*Inria*

**GROUPE RENAULT**

# Reinforcement Learning for Safe Decision-Making in Autonomous Driving

Edouard Leurent[1,2,3],
Odalric-Ambrym Maillard[1],
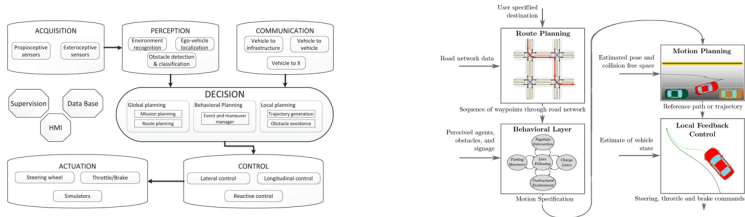Denis Efimov[2]

[1]Inria SequeL,
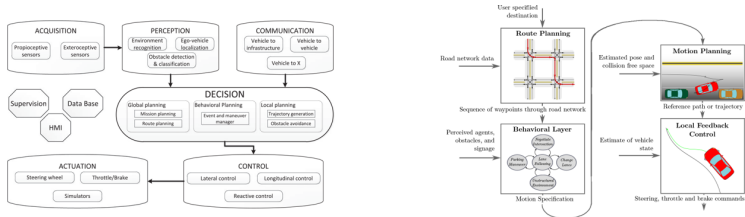[2]Inria Valse,
[3]Renault Group

# 01

## Motivation and Scope

*Inria*
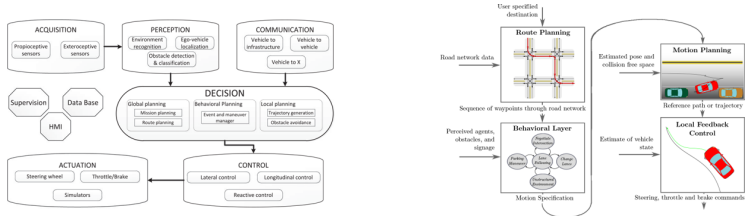
## Classic Autonomous Driving Pipeline

*Ínría*

## Classic Autonomous Driving Pipeline



### (Bold?) Claim

If we remove the humans on the road, the problem becomes easy.

## Classic Autonomous Driving Pipeline



### (Bold?) Claim

If we remove the humans on the road, the problem becomes easy.

- ✓ Even with obstacles, partial observability, disturbances, etc.
- ✓ The problems of Route Planning, Motion Planning, Local Feedback Control are basically solved.

✗ We focus instead on the (arguably) harder challenge:

Behavioural Planning

**What we have**

- In practice, often a hand-crafted rule-based system (FSM).
- Won't scale to complex scenes

**What we want**

- Handle human agents with uncertain behaviours
- Handle the interactions between agents

↳ We turn to learning-based approaches

*Inría*

## Markov Decision Processes

1. Observe state $s \in S$;

*Inria*

**Markov Decision Processes**

1. Observe state $s \in S$;
2. Pick action $a \in A$ according to our policy $\pi(a|s)$;

*Inría*

## Markov Decision Processes

1. Observe state $s \in S$;
2. Pick action $a \in A$ according to our policy $\pi(a|s)$;
3. Transition to a next state $s' \sim \mathbb{P}\left(s'|s, a\right)$;

**Markov Decision Processes**

1. Observe state $s \in S$;
2. Pick action $a \in A$ according to our policy $\pi(a|s)$;
3. Transition to a next state $s' \sim \mathbb{P}\left(s'|s, a\right)$;
4. Receive a reward $r$.

*Inría*

## Markov Decision Processes

1. Observe state $s \in S$;
2. Pick action $a \in A$ according to our policy $\pi(a|s)$;
3. Transition to a next state $s' \sim \mathbb{P}\left(s'|s, a\right)$;
4. Receive a reward $r$.

$$\text{Objective: maximise } V = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t\right]$$

*Inria*

Markov Decision Processes

1. Observe state $s \in S$;
2. Pick action $a \in A$ according to our policy $\pi(a|s)$;
3. Transition to a next state $s' \sim \mathbb{P}\left(s'|s,a\right)$;
4. Receive a reward $r$.

$$\text{Objective: maximise } V = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t\right]$$

- **States**: Ground truth for vehicles, roads, signals, etc.
  - ↪ Continuous
- **Actions**: Semantic decisions: change lane, yield, pass, etc.
  - ↪ Discrete

*Inria*

**Model-free**

1. Directly optimise $\pi(a|s)$ through policy evaluation and policy improvement
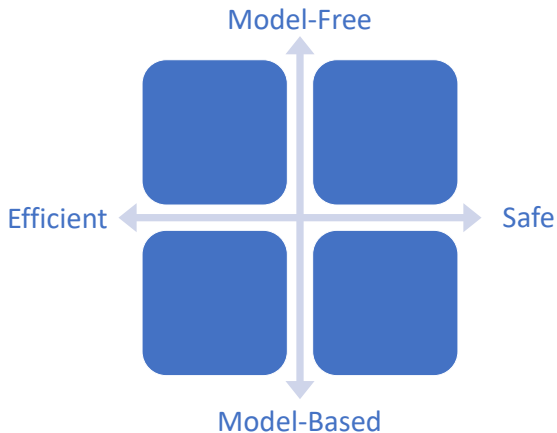
*Inria*

### Model-free

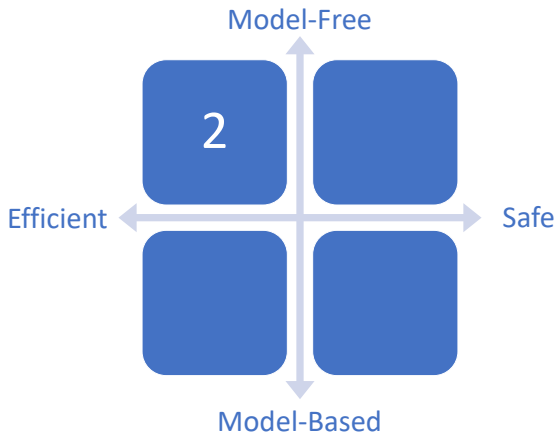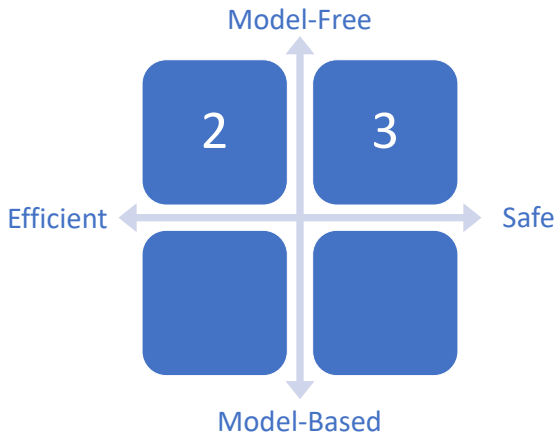1. Directly optimise $\pi(a|s)$ through policy evaluation and policy improvement

### Model-based

1. Learn a model for the dynamics $\hat{T}(s_{t+1}|s_t, a_t)$,
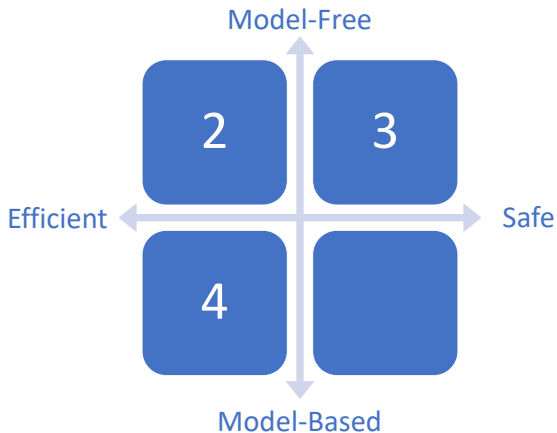2. (*Planning*) Leverage it to compute

$$\max_\pi \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \,\middle|\, a_t \sim \pi(s_t), s_{t+1} \sim \hat{T}(s_t, a_t)\right]$$
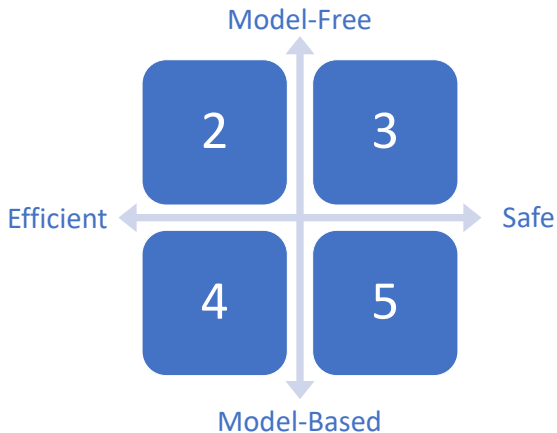
+ Better sample efficiency, interpretability, priors.

