

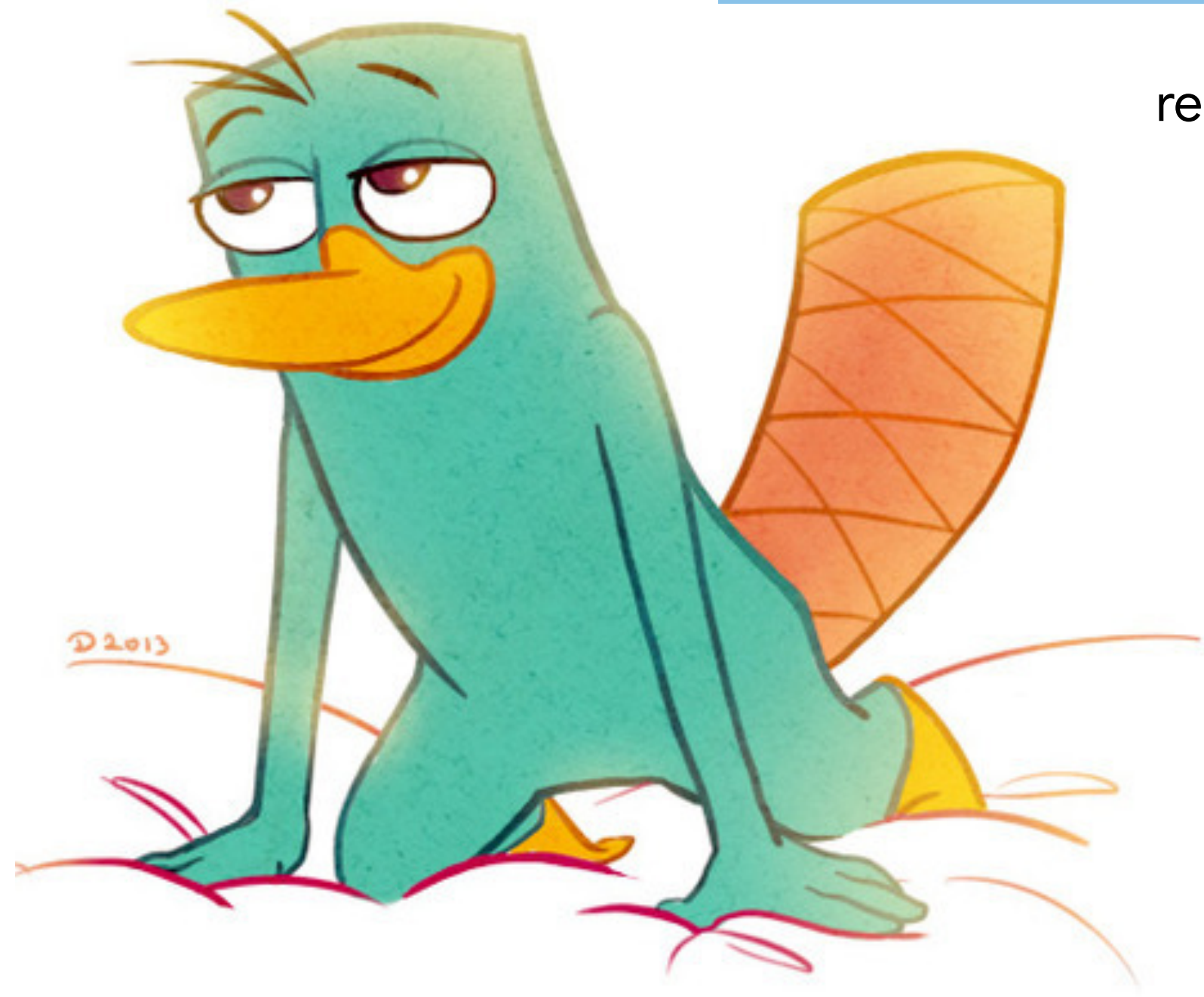
SCALE-FREE ADAPTIVE PLANNING

FOR DETERMINISTIC DYNAMICS & γ -DISCOUNTED REWARDS

PETER BARTLETT, VICTOR GABILLON, JENNIFER HEALEY & MICHAL VALKO

BRAND NEW ADAPTIVE MCTS PLANNER: **Plat γ POOS**

adapts its behavior to an unknown range of rewards



requires **no assumptions** or knowledge of noise

empirically learns much **faster** than UCB approaches

gets the **fast rate** of deterministic planning in low noise for all regimes
 → exponentially faster than OLOP
not a rare case!

