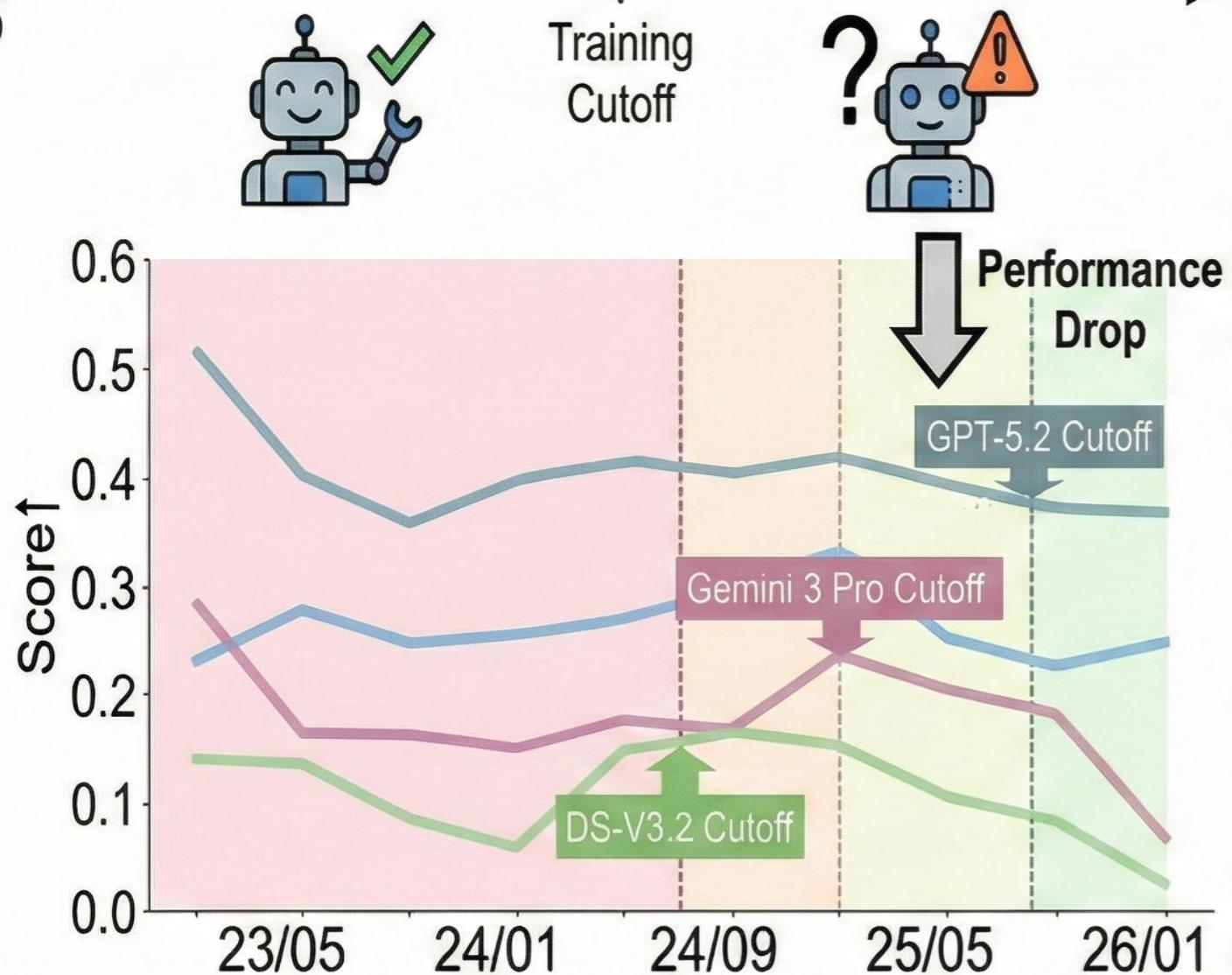
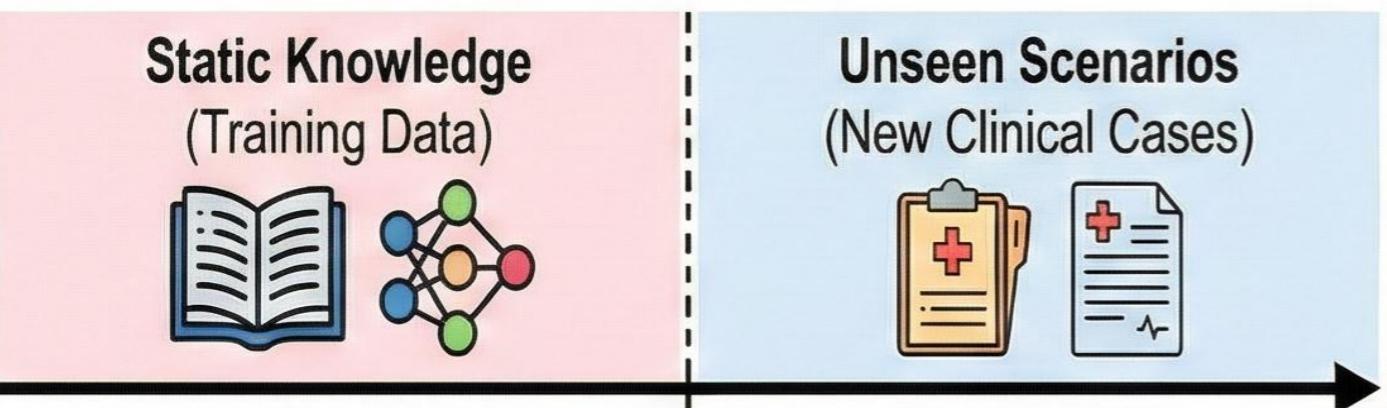
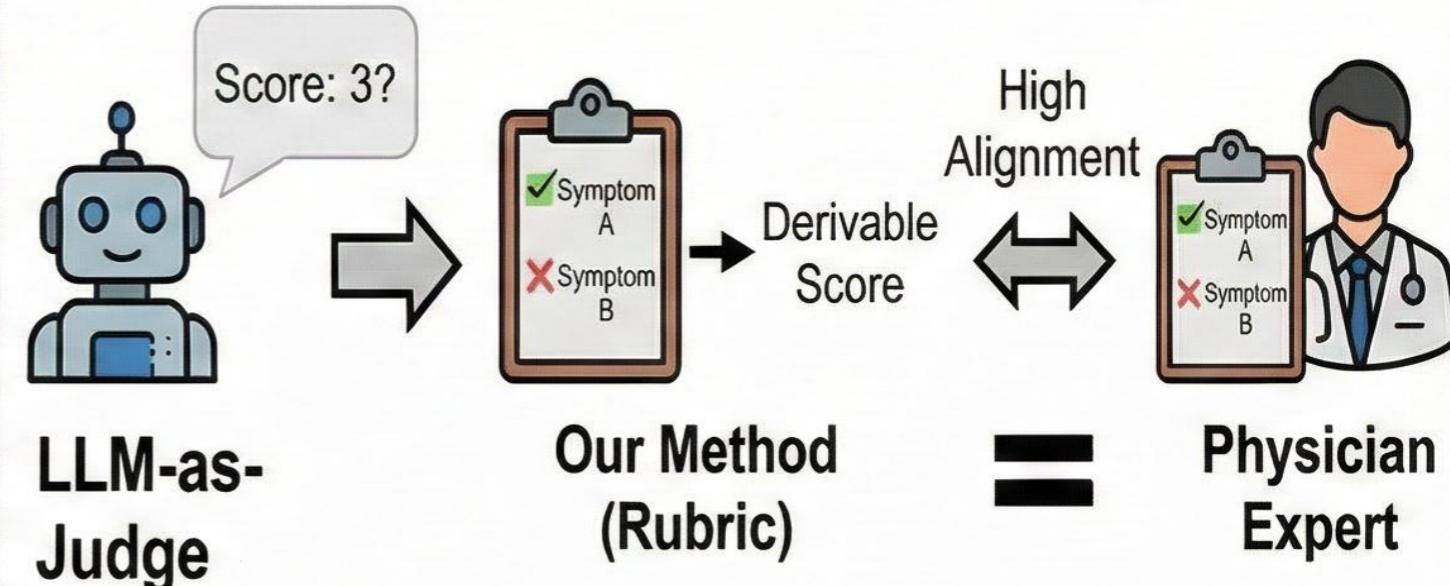


## Panel A: The Challenge – Data Contamination

Schematic



## Panel B: The Challenge – Evaluation Misalignment



### Score Distribution Comparison

