



Aalto University
School of Electrical
Engineering

ELEC-E8125 Reinforcement Learning

Model-based RL

Joni Pajarinen

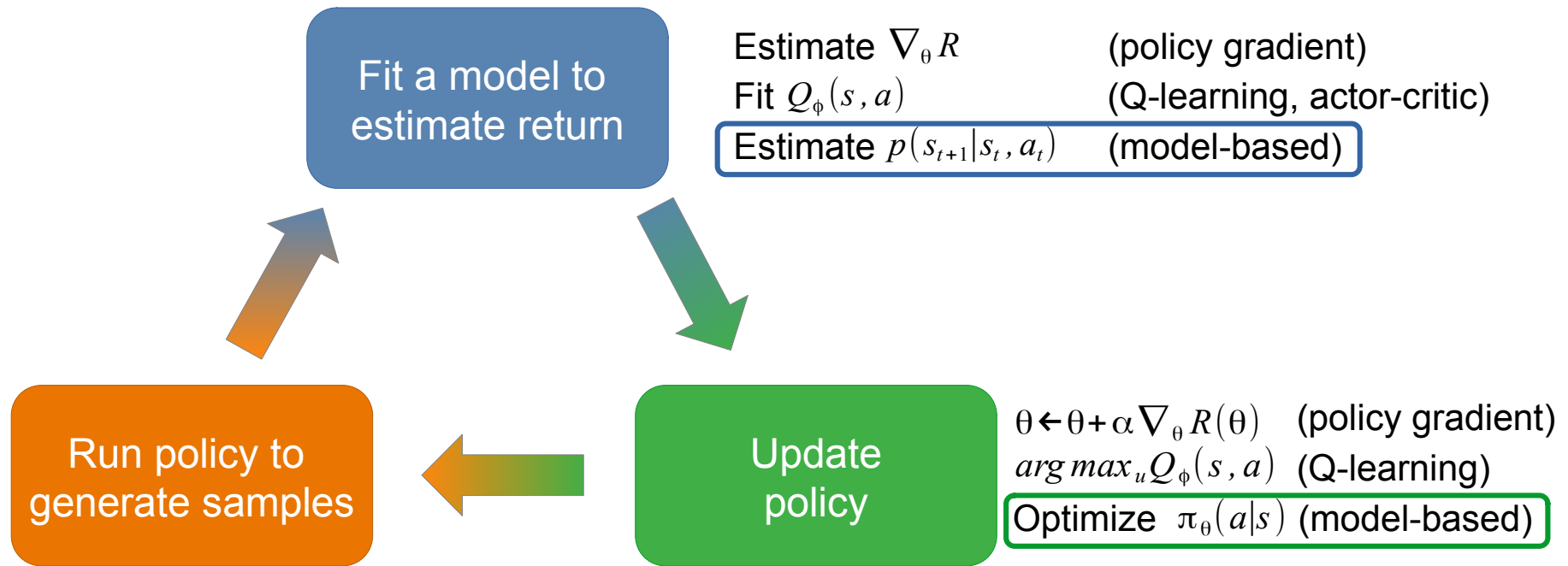
2.11.2021

Learning goals

- Understand basic approaches for model-based reinforcement learning.

Anatomy of reinforcement learning

Model-based



Motivation (partial recap)

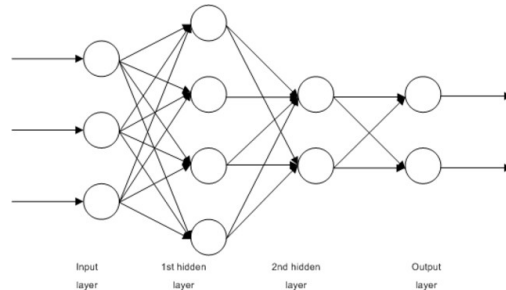
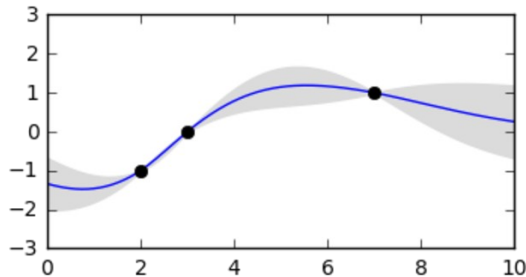
- Reinforcement learning has limited sample efficiency.
- Learned policies are task(reward-function)-specific, learned policies cannot be directly reused.
- Learned dynamics model is reusable and can be used to reason about potential futures.
- Sometimes we know the model, e.g. in games!



Model definition and types

- Dynamics model $s_{t+1} = f(s_t, a_t)$ or $f(s_{t+1} | s_t, a_t)$
- Reward model $r_{t+1} = r(s_t, a_t)$ or $r(r_{t+1} | s_t, a_t)$
- Models are usually learned.
 - Parametric regression (e.g. neural net) common.
- May be also known (e.g. games, simulators)
 - Even physics based models need to be often calibrated.
- Also other possibilities (active research area)
 - Latent variable models, graph neural networks, non-parametric regression models such as Gaussian processes, ...

Which model to use?



$$Y_i = \beta_0 + \beta_1 \phi_1(X_{i1}) + \dots + \beta_p \phi_p(X_{ip})$$

Gaussian process (GP)

- Data-efficient
- Slow with big datasets
- May be too smooth for non-smooth dynamics

Neural networks

- Expressive
- Unpredictable with sparse data (overfit)

Linear models

- May be used locally
- Do not overfit

Domain specific parametric models (e.g. physics parameters) can also be used.
→ Traditional control engineering approach of model identification + control.

