# Automatically Labeling Clinical Trial Outcomes: A Large-Scale Benchmark for Drug Development

Chufan Gao*, Jathurshan Pradeepkumar*, Trisha Das*, Shivashankar Thati, and Jimeng Sun
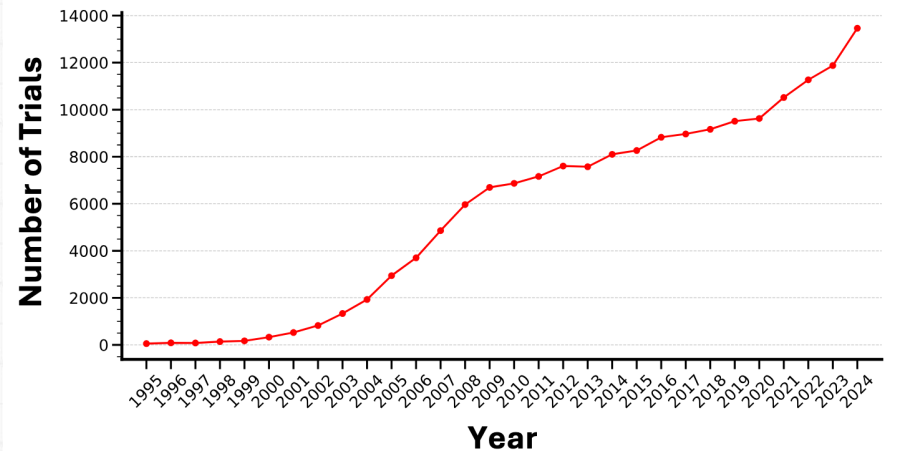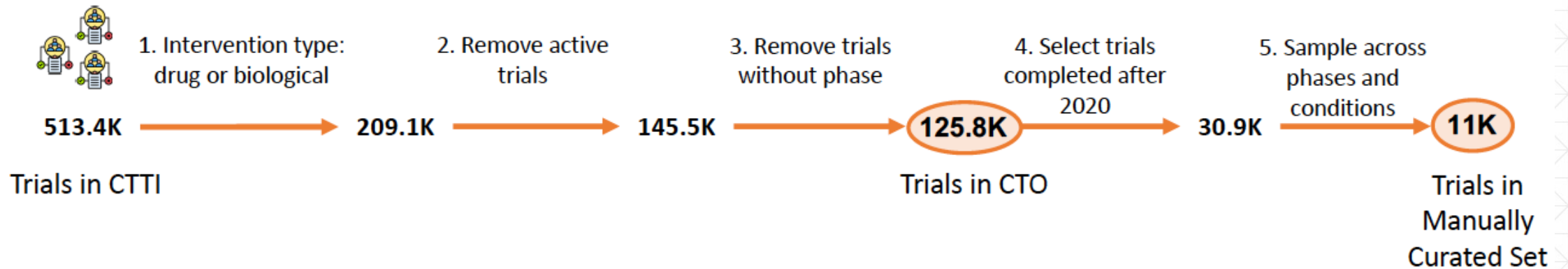
Project Page
Paper
Code

# Background

- Drug discovery and development are expensive, with clinical trial results vital for regulatory approval and patient care.

- Large-scale, high-quality clinical trial outcome data remains limited.
  - Hindering the development of predictive models

- Dynamic & Rapidly Growing Data
  - Clinical trial data grows rapidly and is affected by diverse external factors (e.g., COVID-19, regulatory changes).
  - Frequent label updates are needed, but **manual labeling** is **impractical** at scale.

