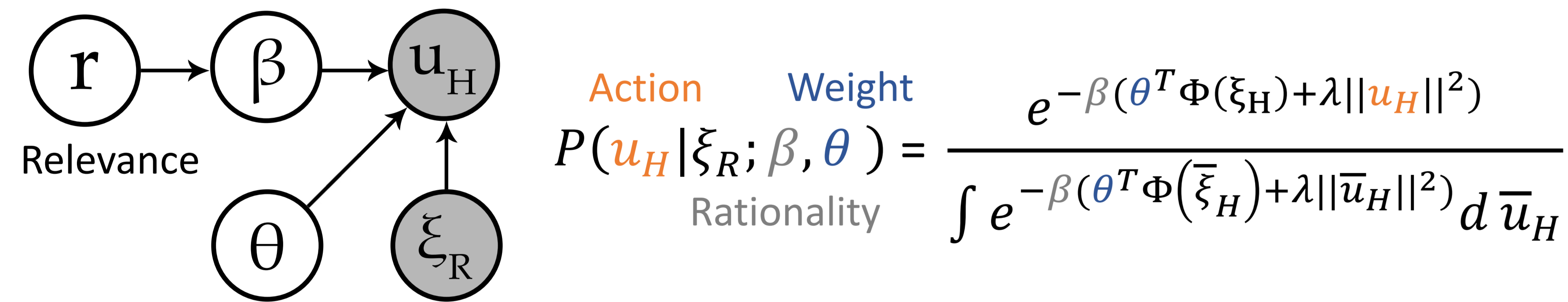


What if what  $H$  wants is outside the robot's hypothesis space  $\Theta$ ?

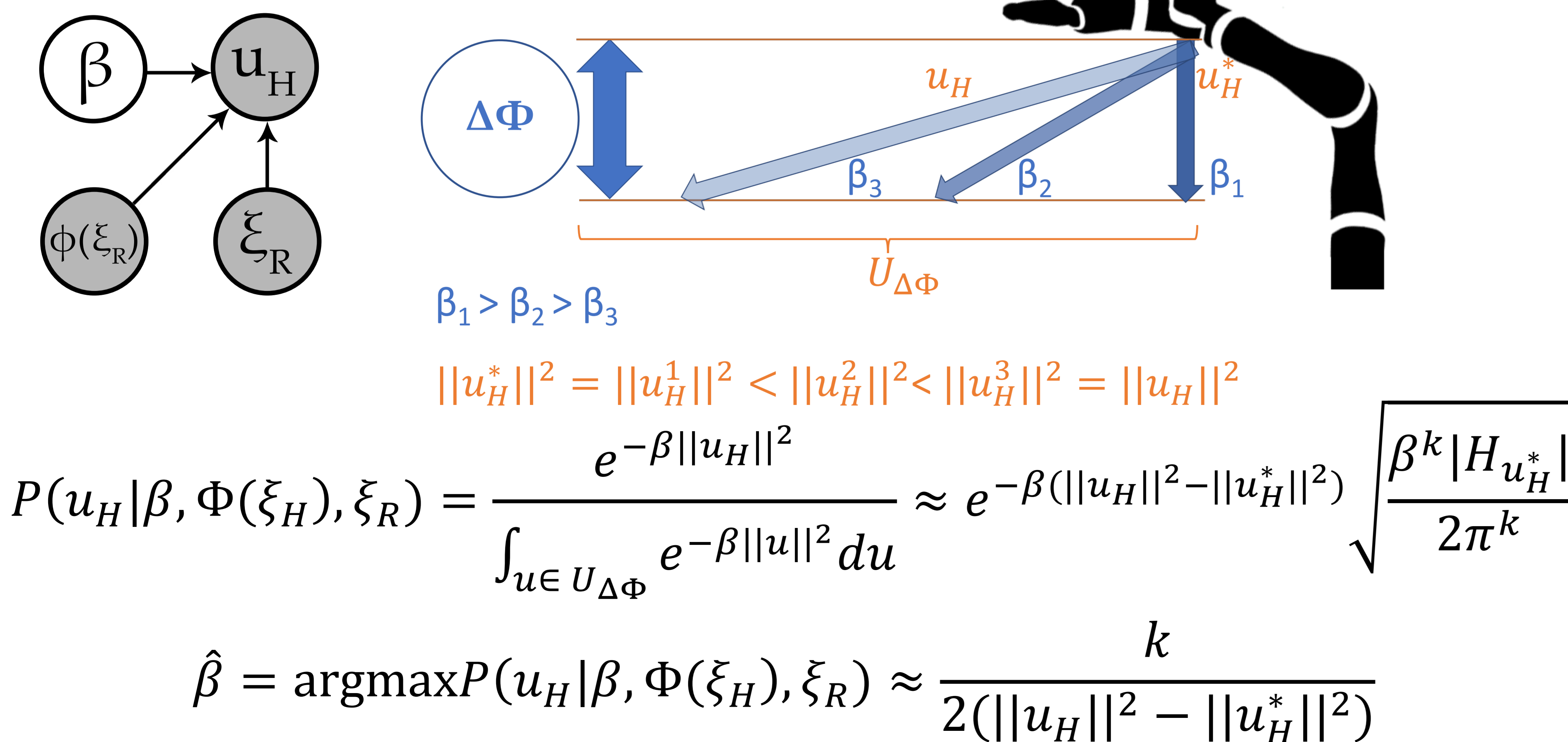
Key Insight: If the human *seems* to be *suboptimal* for any hypothesis, chances are we don't have the *right* hypothesis space.

