

## Motivation

- Reduce expensive reward calculation;
- Leverage wasted historical and expert data;
- Accelerate planning without performance loss.

## Key idea

- Efficient incremental update method for the Multiple Importance Sampling estimator;
- Experience-based value estimation utilizing expert demonstrations without planning;
- An MCTS-inspired online algorithm (IR-PFT) that accelerates computations by reusing data from previous planning sessions.

## Problem Formulation

$$V^\pi(b_k) = \mathbb{E}_{b_{k+1:k+d}} [G_k | b_k, \pi],$$

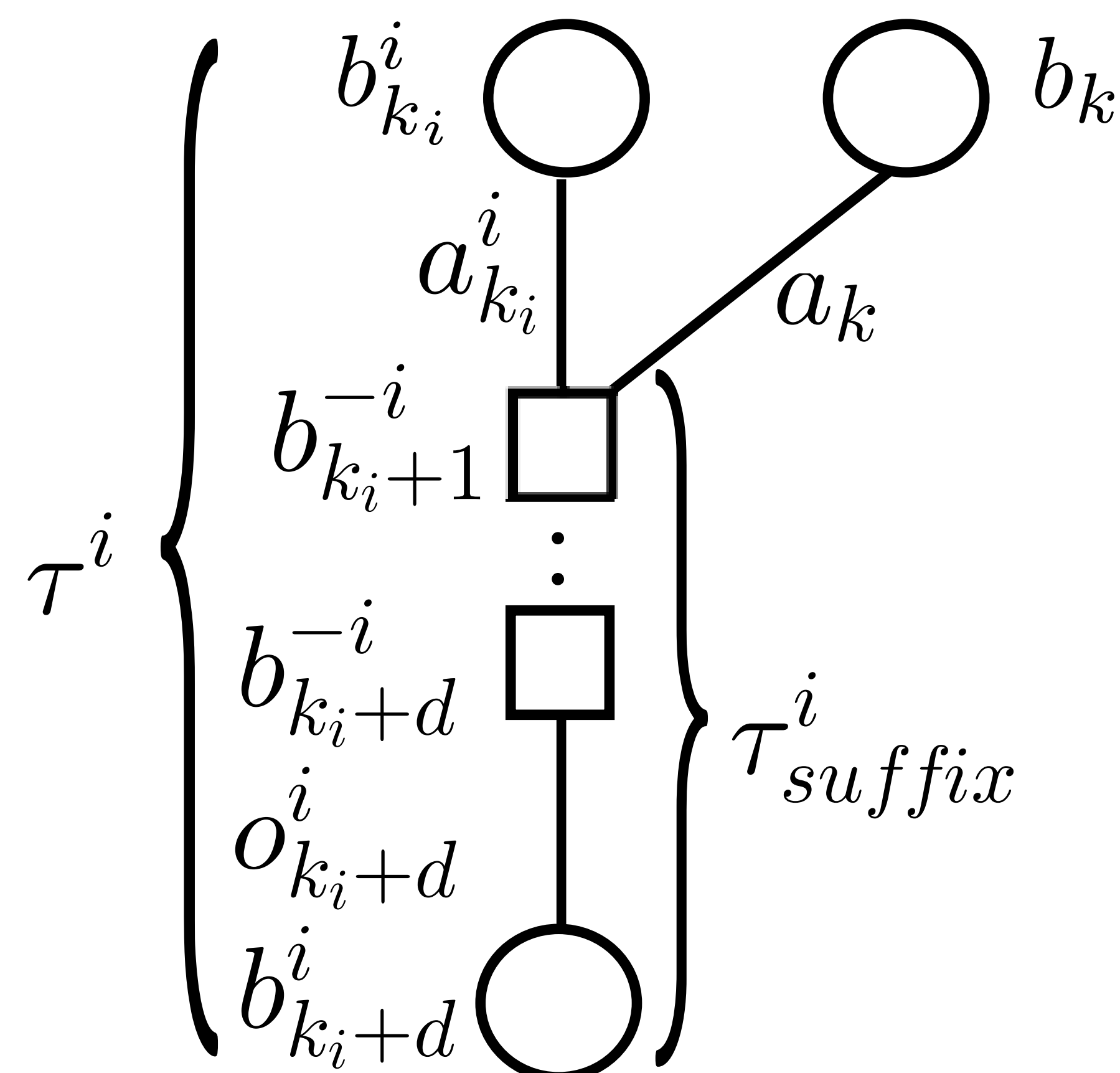
$$\pi \triangleq \pi_{k:k+d-1},$$

$$G_k = \sum_{i=k}^{k+d-1} r(b_i, \pi_i(b_i), b_{i+1}),$$

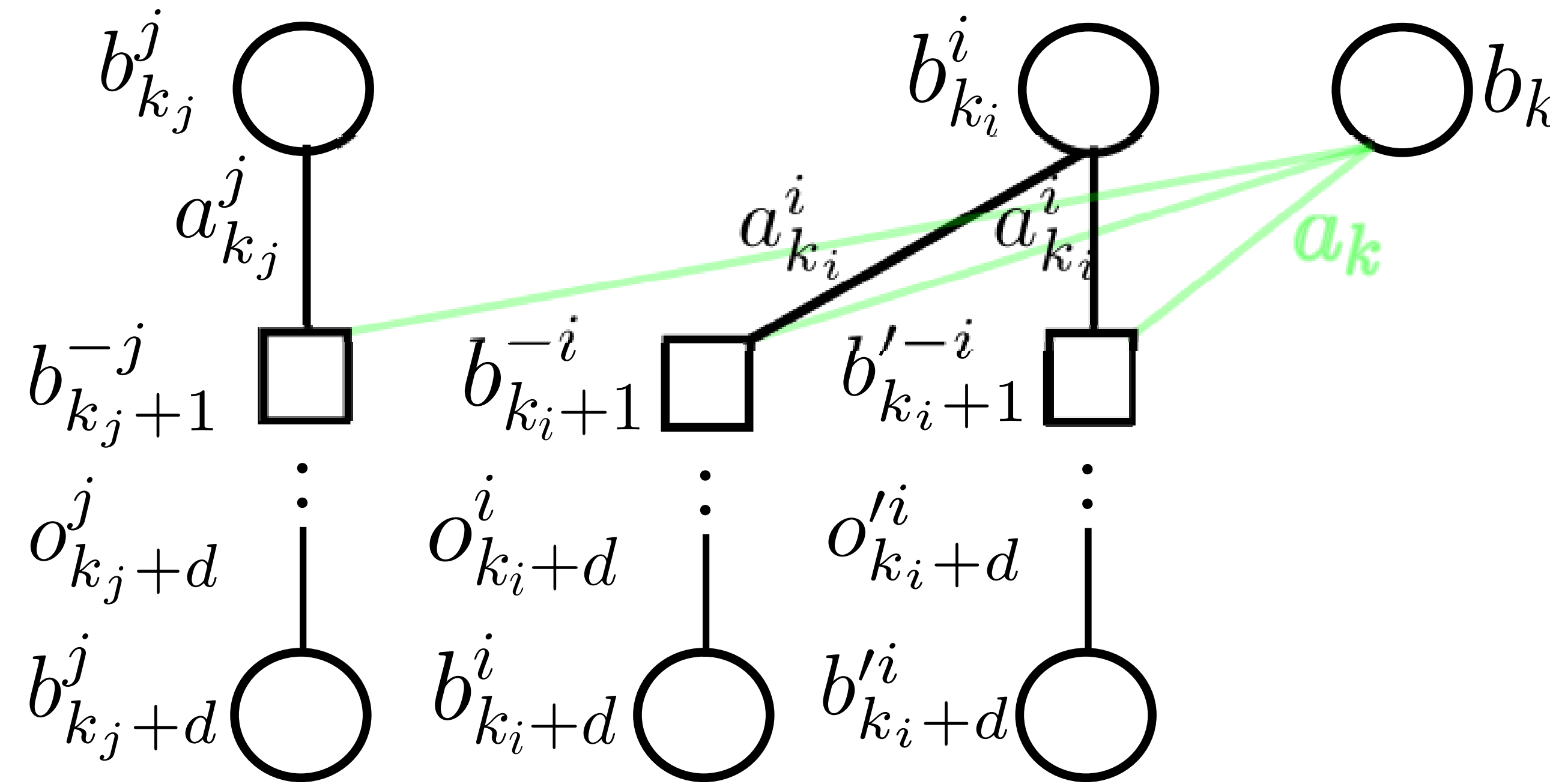
$$Q^\pi(b_k, a) = \mathbb{E}_{b_{k+1}} [r(b_k, a, b_{k+1}) + V^\pi(b_{k+1})].$$

$$\hat{Q}_{MIS}^\pi(b_k, a_k) \triangleq \sum_{m=1}^M \sum_{l=1}^{n_m} \frac{P(b_{k_m+1}^{-l,m} | b_k, a_k)}{\sum_{j=1}^M n_j \cdot P(b_{k_m+1}^{-l,m} | b_{k_j}^j, a_{k_j}^j)} \cdot \hat{G}_k^{l,m}.$$

