
74  AutoPrognosis 2.0: Brier score achieved: 0.007
75  AutoPrognosis 2.0: Investigating model:
76  AutoPrognosis 2.0: Study finished.
77  ["1.0 * data_cleanup->cox_ph({'cox_ph': {}})", "1.0 * data_cleanup
        ->cox_ph({'cox_ph': {}})", "1.0 * data_cleanup->cox_ph({'cox_ph
        ': {}})"]
78  AutoPrognosis 2.0: C-index achieved: 0.903
79  AutoPrognosis 2.0: Brier score achieved: 0.007
80  AutoPrognosis 2.0: Study finished.
81  Survival analysis study completed!
82  Final model summary:
83  Metrics (train data): "c_index": "0.925 +/- 0.029",
84  Metrics (test data): "c_index": "0.848 +/- 0.005"
```

## E.4 EXAMPLE: MODEL EVALUATION & FEATURE IMPORTANCE

Based on the results, CliMB-DC will provide feature importance explanations for the users to better understand the prediction model, ensuring the interpretability in healthcare.

**Summary of Progress to Date:** CliMB-DC has fitted a survival model for the given dataset and the task.

. . .

**Assistant:** Would you like to see feature importance plots for your survival analysis task?

**User:** Yes, please

**Assistant:** Tool `permutation_explainer` completed successfully ✓
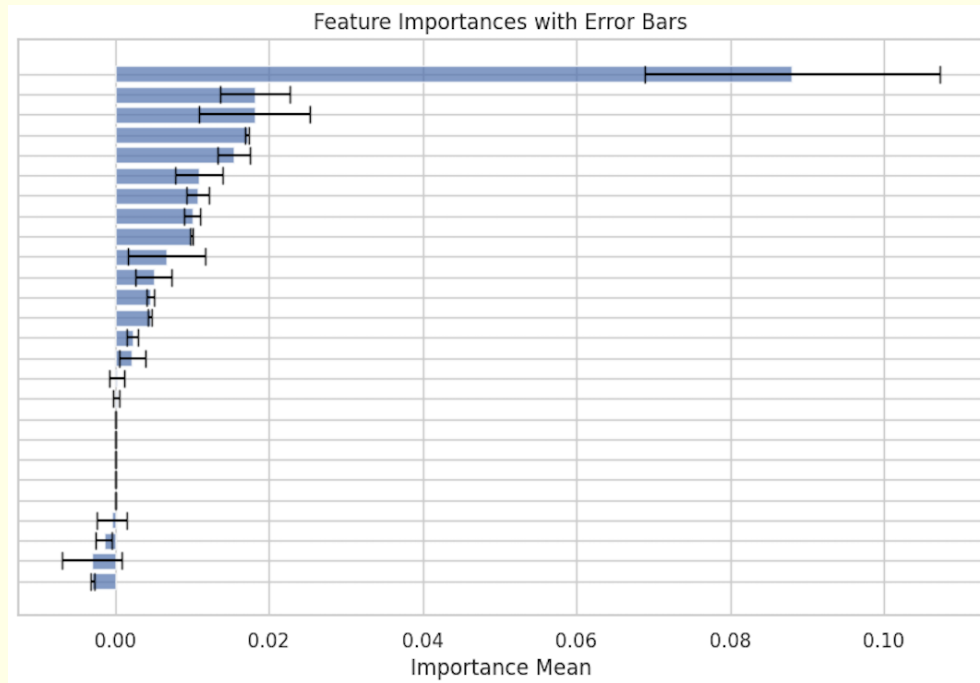**Tool logs:**