Spectrum of model-based RL

how to act in
current situation
(choose action)

Time of planning

learn to act in
any situation
(learn policy)

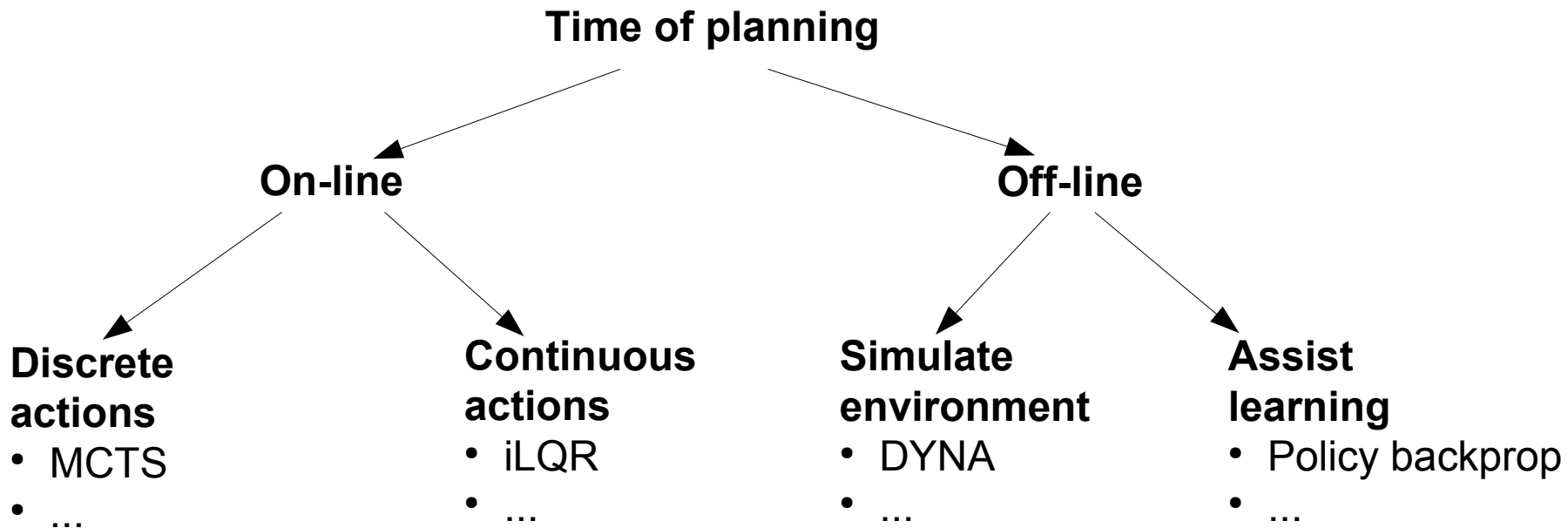
On-line

- Act on current state
- Act without learning
- Better in unfamiliar situations

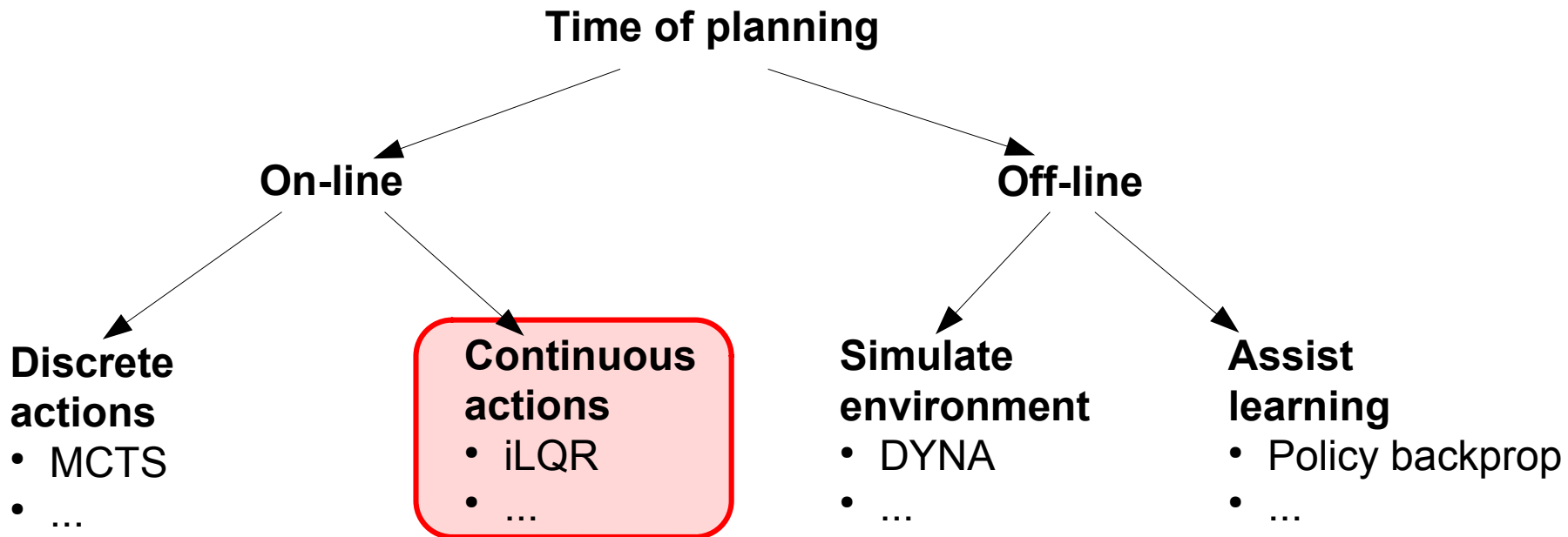
Off-line

- Fast online computation
- Predictable within familiar situations

Spectrum of model-based RL



Spectrum of model-based RL



Continuous on-line planning: iLQR + learned model

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(s, a, s')_i\}$

Repeat

Fit dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$

Use model to plan (e.g. iLQR) actions

Execute first planned action, observe resulting state s'

Update dataset $D \leftarrow D \cup \{(s, a, s')\}$

Continuous on-line planning: iLQR + learned model

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(s, a, s')_i\}$

Repeat

Fit dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$

Use model to plan (e.g. iLQR) actions