adapts also to the **global smoothness** ρ and beyond the base smoothness provided by γ



OUR LOVELY PLAT γ POOS

Input: n, A

Initialization: open the root node \emptyset , h_{\max} times
 $h_{\max} \leftarrow \lfloor \frac{n}{2(\log_2 n + 1)^2} \rfloor$, $p_{\max} \leftarrow \lfloor \log_2(h_{\max}) \rfloor$

For $h = 1$ to h_{\max} ◀ **exploration** ▶

For $p = \lfloor \log_2(h_{\max} / \lceil h^2 \gamma^{2h} \rceil) \rfloor$ down to 0
 open $\lceil h 2^p \gamma^{2h} \rceil$ times the at most $\lfloor \frac{h_{\max}}{h \lceil h 2^p \gamma^{2h} \rceil} \rfloor$
 non-opened $a^{h,i} \in A^h$ with highest values $\hat{u}(a^{h,i})$ and given $T_{a^{h,i}} \geq \lceil (h-1) 2^p \gamma^{2(h-1)} \rceil$

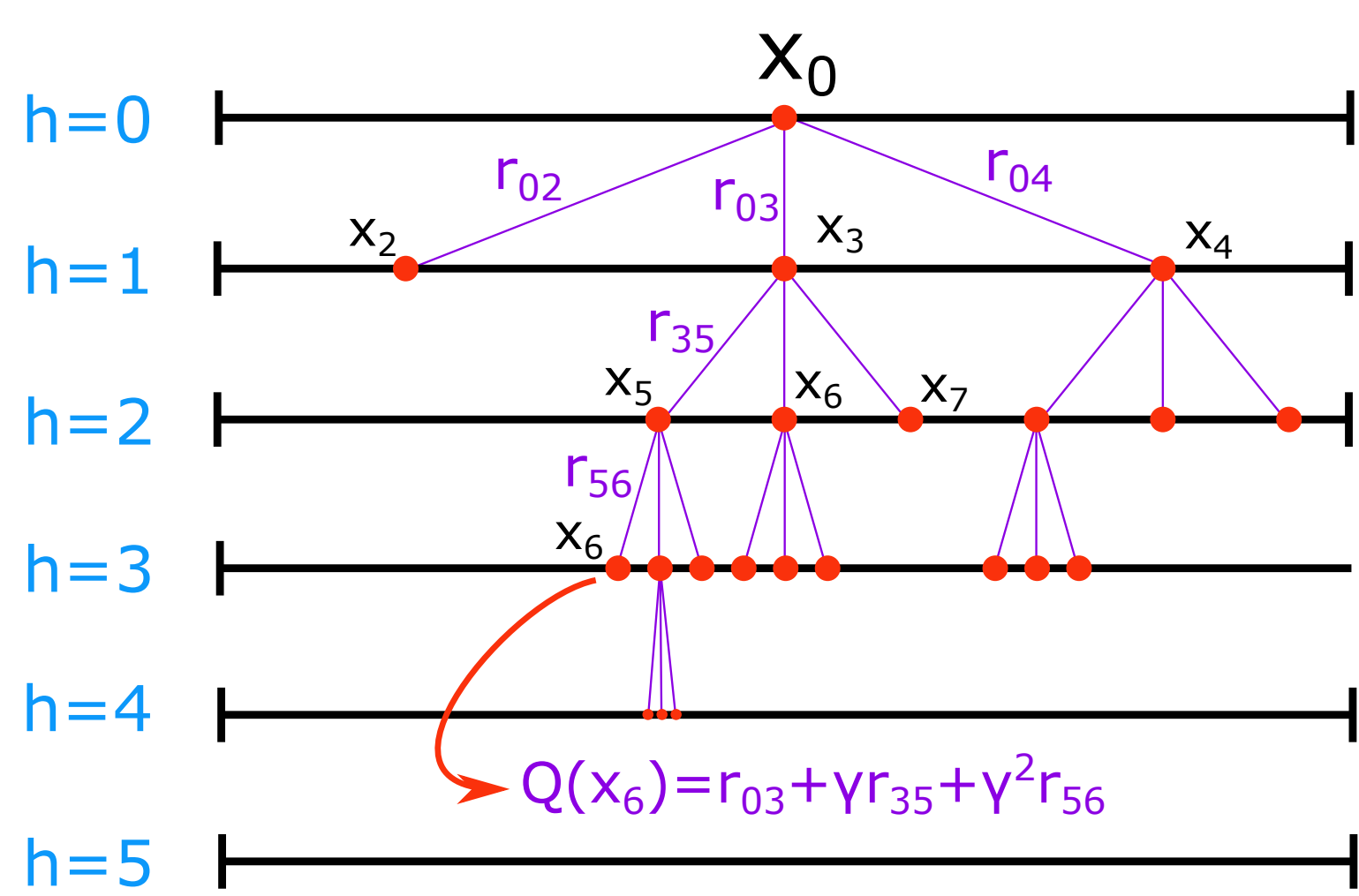
For $p \in [0 : p_{\max}]$ ◀ **cross-validation** ▶

