

## Background & objective

### Background

Early identification of individuals at higher risk for cancer can improve outcomes and help target screening resources. Existing machine learning approaches require substantial data preprocessing, manual feature engineering, and task-specific model training, and their interpretability is often limited.

### Objective

We evaluated whether a large language model-based chain-of-agents framework can estimate 1-year cancer risk directly from raw longitudinal EHR data across multiple cancer types.

## Study population

### Lung cancer (unmatched)

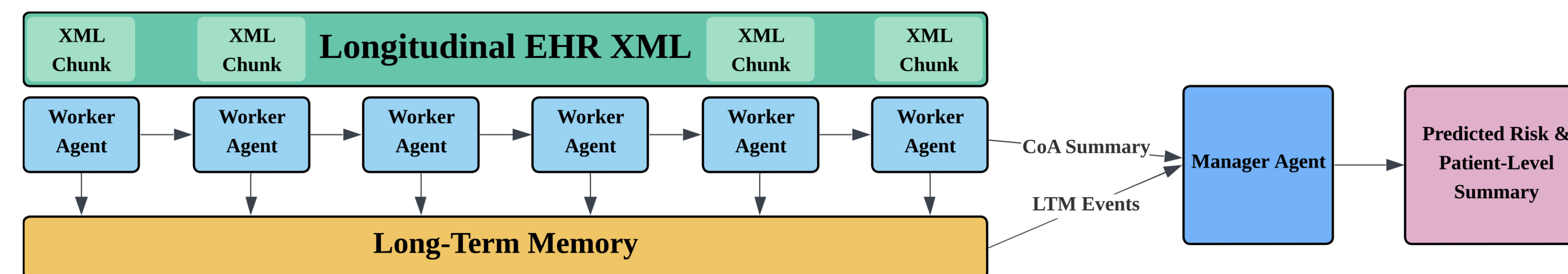
- From 120 million patients from Truveta Data
- 500 lung cancer cases
- 125,000 controls (randomly sampled, approximate cumulative incidence)
- 1-year risk prediction

### Multi-cancer evaluation (matched)

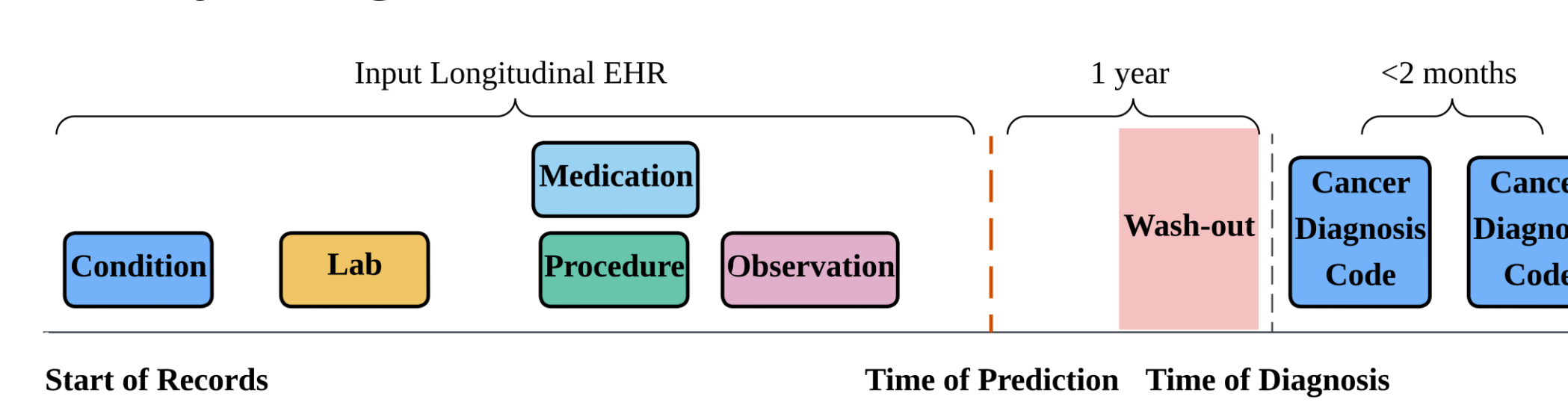
- 15 cancer types
- Each cancer type: 500 cases, 500 matched controls (1:1 matching on sex and 10-year age group)
- 1-year risk prediction