Execute first planned action, observe resulting state s'

Update dataset $D \leftarrow D \cup \{(s, a, s')\}$

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly \rightarrow model may be fitted only periodically (every n steps).
 - Planning costly for long horizons.
- Robust to moderate model errors.
- Choice of regression model is an important design parameter.

Continuous on-line planning: iLQR + learned model

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(s, a, s')_i\}$

Repeat

Fit dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$

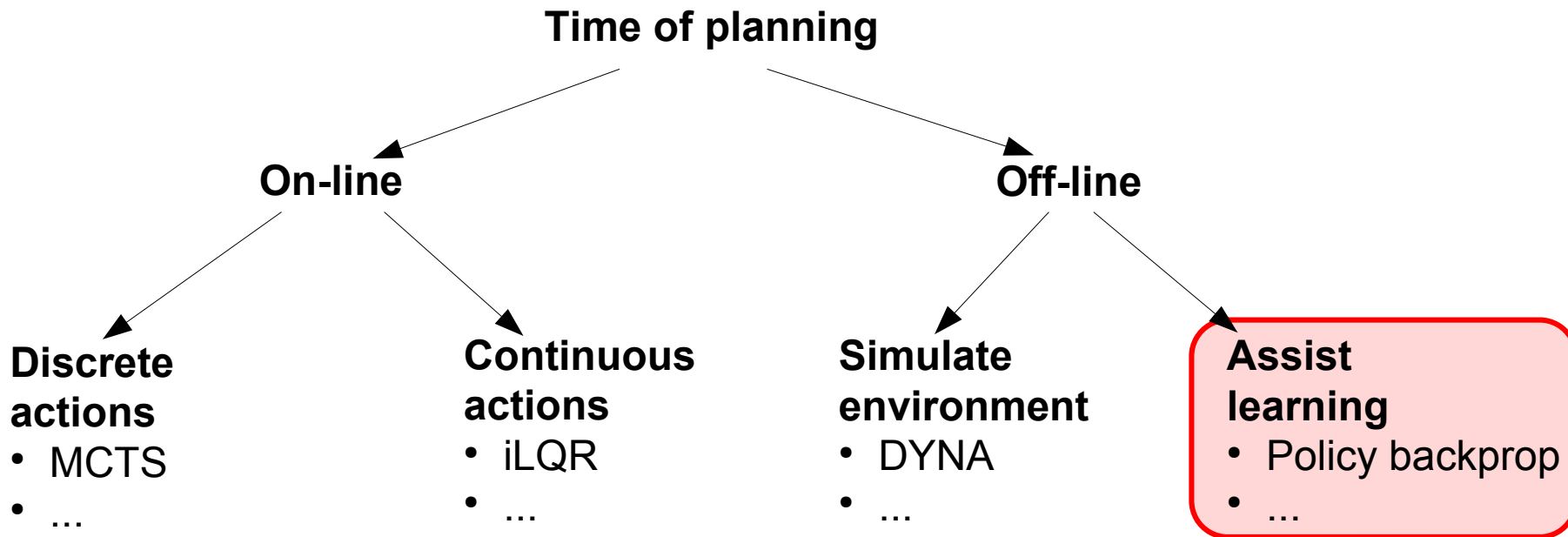
Use model to plan (e.g. iLQR) actions

Execute first planned action, observe resulting state s'

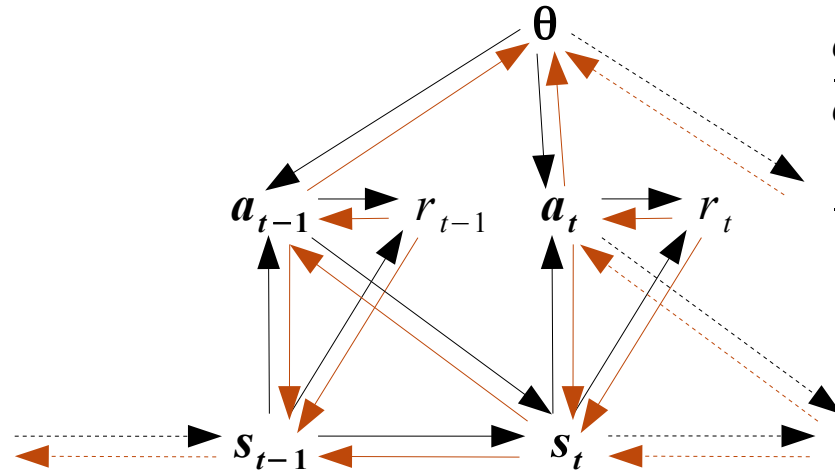
Update dataset $D \leftarrow D \cup \{(s, a, s')\}$

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly \rightarrow model may be fitted only periodically (every n steps).
 - **Planning costly for long horizons.**
- Robust to moderate model errors.
- Choice of regression model is an important design parameter.