*Inría*

# 02

## Efficient Model-Free

*Inría*

Definition (Optimal State-action Value Function $Q^*$)

$$Q^*(s, a) = \max_\pi \mathbb{E}_\pi \left[ \sum_{t=0}^\infty \gamma^t R(s_t, a_t) \,\middle|\, s_0 = s, a_0 = a \right]$$

How to learn $Q^*$ ?

Proposition (Bellman Optimality Equation)

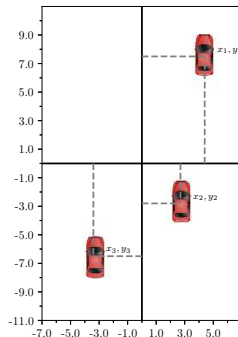$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s'} \max_{a'} Q^*(s', a')$$

↳ Represent $Q^*$ with function approximation
  (e.g. a neural network in DQN)

↳ Apply fixed-point iteration over samples $(s, a, s')$ until
  convergence

*Inria*

## The list of features representation

A joint state $s$ of $N + 1$ observed vehicles

$$s = (s_i)_{i \in [0,N]}$$

$$s_i = \begin{bmatrix} x_i & y_i & v_i^x & v_i^y & \cos \psi_i & \sin \psi_i \end{bmatrix}^T$$

Ínría

Issues related to function approximation

1. Variable size
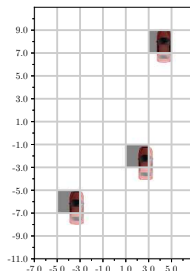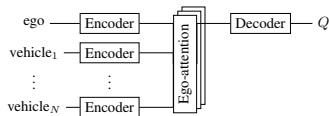   ↳ usual models accept fixed-size inputs
2. Sensitivity to the ordering
   ↳ we want the policy to be permutation-invariant:

$$\forall \tau \in \mathfrak{S}_N, \quad \pi(\cdot | (s_0, s_1, \ldots, s_N)) = \pi(\cdot | (s_0, s_{\tau(1)}, \ldots, s_{\tau(N)}))$$
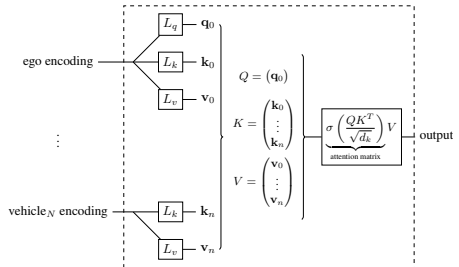
*Inría*

## Occupancy grid representation



✓ Fixed-size

✓ Does not depend on an ordering
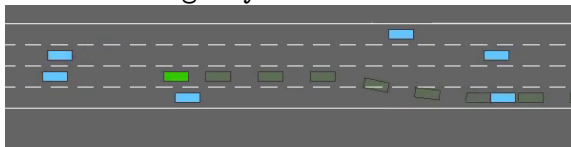
✗ Suffers from an accuracy / size tradeoff

*Inria*

Model architecture

Ego-attention block

✓ Inputs can have a variable size
✓ Based on a dot product
   ↳ permutation-invariant
✓ Compact size with no accuracy loss

*Inria*

## The `highway-env` environment



| Agent | FCN/List | CNN/Grid | Ego-Attention |
|---|---|---|---|
| Input sizes | [15, 7] | [32, 32, 7] | [ · , 7] |
| Layers sizes | [128, 128] | Convolutional layers: 3 | Encoder: [64, 64] |
| | | Kernel Size: 2 | Attention: 2 heads |
| | | Stride: 2 | $d_k = 32$ |
| | | Head: [20] | Decoder: [64, 64] |
| Number of parameters | 3.0e4 | 3.2e4 | 3.4e4 |
| Variable input size | No | No | Yes |
| Permutation invariant | No | Yes | Yes |

*Inria*

Head specialisation



Distance

*Inría*

Sensitivity to uncertainty



A full episode

*Inría*

# 03

Safe Model-Free

*Inria*

Reinforcement learning relies on a single reward function $R$

*Inría*

Reinforcement learning relies on a single reward function $R$

✓ A convenient formulation, but;

Reinforcement learning relies on a single reward function $R$

- ✓ A convenient formulation, but;
- ✗ $R$ is not always easy to design.

*Inria*

Reinforcement learning relies on a single reward function $R$

- ✓ A convenient formulation, but;
- ✗ $R$ is not always easy to design.

Conflicting Objectives

Complex tasks require multiple contradictory aspects. Typically:

<div align="center">

Task completion    vs    Safety

</div>

*Inría*

Reinforcement learning relies on a single reward function $R$

- ✓ A convenient formulation, but;
- ✗ $R$ is not always easy to design.

Conflicting Objectives

Complex tasks require multiple contradictory aspects. Typically:

<div align="center">

Task completion    vs    Safety

</div>

For example...

*Inría*

## Two-Way Road

The agent is driving on a two-way road with a car in front of it,

- it can stay behind (safe/slow);
- it can overtake (unsafe/fast).

Reinforcement learning relies on a single reward function $R$

  ✓ A convenient formulation, but;

  ✗ $R$ is not always easy to design.

Conflicting Objectives

Complex tasks require multiple contradictory aspects. Typically:

$$\text{Task completion} \quad \text{vs} \quad \text{Safety}$$

For example...

For a fixed reward function $R$,

  ↳ $\pi^*$ is only guaranteed to lie on a Pareto front $\Pi^*$

  ↳ no control over the $\frac{\text{Task Completion}}{\text{Safety}}$ trade-off

*Inría*

Task Completion $G_1 = \sum \gamma^t R_1^t$

Pareto-optimal curve $\Pi^*$

$$\underset{\pi}{\text{argmax}} \sum \gamma^t R_t(R_1, R_2)$$

Safety $G_2 = \sum \gamma^t R_2^t$

*Inria*

Task Completion $G_r$

Pareto-optimal curve $\Pi^*$

$\underset{\pi}{\text{argmax}} \sum \gamma^t R_r^t$

$s.t. \sum \gamma^t R_c^t < \beta$

$\pi^*$

Risk $G_c$

$\beta$

*Inria*

Markov Decision Process

An MDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, \gamma)$ with:

- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$

Objective

Maximise rewards

$$\max_{\pi \in \mathcal{M}(\mathcal{A})^{\mathcal{S}}} \quad \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s\right]$$

*Inría*

## Constrained Markov Decision Process

A CMDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, R_c, \gamma, \beta)$ with:

- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Costs $R_c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Budget $\beta$

### Objective

Maximise rewards while keeping costs under a fixed budget

$$\max_{\pi \in \mathcal{M}(\mathcal{A})^{\mathcal{S}}} \quad \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s\right]$$
$$\text{s.t.} \quad \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_c(s_t, a_t) \mid s_0 = s\right] \leq \beta$$

*Inria*

Task Completion $G_r$

Pareto-optimal curve $\Pi^*$

$$\underset{\pi}{\mathrm{argmax}} \sum \gamma^t R_r^t$$
$$s.t. \sum \gamma^t R_c^t < \beta$$

$\pi^*$

Risk $G_c$

$\beta$

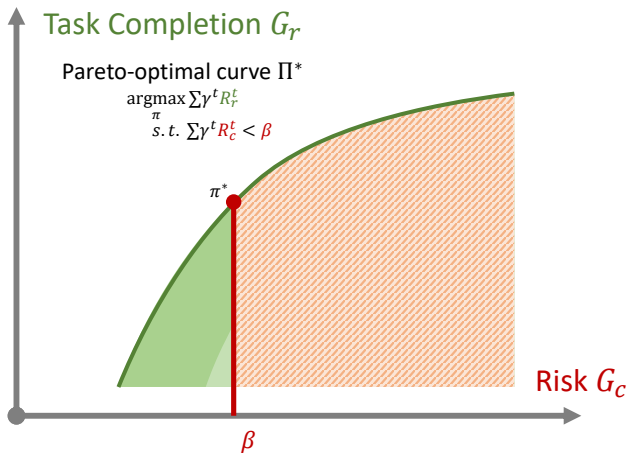*Inria*

Budgeted Markov Decision Process

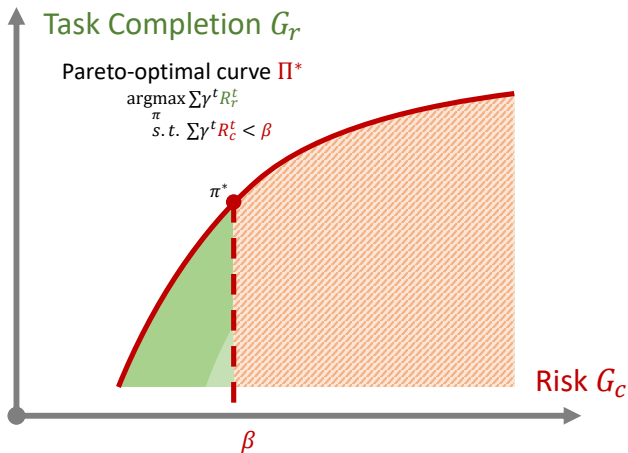A BMDP is a tuple $(\mathcal{S}, \mathcal{A}, P, R_r, R_c, \gamma, \mathcal{B})$ with:

- Rewards $R_r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Costs $R_c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
- Budget space $\mathcal{B}$

Objective

Maximise rewards while keeping costs under an adjustable budget. $\forall \beta \in \mathcal{B}$,

$$\max_{\pi \in \mathcal{M}(\mathcal{A} \times \mathcal{B})^{\mathcal{S} \times \mathcal{B}}} \quad \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_r(s_t, a_t) \mid s_0 = s, \beta_0 = \beta\right]$$
$$\text{s.t.} \quad \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_c(s_t, a_t) \mid s_0 = s, \beta_0 = \beta\right] \leq \beta$$

*Inria*

Budgeted policies $\pi$

- Take a budget $\beta$ as an additional input
- Output a next budget $\beta'$
- $\pi : \underbrace{(s, \beta)}_{\overline{s}} \to \underbrace{(a, \beta')}_{\overline{a}}$

$\hookrightarrow$ Augment the spaces with the budget $\beta$

*Inría*

**Definition (Augmented spaces)**

- States $\overline{\mathcal{S}} = \mathcal{S} \times \mathcal{B}$.
- Actions $\overline{\mathcal{A}} = \mathcal{A} \times \mathcal{B}$.
- Dynamics $\overline{P}$

  state $(s, \beta)$, action $(a, \beta_a) \to$ next state $\begin{cases} s' \sim P(s'|s, a) \\ \beta' = \beta_a \end{cases}$

**Definition (Augmented signals)**

1. Rewards $R = (R_r, R_c)$
2. Returns $G^\pi = (G_r^\pi, G_c^\pi) \overset{\text{def}}{=} \sum_{t=0}^{\infty} \gamma^t R(\overline{s}_t, \overline{a}_t)$
3. Value $V^\pi(\overline{s}) = (V_r^\pi, V_c^\pi) \overset{\text{def}}{=} \mathbb{E}\left[ G^\pi \mid \overline{s_0} = \overline{s} \right]$
4. Q-Value $Q^\pi(\overline{s}, \overline{a}) = (Q_r^\pi, Q_c^\pi) \overset{\text{def}}{=} \mathbb{E}\left[ G^\pi \mid \overline{s_0} = \overline{s}, \overline{a_0} = \overline{a} \right]$

*Inría*

## Definition (Budgeted Optimality)

In that order, we want to:

(i) Respect the budget $\beta$:

$$\Pi_a(\bar{s}) \overset{\text{def}}{=} \{\pi \in \Pi : V_c^\pi(s, \beta) \leq \beta\}$$

(ii) Maximise the rewards:

$$V_r^*(\bar{s}) \overset{\text{def}}{=} \max_{\pi \in \Pi_a(\bar{s})} V_r^\pi(\bar{s}) \qquad \Pi_r(\bar{s}) \overset{\text{def}}{=} \arg\max_{\pi \in \Pi_a(\bar{s})} V_r^\pi(\bar{s})$$

(iii) Minimise the costs:

$$V_c^*(\bar{s}) \overset{\text{def}}{=} \min_{\pi \in \Pi_r(\bar{s})} V_c^\pi(\bar{s}), \qquad \Pi^*(\bar{s}) \overset{\text{def}}{=} \arg\min_{\pi \in \Pi_r(\bar{s})} V_c^\pi(\bar{s})$$

We define the budgeted action-value function $Q^*$ similarly

*Inria*

### Theorem (Budgeted Bellman Optimality Equation)

$Q^*$ *verifies the following equation:*

$$Q^*(\overline{s}, \overline{a}) = \mathcal{T}Q^*(\overline{s}, \overline{a})$$
$$\overset{def}{=} R(\overline{s}, \overline{a}) + \gamma \sum_{\overline{s}' \in \overline{\mathcal{S}}} \overline{P}(\overline{s}'|\overline{s}, \overline{a}) \sum_{\overline{a}' \in \overline{\mathcal{A}}} \pi_{greedy}(\overline{a}'|\overline{s}'; Q^*) Q^*(\overline{s}', \overline{a}')$$

where the greedy policy $\pi_{\text{greedy}}$ is defined by:

$$\pi_{\text{greedy}}(\overline{a}|\overline{s}; Q) \in \arg\min_{\rho \in \Pi_r^Q} \mathbb{E}_{\overline{a} \sim \rho} Q_c(\overline{s}, \overline{a}),$$

$$\text{where} \quad \Pi_r^Q \overset{\text{def}}{=} \arg\max_{\rho \in \mathcal{M}(\overline{\mathcal{A}})} \mathbb{E}_{\overline{a} \sim \rho} Q_r(\overline{s}, \overline{a})$$

$$\text{s.t.} \quad \mathbb{E}_{\overline{a} \sim \rho} Q_c(\overline{s}, \overline{a}) \leq \beta$$

*Inría*

**Proposition (Optimality of the policy)**

$\pi_{greedy}(\cdot \; ; Q^*)$ is *simultaneously optimal* in all states $\overline{s} \in \overline{\mathcal{S}}$:

$$\pi_{greedy}(\cdot \; ; Q^*) \in \Pi^*(\overline{s})$$

In particular, $V^{\pi_{greedy}(\cdot;Q^*)} = V^*$ and $Q^{\pi_{greedy}(\cdot;Q^*)} = Q^*$.

**Proposition (Solving the non-linear program)**

$\pi_{greedy}$ can be computed *efficiently*, as a mixture $\pi_{hull}$ of two points that lie on the convex hull of $Q$.

$$\pi_{greedy} = \pi_{hull}$$

*Inria*

Recall what we've shown so far:

$$\mathcal{T} \xrightarrow[\text{fixed−point}]{} Q^* \xrightarrow[\text{tractable}]{} \pi_{\text{hull}}(Q^*) \xrightarrow[\text{equal}]{} \pi_{\text{greedy}}(Q^*) \xrightarrow[\text{optimal}]{}$$

*Inria*

Recall what we've shown so far:

$$\mathcal{T} \xrightarrow[fixed-point]{} Q^* \xrightarrow[tractable]{} \pi_{\text{hull}}(Q^*) \xrightarrow[equal]{} \pi_{\text{greedy}}(Q^*) \xrightarrow[optimal]{}$$

We're almost there!
All that is left is to perform Fixed-Point Iteration to compute $Q^*$.

*Inria*

Recall what we've shown so far:

$$\mathcal{T} \xrightarrow[\text{fixed-point}]{} Q^* \xrightarrow[\text{tractable}]{} \pi_{\text{hull}}(Q^*) \xrightarrow[\text{equal}]{} \pi_{\text{greedy}}(Q^*) \xrightarrow[\text{optimal}]{}$$

We're almost there!

All that is left is to perform Fixed-Point Iteration to compute $Q^*$.

---

Theorem (Non-Contractivity)

*For any BMDP $(\mathcal{S}, \mathcal{A}, P, R_r, R_c, \gamma)$ with $|\mathcal{A}| \geq 2$, $\mathcal{T}$ is **not** a contraction.*

$$\forall \varepsilon > 0, \exists Q^1, Q^2 \in (\mathbb{R}^2)^{\overline{\mathcal{SA}}} : \|\mathcal{T}Q^1 - \mathcal{T}Q^2\|_\infty \geq \frac{1}{\varepsilon} \|Q^1 - Q^2\|_\infty$$

---

✗ We cannot guarantee the convergence of $\mathcal{T}^n(Q_0)$ to $Q^*$

*Inria*

Thankfully,

### Theorem (Contractivity on smooth Q-functions)

$\mathcal{T}$ *is a contraction* when restricted to the subset $\mathcal{L}_\gamma$ of Q-functions such that "$Q_r$ is L-Lipschitz with respect to $Q_c$", with $L < \frac{1}{\gamma} - 1$.

$$\mathcal{L}_\gamma = \left\{ \begin{array}{l} Q \in (\mathbb{R}^2)^{\overline{\mathcal{S}\mathcal{A}}} \text{ s.t. } \exists L < \frac{1}{\gamma} - 1 : \forall \overline{s} \in \overline{\mathcal{S}}, \overline{a}_1, \overline{a}_2 \in \overline{\mathcal{A}}, \\ |Q_r(\overline{s}, \overline{a}_1) - Q_r(\overline{s}, \overline{a}_2)| \leq L|Q_c(\overline{s}, \overline{a}_1) - Q_c(\overline{s}, \overline{a}_2)| \end{array} \right\}$$
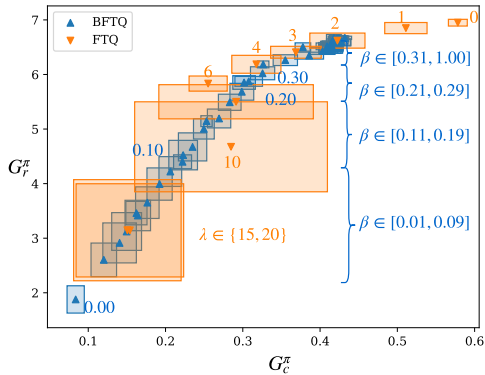
✓ We guarantee convergence under some (strong) assumptions
✓ We observe empirical convergence