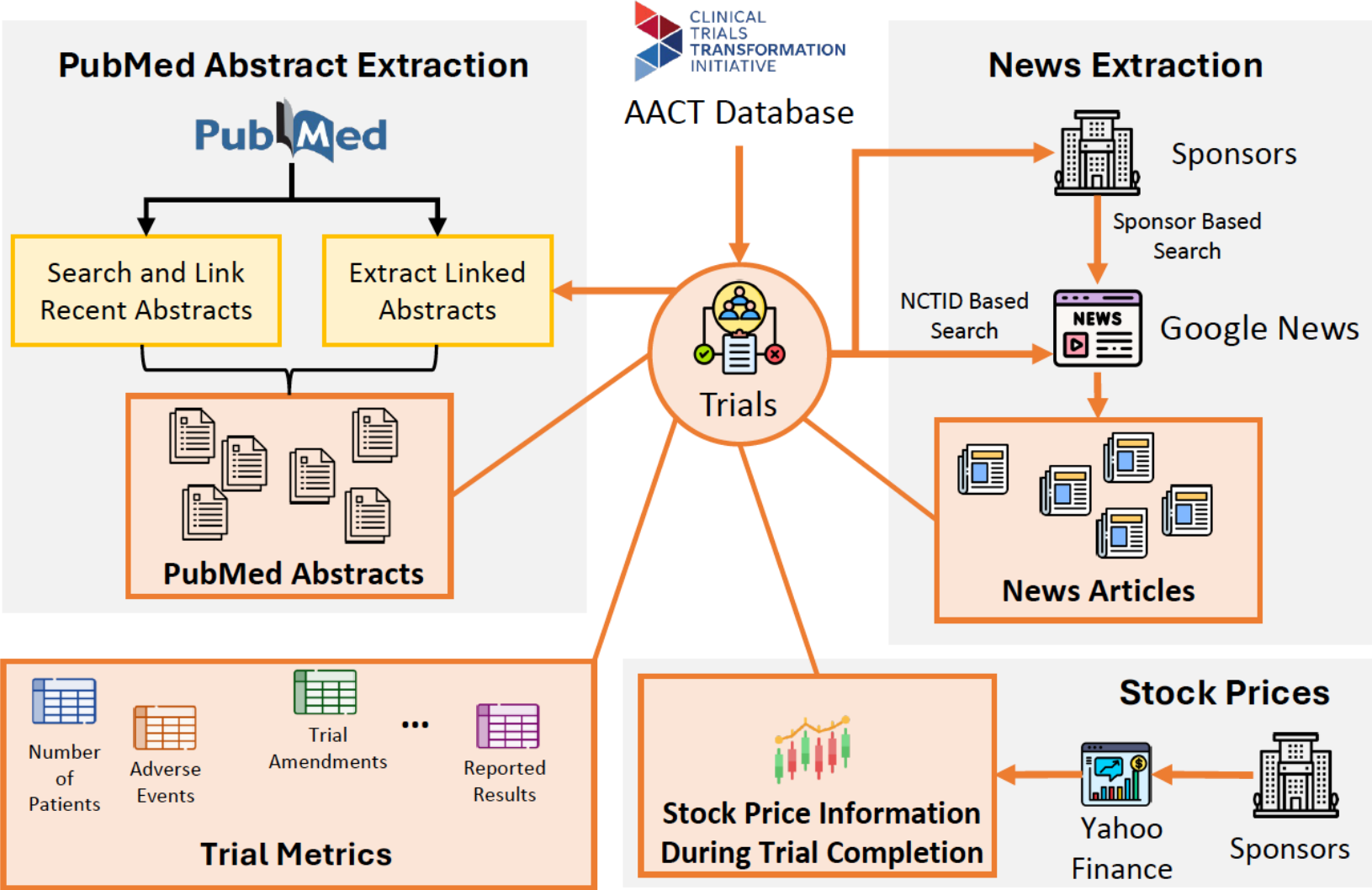### A: Trial Distribution by Completed Year
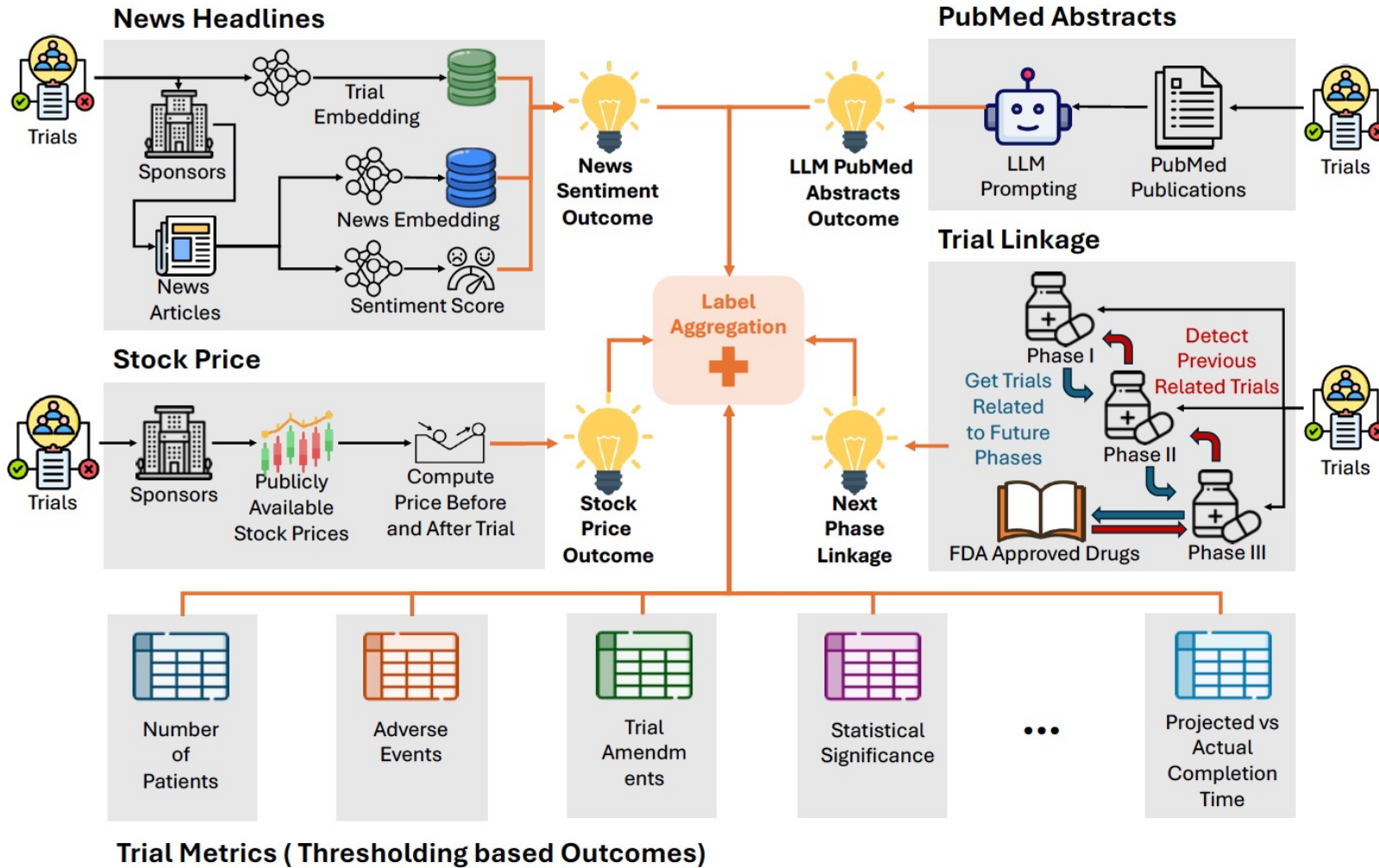
# Clinical Trial Outcome Benchmark (CTO)