Spectrum of model-based RL



Combining parametric policy with learned dynamics by backpropagation

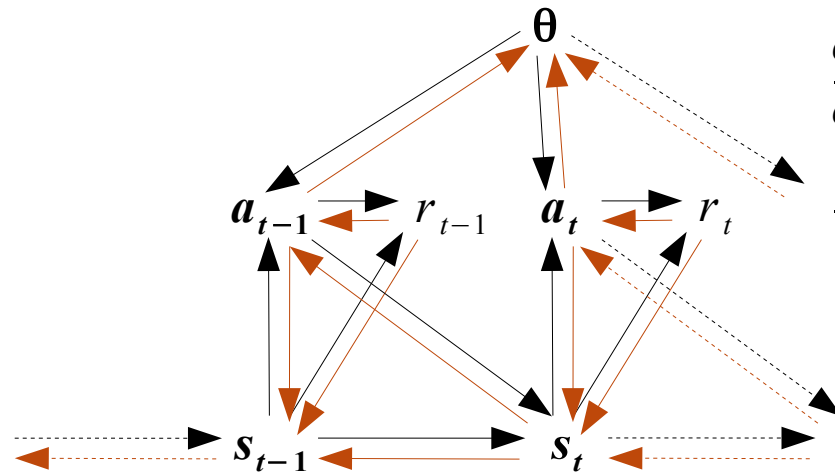


$$\frac{\partial r_t}{\partial \theta} = \frac{\partial r_t}{\partial a_t} \frac{\partial a_t}{\partial \theta} + \frac{\partial r_t}{\partial s_t} \frac{\partial s_t}{\partial \theta}$$

$$\frac{\partial s_t}{\partial \theta} = \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial \theta} + \frac{\partial s_t}{\partial a_{t-1}} \frac{\partial a_{t-1}}{\partial \theta}$$

policy	reward	dynamics
$\nabla_{\theta} \pi(s_t, a_t)$	$\nabla_a r(s_t, a_t)$	$\nabla_s f(s_{t-1}, a_{t-1})$
	$\nabla_s r(s_t, a_t)$	$\nabla_a f(s_{t-1}, a_{t-1})$

Combining parametric policy with learned dynamics by backpropagation



$$\frac{\partial r_t}{\partial \theta} = \frac{\partial r_t}{\partial \mathbf{a}_t} \frac{\partial \mathbf{a}_t}{\partial \theta} + \frac{\partial r_t}{\partial \mathbf{s}_t} \frac{\partial \mathbf{s}_t}{\partial \theta}$$

$$\frac{\partial \mathbf{s}_t}{\partial \theta} = \frac{\partial \mathbf{s}_t}{\partial \mathbf{s}_{t-1}} \frac{\partial \mathbf{s}_{t-1}}{\partial \theta} + \frac{\partial \mathbf{s}_t}{\partial \mathbf{a}_{t-1}} \frac{\partial \mathbf{a}_{t-1}}{\partial \theta}$$

Run base policy to collect data $D \leftarrow \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')_i\}$

Repeat

Fit dynamics model $f_\phi(\mathbf{s}, \mathbf{a})$ to minimize $\sum_i \|f_\phi(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{s}'_i\|^2$

Calculate policy gradient update by backpropagating through dynamics

Execute updated policy (1 or more steps), collect data

Update dataset $D \leftarrow D \cup \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')\}$

Continuous on-line planning: iLQR + learned model

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(s, a, s')_i\}$

Repeat

Fit dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$

Use model to plan (e.g. iLQR) actions

Execute first planned action, observe resulting state s'

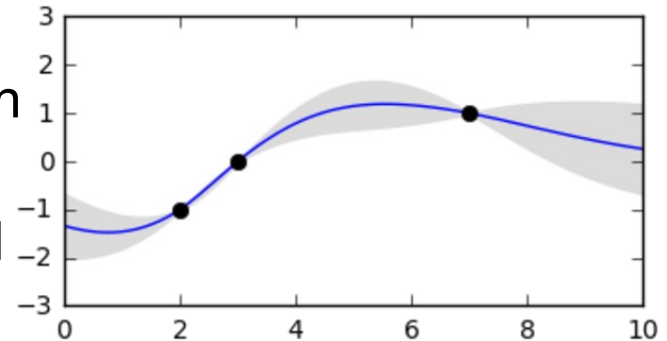
Update dataset $D \leftarrow D \cup \{(s, a, s')\}$

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly \rightarrow model may be fitted only periodically (every n steps).
 - Planning costly for long horizons.
- Robust to moderate model errors.
- **Choice of regression model is an important design parameter.**