## TrajOnco framework

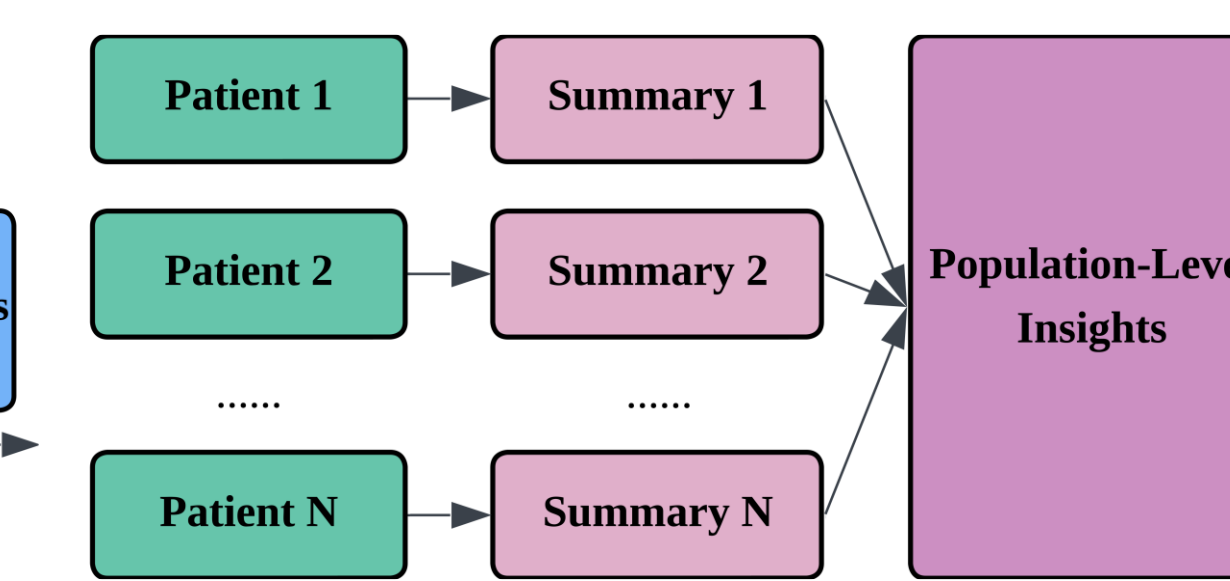
### Architecture overview



### Study design



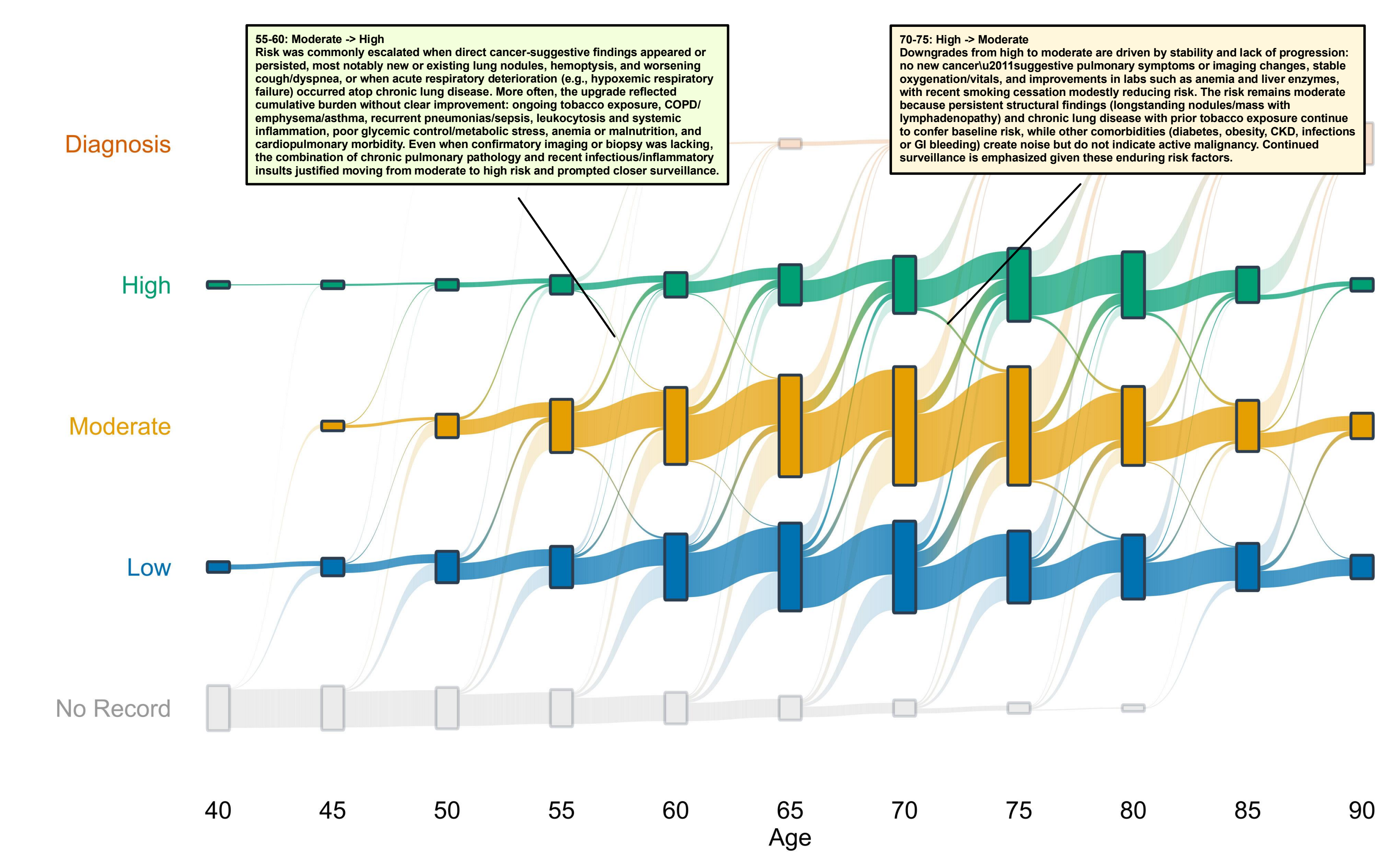
### Insights generation



- Zero-shot chain-of-agents (CoA) processes raw longitudinal EHR in XML format with sequential worker agents and a long-term memory
- Manager agent synthesizes patient-level summary, risk score (1-10), and evidence-linked rationale
- Interpretability can be achieved in patient-level and population-level