## Relevance to $R$ 's hypothesis space dictates apparent rationality

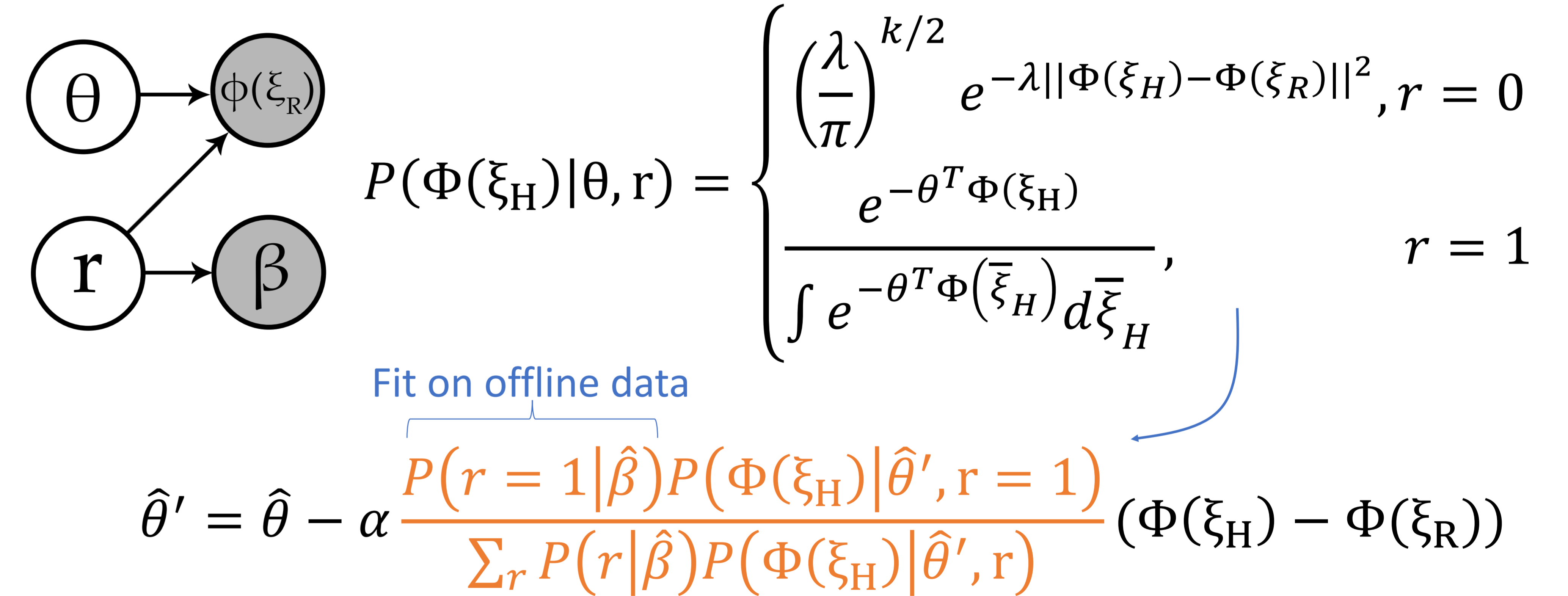


## Real-time approximation

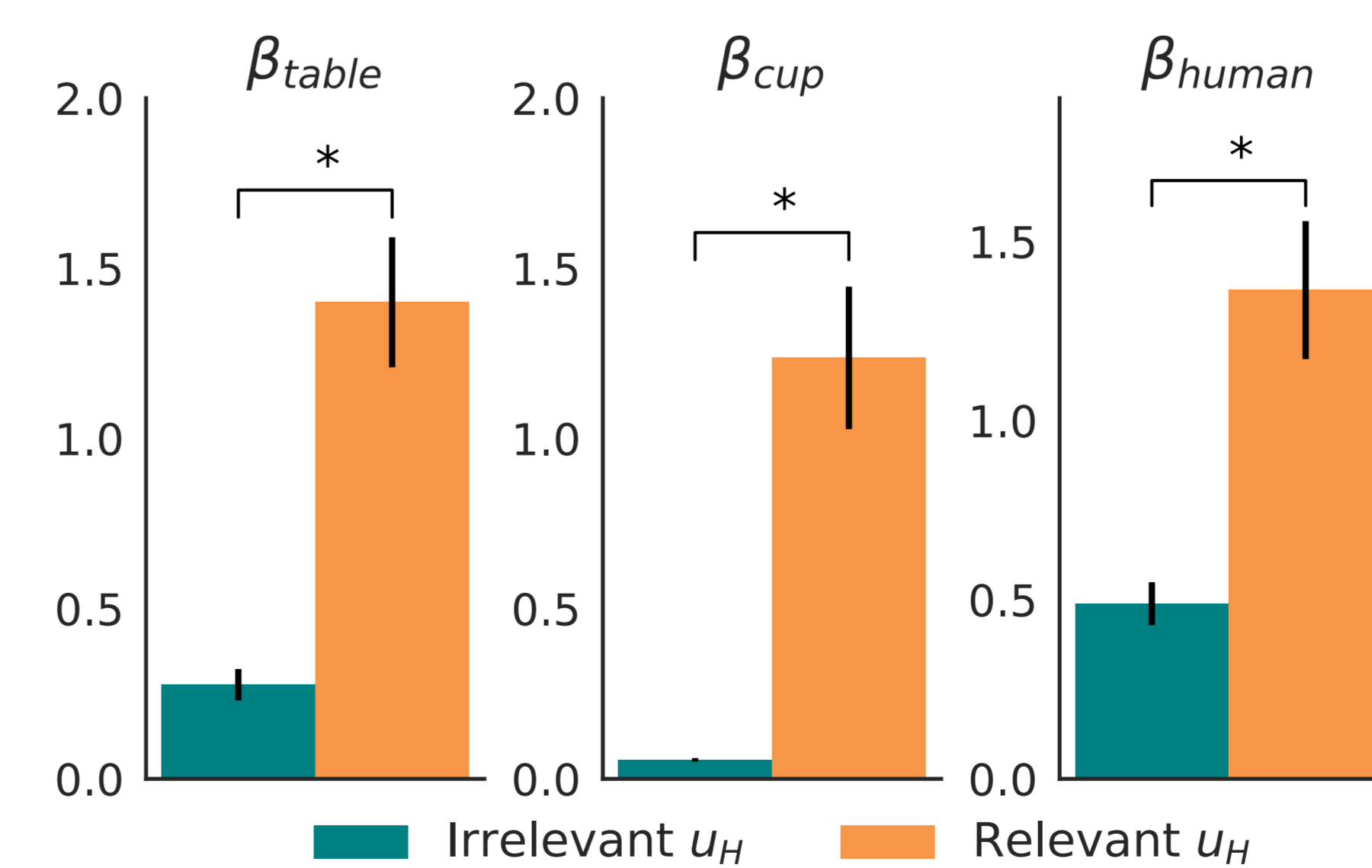