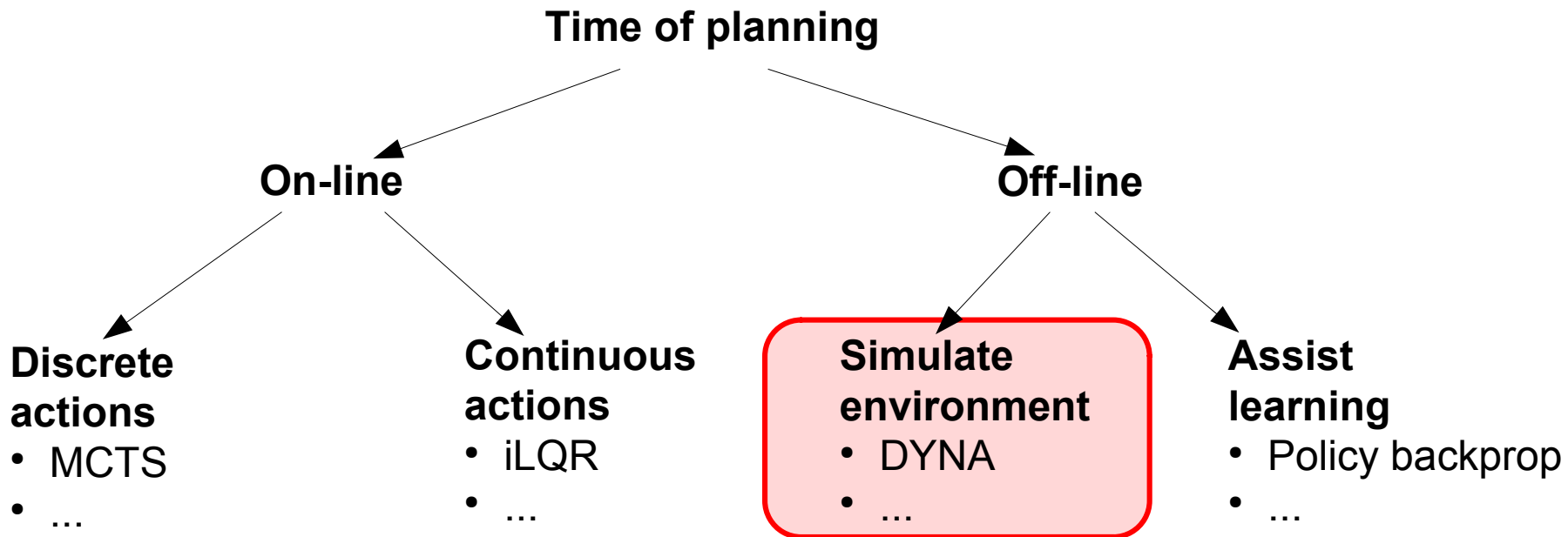
Example

PILCO (Deisenroth&Rasmussen, 2011)

- Dynamics learning: Use Gaussian process models to include model uncertainty. Known quadratic reward.
- Simulation: Simulate trajectory with learned model, including uncertainty.
- Policy: Radial basis function.
- Policy update: Calculate analytically policy gradient using learned dynamics and optimize with quasi-Newton optimizer (BFGS).
- GP → Very sample efficient. Cannot handle large dataset.



Spectrum of model-based RL



Simulate environment to generate additional data: DYNA

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

(a) $S \leftarrow$ current (nonterminal) state

(b) $A \leftarrow \varepsilon$ -greedy(S, Q)

(c) Take action A ; observe resultant reward, R , and state, S'

(d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

(e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)