```
1   Loading the data...
2   Data loaded with shape: 4800 rows and 28 columns.
3   Loading model from file: path/model_all.p
4   Running the permutation explainer, this can take a while...
5
6         This tool takes approximately 1 minute per column. If this
            is too slow, consider reducing the number of columns.
7         Reducing the number of column can be done with the feature
            selection tool. You can cancel this tool with the
8         'Restart from last reasoning step' button. Then use the
            feature selection tool to reduce the number of columns.
9   Permutation explainer completed!
10                          importances_mean   importances_std
11  age_at_recruitment           0.088055           0.019220
12  copd_emphysema_bronchitis    0.018188           0.004543
13  packyears                    0.018143           0.007201
14  pneumonia_self_reported      0.017176           0.000252
15  age_started_smoking          0.015429           0.002090
16  personal_hx_cancer           0.010980           0.003070
17  ...
```

**Assistant:** The feature importance plot is here[a]:

Feature Importances with Error Bars



[a]We omit the feature names in the figure.

## F  MORE CLIMB-DC INTERACTION EXAMPLES

In this section, we present more examples highlighting the advantages of CliMB-DC in executing effective data preprocessing.

### F.1  EXAMPLE: DATA EXPLORATION

For the input data, CliMB-DC provides users with a detailed data analysis as well as visualization using tools `EDA`, `descriptive statistics`, and `lifelines.KaplanMeierFitter`.

**Summary of Progress to Date**: The data file has been uploaded.

. . .

**Assistant:** Tool `EDA` completed successfully ✓

**Tool logs:**

```
1  Getting dataset basic info...
2  Getting descriptive statistics for numerical features...
3  Getting detailed information on categorical variables...
4  Performing missing values analysis...
5  Identified numeric columns that should be considered categorical:
6  ['asbestosis_self_reported', 'pneumonia_self_reported', '
      chronic_obstructive_airways_disease_copd_self_reported', '
      emphysema_chronic_bronchitis', 'copd_emphysema_bronchitis', '
      emphysema', 'asthma', 'eczema_allergicrhinitis_hayfever', '
      number_of_self_reported_cancers', 'personal_hx_cancer', '
```

```
        family_hx_lung_cancer_mother', 'family_hx_lung_cancer_father',
        'family_hx_lung_cancer_siblings', 'asbestos', 'lung_cancer']
7  Performing correlation analysis...
8  Performing potential outliers identification...
9  Performing duplicate records analysis...
```
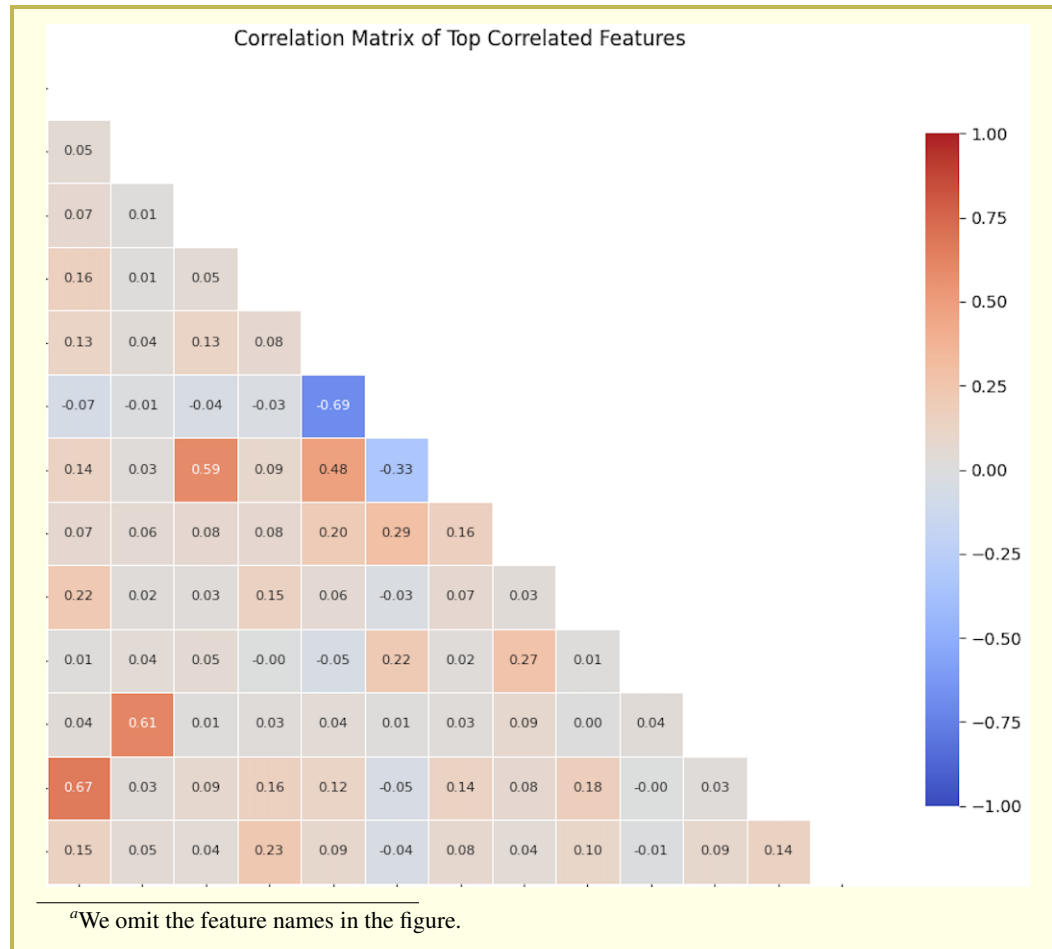
**Tool output:**