evaluate $(t+1) \gamma^{2t} h_{\max} (1-\gamma^2)^2$ times the actions at round t , a_t^p , of the candidates:
 $a^p \leftarrow \arg \max_{a \in A^p : \forall t \in [2:h(a)], T_{a^{h,t}} \geq \lceil (t-1) 2^p \gamma^{2(t-1)} \rceil} \hat{u}(a)$

Output $a^n \leftarrow \arg \max_{\{a^p, p \in [0:p_{\max}]\}} \hat{u}(a^p)$

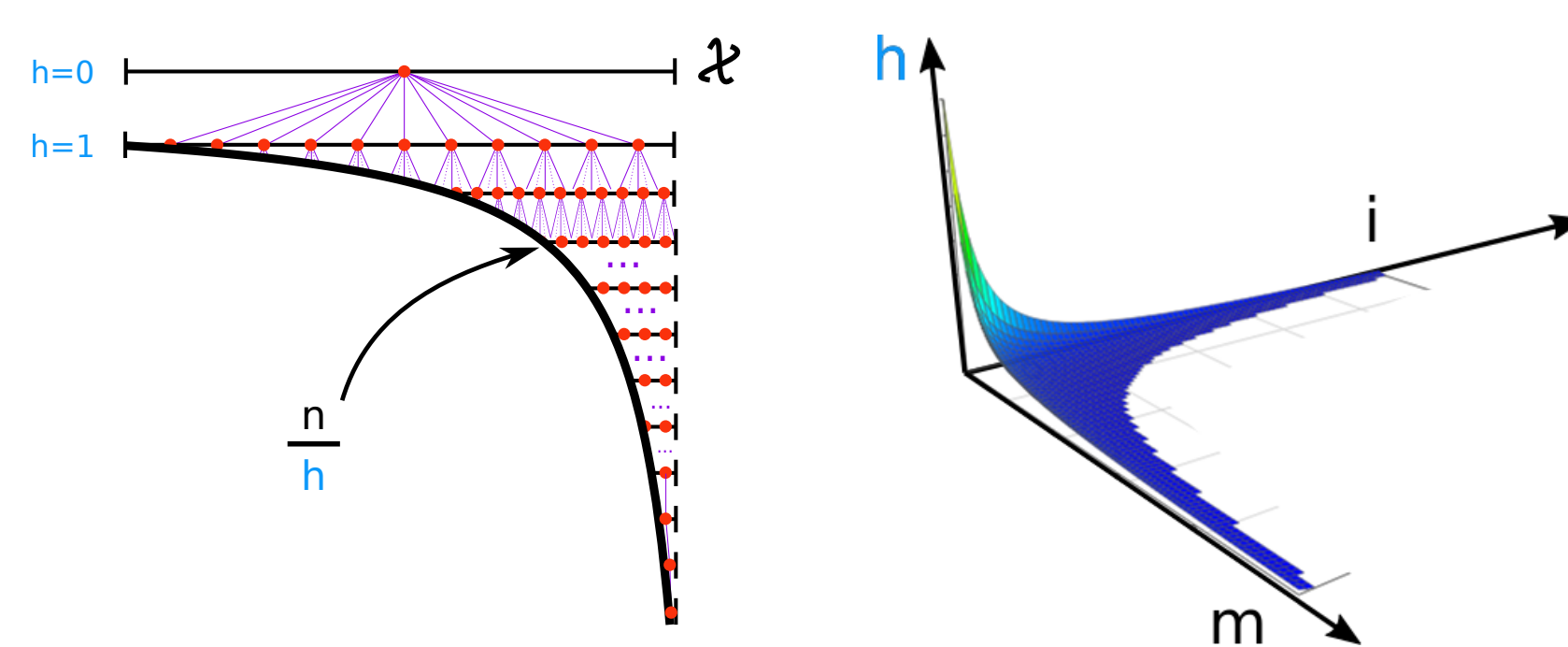
- implements **Zipf** exploration for MCTS Stroqu00L
- explicitly pulls an action at depth $h+1$, γ times less than action at depth h , ($Q^*(x, a) = r(x, a) + \sup_{\pi} \sum \gamma^t r(x_t, \pi(x_t))$),
- does not use UCB & no use of R_{\max} and b