*Inría*

## Lagrangian Relaxation Baseline

Consider the dual problem so as to replace the hard constraint by a soft constraint penalised by a Lagrangian multiplier $\lambda$:

$$\max_{\pi} \mathbb{E} \sum_t \gamma^t R_r(s, a) - \lambda \gamma^t R_c(s, a)$$

- Train many policies $\pi_k$ with penalties $\lambda_k$ and recover the cost budgets $\beta_k$
- Very data/memory-heavy

*Inria*

# 04

Efficient
Model-Based

*Inria*

## Model estimation

Learn a model for the dynamics $\hat{T}(s_{t+1}|s_t, a_t)$. For instance:

1. Least-square estimate: $\min_{\hat{T}} \sum_t \|s_{t+1} - \hat{T}(s_t, a_t)\|_2^2$
2. Maximum Likelihood estimate: $\max_{\hat{T}} \sum_t \hat{T}(s_{t+1}|s_t, a_t)$
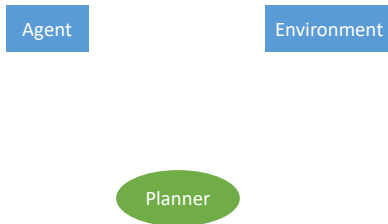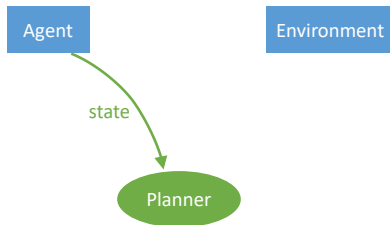
## Planning

Leverage $\hat{T}$ to compute

$$\max_\pi \mathbb{E}\left[\sum_{t=0}^\infty \gamma^t r(s_t, a_t) \;\middle|\; a_t \sim \pi(s_t), s_{t+1} \sim \hat{T}(s_t, a_t)\right]$$
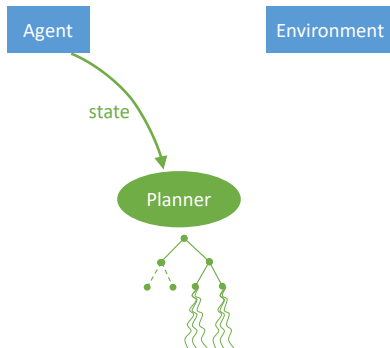
How?

*Inria*

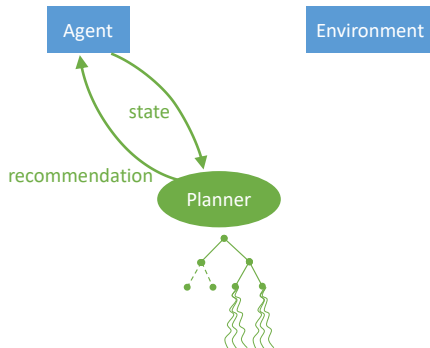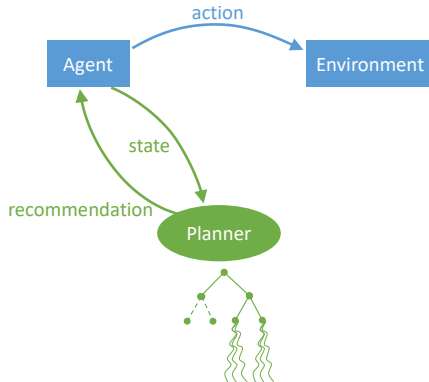We can use $\hat{T}$ as a generative model:

Agent

Environment

Planner

We can use $\hat{T}$ as a generative model:
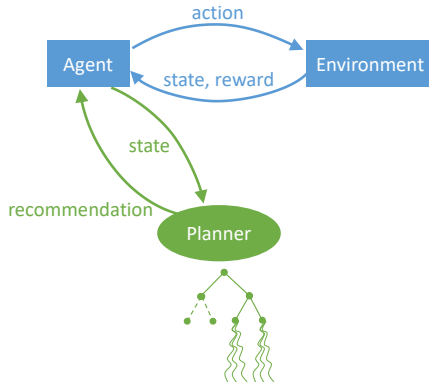
We can use $\hat{T}$ as a generative model:

We can use $\hat{T}$ as a generative model:

We can use $\hat{T}$ as a generative model:

We can use $\hat{T}$ as a generative model:

**Online *Planning***

- **fixed budget**: the model can only be queried $n$ times

    Objective: minimize $\mathbb{E} \underbrace{V^* - V(n)}_{\text{Simple Regret } r_n}$

    An exploration-exploitation problem.

**Optimism in the Face of Uncertainty**

Given a set of options $a \in A$ with uncertain outcomes, try the one with the highest possible outcome.

*Inria*

## Optimism in the Face of Uncertainty

Given a set of options $a \in A$ with uncertain outcomes, try the one with the highest possible outcome.

- Either you performed well;

*Inria*

### Optimism in the Face of Uncertainty

Given a set of options $a \in A$ with uncertain outcomes, try the one with the highest possible outcome.

- Either you performed well;
- or you learned something.
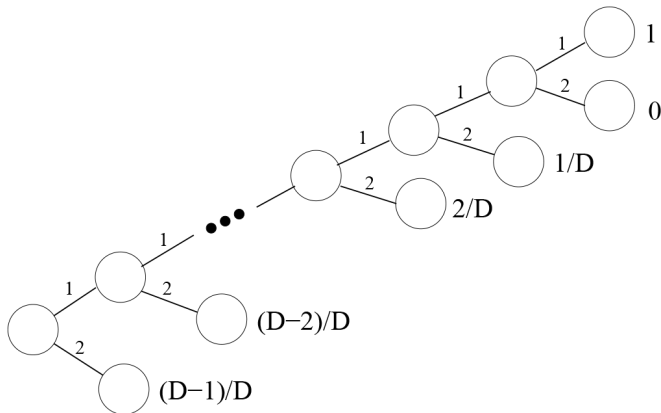
*Inria*

## Optimism in the Face of Uncertainty

Given a set of options $a \in A$ with uncertain outcomes, try the one with the highest possible outcome.

- Either you performed well;
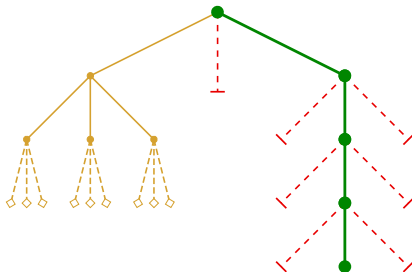- or you learned something.

## Instances
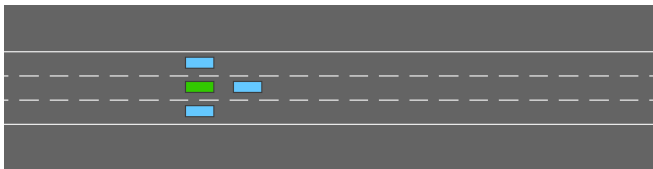
- Monte-carlo tree search (MCTS) (Coulom, 2006): `CrazyStone`
- Reframed in the bandit setting as UCT (Kocsis and Szepesvári, 2006), still very popular (e.g. `Alpha Go`).
- Proved asymptotic consistency, but no regret bound.

*Inria*

It was analysed in (Coquelin and Munos, 2007)]



The sample complexity of is lower-bounded by $O(\exp(\exp(D)))$.

Not just a theoretical counter-example.

OPD: Optimistic Planning for Deterministic systems

- Introduced by (Hren and Munos, 2008)
- Another optimistic algorithm
- Only for deterministic MDPs

Theorem (OPD sample complexity)

$$\mathbb{E}\, r_n = \mathcal{O}\left( n^{-\frac{\log 1/\gamma}{\log \kappa}} \right), \text{if } \kappa > 1$$

*Inria*

**OPD: Optimistic Planning for Deterministic systems**
- Introduced by (Hren and Munos, 2008)
- Another optimistic algorithm
- Only for deterministic MDPs

**Theorem (OPD sample complexity)**

$$\mathbb{E}\, r_n = \mathcal{O}\left(n^{-\frac{\log 1/\gamma}{\log \kappa}}\right), \textit{if } \kappa > 1$$

**OLOP: Open-Loop Optimistic Planning**
- Introduced by (Bubeck and Munos, 2010)
- Extends OPD to the stochastic setting
- Only considers open-loop policies, i.e. sequences of actions

