

a

Exploration with
outcome prediction

Exploitation with
policy gradient

Uncertainty

Confidence

π^Q

π^{PG}

Convex policy