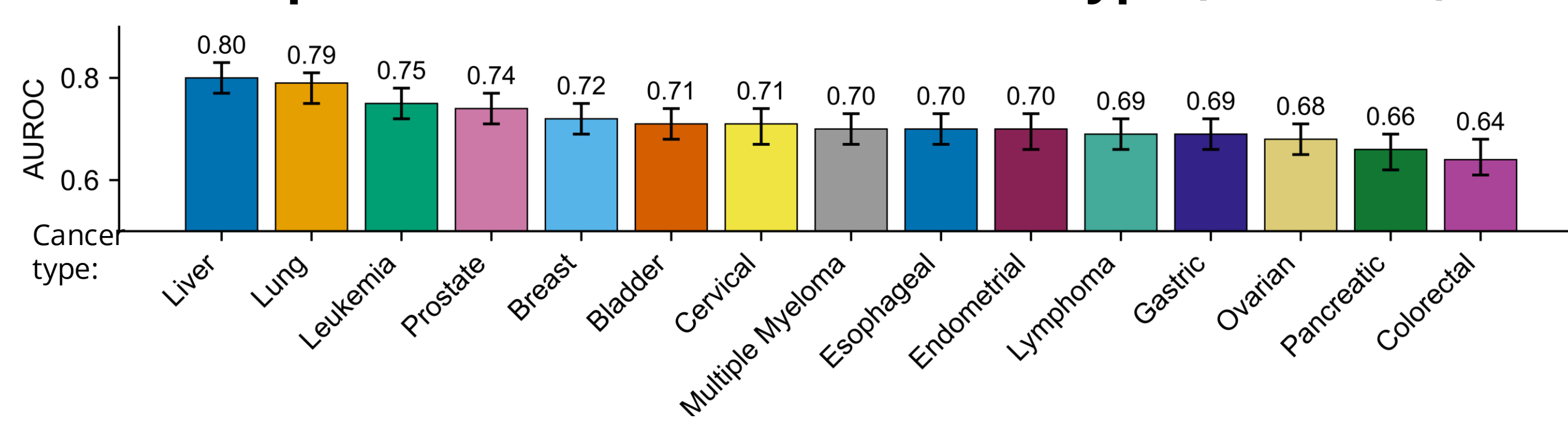
## Interpretation

### Risk transition summarization

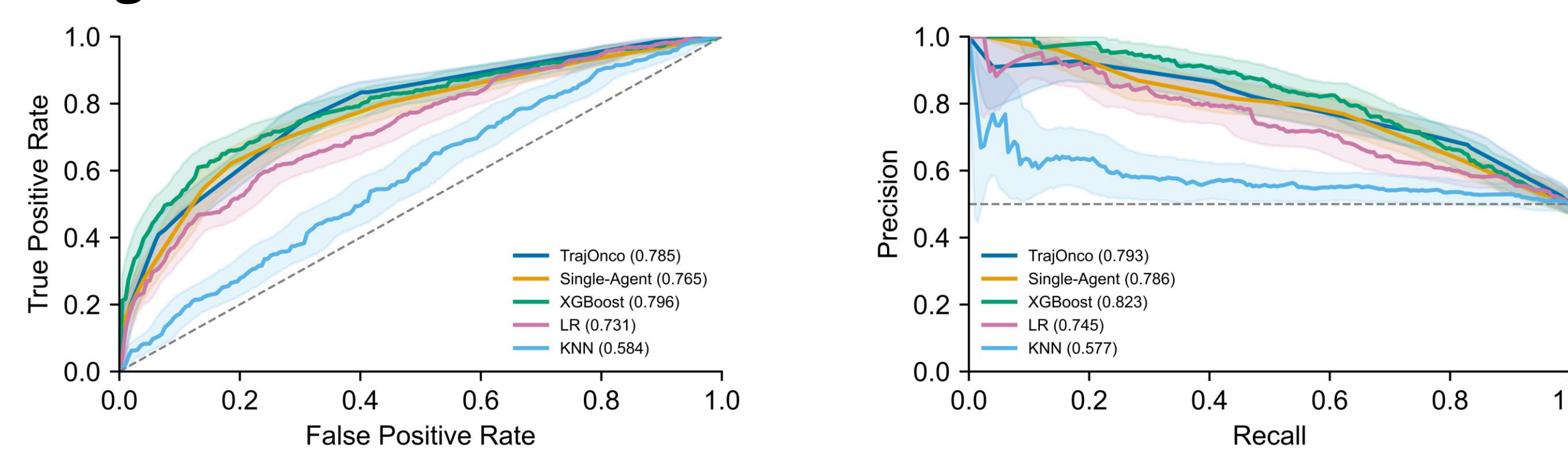


