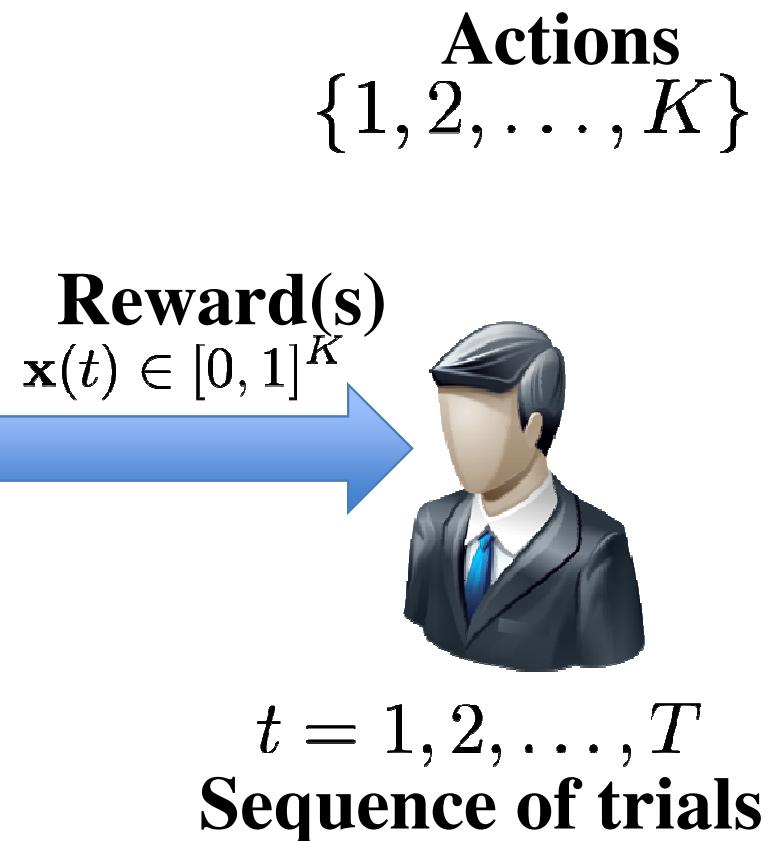


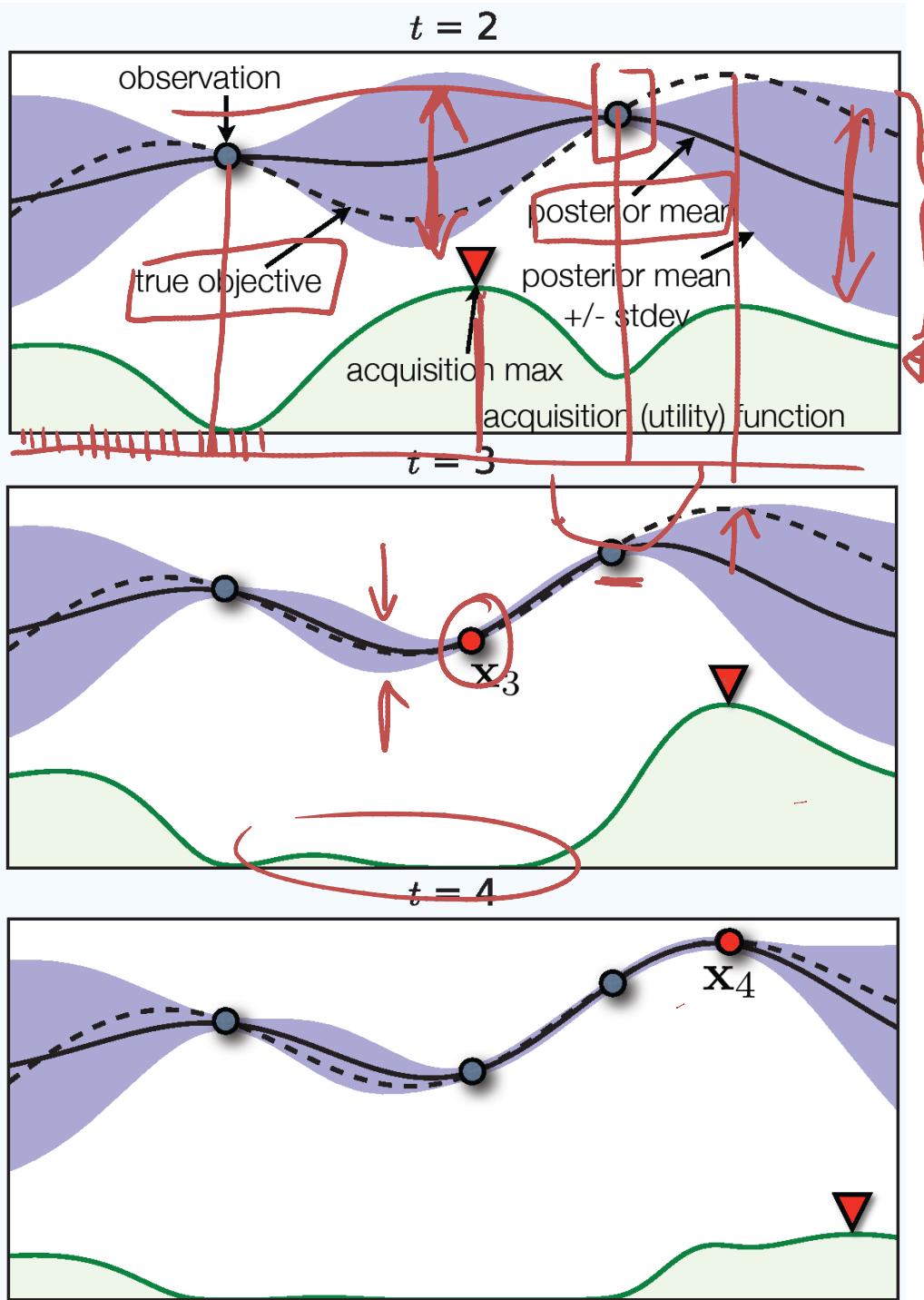
Multi-armed bandit problem



Multi-armed bandit problem



- Trade-off between Exploration and Exploitation
- Regret = Player reward – Reward of best action



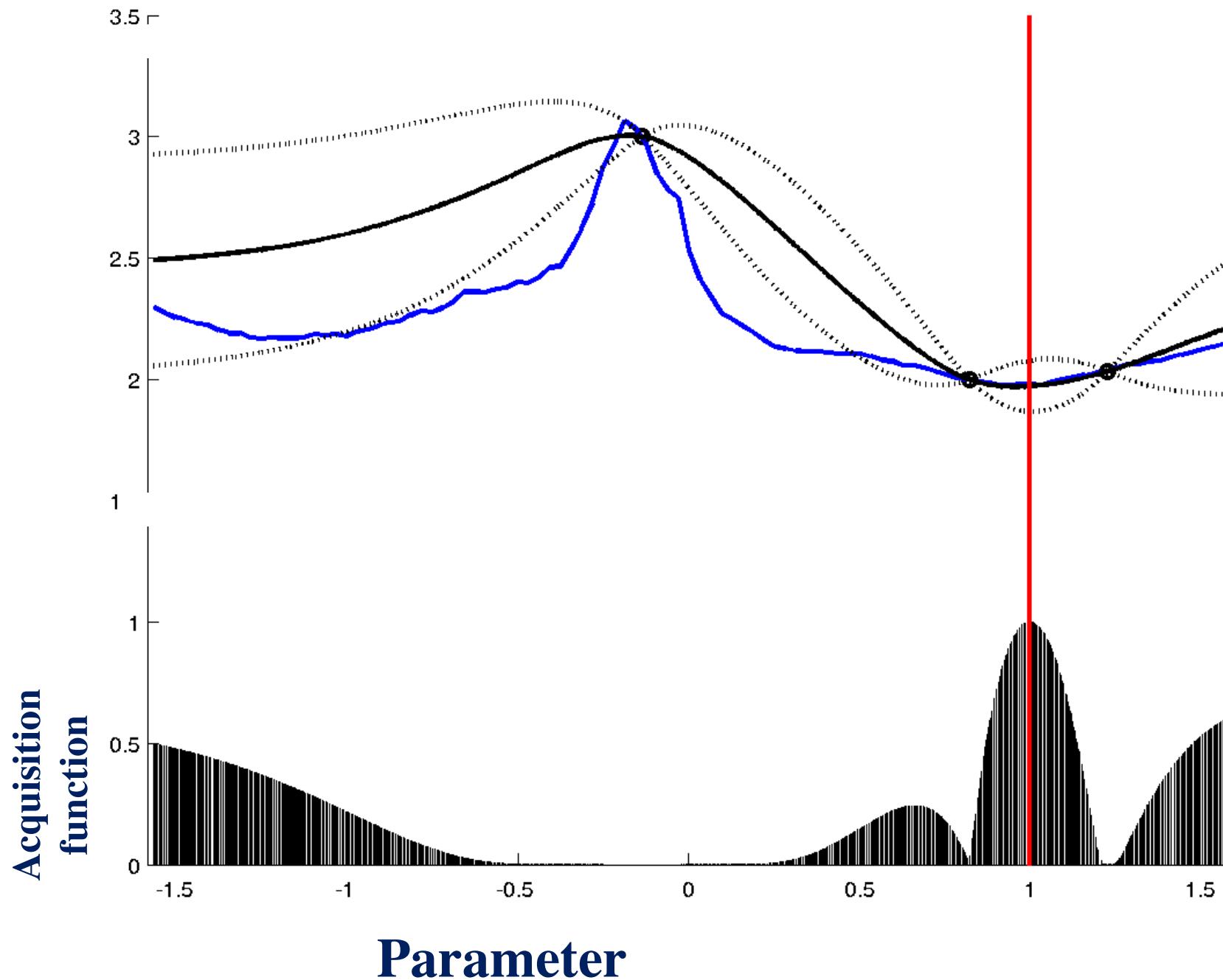
Bayesian optimization

- 1: **for** $t = 1, 2, \dots$ **do**
- 2: Find x_t by combining attributes of the posterior distribution in a utility function u and maximizing:

$$x_t = \operatorname{argmax}_x u(x | \mathcal{D}_{1:t-1}).$$
- 3: Sample the objective function:

$$y_t = f(x_t) + \varepsilon_t.$$
- 4: Augment the data $\mathcal{D}_{1:t} = \{\mathcal{D}_{1:t-1}, (x_t, y_t)\}$ and update the GP.
- 5: **end for**

$$\{(x_1, y_1), (x_2, y_2)\} \\ u(x | \text{data})$$



Exploration-exploitation tradeoff

Recall the expressions for GP prediction:

$$\begin{aligned} P(y_{t+1} | \mathcal{D}_{1:t}, \underline{\mathbf{x}}_{t+1}) &= \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{\text{noise}}^2) \\ \text{Prediction } \mu_t(\mathbf{x}_{t+1}) &= \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{y}_{1:t} \\ \sigma_t^2(\mathbf{x}_{t+1}) &= k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{k} \end{aligned}$$