(f) Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

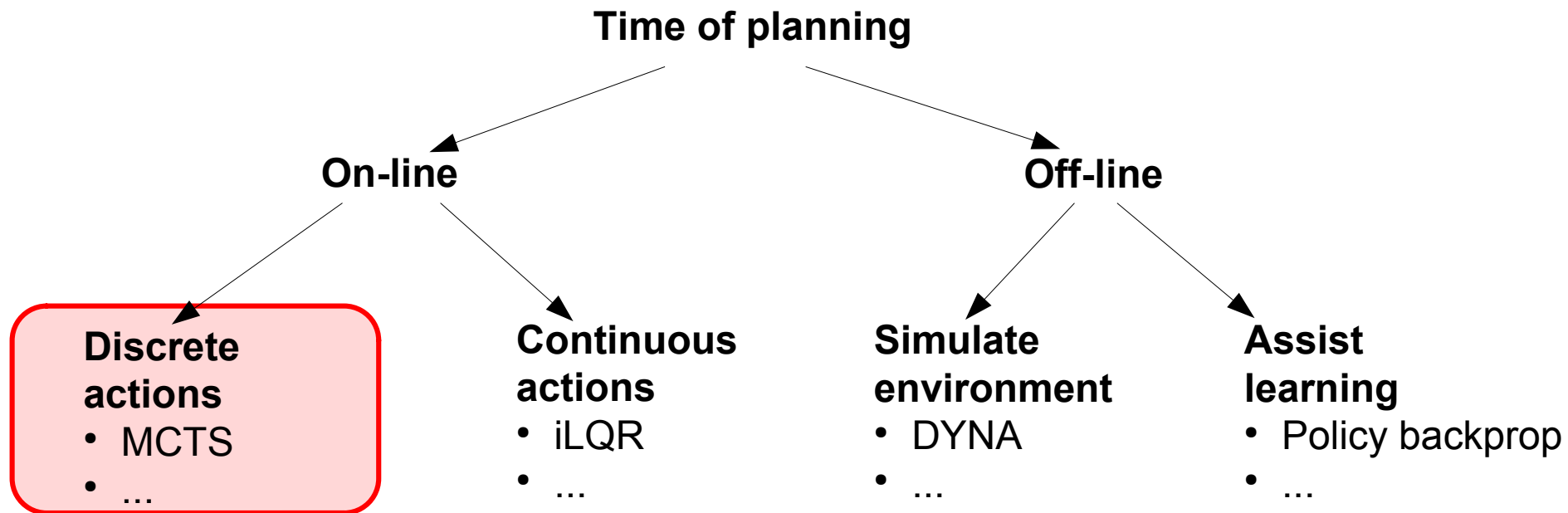
Update using experience

Learn dynamics model

Generate data by simulating dynamics

Update using simulated experience

Spectrum of model-based RL



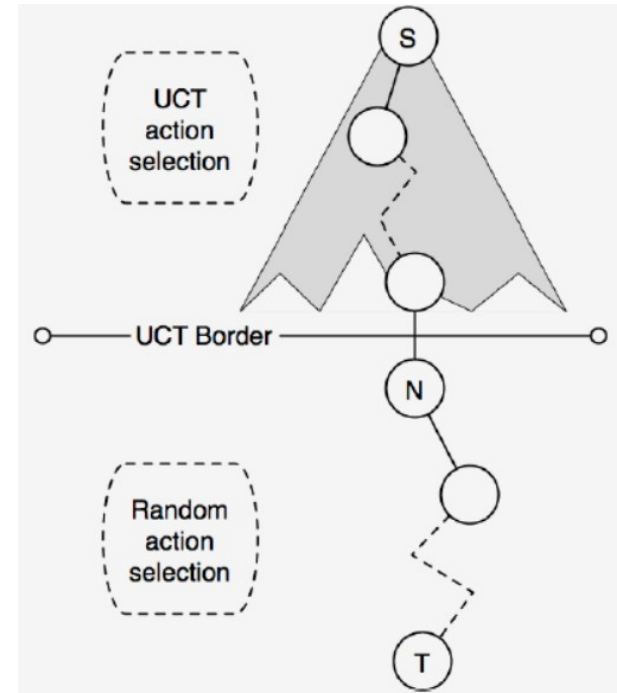
Monte Carlo tree search

- Search method for optimal decision making.
- State-of-the-art for playing games (e.g. Alpha Go).

- Iteratively builds a search tree.
- Phases:
 - Selection: Choose a promising node to expand.
 - Expansion: Add a new node.
 - Simulation: Simulate value for new node.
 - Backup: Back-up value to root (update values for parents).

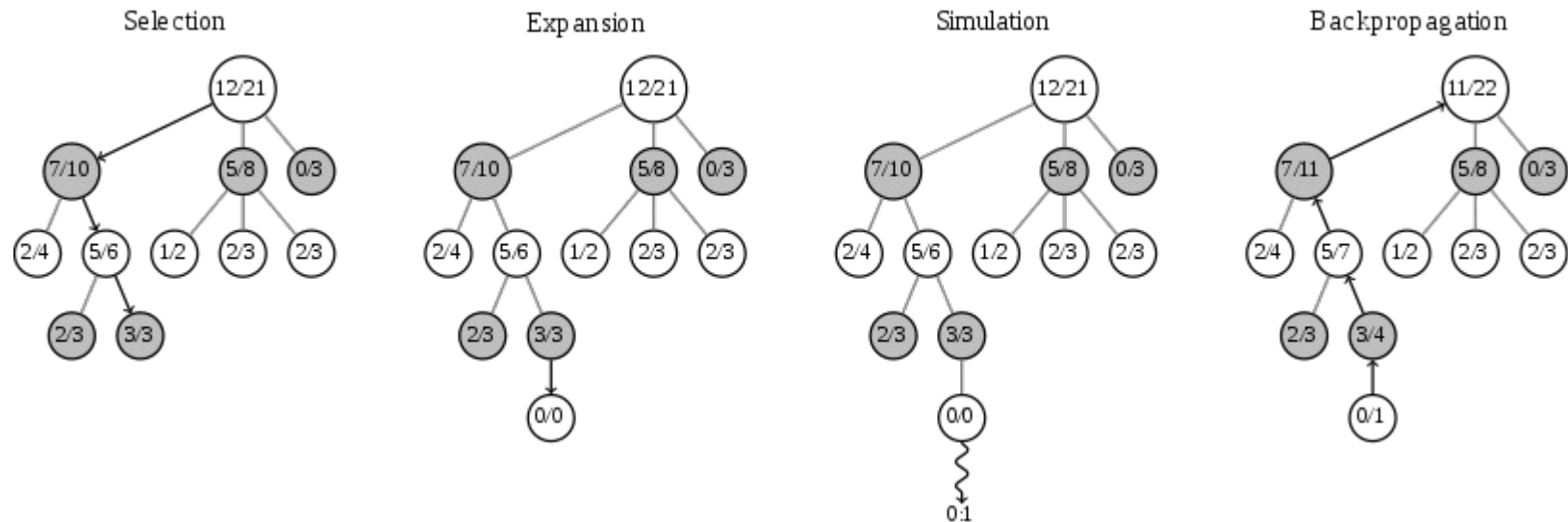
MCTS operation

- From start node S choose actions to walk down tree until reaching a leaf node.
- Choose an action and create a child node N for that action.
- Perform a **random** roll-out (take random actions) until end of episode (or for a fixed horizon).
- Record returns as value for N and back up value to root.