TREE SEARCH FOR THE WIN!



Is this zero order optimization?

ZIPF: Sequ00L AND Stroqu00L



A simple parameter-free and adaptive approach to optimization under a minimal local smoothness assumption, Bartlett, Gabillon and Valko, Algorithmic Learning Theory, 2019

MCTS SETTING

MDP with starting state $x_0 \in X$, action space A

n interactions: At time t playing a_t in x_t leads to

Deterministic dynamics $g: x_{t+1} \triangleq g(x_t, a_t)$,

Reward: $r_t(x_t, a_t) + \varepsilon_t$ with ε_t being the noise

Objective: Recommend action $a(n)$ minimizing

$$r_n \triangleq \max_{a \in A} Q^*(x, a) - Q^*(x, a(n)) \quad \text{simple regret}$$

where $Q^*(x, a) \triangleq r(x, a) + \sup_{\pi} \sum \gamma^t r(x_t, \pi(x_t))$

Assumption: $r_t \in [0, R_{\max}]$ and $|\varepsilon_t| \leq b$

Approach: Explore without parameters R_{\max} & b

OLOP (BUBECK AND MUNOS, 2010)

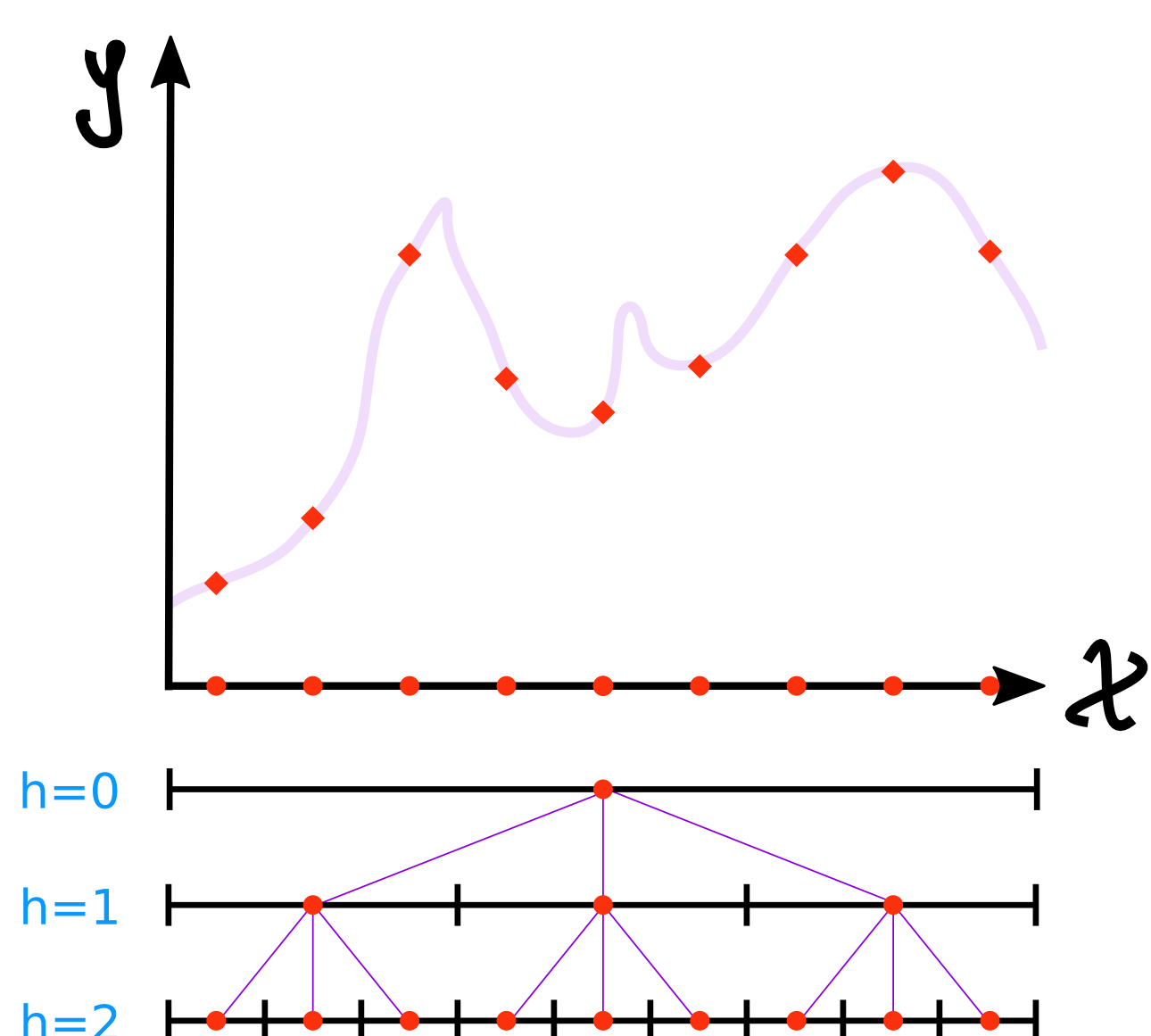
OLOP implements Optimistic Planning using Upper Confidence Bound (UCB) on the Q value of a sequence of q actions a_1, \dots, a_q :