## Results

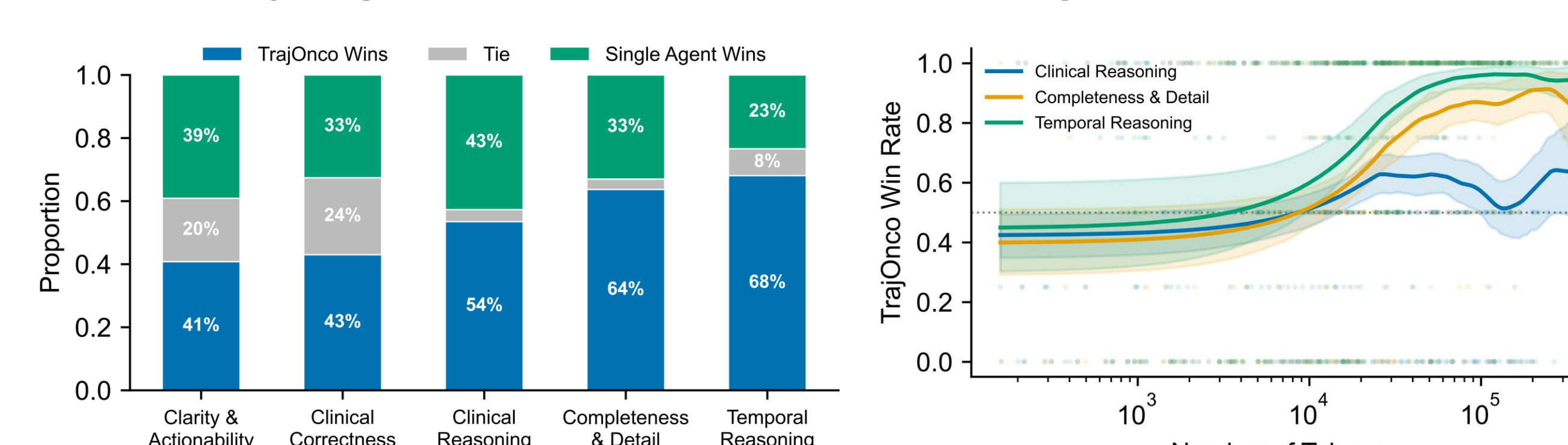
### Zero-shot performance across 15 cancer type (matched)



### Comparison with trained machine learning