### a) Apparent rationality estimation



## b) Relevance-aware approximate MAP estimate:

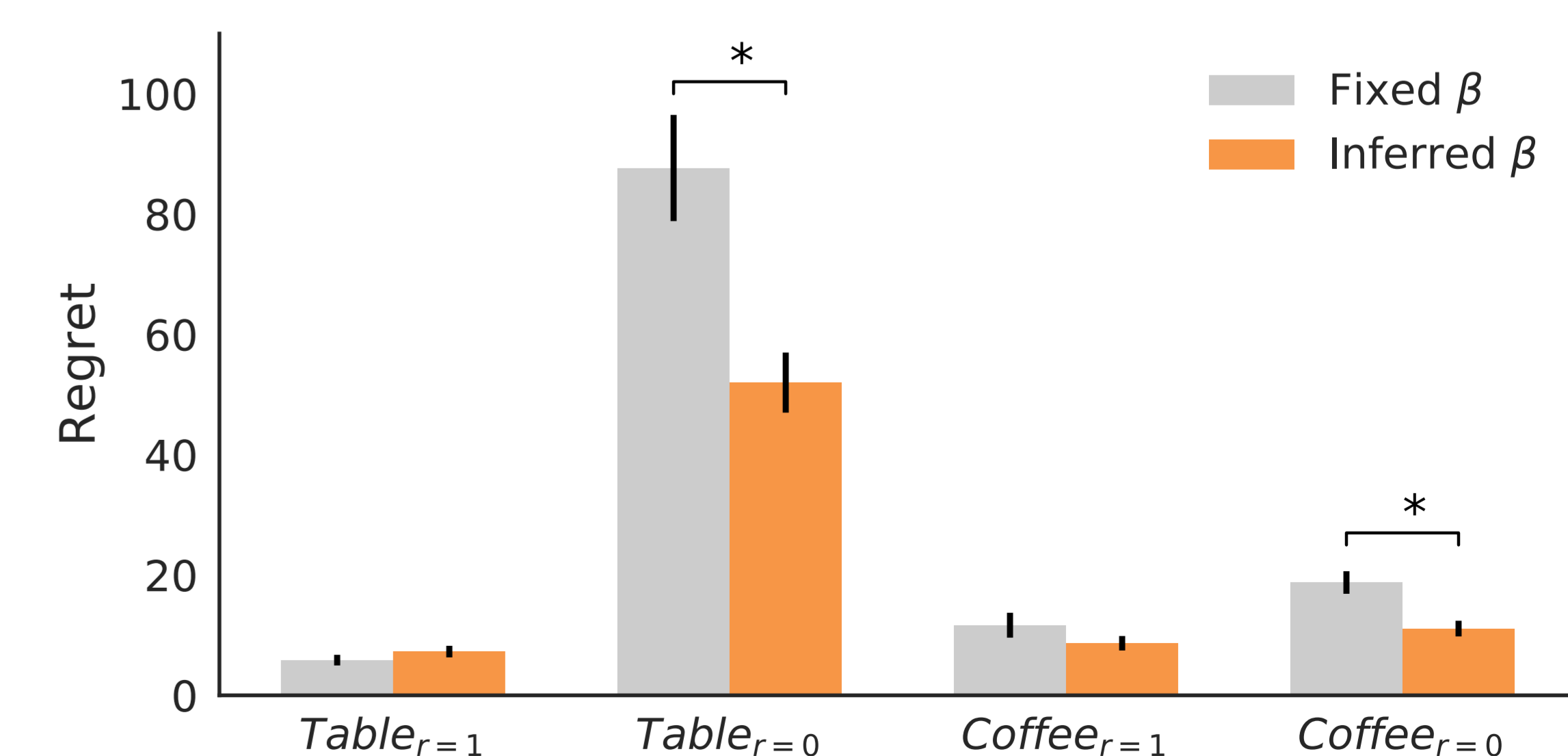


## Offline $\beta$ -estimation



If  $r = 1$ , action appears rational ( $\beta \uparrow$ ).  
If  $r = 0$ , action appears irrational ( $\beta \downarrow$ ).

## User study



If  $r=0$ , **relevance-aware** reduces unintended learning, while keeping good accuracy if  $r=1$ .