- Clinical Trial Outcome (CTO) benchmark, a fully reproducible, regularly updated, large-scale repository encompassing approximately ~125K drug and biologics trials.

  - A comprehensive trial knowledge base

  - Automated labeling framework based on aggregation of weak labeling function
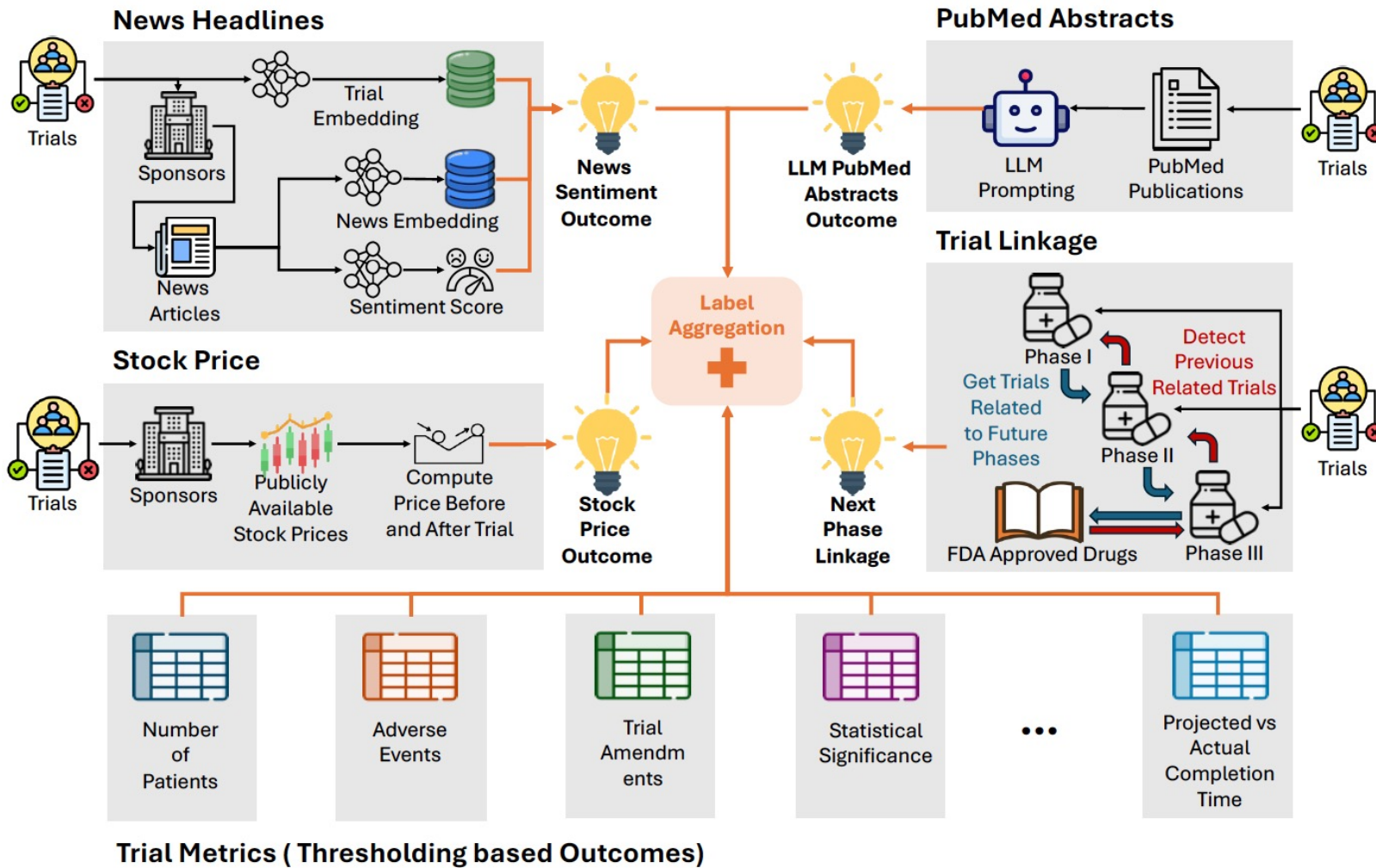
  - Manually curated around 11K trials.