MCTS: Example in game playing

- Value: number of won/simulated games.



Node selection in MCTS

- Node selection in search has to balance between exploration and exploitation (note difference to RL, here exploration & exploitation only using simulation).
- Idea: Explore when uncertain of outcome.
- Upper confidence bound 1 (UCB1) on trees (UCT).
 - A bound for value of a node (Kocsis&Szepesvari, 2006).

$$Q^+(s, a) = Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

Positive exploration constant

Visitation count

MCTS simulation phase

- Perform one or several roll-outs from leaf node using random action selection.
- Stop at terminal state or until a discount horizon is reached.
- Estimate value of state as mean return of the N simulations:
$$V(s) = \frac{1}{N} \sum_i G_i$$

MCTS backpropagation

- After simulation phase backpropagate values to the root node
- Estimate value of state as mean return of the N simulations:

$$V(s) = \sum_a \frac{N_a}{N} Q(s, a)$$

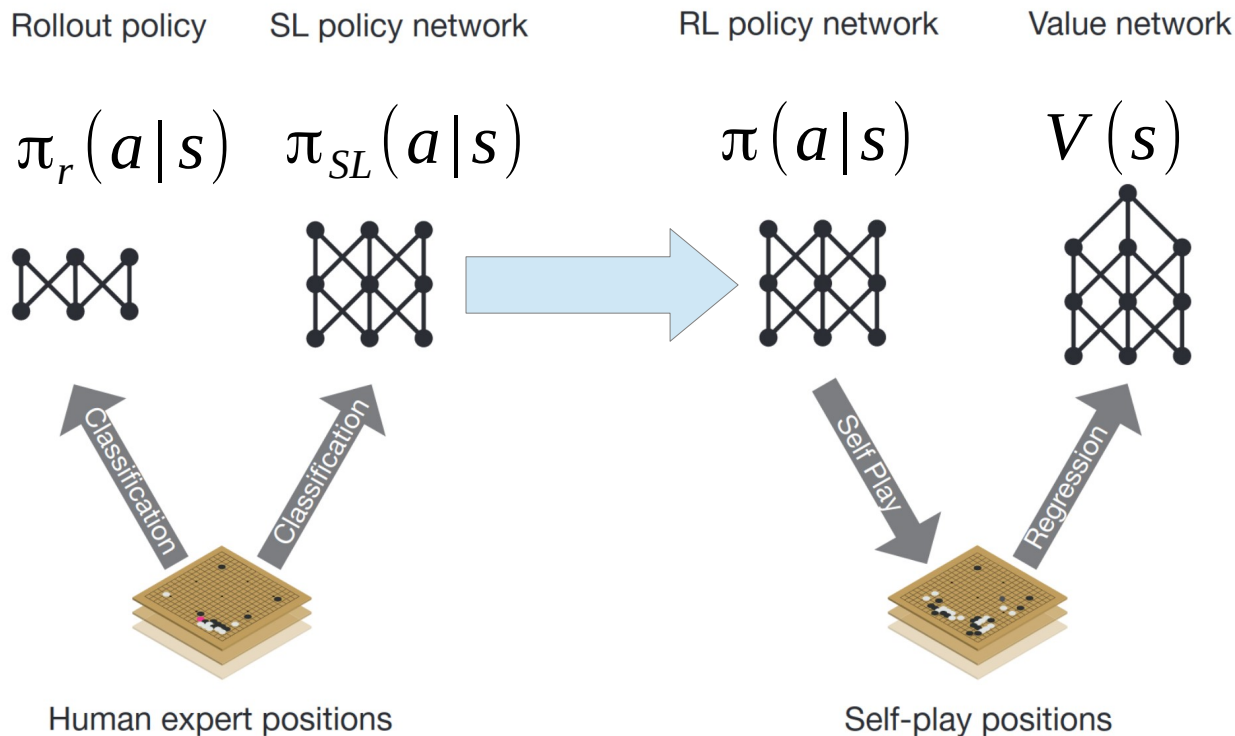
$$Q(s, a) = E_{s' \sim p(\cdot | s, a)} [V(s')]$$

MCTS extensions

- AlphaGo (2016)
 - Learn initial policy from expert demonstrations
 - Update policy using self-play and MCTS
- AlphaZero (2017, 2018)
 - No expert demonstrations needed
- MuZero (2020)
 - Similar to AlphaZero but interleaves model learning and MCTS
 - Does not require a known model