*Inría*

A direct application of Optimism in the Face of Uncertainty

1. We want

$$\max_a V(a)$$

A direct application of Optimism in the Face of Uncertainty

1. We want

$$\max_a V(a)$$

2. Form upper confidence-bounds of sequence values:

$$V(a) \leq U_a \quad \text{w.h.p}$$

*Inria*

A direct application of Optimism in the Face of Uncertainty
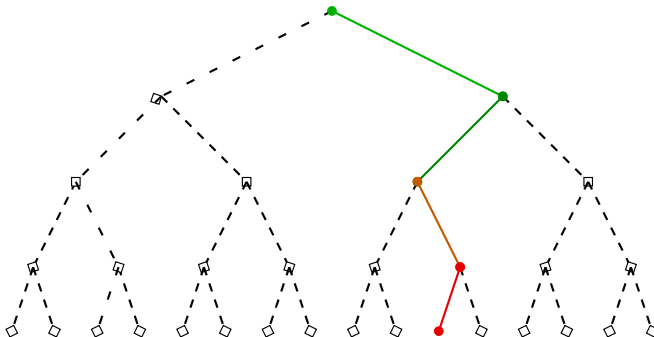
1. We want

$$\max_a V(a)$$
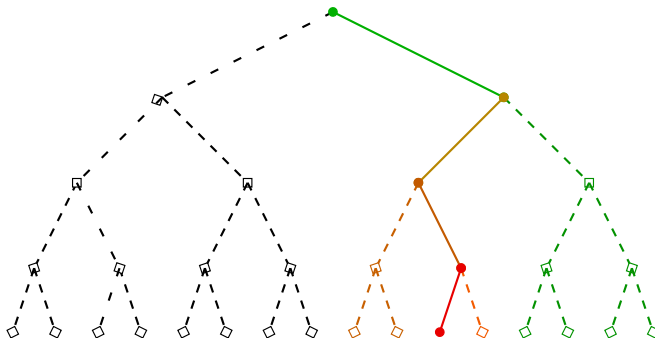
2. Form upper confidence-bounds of sequence values:

$$V(a) \leq U_a \quad \text{w.h.p}$$

3. Sample the sequence with highest UCB:

$$\arg\max_a U_a$$

*Inria*

Upper-bounding the value of sequences

$$V(a) = \overbrace{\sum_{t=1}^{h} \gamma^t \mu_{a_{1:t}}}^{\text{follow the sequence}} + \overbrace{\sum_{t \geq h+1} \gamma^t \mu_{a^*_{1:t}}}^{\text{act optimally}}$$

*Inria*

Upper-bounding the value of sequences

$$V(a) = \overbrace{\sum_{t=1}^{h} \gamma^t \underbrace{\mu_{a_{1:t}}}_{\leq U^\mu}}^{\text{follow the sequence}} + \overbrace{\sum_{t \geq h+1} \gamma^t \underbrace{\mu_{a^*_{1:t}}}_{\leq 1}}^{\text{act optimally}}$$

*Inria*

## OLOP main tool: the Chernoff-Hoeffding deviation inequality

$$\underbrace{U_a^\mu(m)}_{\text{Upper bound}} \overset{\text{def}}{=} \underbrace{\hat{\mu}_a(m)}_{\text{Empirical mean}} + \underbrace{\sqrt{\frac{2\log M}{T_a(m)}}}_{\text{Confidence interval}}$$

*Inría*

`OLOP` main tool: the Chernoff-Hoeffding deviation inequality

$$\underbrace{U_a^\mu(m)}_{\text{Upper bound}} \stackrel{\text{def}}{=} \underbrace{\hat{\mu}_a(m)}_{\text{Empirical mean}} + \underbrace{\sqrt{\frac{2\log M}{T_a(m)}}}_{\text{Confidence interval}}$$

`OPD`: upper-bound all the future rewards by 1

$$U_a(m) \stackrel{\text{def}}{=} \sum_{t=1}^{h} \underbrace{\gamma^t U_{a_{1:t}}^\mu(m)}_{\text{Past rewards}} + \underbrace{\frac{\gamma^{h+1}}{1-\gamma}}_{\text{Future rewards}}$$

*Inria*

`OLOP` main tool: the Chernoff-Hoeffding deviation inequality

$$\underbrace{U_a^\mu(m)}_{\text{Upper bound}} \overset{\text{def}}{=} \underbrace{\hat{\mu}_a(m)}_{\text{Empirical mean}} + \underbrace{\sqrt{\frac{2\log M}{T_a(m)}}}_{\text{Confidence interval}}$$

`OPD`: upper-bound all the future rewards by 1

$$U_a(m) \overset{\text{def}}{=} \sum_{t=1}^{h} \underbrace{\gamma^t U_{a_{1:t}}^\mu(m)}_{\text{Past rewards}} + \underbrace{\frac{\gamma^{h+1}}{1-\gamma}}_{\text{Future rewards}}$$

*Bounds sharpening*

$$B_a(m) \overset{\text{def}}{=} \inf_{1 \le t \le L} U_{a_{1:t}}(m)$$

*Inría*

## Theorem (OLOP Sample complexity)

*OLOP satisfies:*

$$\mathbb{E}\, r_n = \begin{cases} \widetilde{\mathcal{O}}\left(n^{-\frac{\log 1/\gamma}{\log \kappa'}}\right), & \text{if } \gamma\sqrt{\kappa'} > 1 \\ \widetilde{\mathcal{O}}\left(n^{-\frac{1}{2}}\right), & \text{if } \gamma\sqrt{\kappa'} \leq 1 \end{cases}$$

*"Remarkably, in the case $\kappa\gamma^2 > 1$, we obtain the same rate for the simple regret as Hren and Munos (2008). Thus, in this case, we can say that planning in stochastic environments is not harder than planning in deterministic environments".*

*Inría*

Highway

Our objective: understand and bridge this gap.

Make OLOP *practical*.

*Inria*

## Explanation: inconsistency

- Unintended behaviour happens when $U_a^\mu(m) > 1, \forall a$.

$$U_a^\mu(m) = \underbrace{\hat{\mu}_a(m)}_{\in[0,1]} + \underbrace{\sqrt{\frac{2\log M}{T_a(m)}}}_{>0}$$

*Inría*

Explanation: inconsistency

- Unintended behaviour happens when $U_a^\mu(m) > 1, \forall a$.

$$U_a^\mu(m) = \underbrace{\hat{\mu}_a(m)}_{\in[0,1]} + \underbrace{\sqrt{\frac{2\log M}{T_a(m)}}}_{>0}$$

- Then the sequence $(U_{a_{1:t}}(m))_t$ is increasing

$$U_{a_{1:1}}(m) = \gamma U_{a_1}^\mu(m) + \gamma^2 1 \qquad\qquad + \gamma^3 1 + \dots$$
$$U_{a_{1:2}}(m) = \gamma U_{a_1}^\mu(m) + \gamma^2 \underbrace{U_{a_2}^\mu}_{>1} \qquad\qquad + \gamma^3 1 + \dots$$

*Inría*

Explanation: inconsistency

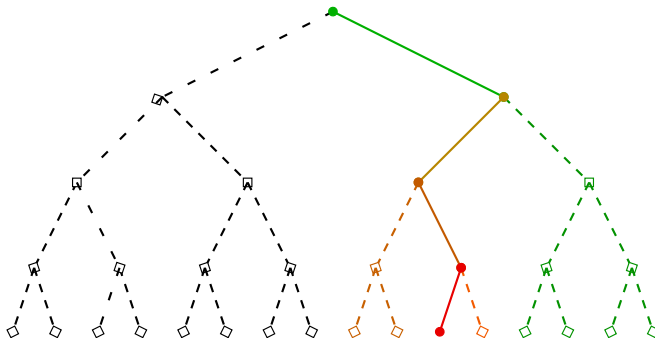- Unintended behaviour happens when $U_a^\mu(m) > 1, \forall a$.

$$U_a^\mu(m) = \underbrace{\hat{\mu}_a(m)}_{\in [0,1]} + \underbrace{\sqrt{\frac{2 \log M}{T_a(m)}}}_{>0}$$
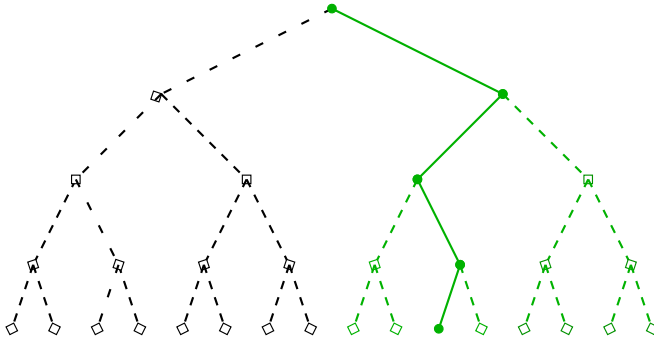
- Then the sequence $(U_{a_{1:t}}(m))_t$ is increasing

$$U_{a_{1:1}}(m) = \gamma U_{a_1}^\mu(m) + \gamma^2 1 \qquad\qquad + \gamma^3 1 + \ldots$$
$$U_{a_{1:2}}(m) = \gamma U_{a_1}^\mu(m) + \gamma^2 \underbrace{U_{a_2}^\mu}_{>1} \qquad + \gamma^3 1 + \ldots$$

- Then $B_a(m) = U_{a_{1:1}}(m)$

*Inria*

What we were promised

*Inría*

What we actually get



`OLOP` behaves as uniform planning!