# Trial Knowledge Base

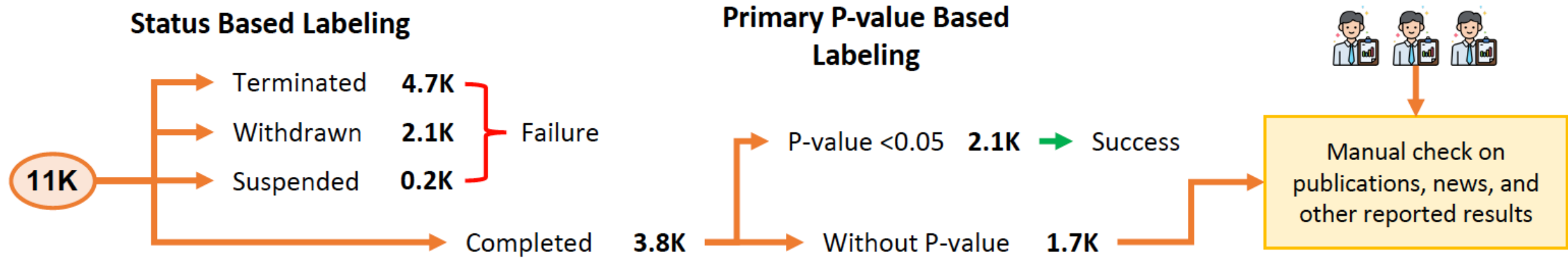# CTO Automated Labeling Framework

# CTO Automated Labeling Framework



| Phase | Aggregation Method | F1 | $\kappa$ |
|---|---|---|---|
| I | $CTO_{MV}$ | 0.726 | 0.490 |
| | $CTO_{DP}$ | 0.870 | 0.700 |
| | $CTO_{RF}$ | **0.913** | **0.790** |
| II | $CTO_{MV}$ | 0.689 | 0.430 |
| | $CTO_{DP}$ | 0.856 | 0.623 |
| | $CTO_{RF}$ | **0.878** | **0.693** |
| III | $CTO_{MV}$ | 0.904 | 0.606 |
| | $CTO_{DP}$ | 0.921 | 0.582 |
| | $CTO_{RF}$ | **0.941** | **0.710** |
| All | $CTO_{MV}$ | 0.793 | 0.529 |
| | $CTO_{DP}$ | 0.884 | 0.646 |
| | $CTO_{RF}$ | **0.909** | **0.729** |