```
1   Dataset Shape: 4800 rows and 31 columns
2   Column Names and Types:
3   index                                           int64
4   age_at_recruitment                            float64
5   sex                                            object
6   body_mass_index                               float64
7   n_cig_per_day                                 float64
8   age_started_smoking                           float64
9   age_stopped_smoking                           float64
10  ...
11
12  Descriptive Statistics for Numerical Features:
13               index  age_at_recruitment  body_mass_index ...
14  count    4800.000000        4800.000000      4800.000000 ...
15  mean    87406.250417          57.417163        27.676107 ...
16  std     49908.726761           8.084722         4.621542 ...
17  min       129.000000          40.000000        16.024943 ...
18  25%     44130.250000          51.167603        24.473217 ...
19  50%     88896.000000          59.135000        27.116986 ...
20  75%    130874.500000          63.719537        30.425950 ...
21  max    173314.000000          72.000000        53.565422 ...
22  skew       -0.033166          -0.447521         0.579946 ...
23  kurt       -1.179482          -0.848562         0.532072 ...
24
25  Identified numeric value columns that should most likely be
        considered categoricals:
26  ['asbestosis_self_reported', 'pneumonia_self_reported', 'asbestos',
        'lung_cancer', ...].
27  This is done by checking whether the column contains only integers
        and has a low number of unique values (<20 or <5% of total
        examples).
28
29  Detailed Information on Categorical Variables:
30  smoking_status - Unique Values: 2
31  Top 5 Values:
32  smoking_status
33  Previous    3603
34  Current     1197
35  ...
36
37  Missing Values Analysis:
38  No missing values found.
39
40  Correlation Analysis: ...
41
42  Outlier Identification for Numerical Features:
43  index - Outliers Count: 0
44  [Lower Bound: -1.08e+05, Upper Bound: 2.83e+05]
45  ...
46
47  Duplicate Records: 0
```

**Assistant:** Here is a correlation gram showing the correlation between features[a]:

Correlation Matrix of Top Correlated Features

*a*We omit the feature names in the figure.

**Assistant:** Would you like to generate descriptive statistics for the training dataset (lung_train.csv)?

**User:** Yes

**Assistant:** Tool `descriptive_statistics` completed successfully ✓

**Tool logs:**

```
1  Creating the descriptive statistics table...
2  Saving the summary table to:
3  /mnt/data-fourtb/Dropbox/Programming/wsl_repos/climb_versions/
       climb_cot/climb/wd/session_2025-01-08_14-33-29/lung_train.
       csv__descriptive_stats.csv
4  Identified numeric columns that should be considered categorical:
5  ['asbestos', 'lung_cancer', ...]
6  Creating plots for the data...
7  Normally distributed features:
8  []
9  Not normally distributed features:
10 ['age_at_recruitment', 'age_started_smoking', 'age_stopped_smoking'
       , 'body_mass_index', 'index', ...]
11 Plotted a bar plot for: 'sex'
12 Plotted a bar plot for: 'smoking_status'
13 Plotted a histogram and box plot for: 'age_at_recruitment'
14 ...
```

```
15  Plotting of 'years_since_stopped_smoking' took too long and was
       skipped.
```

**Tool output:**