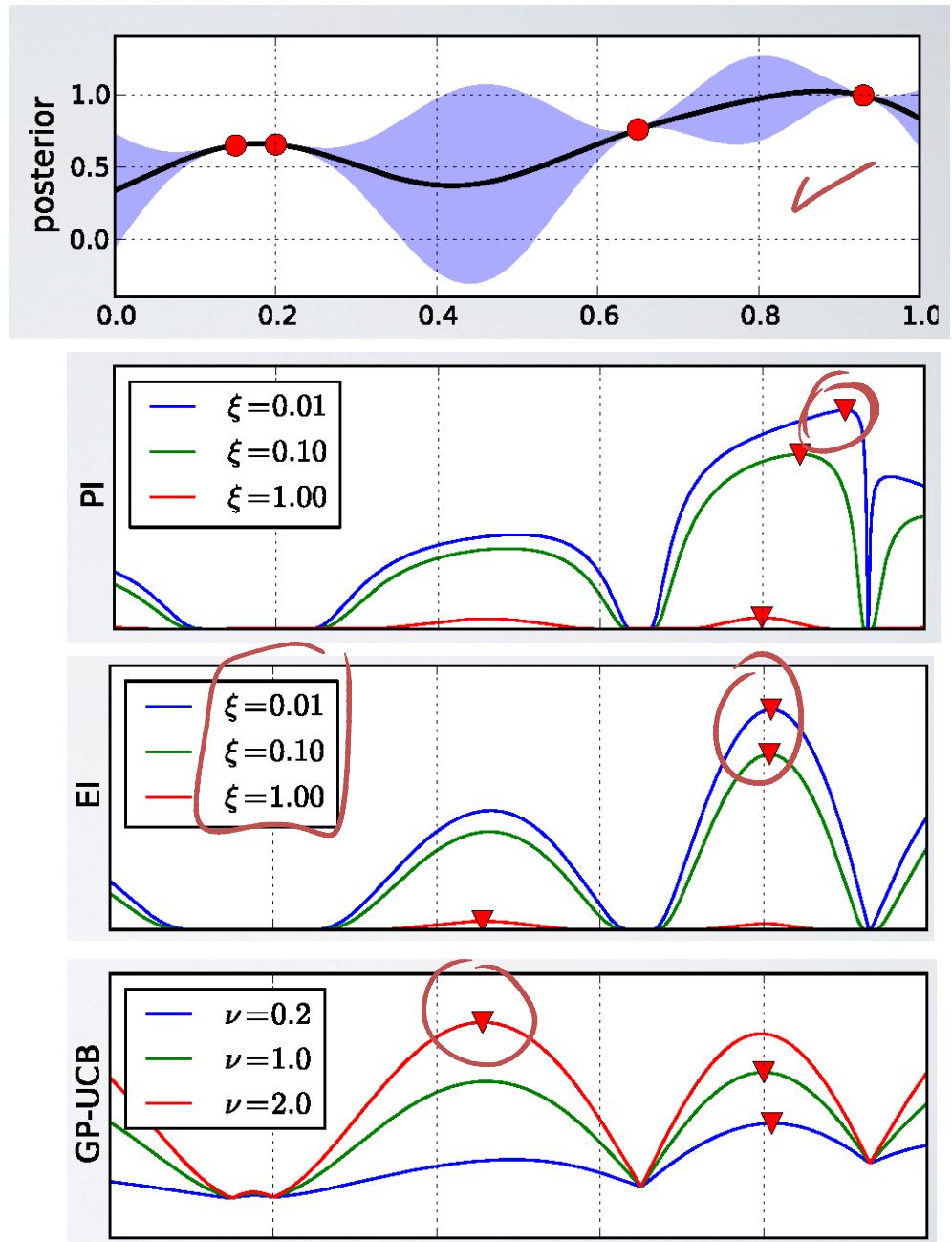
We should choose the next point \mathbf{x} where the mean is high (**exploitation**) and the variance is high (**exploration**).

We could balance this tradeoff with an **acquisition function** as follows:

$$\mu(\mathbf{x}) + \kappa \sigma(\mathbf{x})$$

Acquisition functions

- aka infill, figure of merit
- acquisition function guides the optimization by determining which x_{t+1} to observe next
- uses predictive posterior to combine exploration (high-variance regions) and exploitation (high-mean regions)
- optimize to find sample point (can be done cheaply/approximately)



An acquisition function: Probability of Improvement

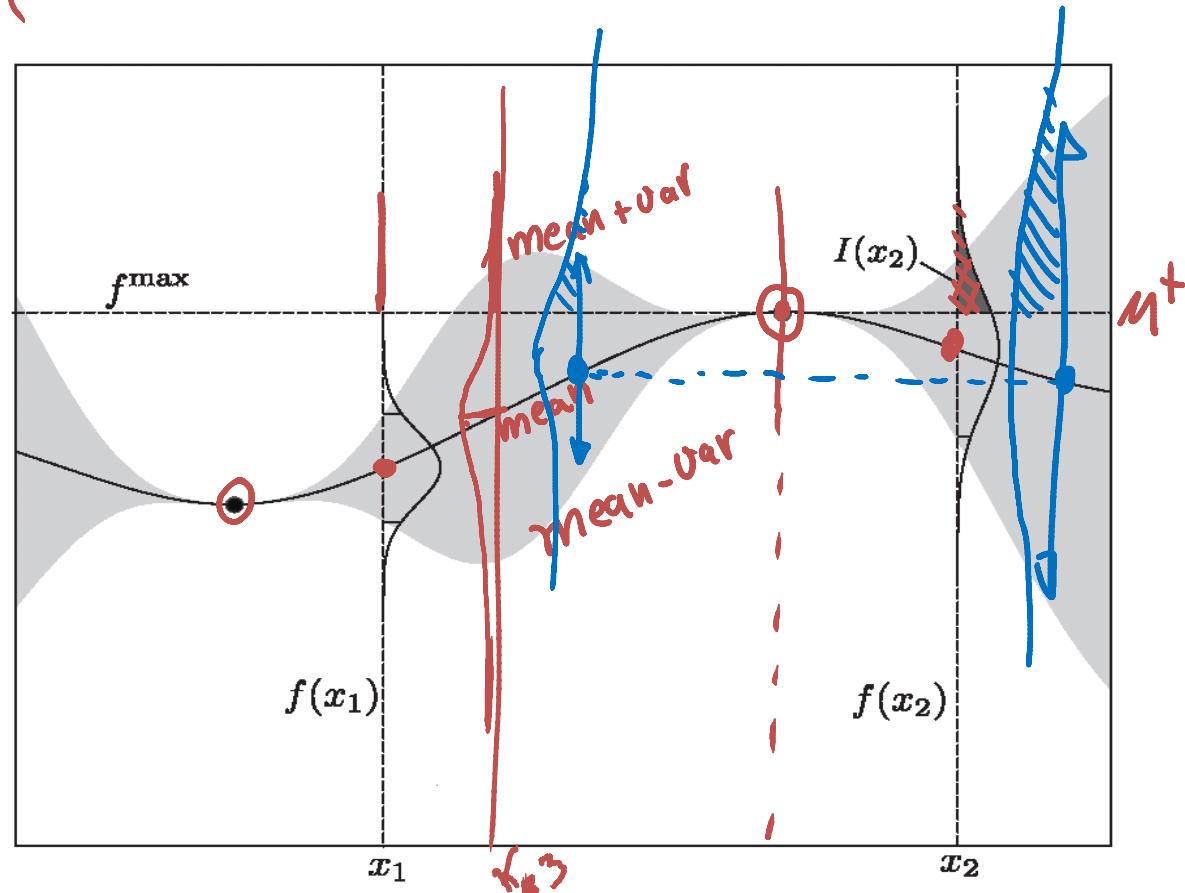
$$\text{PI}(\mathbf{x}) = P(f(\mathbf{x}) \geq \mu^+ + \xi)$$