$$\hat{Q}_t(a_{1:q}) \triangleq \underbrace{\sum_{h=1}^q \left(\gamma^h \hat{r}_h(t) + \frac{\gamma^h b}{\sqrt{T_{a_h}(t)}} \right)}_{\text{estimation of observed reward}} + \underbrace{\frac{R_{\max} \gamma^{q+1}}{1-\gamma}}_{\text{unseen reward}}$$

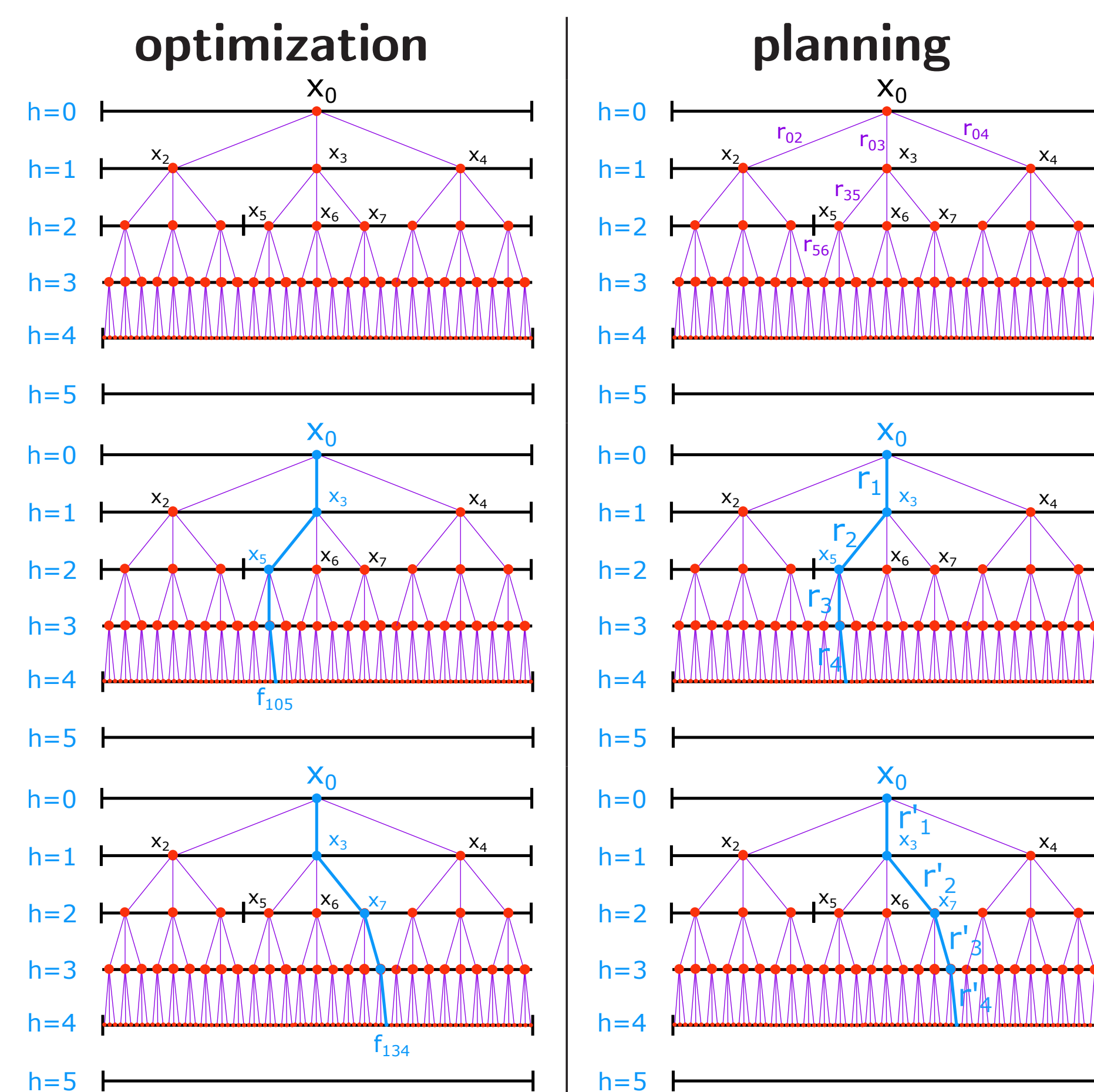
in optimization under a fixed budget n , **excellent strategies ignore R_{\max} or b**