```
1   Variable                        Summary
2   index                                                    88896 (44130
        - 130874)
3   age_at_recruitment                                           59.1
        (51.2 - 63.7)
4   sex
5       Male
            2483/4800 (51.7)
6       Female
            2317/4800 (48.3)
7   ...
```

**Assistant:** The descriptive statistics have been generated successfully. Here is a summary of the key findings:

**Summary of Descriptive Statistics:**
Numerical Variables:

- age_at_recruitment: Median = 59.1 (Q1-Q3: 51.2 - 63.7)
- body_mass_index: Median = 27.1 (Q1-Q3: 24.5 - 30.4)
- ...

Categorical Variables:

- sex: Male = 51.7%, Female = 48.3%
- smoking_status: Previous = 75.1%, Current = 24.9%
- ...

**Key Figures:** Here are some important figures generated during the analysis:
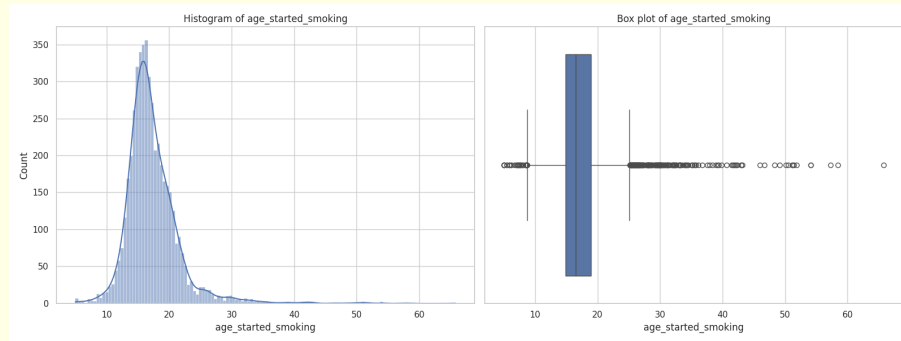
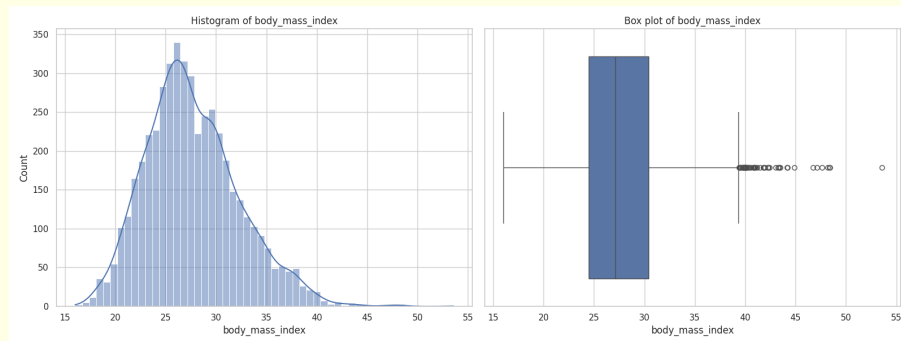- Distribution of age_at_recruitment:



- Distribution of age_started_smoking:
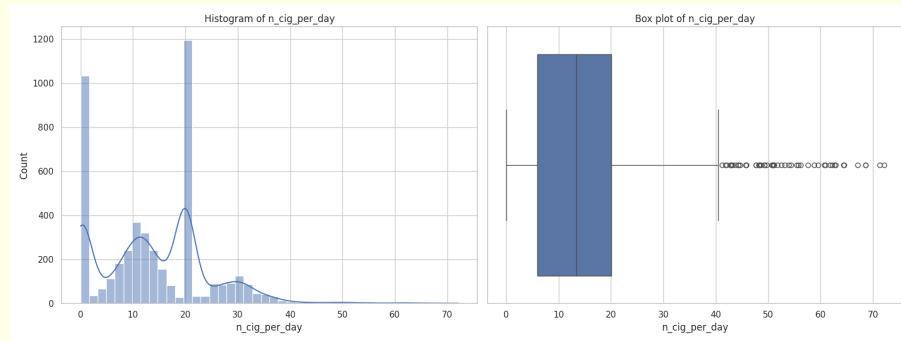


42

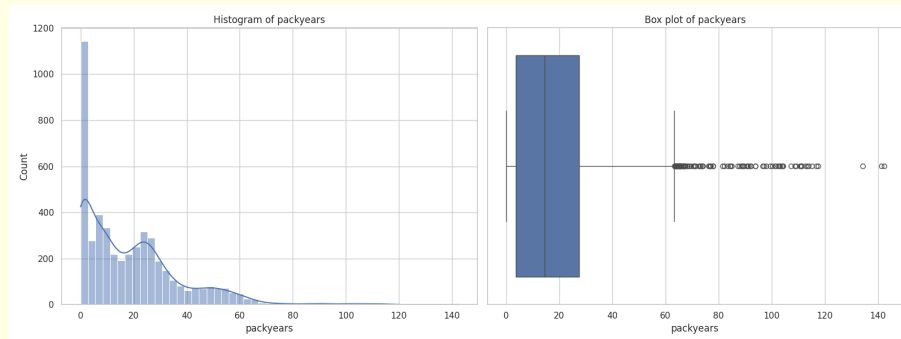- Distribution of age_stopped_smoking:



- Distribution of body_mass_index:



- Distribution of number_of_cigarettes_per_day::



- Distribution of packyears:



- . . .

**Assistant:** Would you like to see a Kaplan-Meier plot for the survival analysis?

**User:** Yes

**Assistant:** I will generate the Kaplan-Meier plot for the training dataset using the `lifelines` library.

**Generated code:**