methods (matched lung cancer cohort)



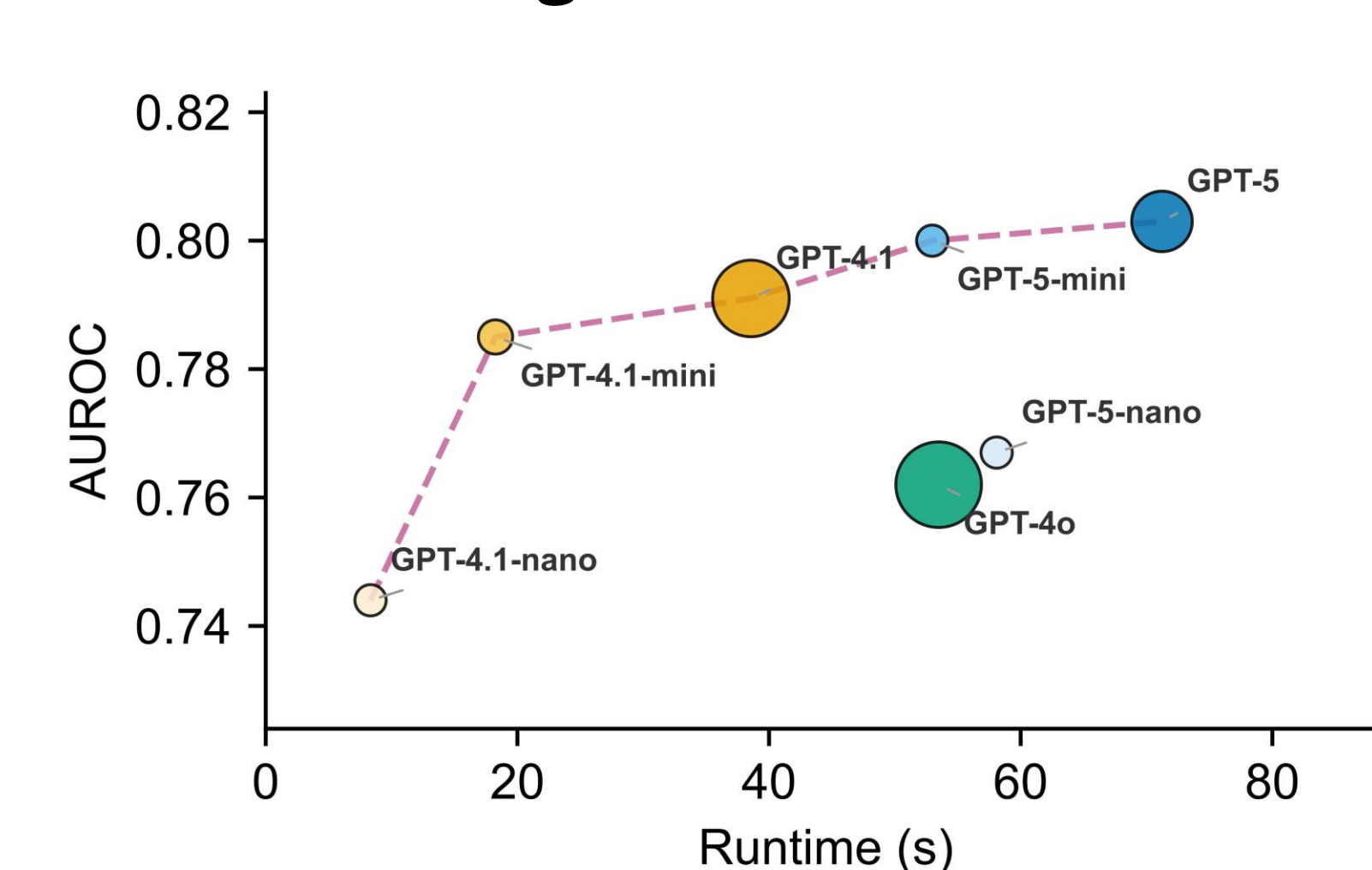
### LLM-as-a-judge evaluation (matched lung cancer cohort)