$$= \Phi\left(\frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}\right)$$

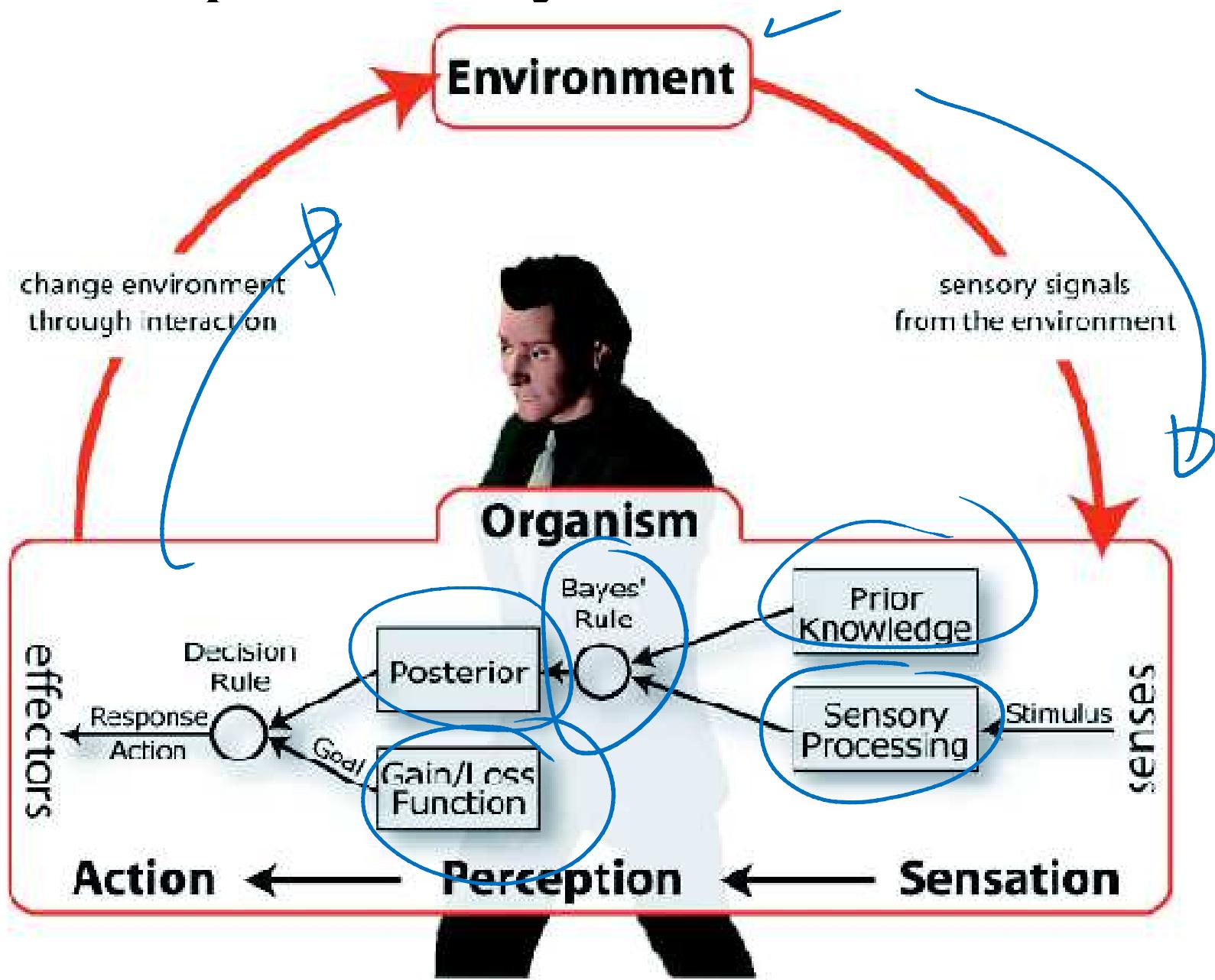
cumulative

μ^+ best observed value

$$f^{\max} = \mu^+ + \xi^{\text{small}}$$



People as Bayesian reasoners



Bayes and decision theory

Utilitarian view: We need models to make the right decisions under uncertainty. Inference and decision making are intertwined.

Learned posterior

$$\left\{ \begin{array}{l} P(x=\text{healthy}|data) = 0.9 \\ P(x=\text{cancer}|data) = 0.1 \end{array} \right. \quad \checkmark$$

Cost/Reward model $u(x,a)$

	$a = \text{no treatment}$	$a = \text{treatment}$
$x = \text{healthy}$	0	-30
$x = \text{cancer}$	-100	-20

We choose the action that maximizes the **expected utility**, or equivalently, which minimizes the **expected cost**.

$$EU(a) = \sum_x u(x,a) P(x|data)$$

$$EU(a=\text{treatment}) = u(\text{healthy, treatment}) P(x=\text{healthy}|d) + u(\text{cancer, treatment}) P(x=\text{cancer}|d) \\ (-30)(0.9) + (-20)(0.1) =$$

$$EU(a=\text{no treatment}) =$$

An expected utility criterion

At iteration $n+1$, choose the point that minimizes the distance to the objective evaluated at the maximum \underline{x}^* :

$$\begin{aligned}\underline{\mathbf{x}}_{n+1} &= \arg \min_{\underline{\mathbf{x}}} \mathbb{E}(\|f_{n+1}(\underline{\mathbf{x}}) - f(\underline{\mathbf{x}}^*)\| \mid \mathcal{D}_n) \\ &= \arg \min_{\underline{\mathbf{x}}} \int \|f_{n+1}(\underline{\mathbf{x}}) - f(\underline{\mathbf{x}}^*)\| p(f_{n+1} \mid \mathcal{D}_n) df_{n+1}\end{aligned}$$


We don't know the true objective at the maximum. To overcome this, Mockus proposed the following acquisition function:

$$\underline{\mathbf{x}} = \arg \max_{\underline{\mathbf{x}}} \mathbb{E}(\max\{0, f_{n+1}(\underline{\mathbf{x}}) - f^{\max}\} \mid \mathcal{D}_n)$$

Expected improvement

$$\mathbf{x} = \arg \max_{\mathbf{x}} \mathbb{E} \left(\max \left\{ 0, f_{n+1}(\mathbf{x}) - \frac{f^{\max}}{\mu^+ + \xi} \right\} \mid \mathcal{D}_n \right)$$

For this acquisition, we can obtain an analytical expression:

$$\text{EI}(\mathbf{x}) = \begin{cases} (\mu(\mathbf{x}) - \mu^+ - \xi)\Phi(Z) + \sigma(\mathbf{x})\phi(Z) & \text{if } \sigma(\mathbf{x}) > 0 \\ 0 & \text{if } \sigma(\mathbf{x}) = 0 \end{cases}$$

$Z = \frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}$

where $\phi(\cdot)$ and $\Phi(\cdot)$ denote the PDF and CDF of the standard Normal

A third criterion: GP-UCB

Define the *regret* and cumulative regret as follows:

$$\begin{aligned} r(\mathbf{x}) &= \underbrace{f(\mathbf{x}^*)}_{\text{optimal}} - \underbrace{f(\mathbf{x})}_{\text{observed}} \\ R_T &= \underbrace{r(\mathbf{x}_1) + \cdots + r(\mathbf{x}_T)}_{\text{cumulative}} \end{aligned}$$

The GP-UCB criterion is as follows:

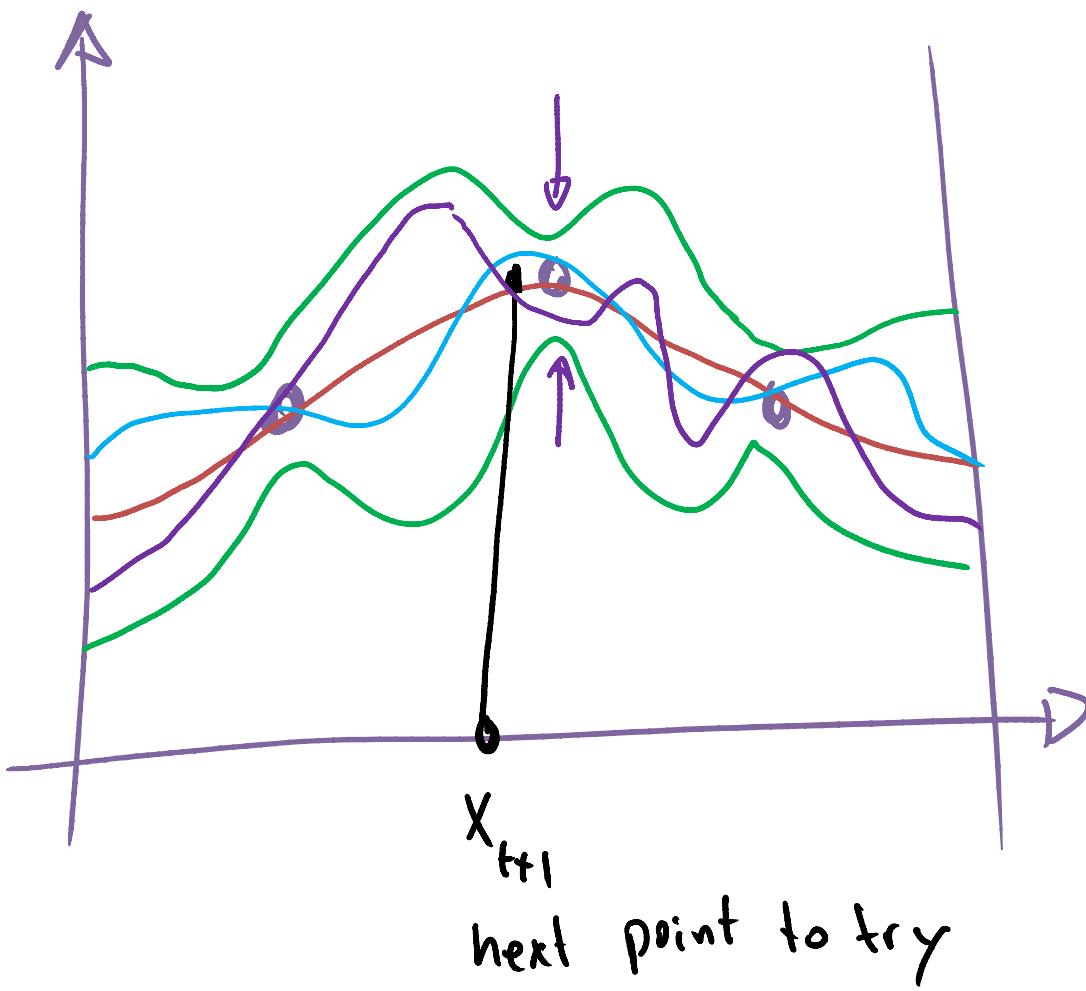
$$\text{GP-UCB}(\mathbf{x}) = \underbrace{\mu(\mathbf{x})}_{\text{mean}} + \sqrt{\underbrace{\nu \beta_t}_{\text{variance}} \sigma(\mathbf{x})}$$

Beta is set using a simple concentration bound:

With $\nu = 1$ and $\beta_t = 2 \log(t^{d/2+2} \pi^2 / 3\delta)$, it can be shown² with high probability that this method is *no regret*, i.e. $\lim_{T \rightarrow \infty} R_T/T = 0$. This in turn implies a lower-bound on the convergence rate for the optimization problem.