BLACK-BOX OPTIMIZATION

use the partitioning to explore f (uniformly)

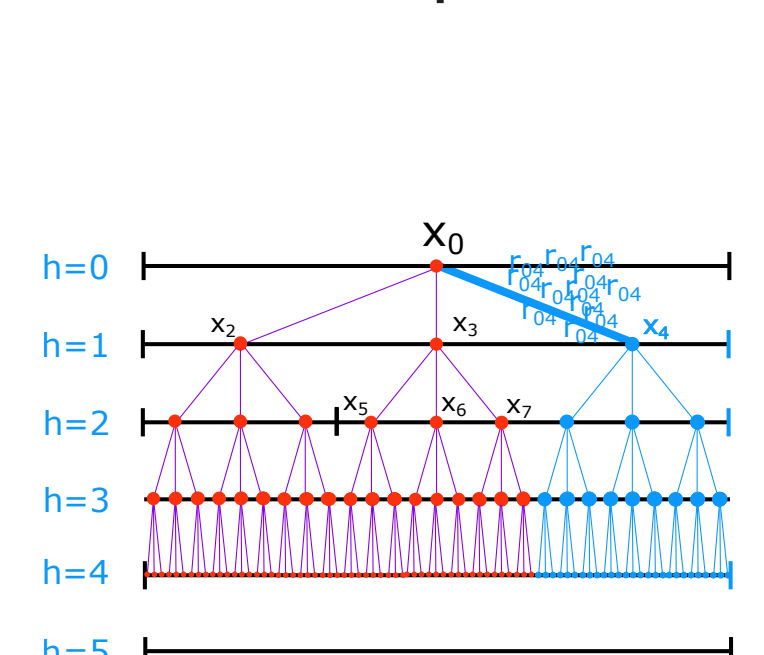


OPTIMIZATION VS. PLANNING



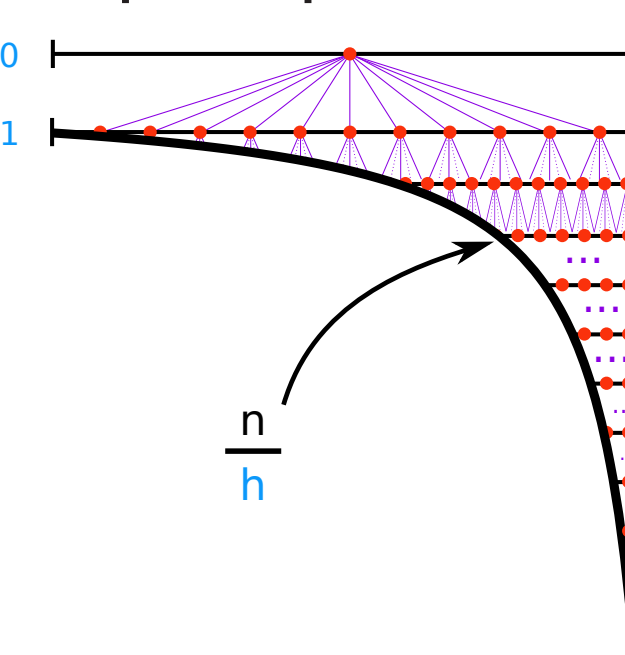
lower regret for planning! (Bubeck+Munos'10)
 thanks to the **reuse of samples**