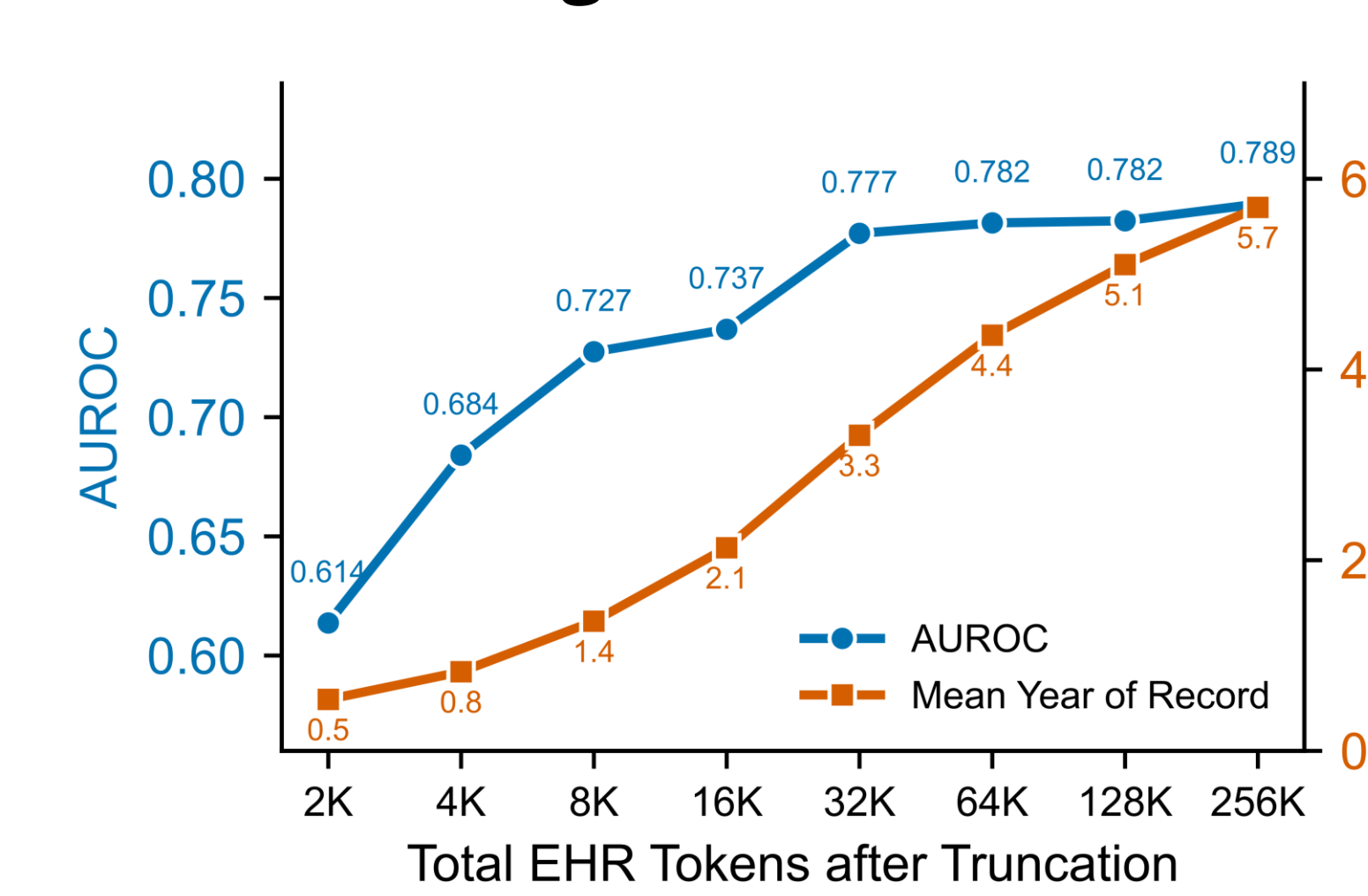
### Performance on unmatched lung cancer cohort (unmatched)