[Srinivas et al, 2010]

A fourth criterion: Thompson sampling



$$\mu^+ = \operatorname{argmax}_{\mathbf{x}_i \in \mathbf{x}_{1:t}} \mu(\mathbf{x}_i)$$

- Probability of Improvement

$$\text{PI}(\mathbf{x}) = \Phi \left(\frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})} \right)$$

Kushner 1964

- Expected Improvement

$$\text{EI}(\mathbf{x}) = (\mu(\mathbf{x}) - \mu^+ - \xi) \Phi(Z) + \sigma(\mathbf{x}) \phi(Z)$$

$$Z = \frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}$$

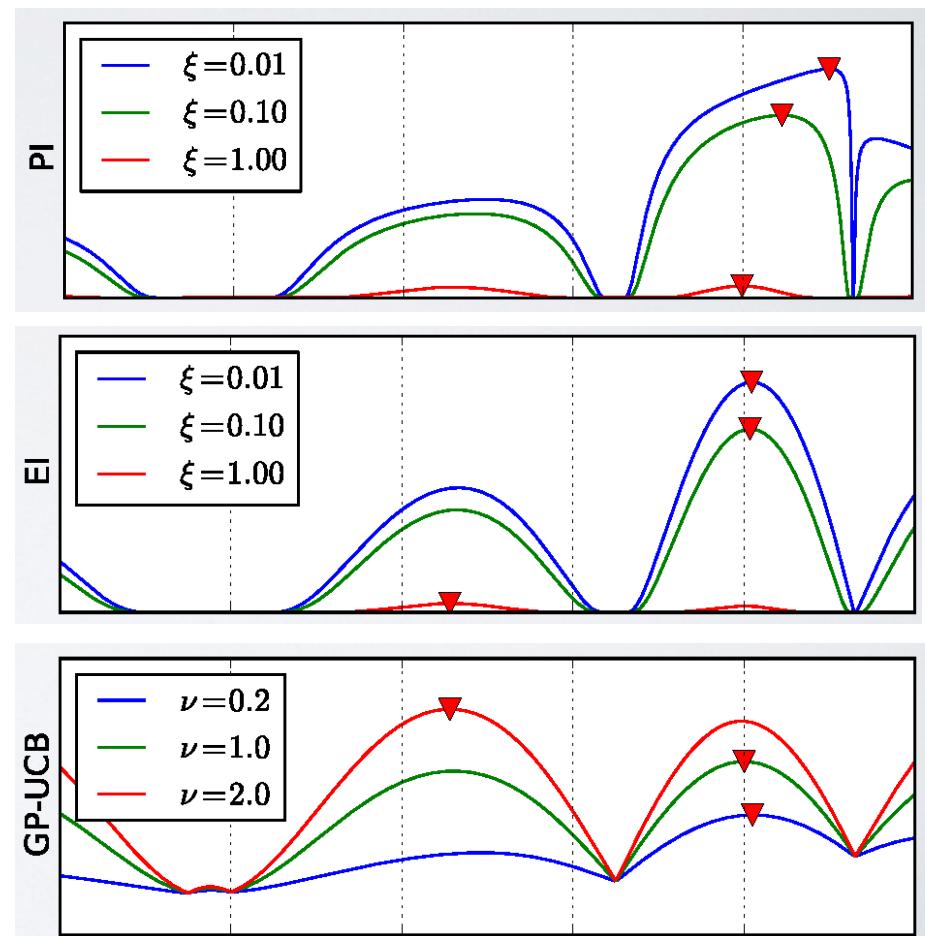
Mockus 1978

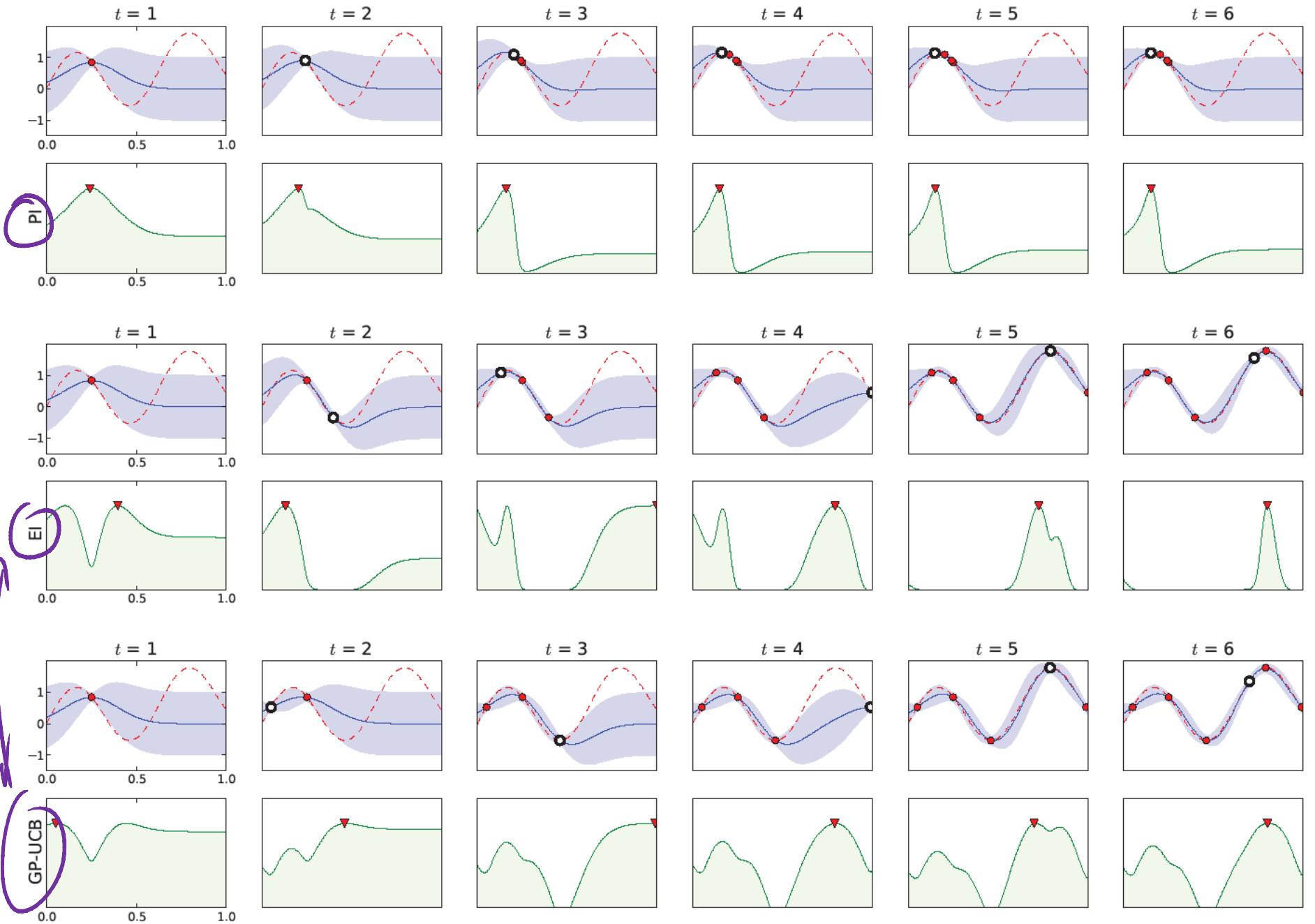
- Upper Confidence Bound

$$\text{GP-UCB}(\mathbf{x}) = \mu(\mathbf{x}) + \sqrt{\nu \tau_t} \sigma(\mathbf{x})$$