Agreement of CTO automated labels with human-labeled TOP[1] dataset

[1] Fu, Tianfan, et al. "Hint: Hierarchical interaction network for clinical-trial-outcome predictions." *Patterns* 3.4 (2022).
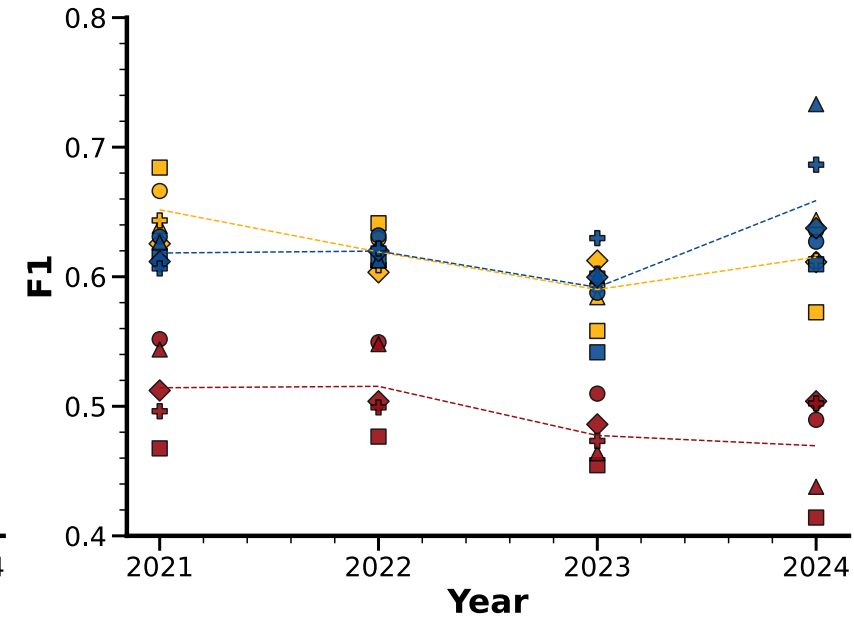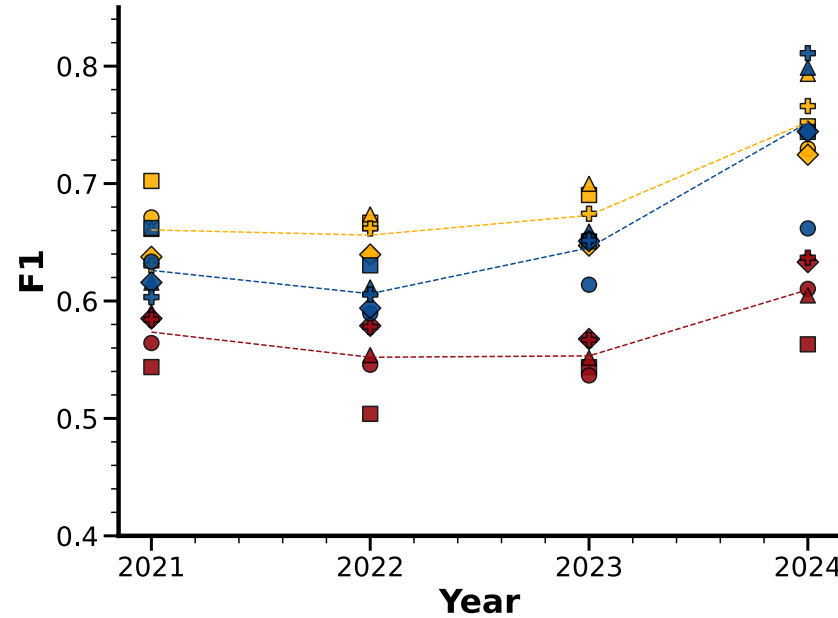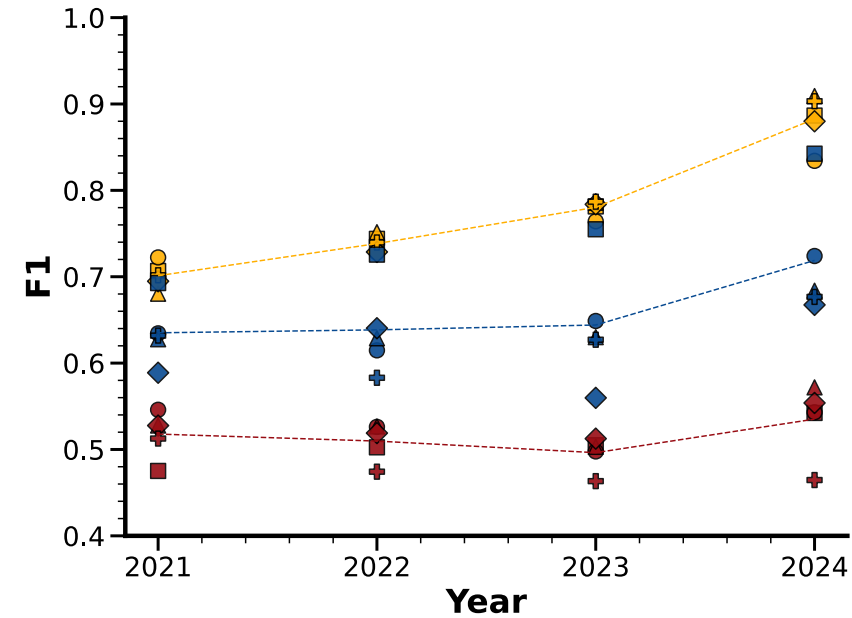
# Manual Curation Process



**Status Based Labeling**

11K →
- Terminated **4.7K** ⎤
- Withdrawn **2.1K** ⎬ Failure
- Suspended **0.2K** ⎦
- Completed **3.8K**