	AUROC	AUPRC	NPV	PPV	Sensitivity	Specificity
<b>Logistic Regression</b>	0.903 (0.892, 0.912)	0.055 (0.046, 0.070)	<b>1.000</b> (0.999, 1.000)	0.014 (0.013, 0.016)	<b>0.907</b> (0.880, 0.929)	0.737 (0.734, 0.739)
<b>XGBoost</b>	<b>0.925</b> (0.916, 0.934)	<b>0.195</b> (0.162, 0.233)	0.999 (0.999, 1.000)	0.017 (0.015, 0.019)	0.885 (0.856, 0.912)	0.785 (0.782, 0.787)
<b>KNN</b>	0.790 (0.772, 0.805)	0.014 (0.012, 0.017)	0.999 (0.998, 0.999)	0.010 (0.009, 0.011)	0.794 (0.754, 0.830)	0.657 (0.654, 0.659)
<b>Single-agent (GPT-4.1-mini)</b>	0.866 (0.849, 0.883)	0.068 (0.051, 0.089)	0.999 (0.999, 0.999)	<b>0.018</b> (0.012, 0.020)	0.741 (0.710, 0.859)	<b>0.836</b> (0.747, 0.839)
<b>TrajOnco (GPT-4.1-mini)</b>	<b>0.871</b> (0.855, 0.885)	<b>0.071</b> (0.053, 0.092)	<b>0.999</b> (0.999, 0.999)	<b>0.017</b> (0.013, 0.019)	<b>0.772</b> (0.742, 0.861)	<b>0.825</b> (0.741, 0.827)