We summon the upper-confidence bound from `kl-UCB` (Cappé et al., 2013):

$$U_a^\mu(m) \stackrel{\text{def}}{=} \max\left\{q \in I : T_a(m)d(\hat{\mu}_a(m), q) \le f(m)\right\}$$

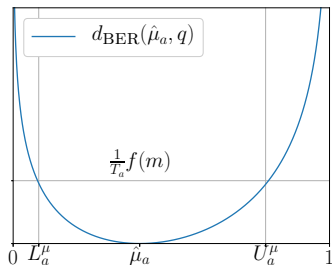*Inría*

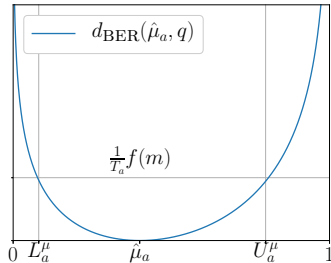We summon the upper-confidence bound from `kl-UCB` (Cappé et al., 2013):

$$U_a^\mu(m) \stackrel{\text{def}}{=} \max\left\{q \in I : T_a(m)d(\hat{\mu}_a(m), q) \leq f(m)\right\}$$

| Algorithm | `OLOP` | `KL-OLOP` |
|-----------|--------|-----------|
| Interval $I$ | $\mathbb{R}$ | $[0, 1]$ |
| Divergence $d$ | $d_{\text{QUAD}}$ | $d_{\text{BER}}$ |
| $f(m)$ | $4 \log M$ | $2 \log M + 2 \log \log M$ |

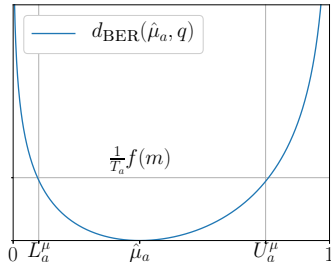$$d_{\text{QUAD}}(p, q) \stackrel{\text{def}}{=} 2(p - q)^2$$
$$d_{\text{BER}}(p, q) \stackrel{\text{def}}{=} p \log \frac{p}{q} + (1 - p) \log \frac{1 - p}{1 - q}$$
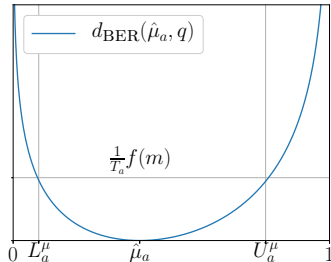
*Inría*

And now,

- $U_a^\mu(m) \in I = [0, 1], \forall a$.

And now,

- $U_a^\mu(m) \in I = [0,1], \forall a$.
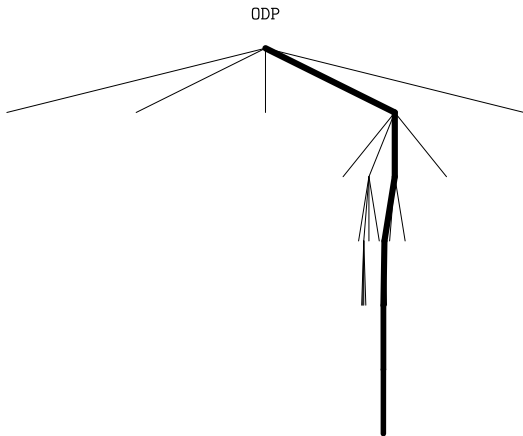- The sequence $(U_{a_{1:t}}(m))_t$ is non-increasing

*Inria*

And now,

- $U_a^\mu(m) \in I = [0,1], \forall a$.
- The sequence $(U_{a_{1:t}}(m))_t$ is non-increasing
- $B_a(m) = U_a(m)$, the *bound sharpening* step is superfluous.

*Inria*

**Theorem (Sample complexity)**
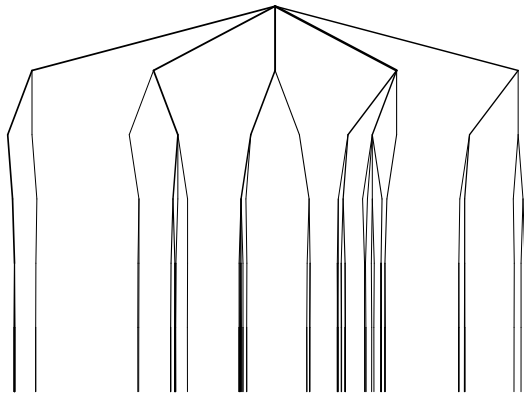
*KL-OLOP enjoys the same regret bounds as OLOP. More precisely, KL-OLOP satisfies:*

$$\mathbb{E} \, r_n = \begin{cases} \widetilde{\mathcal{O}} \left( n^{-\frac{\log 1/\gamma}{\log \kappa'}} \right), & \text{if } \gamma\sqrt{\kappa'} > 1 \\ \widetilde{\mathcal{O}} \left( n^{-\frac{1}{2}} \right), & \text{if } \gamma\sqrt{\kappa'} \leq 1 \end{cases}$$

*Inria*

ODP

OLOP

KL-OLOP

Highway

Stochastic Gridworld

# 05

Safe Model-Based

*Inria*

**Model-based RL** learns the dynamics $\hat{T}$ and optimizes

$$\max_\pi \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \,\middle|\, a_t \sim \pi(s_t), s_{t+1} \sim \hat{T}(s_t, a_t)\right]$$

Definition (Model Bias)

$$T \neq \hat{T}$$

- Video example

*Inria*

1. Build a confidence region $C_\delta$ around the true dynamics $T$

$$\mathbb{P}\left(T \in C_\delta\right) > 1 - \delta$$

2. Plan robustly with respect to this ambiguity

$$\max_\pi \underbrace{\min_{T \in C_\delta} \sum_{t=0}^{\infty} \gamma^t r_t}_{v^r(\pi)}$$

*Inría*

In order to build $C_\delta$, we rely on a structure assumption

Assumption (Structure)

$$\dot{x}(t) = A(\theta)x(t) + Bu(t) + d(t)$$

*with*

$$A(\theta) = \sum_{i=1}^{d} \theta_i \Phi_i$$

Having observed a history of $\dot{x}(t), x(t)$, we obtain a linear regression problem:

$$\min_{\theta} \|\dot{x}(t) - A(\theta)x(t) - Bu(t)\|_2^2$$

*Inria*

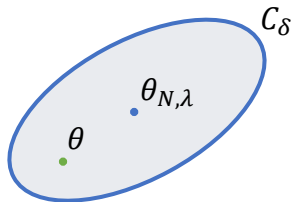**Proposition (Confidence ellipsoid (Abbasi-yadkori, Pál, and Szepesvári, 2011))**

*Under some assumptions on the disturbance $d(t)$, it holds with probability $1 - \delta$ that:*

$$\|\theta - \theta_{Np,\lambda}\|_{G_{Np,\lambda}} \leq \beta_t(\delta)$$

$$where \quad \theta_{Np,\lambda} = G_{Np,\lambda}^{-1} \Phi_{[Np]}^T Y_{[Np]};$$

$$G_{Np,\lambda} = \Phi_{[Np]}^T \Phi_{[Np]} + \lambda I_d.$$

*Ínría*

## Possible trajectories

$$\dot{x}(t) = A(\theta)x(t) + Bu(t) + d(t)$$

There are two sources of uncertainty:

- Parametric uncertainty $A(\theta) \in C_\delta$
- External perturbations $d(t)$



$x(t, \theta(t), d(t))$

$x(0)$

*Inria*

## Interval Prediction

Can we design an interval predictor $[\underline{x}(t), \overline{x}(t)]$ that verifies:

- inclusion property: $\forall t, \underline{x}(t) \leq x(t) \leq \overline{x}(t)$;
- stable dynamics?

We want the predictor to be as tight as possible. How to proceed?



$x(t, \theta(t), d(t))$

$[\underline{x}(t), \overline{x}(t)]$

$x(0)$

*Inria*

Assume that $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$, for some $t \geq 0$.

*Inria*

Assume that $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$, for some $t \geq 0$.

↳ To propagate the interval to $x(t + dt)$, we need to bound $A(\theta)x(t)$.

Assume that $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$, for some $t \geq 0$.

↳ To propagate the interval to $x(t + dt)$, we need to bound $A(\theta)x(t)$.

↳ Why not use interval arithmetics?

*Inría*

Assume that $\underline{x}(t) \le x(t) \le \overline{x}(t)$, for some $t \ge 0$.

↳ To propagate the interval to $x(t + dt)$, we need to bound $A(\theta)x(t)$.

↳ Why not use interval arithmetics?

**Lemma (Image of an interval (Efimov et al., 2012))**

*If $A$ a known matrix, then*

$$A^+\underline{x} - A^-\overline{x} \le Ax \le A^+\overline{x} - A^-\underline{x}.$$

*where $A^+ = \max(A, 0)$ and $A^- = A - A^+$.*

*Inria*

Assume that $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$, for some $t \geq 0$.

↳ To propagate the interval to $x(t + dt)$, we need to bound $A(\theta)x(t)$.

↳ Why not use interval arithmetics?