**Primary P-value Based Labeling**

- P-value <0.05 **2.1K** → Success
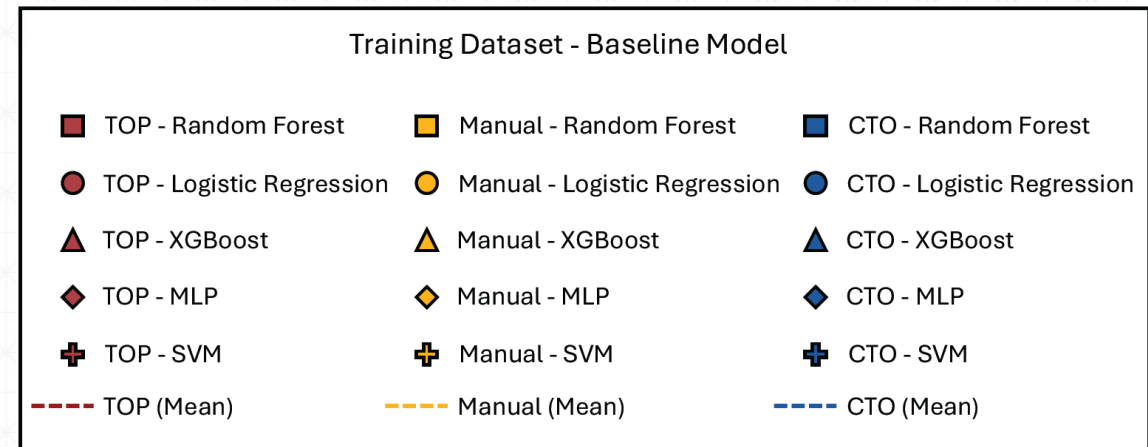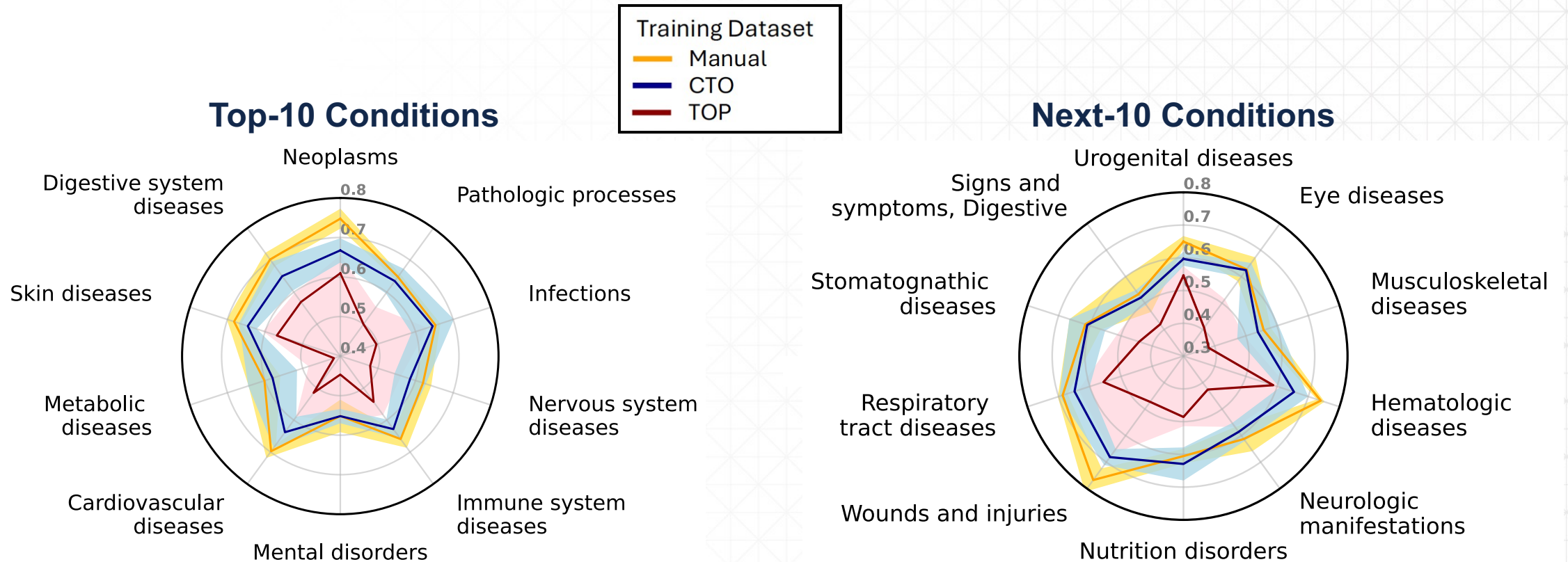- Without P-value **1.7K** → Manual check on publications, news, and other reported results

# Which labels are better?



Manual : Updated human - annotated data
CTO : From our automated labeling framework
TOP : Past benchmark on clinical trial outcomes [1]

[1] Fu, Tianfan, et al. "Hint: Hierarchical interaction network for clinical-trial-outcome predictions." *Patterns* 3.4 (2022).

# Which labels are better?



Top-10 Conditions

Next-10 Conditions

Training Dataset
- Manual
- CTO
- TOP

Manual : Updated human - annotated data
CTO : From our automated labeling framework
TOP : Past benchmark on clinical trial outcomes [1]

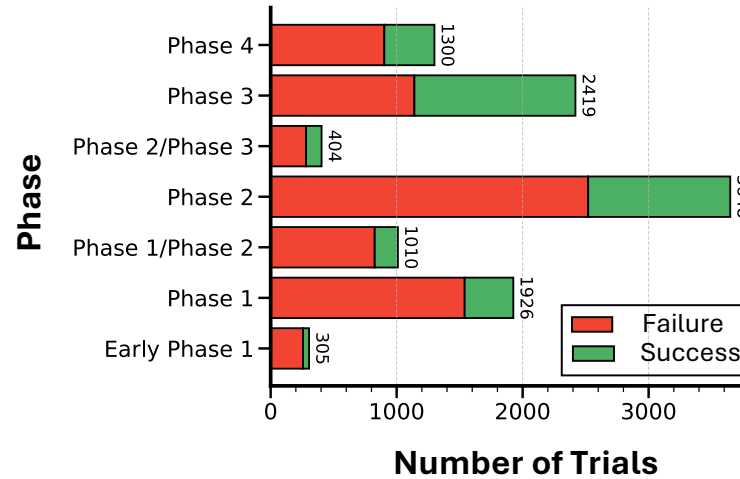[1] Fu, Tianfan, et al. "Hint: Hierarchical interaction network for clinical-trial-outcome predictions." *Patterns* 3.4 (2022).