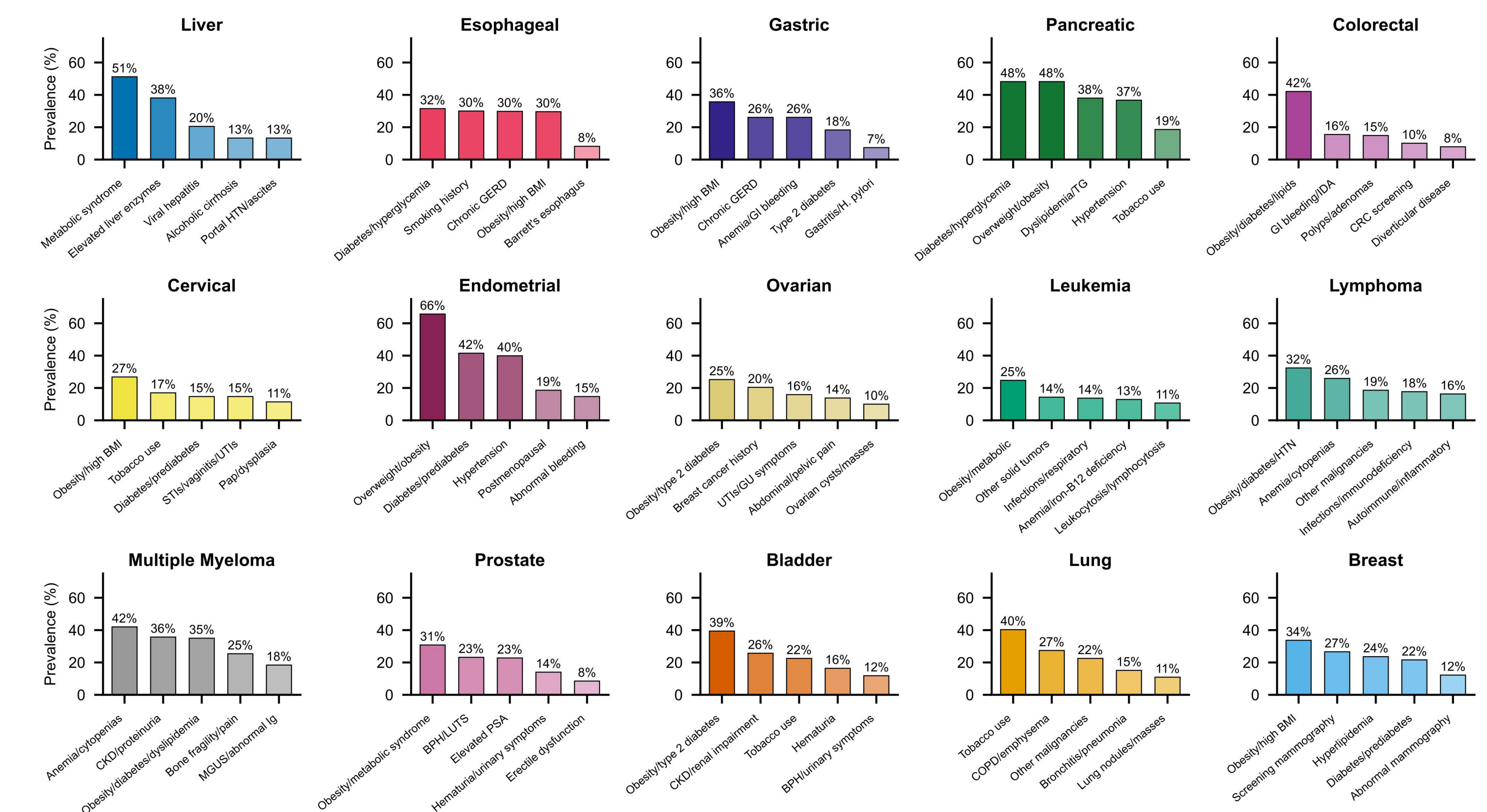
### Base model selection (matched lung cancer cohort)



### Effect of trajectory length (matched lung cancer cohort)



### Top themes for risk estimation



## References

- Zhang et al., Chain of Agents: Large Language Models Collaborating on Long-Context Tasks. NeurIPS 2024
- Zeng et al., Traj-CoA: Patient Trajectory Modeling via Chain-of-Agents for Lung Cancer Risk Prediction. NeurIPS 2025 GenA14Health Workshop
- Zeng et al., TrajOnco: a multi-agent framework for temporal reasoning over longitudinal EHR for multi-cancer early detection. arXiv: 2604.10386