Uniform exploration



not sharing information

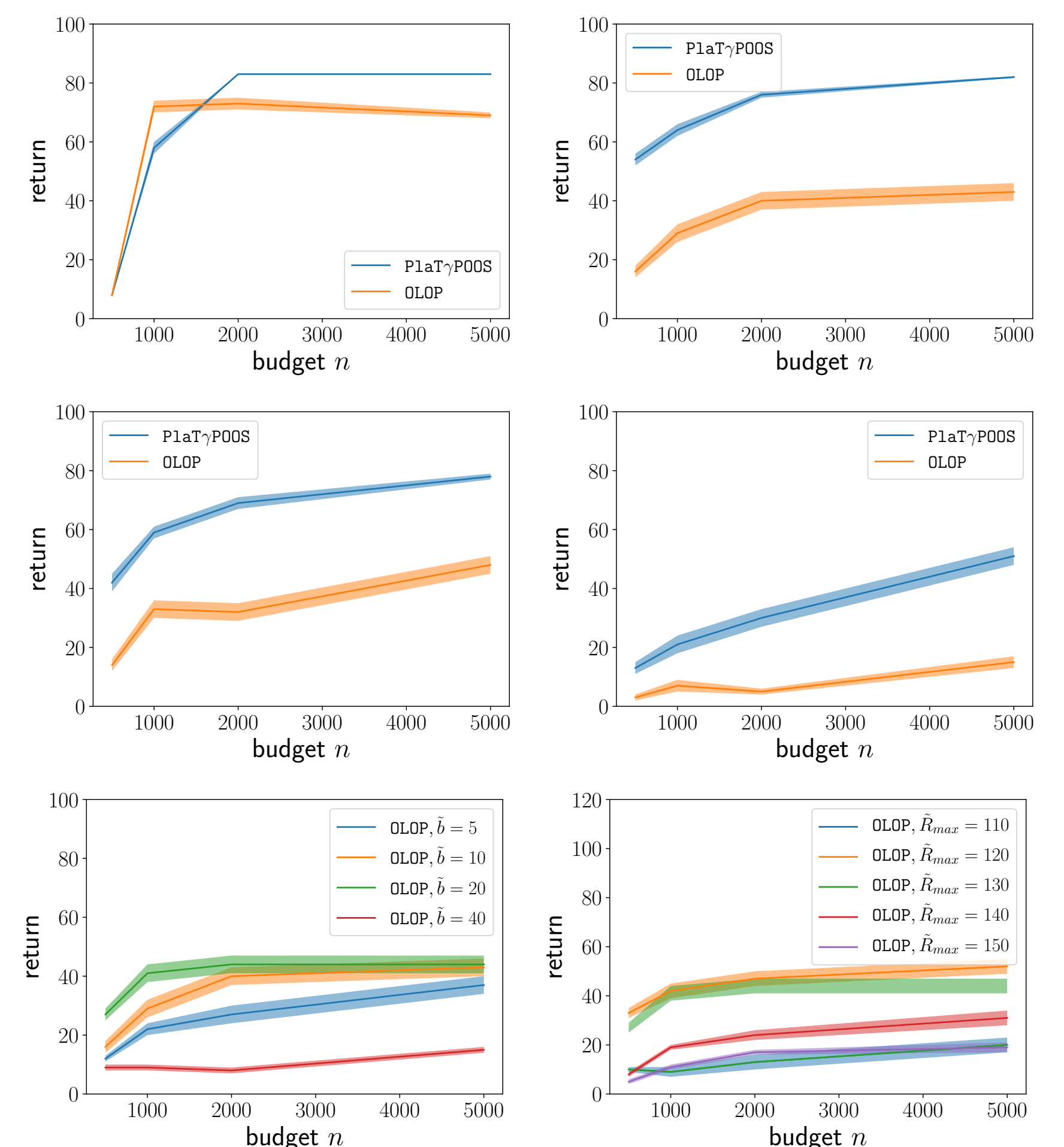
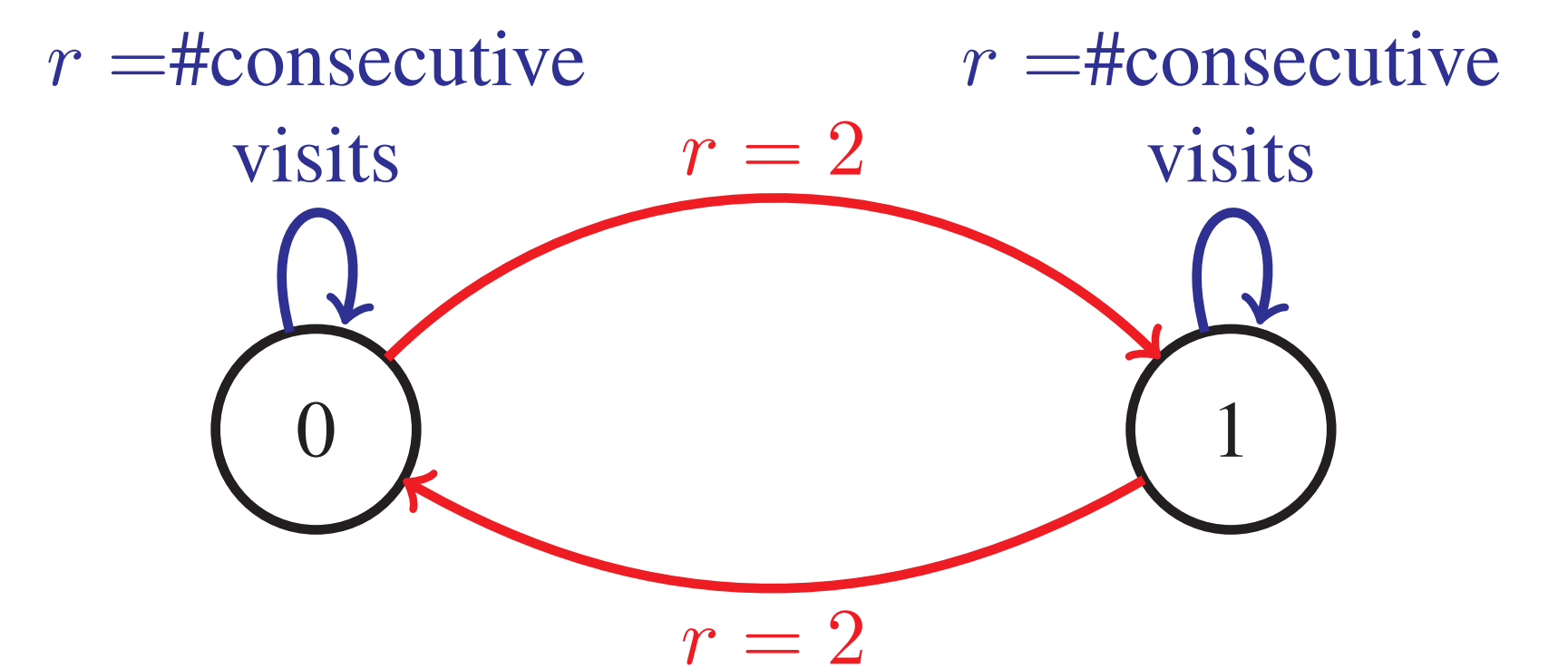
Zipf exploration



Sharing information

Bubeck & Munos: Only for uniform strategies ...
 We figured the amount the samples needed!

NUMERICAL SIMULATIONS



$b = 10$ (top center), $b = 20$, (top right), $b = 50$ (bottom left). Bottom right: true b is set to 10.

Empirical behavior in the figures mimics the behavior of the complexities in the table.

	$\gamma^2 \kappa \leq 1$		$\gamma^2 \kappa \geq 1$	
	High noise (ii)	Low noise (ii)	High noise (iii)	Low noise (iii)
ε	$\left(\frac{n}{b^2}\right)^{-\frac{1}{2}}$	$\rho \sqrt{n}$	$\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\gamma^2 \kappa / \rho^2)}}$	$\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$
ε	$\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$	$\kappa = 1 : \rho^n$ $\kappa > 1 : \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$	$\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$	$\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$