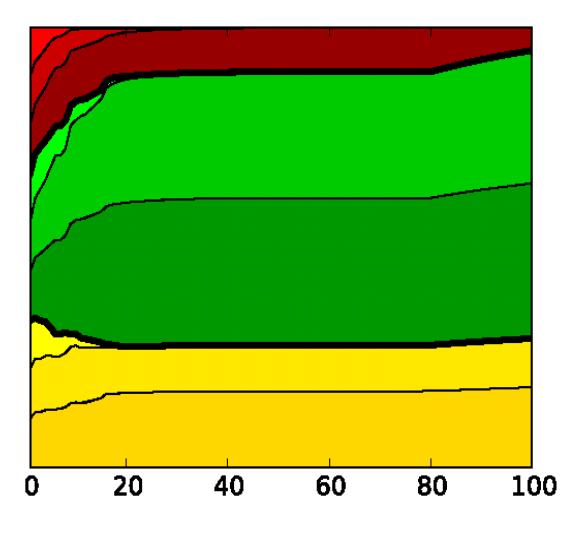
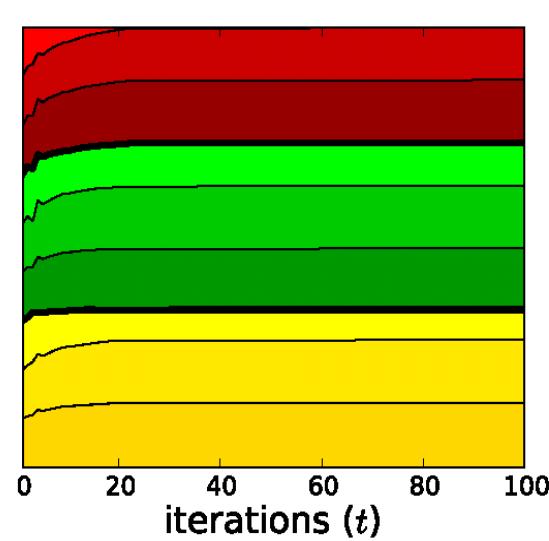
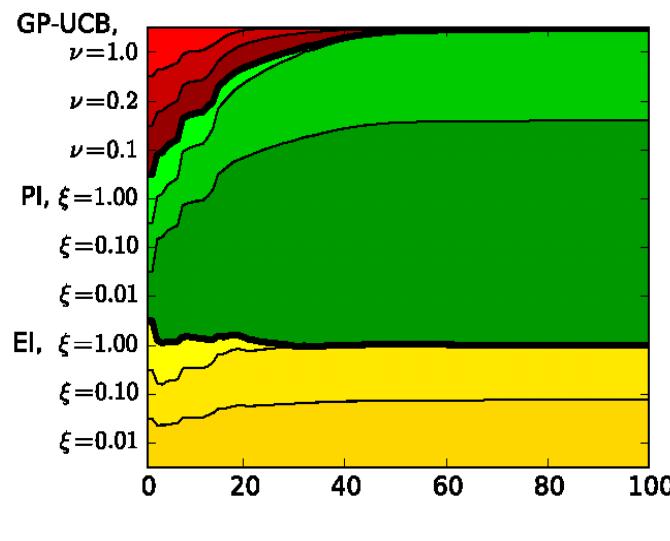
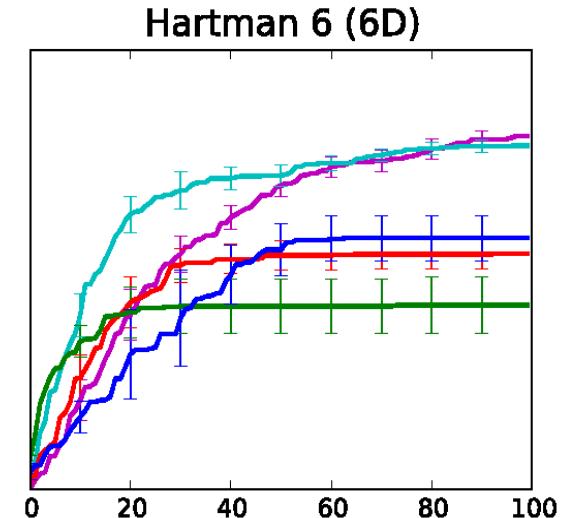
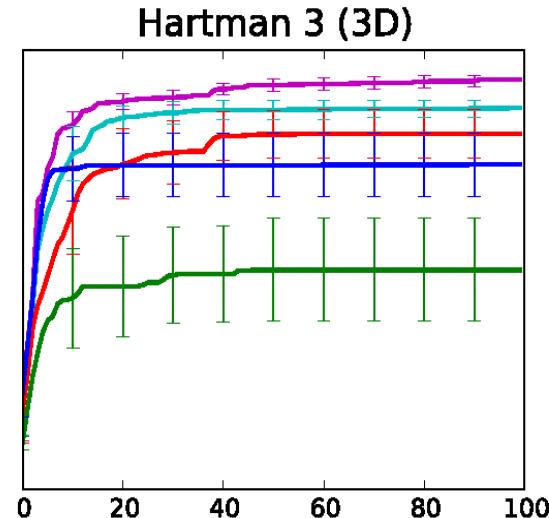
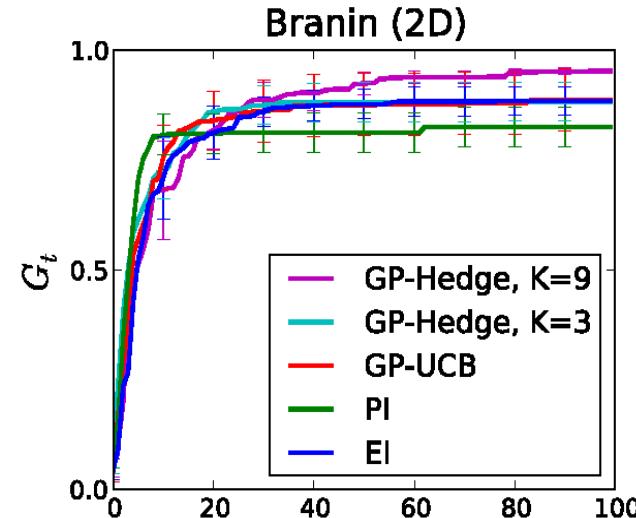
Srinivas *et al.* 2010