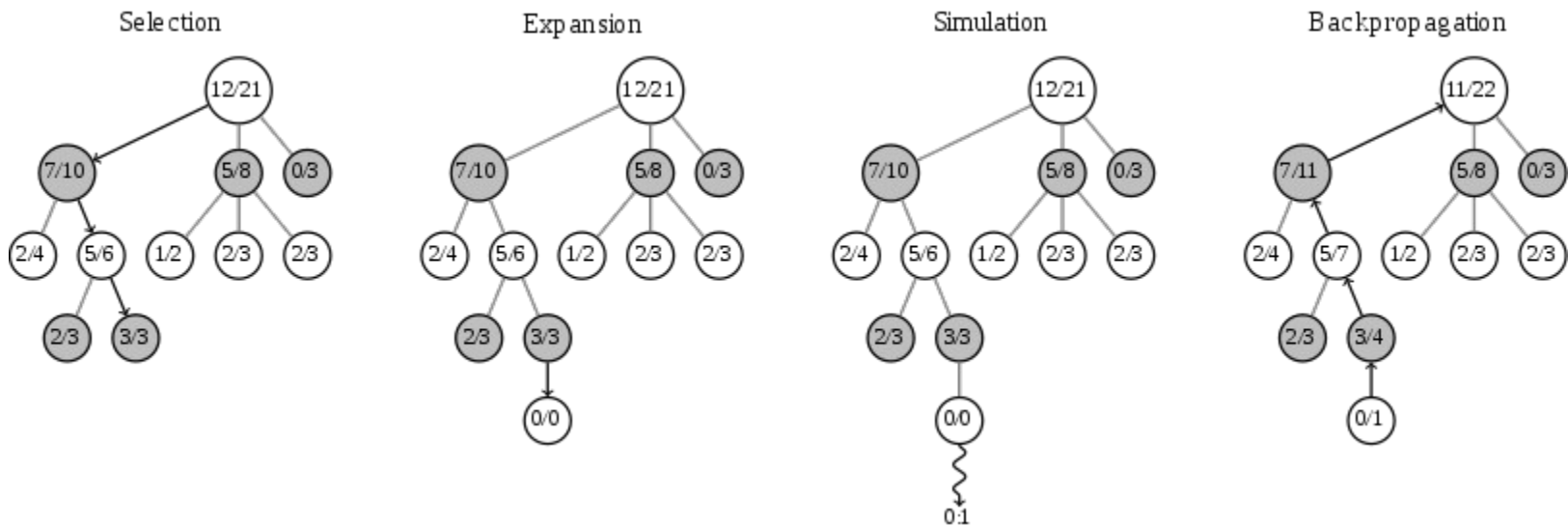
Example: Alpha Go (2016)

- Policy learned initially to imitate human players.
- Updated through policy gradient and self-play.

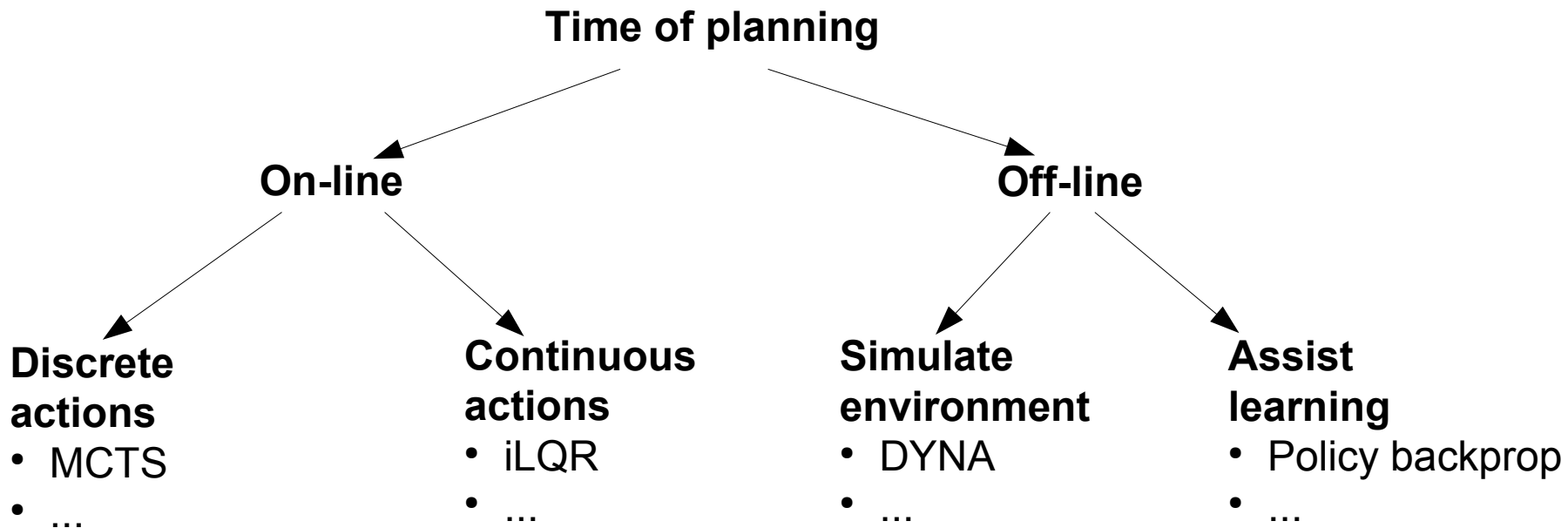


Example: Alpha Go (2016)

- Action chosen by bandit using $Q(s,a)$ and policy
- Leaf-node value: estimated value $V(s)$ plus roll-out value



Spectrum of model-based RL



Summary

- Model-based RL requires typically less data than value-based or policy gradient approaches.
- Learned dynamics can be transferred across tasks.
- Potentially suboptimal: models do not optimize for task performance and policy optimization may be prone to local minima.
- Sometimes models are harder to learn than policy.
- Often require explicit choices (e.g. time horizon).

Next: Guest Lectures

- Next week: Two! guest lectures!
 - Safety and constraints (Gökhan Alcal)
 - Entropy Regularization in Reinforcement Learning (Riad Akrouf)