```python
import pandas as pd
import matplotlib.pyplot as plt
from lifelines import KaplanMeierFitter

# Load the training dataset
df = pd.read_csv("lung_train.csv")

# Define the event and time columns
event_col = "lung_cancer"
time_col = "time_to_event_develop"

# Initialize the Kaplan-Meier fitter
kmf = KaplanMeierFitter()

# Fit the data
kmf.fit(durations=df[time_col], event_observed=df[event_col])

# Plot the Kaplan-Meier curve
plt.figure(figsize=(10, 6))
kmf.plot_survival_function()
plt.title("Kaplan-Meier Survival Curve")
plt.xlabel("Time to Event (days)")
plt.ylabel("Survival Probability")
plt.grid(True)

# Save the plot
plt.savefig("km_plot_training.png")
plt.close()
print("Kaplan-Meier plot for the training dataset has been saved as
      'km_plot_training.png'.")
```
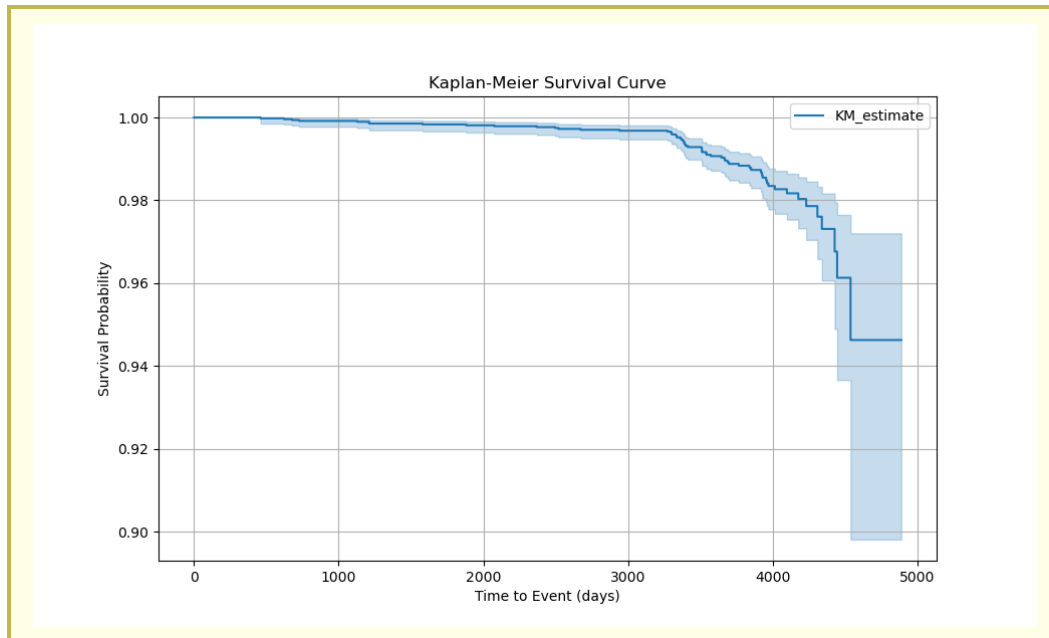
**Code execution output:**
Code execution finished successfully ✓