**Lemma (Product of intervals (Efimov et al., 2012))**

If $A$ is unknown but bounded $\underline{A} \leq A \leq \overline{A}$,

$$\underline{A}^+\underline{x}^+ - \overline{A}^+\underline{x}^- - \underline{A}^-\overline{x}^+ + \overline{A}^-\overline{x}^- \leq Ax$$
$$\leq \overline{A}^+\overline{x}^+ - \underline{A}^+\overline{x}^- - \overline{A}^-\underline{x}^+ + \underline{A}^-\underline{x}^-.$$

*Inria*

Assume that $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$, for some $t \geq 0$.

↳ To propagate the interval to $x(t + dt)$, we need to bound $A(\theta)x(t)$.

↳ Why not use interval arithmetics?

---

**Lemma (Product of intervals (Efimov et al., 2012))**

*If $A$ is unknown but bounded $\underline{A} \leq A \leq \overline{A}$,*

$$\underline{A}^+\underline{x}^+ - \overline{A}^+\underline{x}^- - \underline{A}^-\overline{x}^+ + \overline{A}^-\overline{x}^- \leq Ax$$
$$\leq \overline{A}^+\overline{x}^+ - \underline{A}^+\overline{x}^- - \overline{A}^-\underline{x}^+ + \underline{A}^-\underline{x}^-.$$

---

✓ Since $A(\theta)$ belongs to a known $C_\delta$, we can easily compute such bounds $\underline{A} \leq A(\theta) \leq \overline{A}$

*Inría*

Following this result, define the predictor:

$$
\begin{aligned}
\dot{\underline{x}}(t) &= \underline{A}^+\underline{x}^+(t) - \overline{A}^+\underline{x}^-(t) - \underline{A}^-\overline{x}^+(t) \\
&\quad + \overline{A}^-\overline{x}^-(t) + B^+\underline{d}(t) - B^-\overline{d}(t), \quad\quad (1) \\
\dot{\overline{x}}(t) &= \overline{A}^+\overline{x}^+(t) - \underline{A}^+\overline{x}^-(t) - \overline{A}^-\underline{x}^+(t) \\
&\quad + \underline{A}^-\underline{x}^-(t) + B^+\overline{d}(t) - B^-\underline{d}(t), \\
\underline{x}(0) &= \underline{x}_0,\ \overline{x}(0) = \overline{x}_0,
\end{aligned}
$$

*Inria*

Following this result, define the predictor:

$$
\begin{aligned}
\dot{\underline{x}}(t) &= \underline{A}^+\underline{x}^+(t) - \overline{A}^+\underline{x}^-(t) - \underline{A}^-\overline{x}^+(t) \\
&\quad + \overline{A}^-\overline{x}^-(t) + B^+\underline{d}(t) - B^-\overline{d}(t), \\
\dot{\overline{x}}(t) &= \overline{A}^+\overline{x}^+(t) - \underline{A}^+\overline{x}^-(t) - \overline{A}^-\underline{x}^+(t) \\
&\quad + \underline{A}^-\underline{x}^-(t) + B^+\overline{d}(t) - B^-\underline{d}(t), \\
\underline{x}(0) &= \underline{x}_0, \ \overline{x}(0) = \overline{x}_0,
\end{aligned}
\tag{1}
$$

**Proposition (Inclusion property)**

✓ *The predictor* (1) *satisfies* $\underline{x}(t) \leq x(t) \leq \overline{x}(t)(t)$

*Inría*

Following this result, define the predictor:

$$
\begin{aligned}
\dot{\underline{x}}(t) &= \underline{A}^+ \underline{x}^+(t) - \overline{A}^+ \underline{x}^-(t) - \underline{A}^- \overline{x}^+(t) \\
&\quad + \overline{A}^- \overline{x}^-(t) + B^+ \underline{d}(t) - B^- \overline{d}(t), \\
\dot{\overline{x}}(t) &= \overline{A}^+ \overline{x}^+(t) - \underline{A}^+ \overline{x}^-(t) - \overline{A}^- \underline{x}^+(t) \\
&\quad + \underline{A}^- \underline{x}^-(t) + B^+ \overline{d}(t) - B^- \underline{d}(t), \\
\underline{x}(0) &= \underline{x}_0, \ \overline{x}(0) = \overline{x}_0,
\end{aligned}
\tag{1}
$$

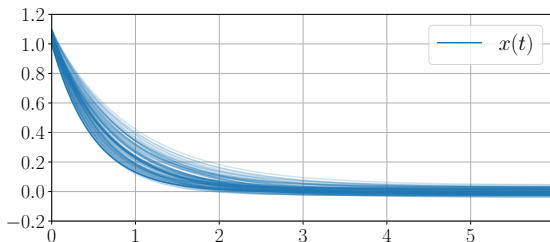**Proposition (Inclusion property)**

✓ *The predictor (1) satisfies* $\underline{x}(t) \leq x(t) \leq \overline{x}(t)(t)$

? But is it stable?

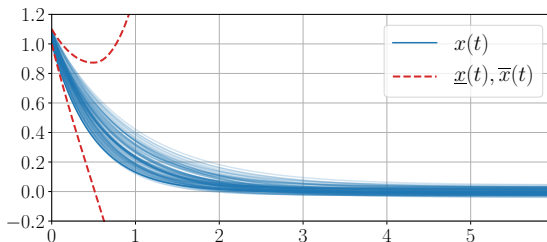*Inria*

Consider the scalar system, for all $t \geq 0$:

$$\dot{x}(t) = -\theta(t)x(t) + d(t), \text{ where } \begin{cases} x(0) \in [\underline{x_0}, \overline{x}_0] = [1.0, 1.1], \\ \theta(t) \in \Theta = [\underline{\theta}, \overline{\theta}] = [1, 2], \\ d(t) \in [\underline{d}, \overline{d}] = [-0.1, 0.1], \end{cases}$$



✓ The system is always stable

*Inria*

Consider the scalar system, for all $t \geq 0$:

$$\dot{x}(t) = -\theta(t)x(t) + d(t), \text{ where } \begin{cases} x(0) \in [\underline{x}_0, \overline{x}_0] = [1.0, 1.1], \\ \theta(t) \in \Theta = [\underline{\theta}, \overline{\theta}] = [1, 2], \\ d(t) \in [\underline{d}, \overline{d}] = [-0.1, 0.1], \end{cases}$$

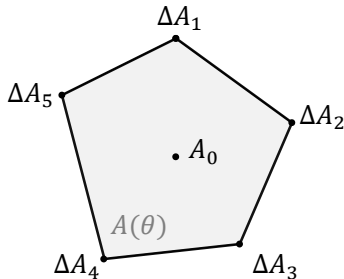

✓ The system is always stable    ✗ The predictor (1) is unstable

## Assumption (Polytopic Structure)

There exist $A_0$ *Metzler* and $\Delta A_0, \cdots, \Delta A_N$ such that:
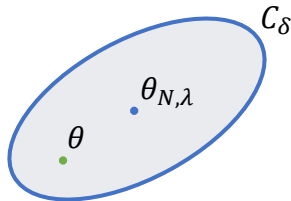
$$A(\theta) = \underbrace{A_0}_{\substack{\text{Nominal} \\ \text{dynamics}}} + \sum_{i=1}^{N} \lambda_i(\theta)\Delta A_i, \quad \sum_{i=1}^{N} \underbrace{\lambda_i(\theta)}_{\geq 0} = 1; \quad \forall \theta \in \Theta$$

*Inria*

## Assumption (Polytopic Structure)

There exist $A_0$ *Metzler* and $\Delta A_0, \cdots, \Delta A_N$ such that:

$$A(\theta) = \underbrace{A_0}_{\substack{\text{Nominal} \\ \text{dynamics}}} + \sum_{i=1}^{N} \lambda_i(\theta) \Delta A_i, \quad \sum_{i=1}^{N} \underbrace{\lambda_i(\theta)}_{\geq 0} = 1; \quad \forall \theta \in \Theta$$

## Assumption (Polytopic Structure)

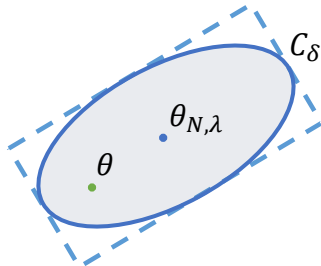There exist $A_0$ *Metzler* and $\Delta A_0, \cdots, \Delta A_N$ such that:

$$A(\theta) = \underbrace{A_0}_{\substack{\text{Nominal} \\ \text{dynamics}}} + \sum_{i=1}^{N} \lambda_i(\theta)\Delta A_i, \quad \sum_{i=1}^{N} \underbrace{\lambda_i(\theta)}_{\geq 0} = 1; \quad \forall \theta \in \Theta$$

Denote

$$\Delta A_+ = \sum_{i=1}^{N} \Delta A_i^+, \ \Delta A_- = \sum_{i=1}^{N} \Delta A_i^-,$$