Acquisition functions

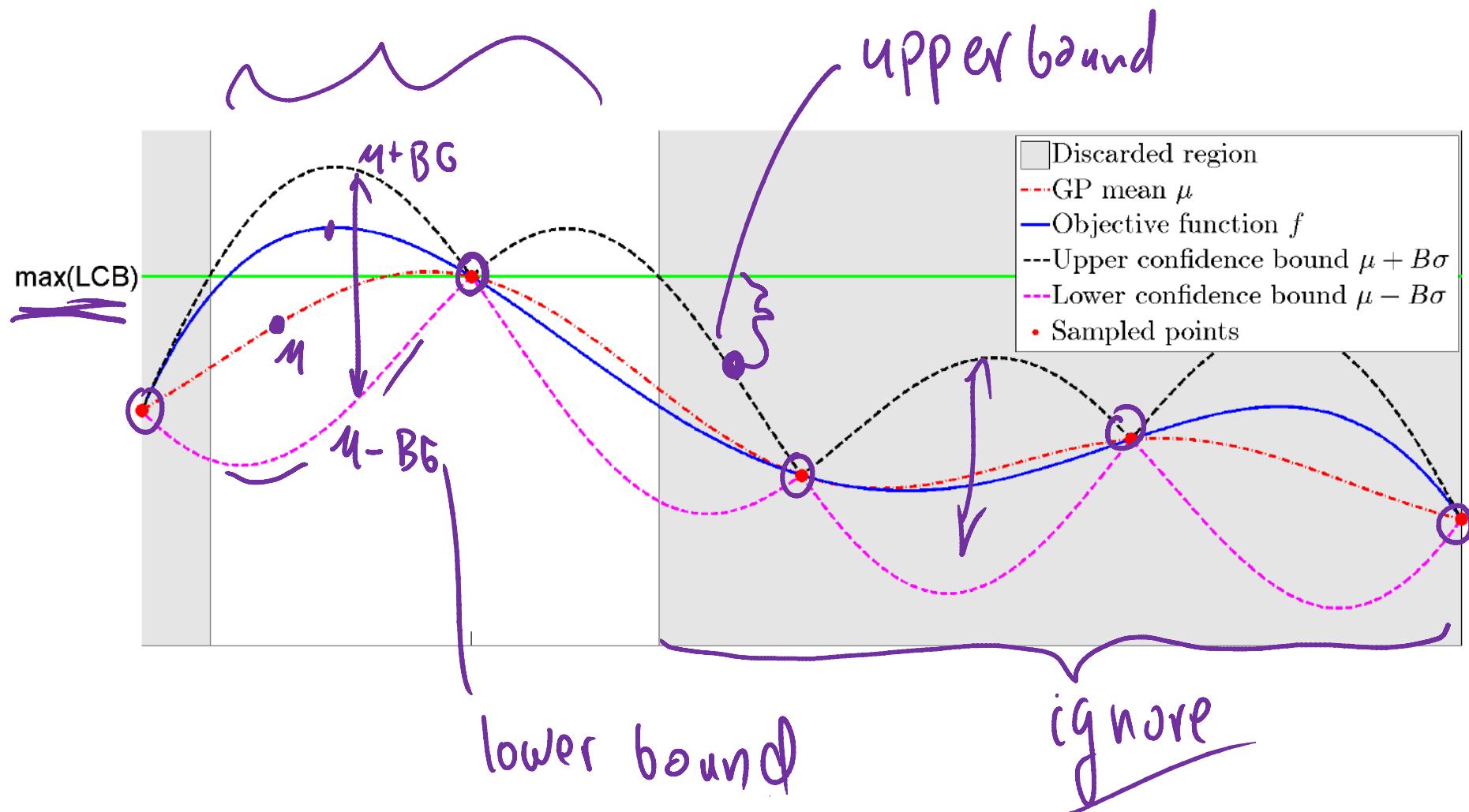




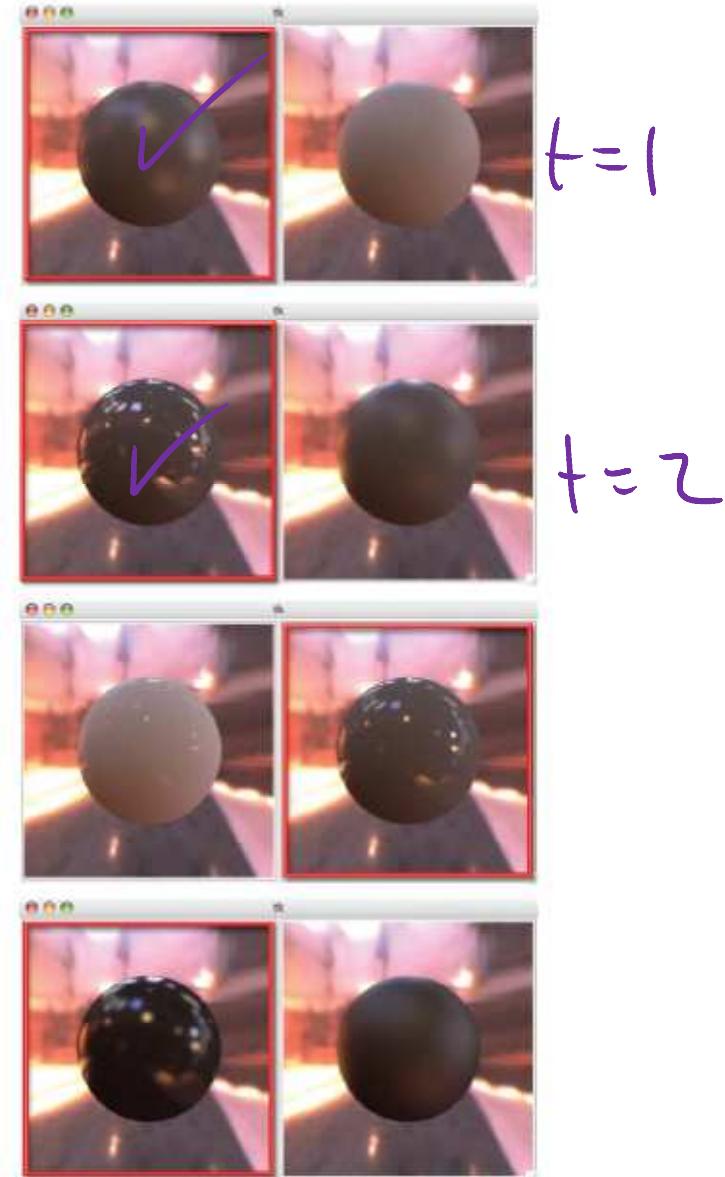
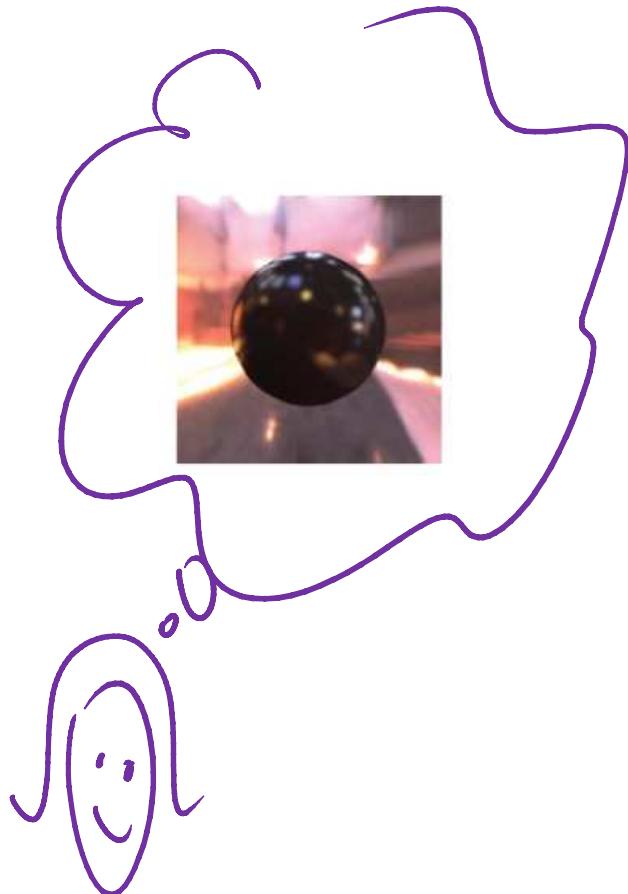
Portfolios of acquisition functions help



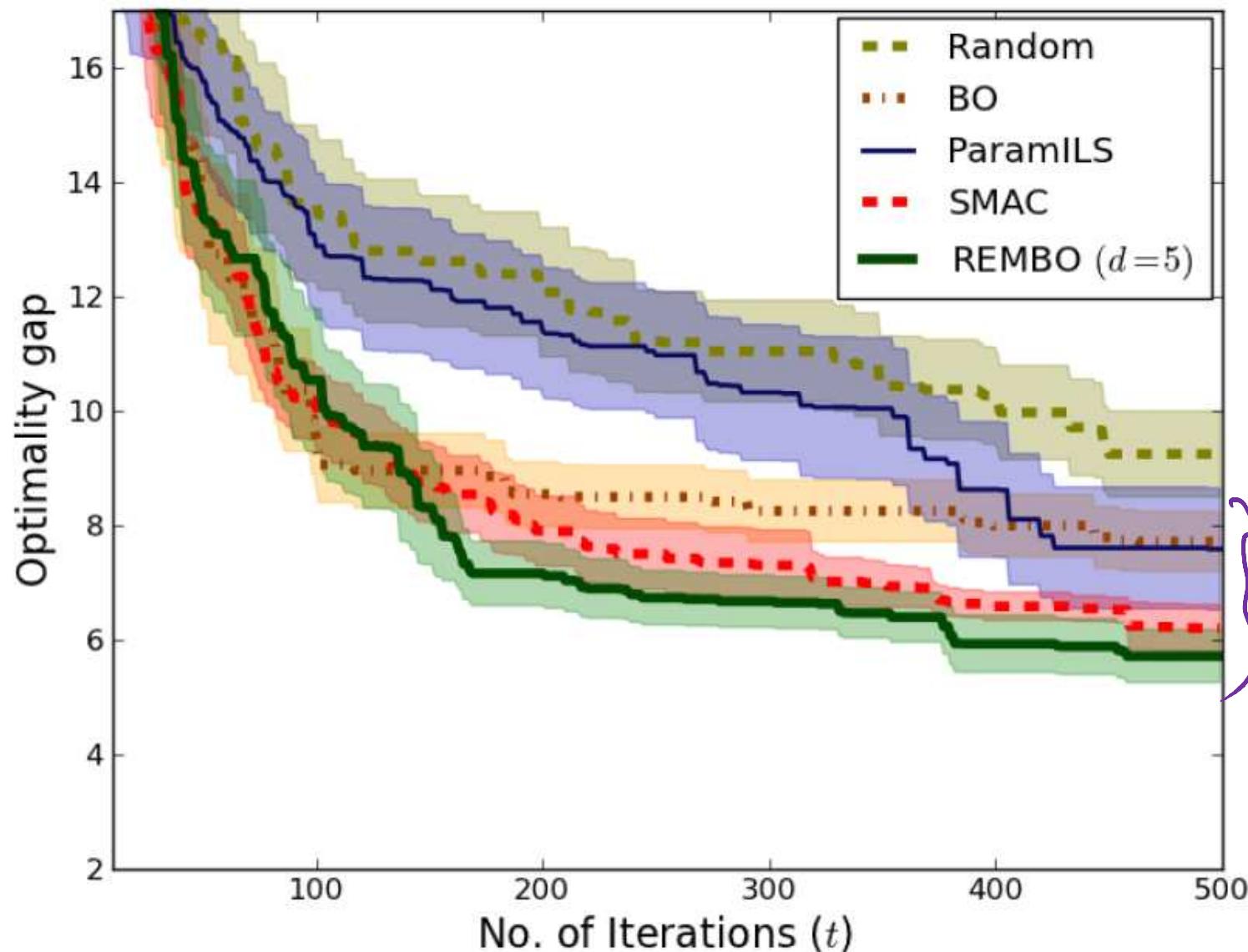
Why Bayesian Optimization works



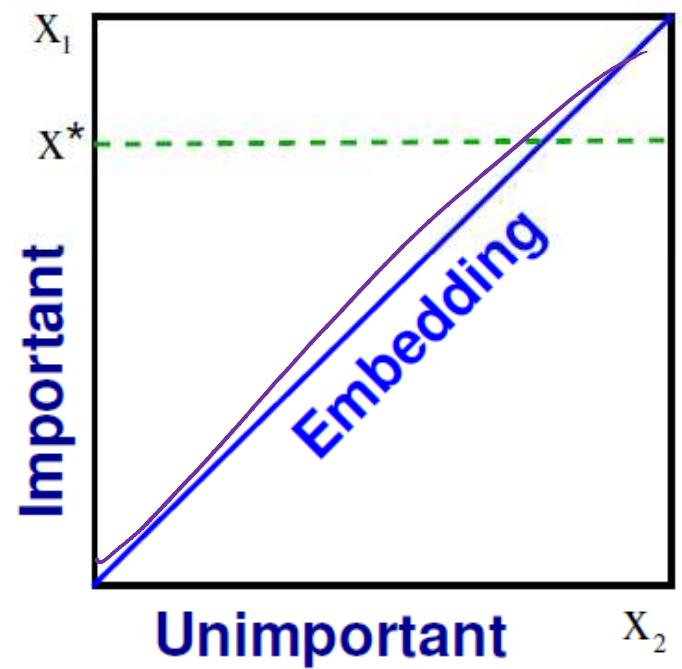
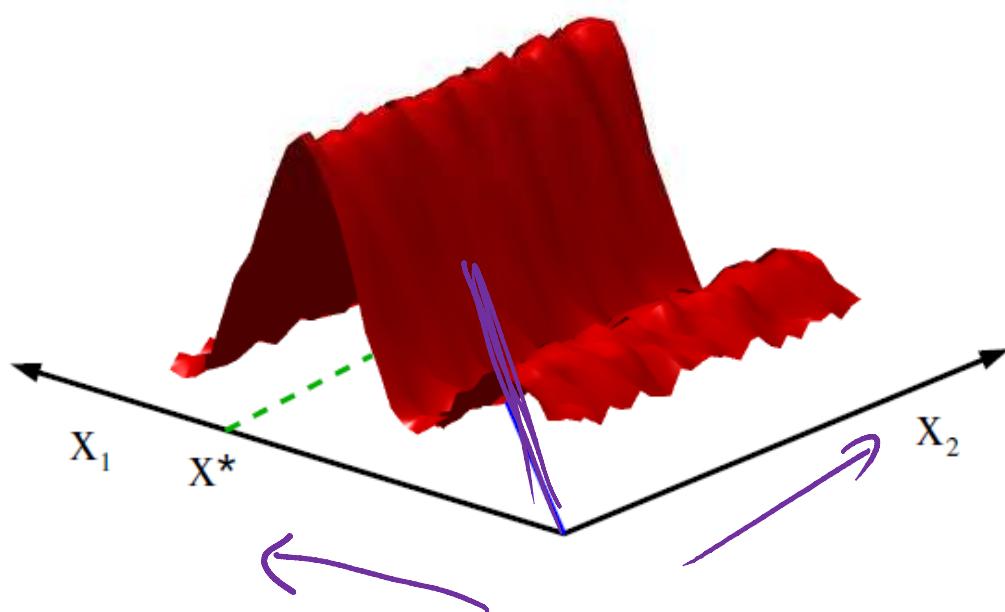
Intelligent user interfaces