```
Kaplan-Meier plot for the training dataset has been saved as '
    km_plot_training.png'.
```

**Assistant:** Here is the Kaplan-Meier plot for the training dataset:

## F.2 EXAMPLE: DATA AGGREGATION

When input data contains multiple measurements, CliMB-DC will auto-check such issue and interact with users to address it.

---

**Summary of Progress to Date** The data file has been uploaded, and relevant background information has been provided to CliMB-DC.

. . .

---

Let's assess the suitability of your data for the task and whether the `AutoPrognosis` tools support it.

**Q1: Is the data suitable for the task?**

- The data is structured with each row representing a patient's information at different follow-up points.
- We need to predict the "Event_status" variable using the "Years_to_event" as the time variable, which indicates a survival analysis problem.

Potential issue: The data has multiple rows per patient, which might not be directly suitable for the `AutoPrognosis` tools that expect one row per patient.

**Q2: Does the AutoPrognosis set of tools support the task?**

- `AutoPrognosis` supports survival analysis through the `autoprognosis_survival_train_test` tool.
- This tool requires the target (event) variable and the time variable, which we have ("Event_status" and "Years_to_event").

**Conclusion**:

- The task can be performed using `AutoPrognosis`, but the data needs to be transformed to have one row per patient.

To transform the data, we can aggregate the follow-up points for each patient. This could involve taking the last follow-up point for each patient or aggregating information in a meaningful way. Would you like to proceed with transforming the data to fit the task requirements?

---

**User:** I would like to take the last follow up point.

**Assistant:** Let's proceed with transforming the data to have one row per patient, taking the last follow-up point for each patient.

**Generated code:**

```python
import pandas as pd

# Load the data
file_path = "PBC_Data_corrected_v2.csv"
data = pd.read_csv(file_path)

# Sort the data by patient ID and follow-up time (assuming columns
    are 'patient_id' and 'Years_to_event')
# Replace 'patient_id' with the actual column name for patient
    identification if different
data_sorted = data.sort_values(by=['patient_id', 'Years_to_event'])

# Take the last follow-up point for each patient
data_last_followup = data_sorted.drop_duplicates(subset=['
    patient_id'], keep='last')

# Save the transformed data
output_file_path = "PBC_Data_last_followup.csv"
data_last_followup.to_csv(output_file_path, index=False)

print(f"Transformed data saved to {output_file_path}")
```

**Code execution output:**
Code execution failed ×

```
KeyError: 'patient_id'
```

**Assistant:** It appears that the column name for patient identification is not "patient_id". Could you please provide the correct column name for patient identification?

**User:** Yes, that would be "id"

**Assistant:** Let's proceed with transforming the data using the correct column name "id" for patient identification.
**Generated code:**

```python
import pandas as pd

# Load the data
file_path = "PBC_Data_corrected_v2.csv"
data = pd.read_csv(file_path)

# Sort the data by patient ID and follow-up time
data_sorted = data.sort_values(by=['id', 'Years_to_event'])

# Take the last follow-up point for each patient
data_last_followup = data_sorted.drop_duplicates(subset=['id'],
    keep='last')

# Save the transformed data
output_file_path = "PBC_Data_last_followup.csv"
data_last_followup.to_csv(output_file_path, index=False)
```

```
16
17  print(f"Transformed data saved to {output_file_path}")
```

**Code execution output:**
Code execution finished successfully ✓

```
1  Transformed data saved to PBC_Data_last_followup.csv
```

**Assistant:** The data has been successfully transformed to have one row per patient, taking the last follow-up point for each patient. The transformed data has been saved as PBC_Data_last_followup.csv.
Based on the assessment:

- The data is now suitable for the task.

- The `AutoPrognosis` tools support survival analysis, which is the type of study you want to perform.

You can proceed with the survival analysis using the `autoprognosis_survival_train_test` tool.