# Manual Curated Trial Outcome Benchmark



A: Top 10 Disease Distribution

B: Phase Distribution

C: Distribution by Completed Year

# CTO Statistic
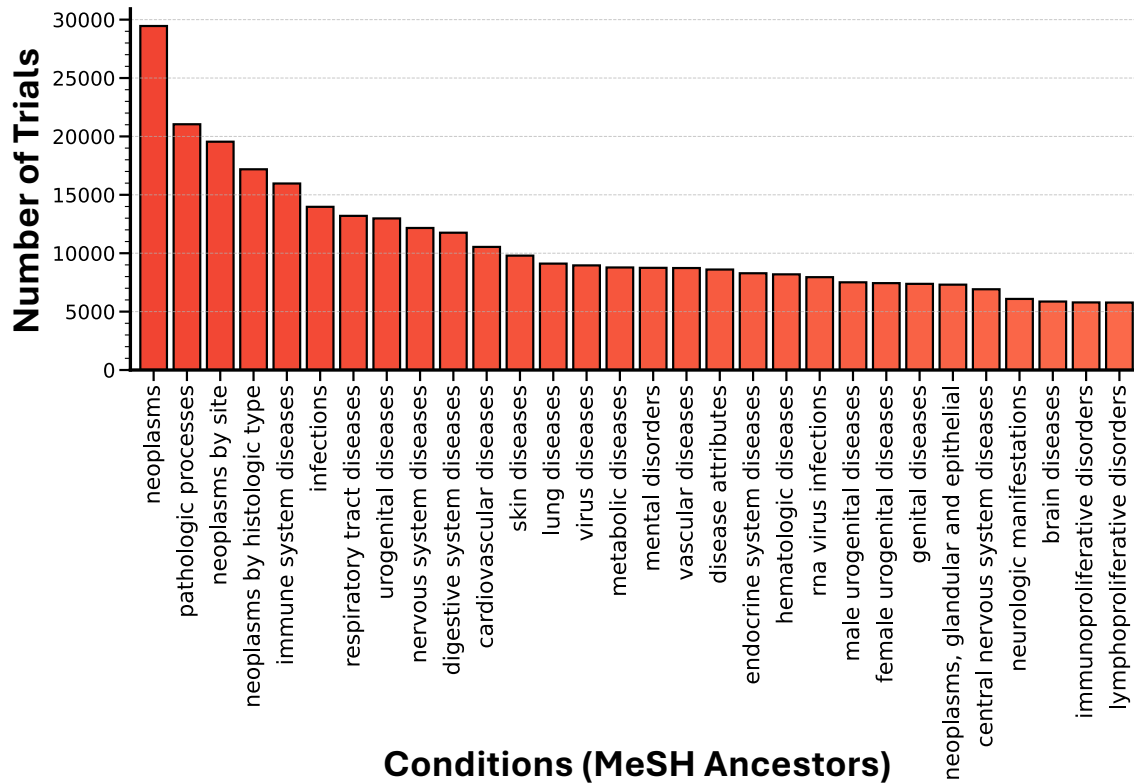
## B: Frequency of Top 30



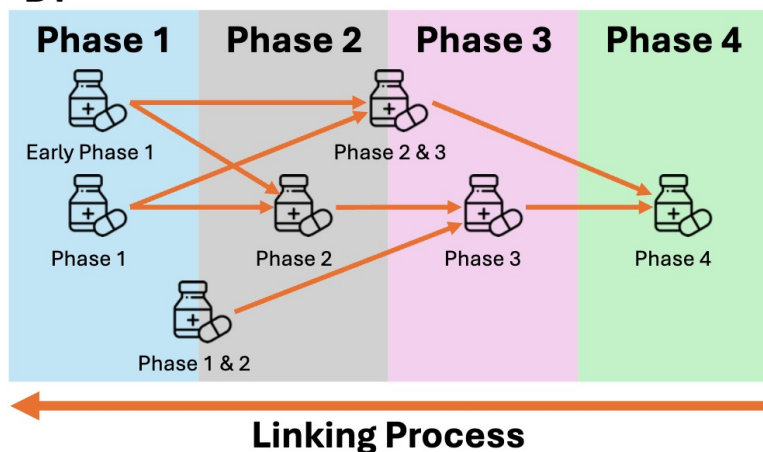## Trial Distribution by Completed Year