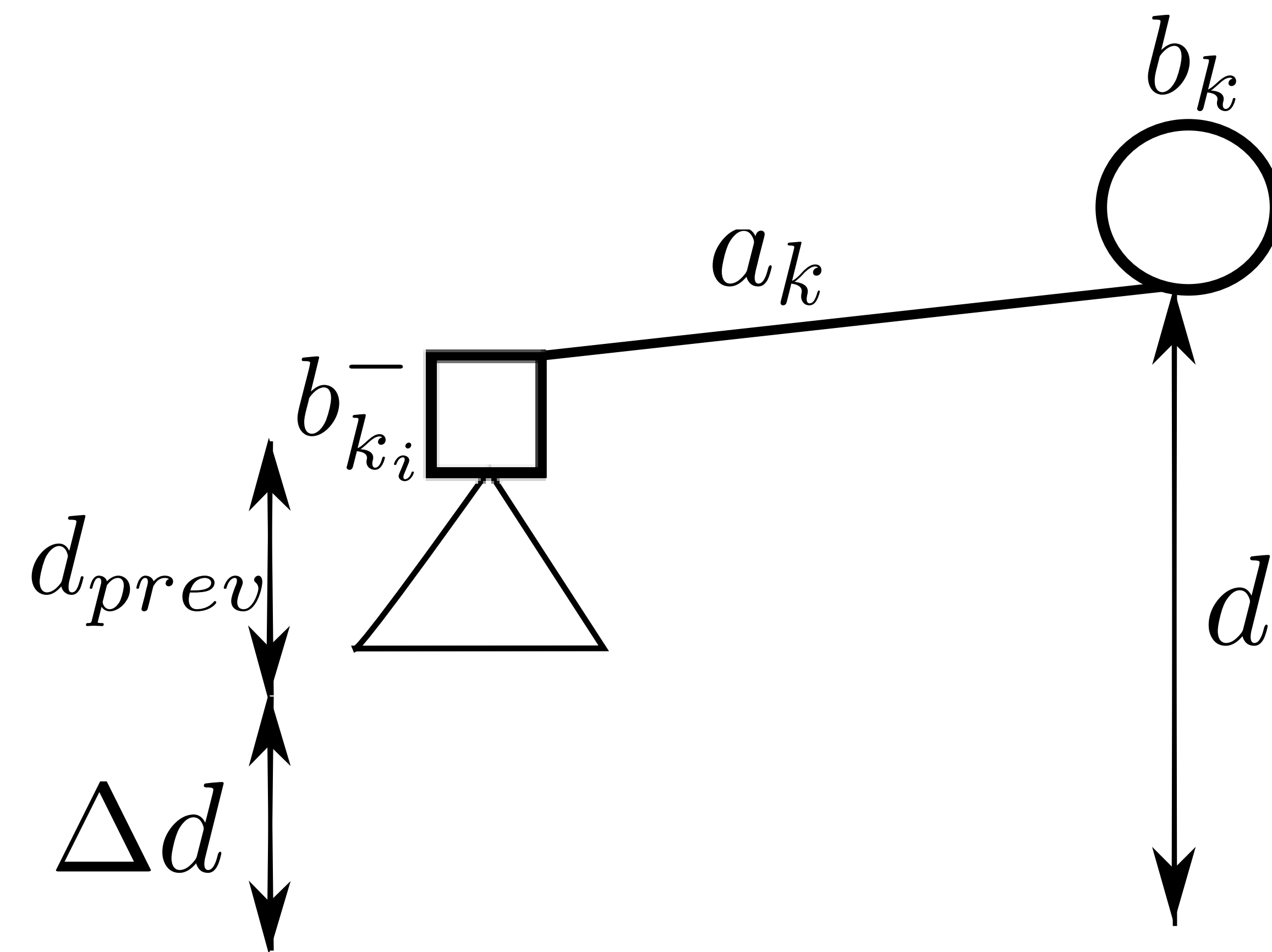
## Trajectory Reuse



## Multiple Trajectory Reuse



## Horizon Alignment

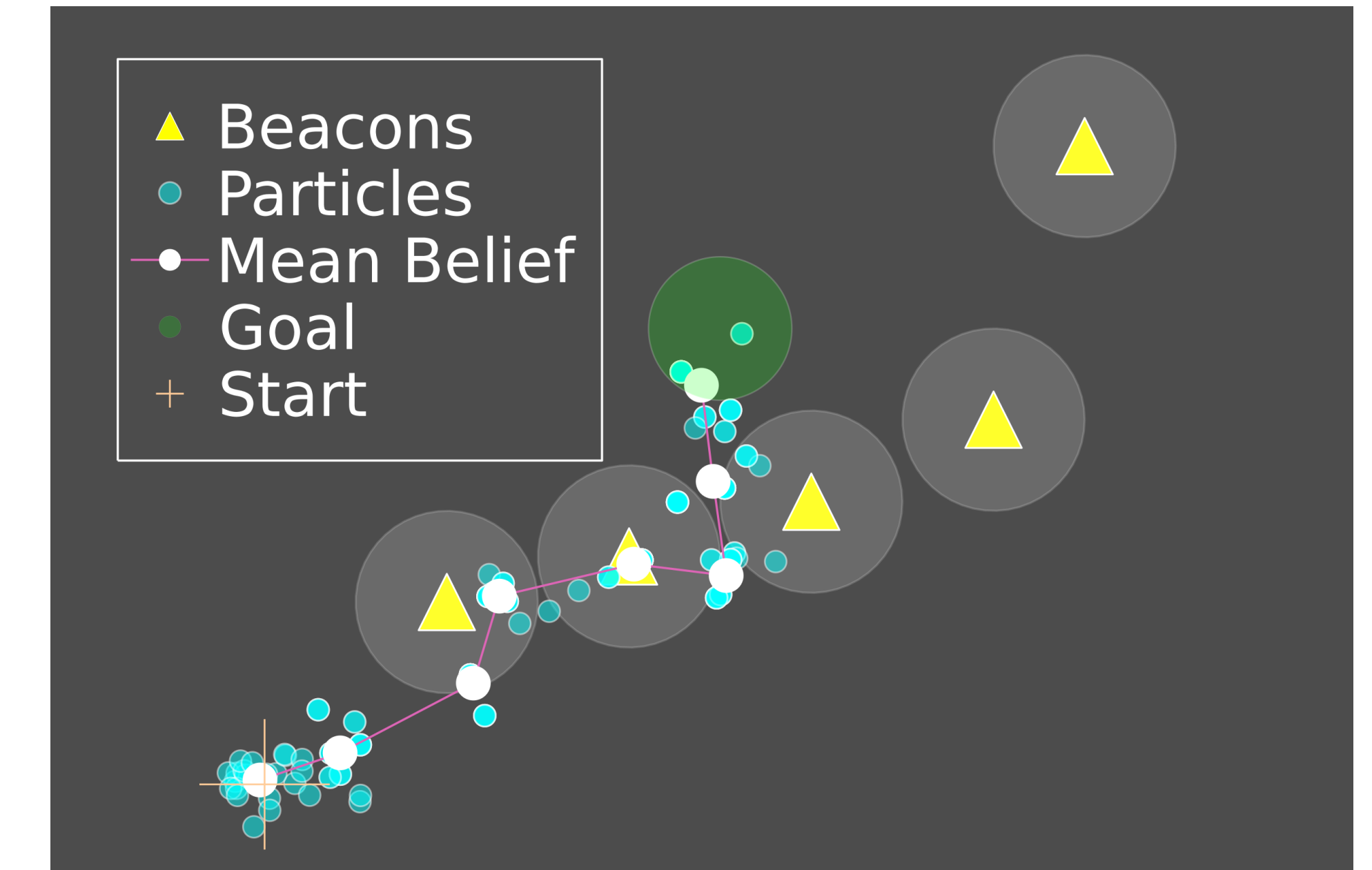


## IR-PFT Decision Flow

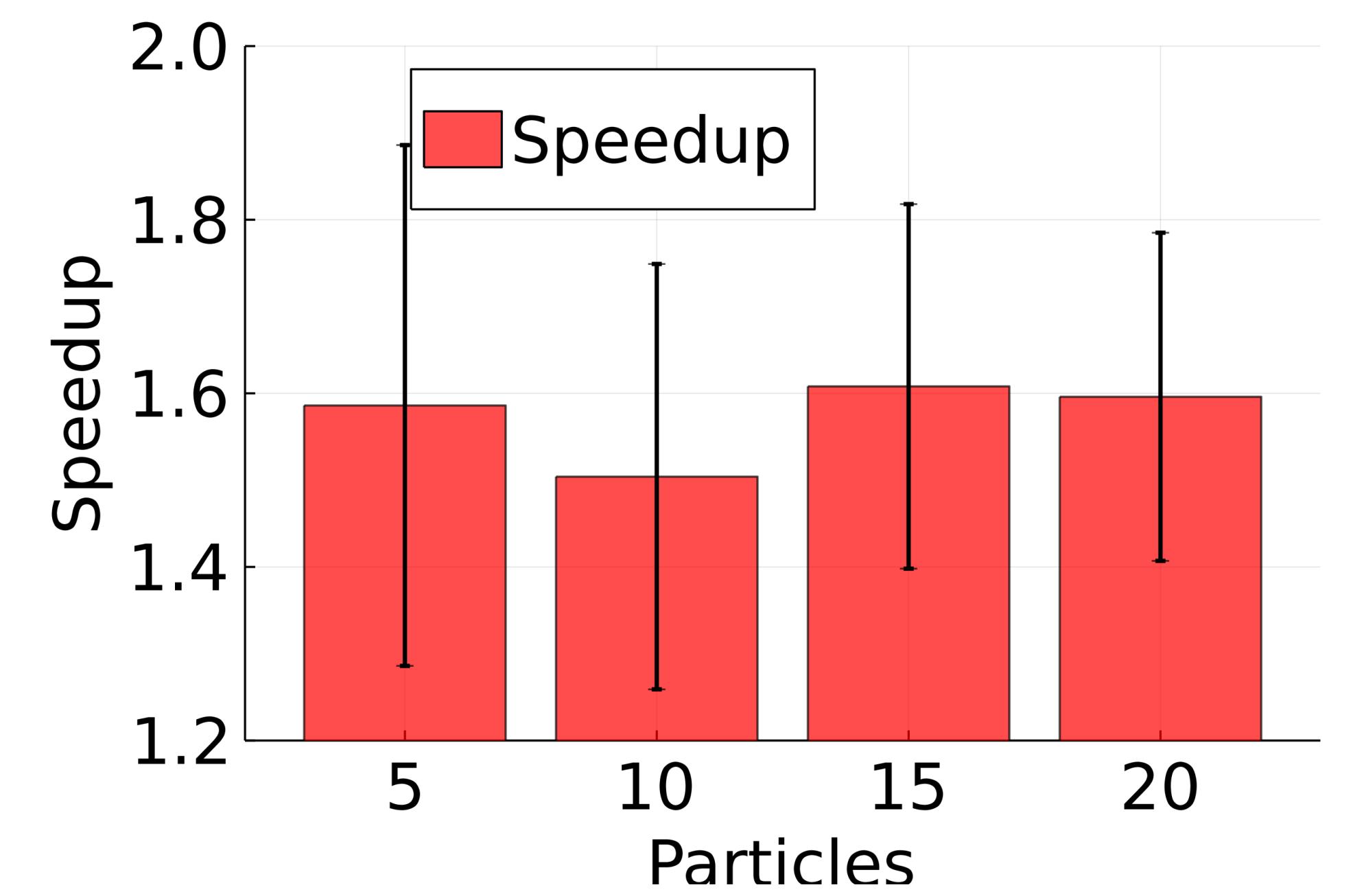
1. **Evaluate Condition:** Is current node the root and are prior candidates available?
2. **If YES (Reuse Path):**
  - **Select:** Retrieve best candidate  $b'^{-}$  from prior dataset  $D$ .
  - **Align:** Resolve horizon discrepancies (*Fill Horizon*).
  - **Update:** Calculate  $Q(b, a)$  efficiently via Incremental MIS.
3. **If NO (Simulate Path):**
  - **Propagate:** Generate new belief via Particle Filter.
  - **Rollout:** Execute standard MCTS trajectory.

## Results

### Light Dark



### Speedup



### Reward Comparison