We define the predictor

$$
\begin{aligned}
\dot{\underline{x}}(t) &= A_0 \underline{x}(t) - \Delta A_+ \underline{x}^-(t) - \Delta A_- \overline{x}^+(t) \\
&\quad + B^+ \underline{d}(t) - B^- \overline{d}(t), \\
\dot{\overline{x}}(t) &= A_0 \overline{x}(t) + \Delta A_+ \overline{x}^+(t) + \Delta A_- \underline{x}^-(t) \\
&\quad + B^+ \overline{d}(t) - B^- \underline{d}(t), \\
\underline{x}(0) &= \underline{x}_0, \ \overline{x}(0) = \overline{x}_0
\end{aligned}
\tag{2}
$$

**Theorem (Inclusion property)**

*The predictor* (2) *ensures* $\underline{x}(t) \leq x(t) \leq \overline{x}(t)$.
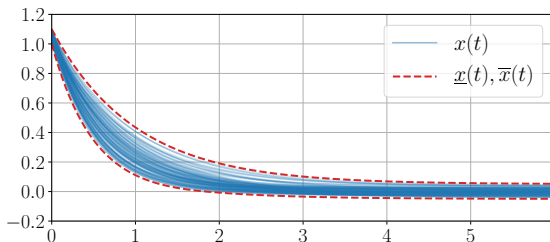
*Inría*

**Theorem (Stability)**

*If there exist diagonal matrices $P$, $Q$, $Q_+$, $Q_-$, $Z_+$, $Z_-$, $\Psi_+$, $\Psi_-$, $\Psi$, $\Gamma \in \mathbb{R}^{2n \times 2n}$ such that the following LMIs are satisfied:*

$$P + \min\{Z_+, Z_-\} > 0, \ \Upsilon \preceq 0, \ \Gamma > 0,$$
$$Q + \min\{Q_+, Q_-\} + 2\min\{\Psi_+, \Psi_-\} > 0,$$

*where $\Upsilon = \Upsilon(A_0, \Delta A_-, \Delta A_+, \Psi_-, \Psi_+, \Psi)$,*
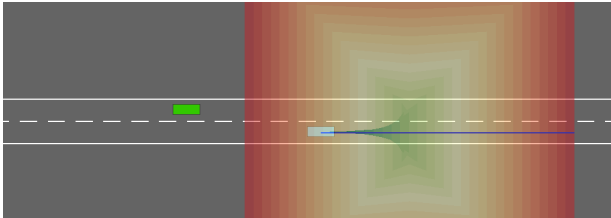*then the predictor* (2) *is input-to-state stable with respect to the inputs $\underline{d}$, $\overline{d}$.*

*Inria*

Recall the scalar system:

$$\dot{x}(t) = -\theta(t)x(t) + d(t), \text{ where } \begin{cases} x(0) \in [\underline{x_0}, \overline{x}_0] = [1.0, 1.1], \\ \theta(t) \in \Theta = [\underline{\theta}, \overline{\theta}] = [1, 2], \\ d(t) \in [\underline{d}, \overline{d}] = [-0.1, 0.1], \end{cases}$$
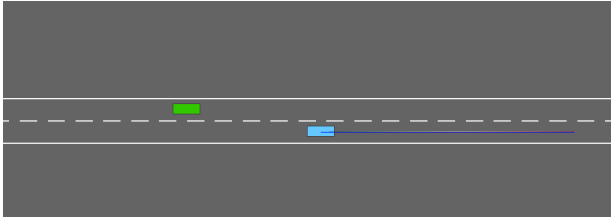


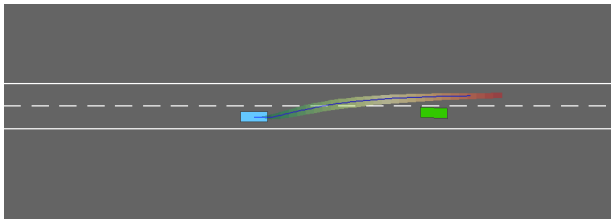✓ The system is always stable    ✓ The predictor (2) is stable
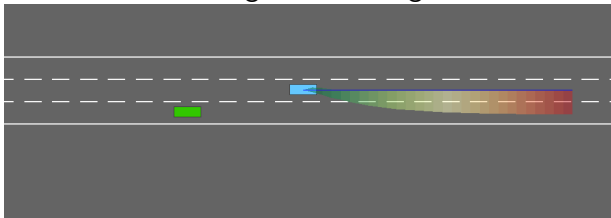
*Inría*

The naive predictor (1) quickly diverges



The proposed predictor (2) remains stable

*Inría*

Prediction during a lane change maneuver



Prediction with uncertainty in the followed lane $L_i$

**Approximate** the robust objective by a tractable surrogate.

Definition (Robust objective $v^r$)

$$v^r(\pi) \stackrel{\text{def}}{=} \min_{A(\theta) \in C_\delta} \sum_{t=0}^{H} \gamma^t R(x_t, \pi(x_t)) \tag{3}$$

Definition (Surrogate objective $\hat{v}^r$)

$$\hat{v}^r(\pi) \stackrel{\text{def}}{=} \sum_{t=0}^{H} \gamma^t \min_{[x \in \underline{x}(t), \overline{x}(t)]} R(x, \pi(x)) \tag{4}$$
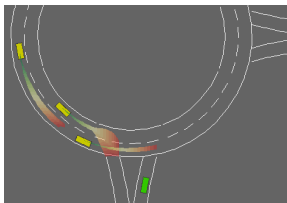
*Inria*

The approximate performance of a policy is guaranteed on the true environment.

Proposition (Lower bound)

*The surrogate objective $\hat{v}^r$ is a lower bound of the true objective $v^r$:*

$$\forall \pi, \; \hat{v}^r(\pi) \leq v^r(\pi) \tag{5}$$

*Inría*

| Ambiguity | Agent | Worst-case | Mean $\pm$ std |
|-----------|-------|------------|----------------|
| None | Oracle | 9.83 | $10.84 \pm 0.16$ |
| Continuous | Nominal | 1.99 | $9.95 \pm 2.38$ |
| | Robust | 7.88 | $10.73 \pm 0.61$ |

*Inría*

Our linear structure assumption is wrong?

Model Adequacy: you can detect it with statistical tests

*Inria*

Our linear structure assumption is wrong?

Model Adequacy: you can detect it with statistical tests

Solution: Multi-Model Prediction

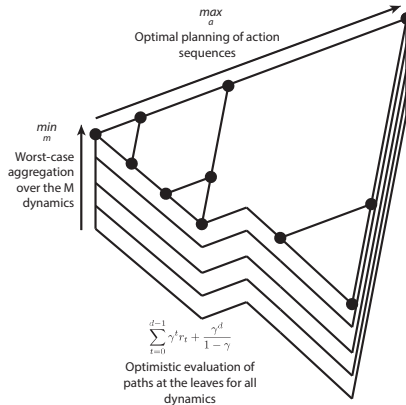Use many linear models with different features. For instance:

- Lane-dependent features
- Neural network features
- Random features

↳ Maintain a set of admissible experts

↳ Perform robust aggregation

*Inria*

## Assumption (Discrete Ambiguity Set)

$$T \in \{T_1, \cdots, T_m\}$$

## Definition (Robust sequence value upper-bound)

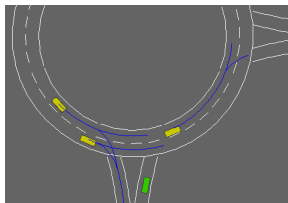Given node $i \in \mathcal{T}$, define the robust B-value:

$$B_i^r(n) \stackrel{\text{def}}{=} \begin{cases} \min_{m \in [1,M]} \sum_{t=0}^{d-1} \gamma^t r_t + \frac{\gamma^d}{1-\gamma} & \text{if } i \in \mathcal{L}_n \text{ ;} \\ \max_{a \in \mathcal{A}} b_{ia}^r(n) & \text{if } i \in \mathcal{T}_n \setminus \mathcal{L}_n \end{cases}$$

## Theorem (Regret bound)

*The corresponding planning algorithm enjoys a simple regret of:*

$$\text{If } \kappa > 1, \qquad \mathcal{R}_n = O\left( n^{-\frac{\log 1/\gamma}{\log \kappa}} \right) \tag{6}$$

*Inria*

| Ambiguity | Agent | Worst-case | Mean $\pm$ std |
|-----------|-------|------------|----------------|
| None | Oracle | 9.83 | $10.84 \pm 0.16$ |
| Continuous | Nominal | 1.99 | $9.95 \pm 2.38$ |
| | Robust | 7.88 | $10.73 \pm 0.61$ |
| Discrete | Nominal | 2.09 | $8.85 \pm 3.53$ |
| | Robust | 8.99 | $10.78 \pm 0.34$ |

*Inría*

Decision-making among interacting drivers with behavioural uncertainty

**Model-free**

1. Self-attention model for permutation invariance and variable size

2. Budgeted reinforcement learning to constrain the expected risk

**Model-based**

3. Efficient tree-based planning with tight statistical bounds

4. Tackle the issue of model bias

   ↳ Build a confidence region around the true model
   ↳ Design a stable interval predictor
   ↳ Perform robust control with respect to this uncertainty

*Inría*

Thank You!

*Inria*