## D: Phase Distribution
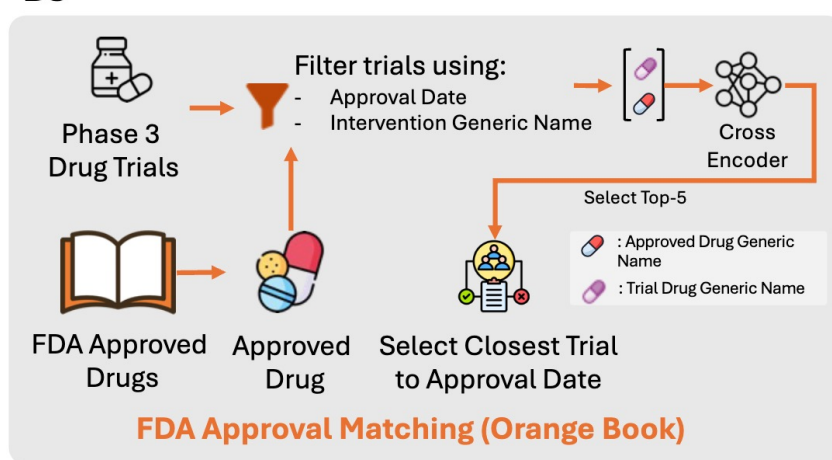
## C: Agreement between Labeling Functions in CTO

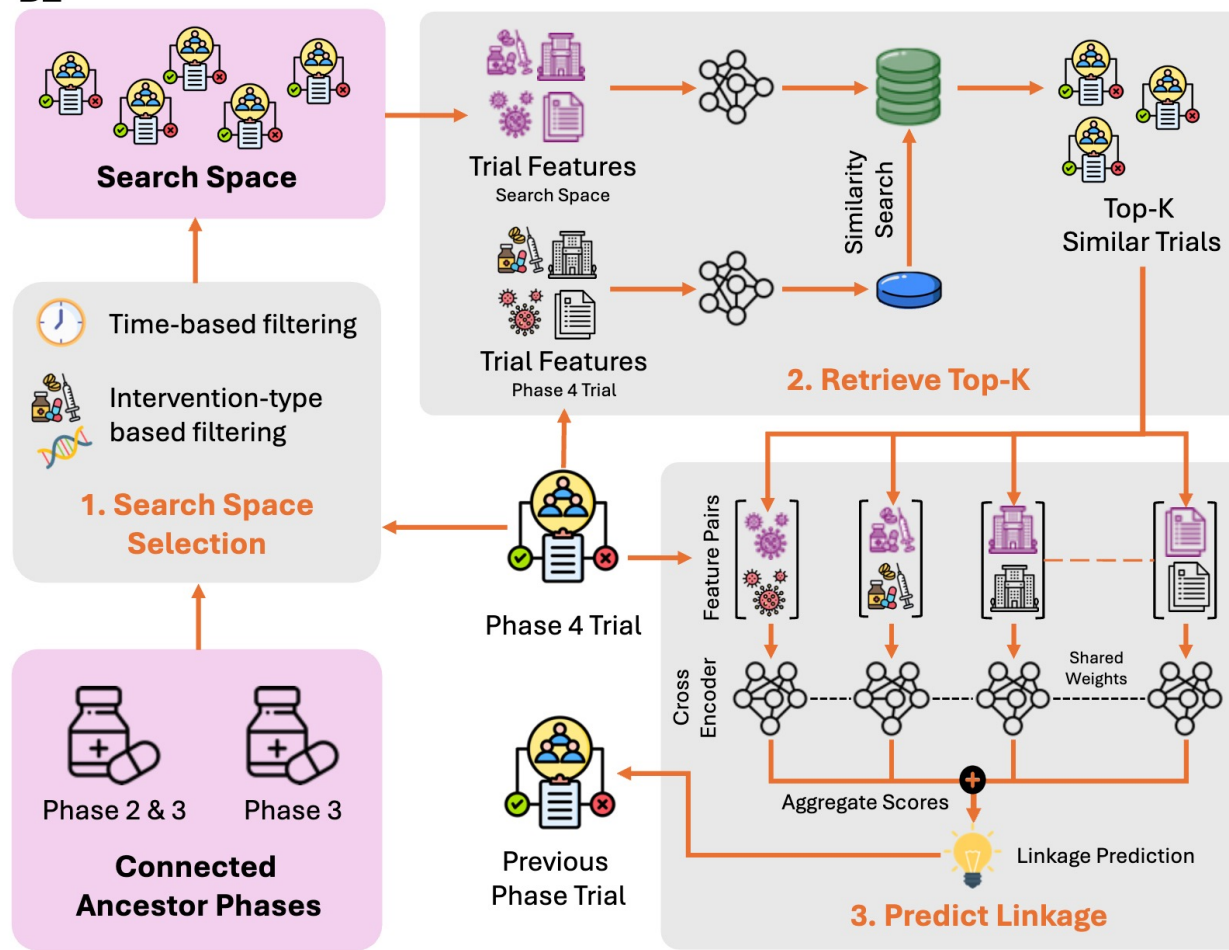# Additional Method / Weak Labeling Slides

## B. Trial Linkage Algorithm

# Additional Method / Weak Labeling Slides