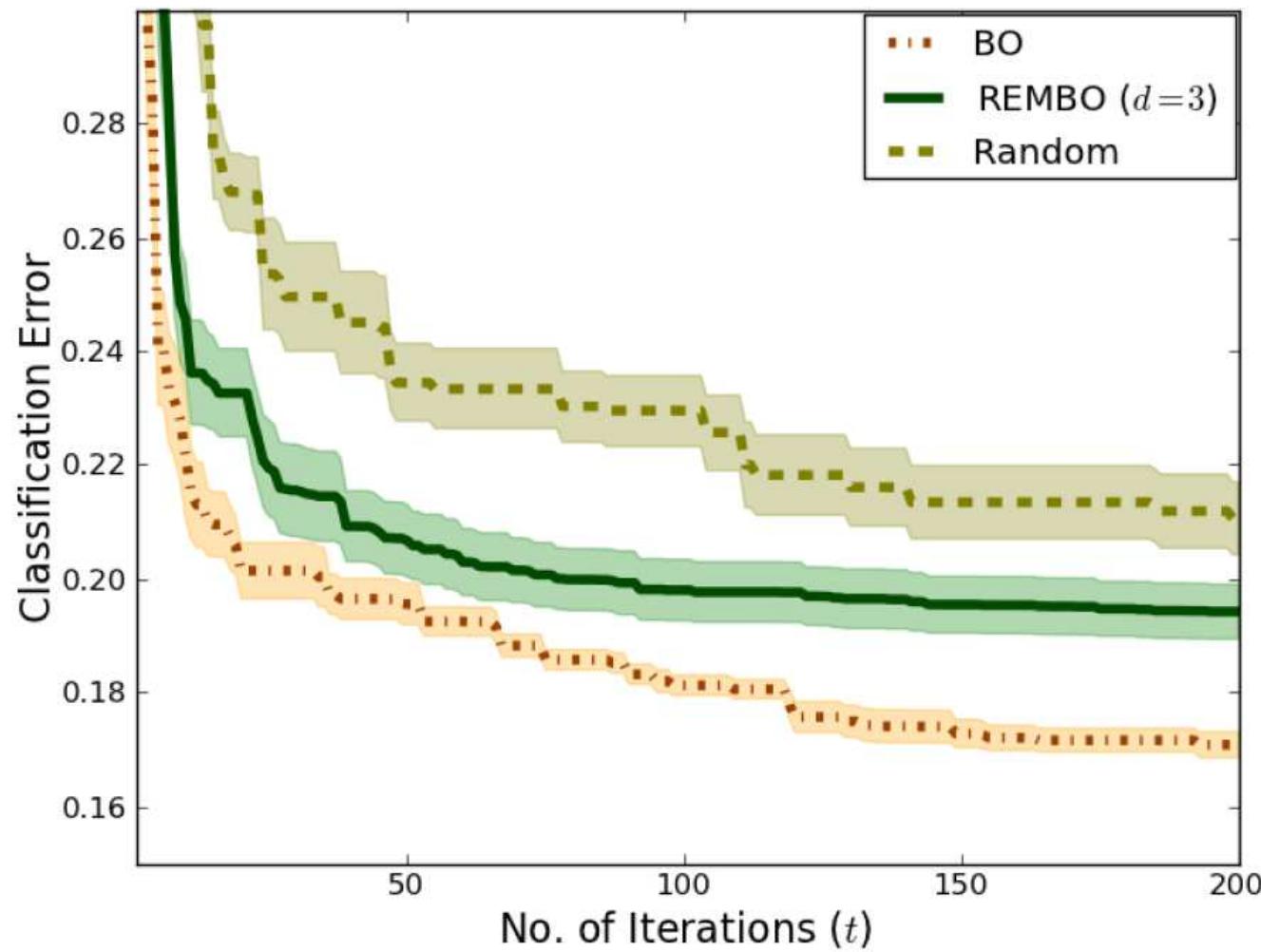
Example: Tuning NP hard problem solvers



Why random tuning works sometimes



Example: Tuning random forests



Example: Tuning hybrid Monte Carlo

Table 4.1: Mean squared test error for the robot arm data set.

Method	Mean Squared Error
Rios Insua and Muller's (1998) MLP with reversible-jump MCMC	0.00620
Mackay's (1992) Gaussian approximation with highest evidence	0.00573
Neal's (1996) HMC	0.00554
Neal's (1996) HMC with ARD	0.00549
Reversible-jump MCMC with Bayesian model by Andrieu et al.	0.00502
Adaptive HMC (Median Error)	0.00499
Adaptive HMC (Mean Error)	0.00498 ± 0.00012

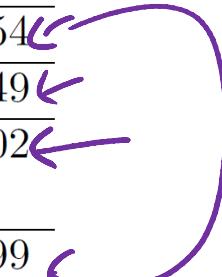
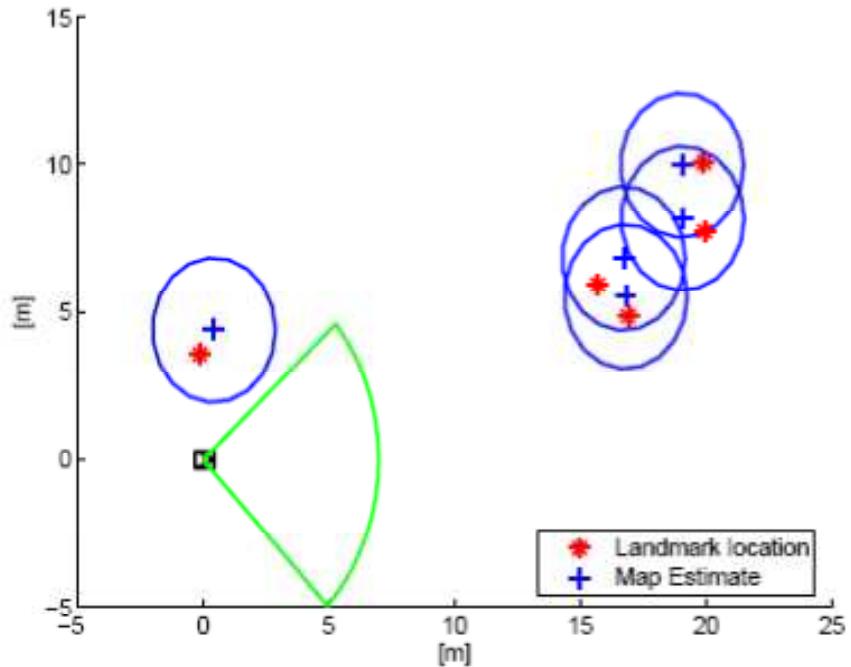


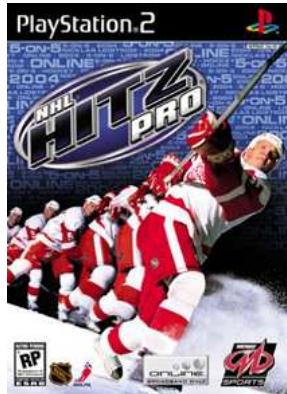
Table 4.2: Classification error on the validation set of the Dexter data set.

Method	Classification Error
New-Bayes-nn-sel	0.0800
Adaptive HMC (Mean error)	0.0730 ± 0.0096
Adaptive HMC (Median error)	0.0700
Adaptive HMC + Majority Voting	0.0667



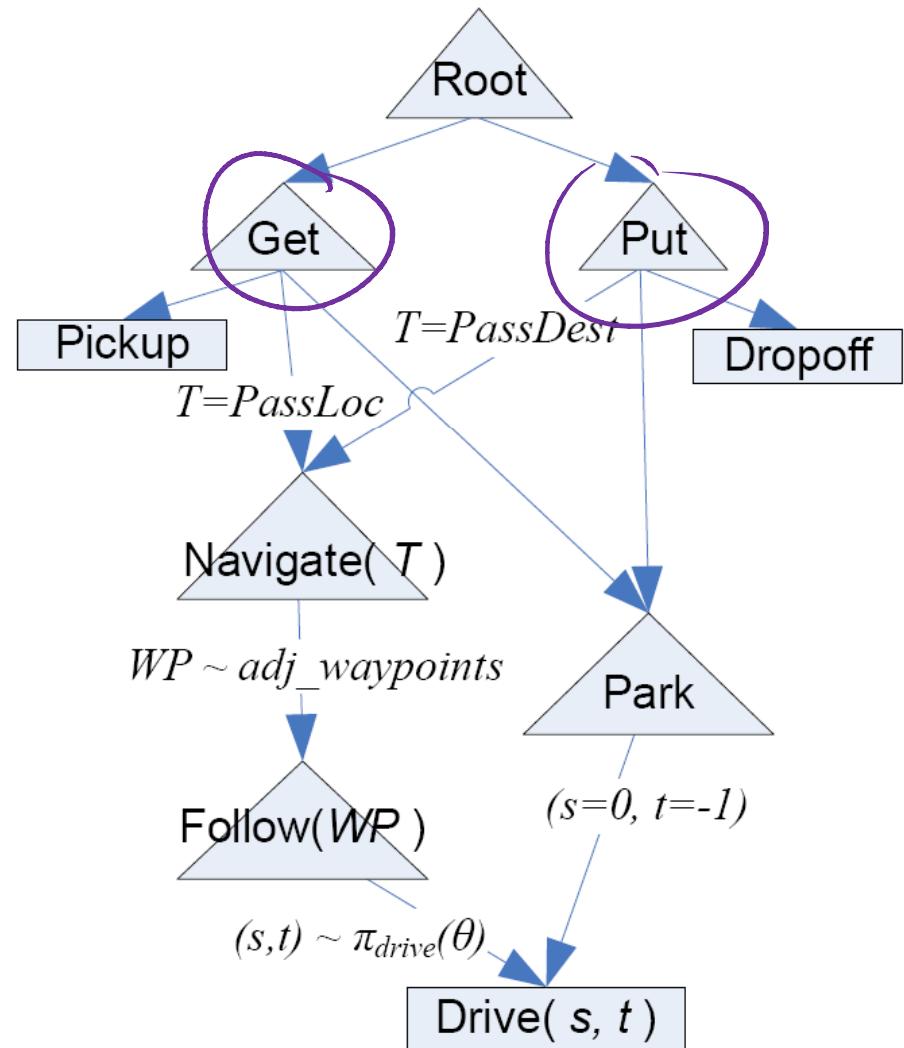


The games industry, rich in sophisticated large-scale simulators, is a great environment for the design and study of automatic decision making systems.



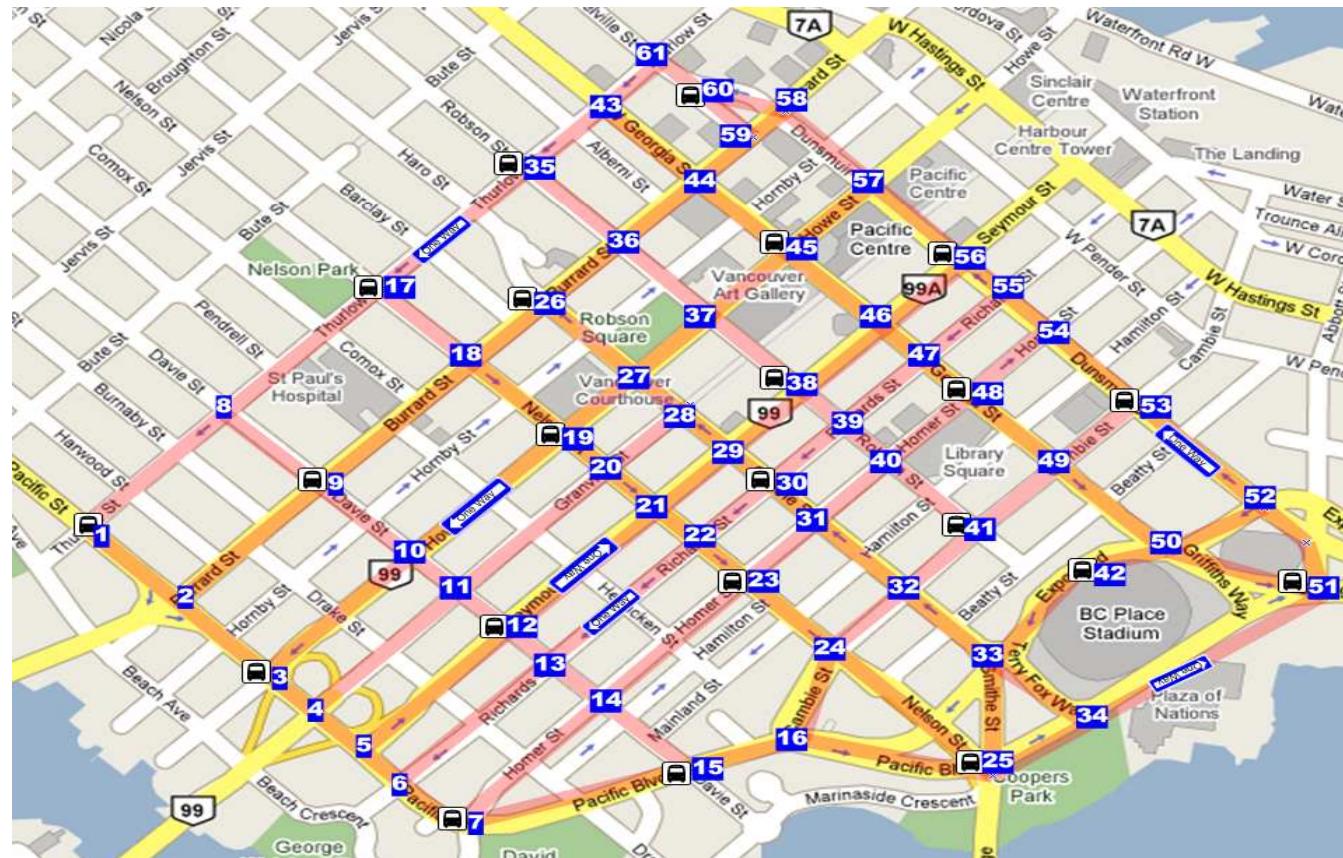
Hierarchical policy example

- **High-level** model-based learning for deciding when to navigate, park, pickup and dropoff passengers.
- **Mid-level** active path learning for navigating a topological map.
- **Low-level** active policy optimizer to learn control of continuous non-linear vehicle dynamics.



Active Path Finding in Middle Level

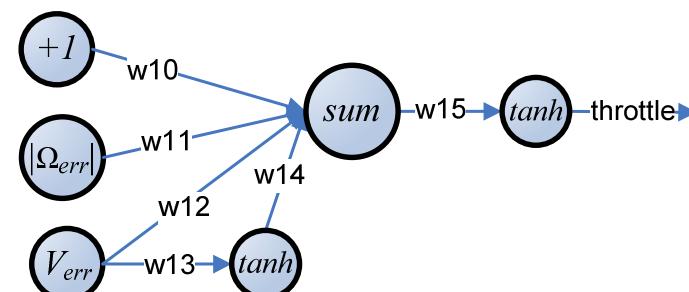
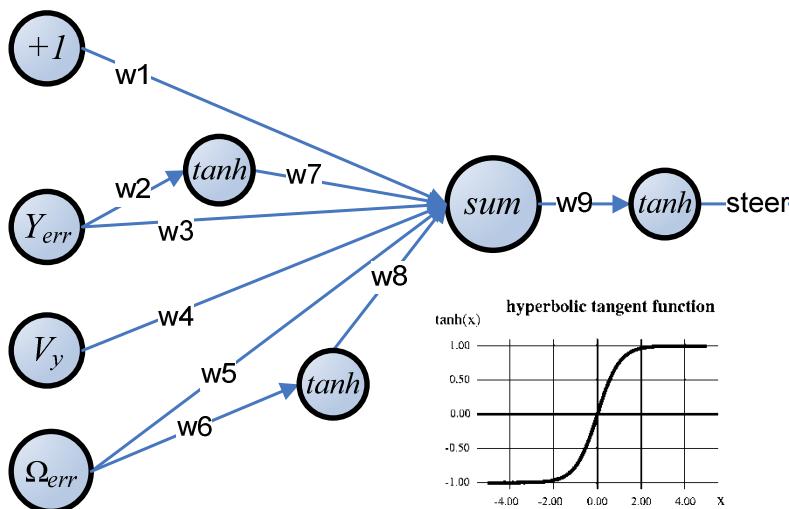
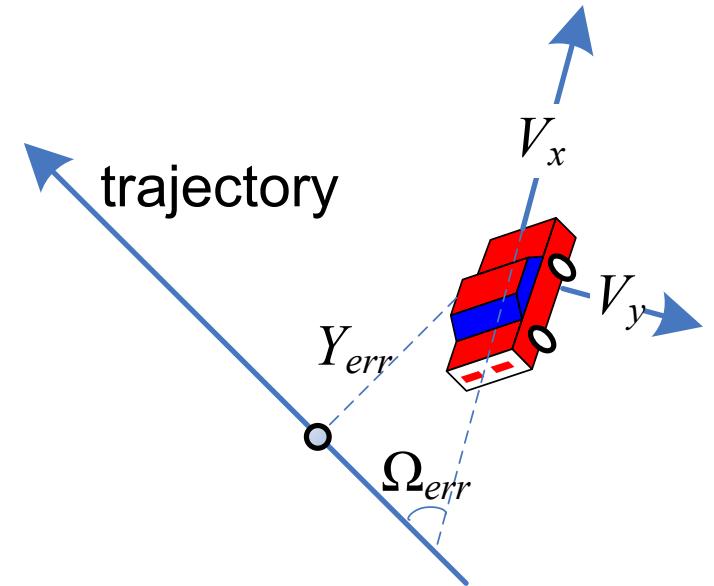
- Mid-level *Navigate* policy generates sequence of waypoints on a topological map to navigate from a location to a destination. $V(\theta)$ value function represents the path length from the current node, to the target.



Low-Level: Trajectory following



TORCS: 3D game engine that implements complex vehicle dynamics complete with manual and automatic transmission, engine, clutch, tire, and suspension models.



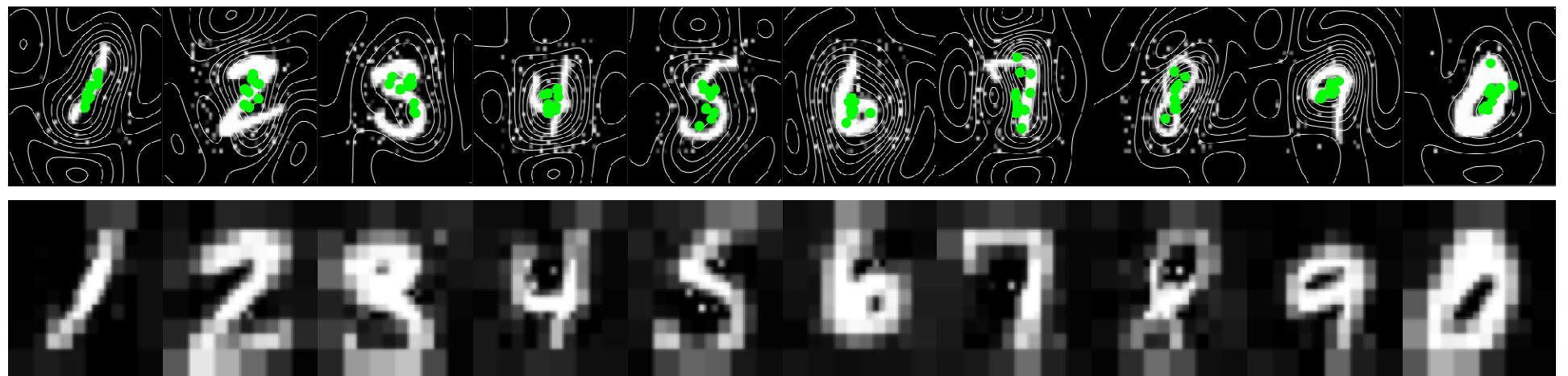
Hierarchical systems apply to many robot tasks – key to build large systems



We used TORCS: A 3D game engine that implements complex vehicle dynamics complete with manual and automatic transmission, engine, clutch, tire, and suspension models.

Gaze planning

Digits Experiment:



Face Experiment:

