# On Time with Minimal Expected Cost !
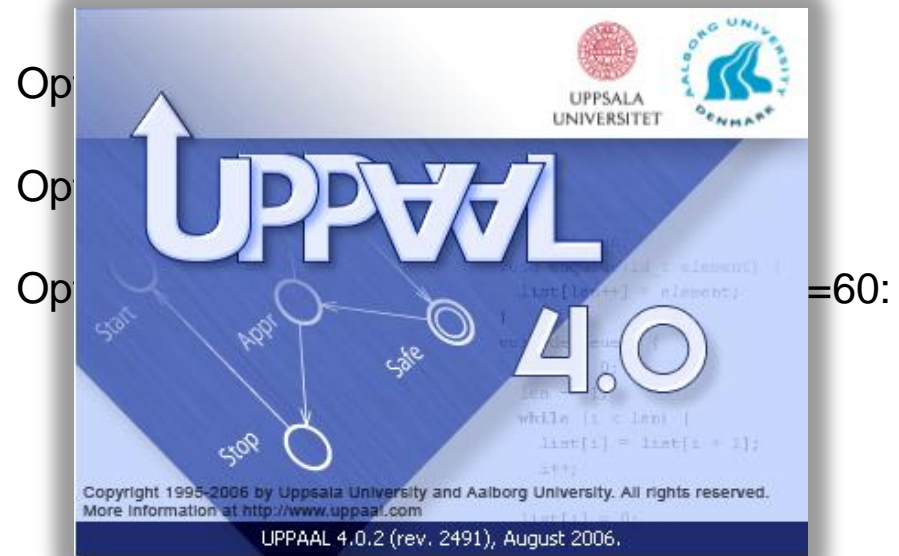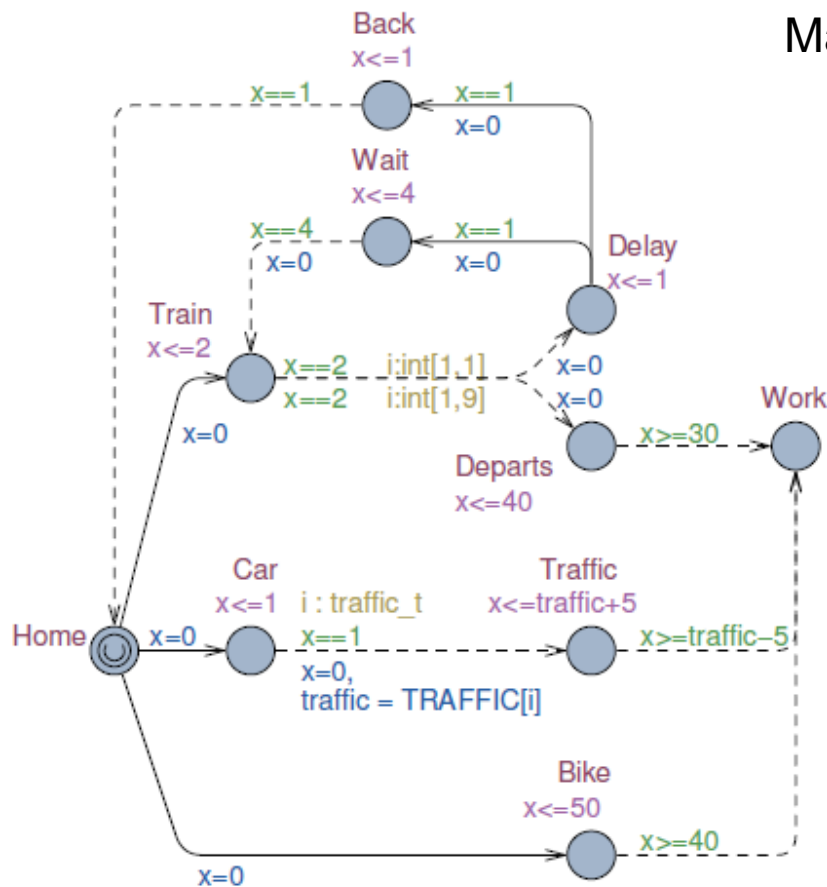
**Alexandre David, Peter G Jensen,**
**Kim G Larsen, Axel Legay, Didier Lime,**
**Mathias G Sørensen, Jakob H Taankvist**

CISS – Aalborg University

DENMARK

# Motivation



2-Player Game (Antagonistic opponent)
Markov Decision Process (probabilistic opponent)

Op

Op

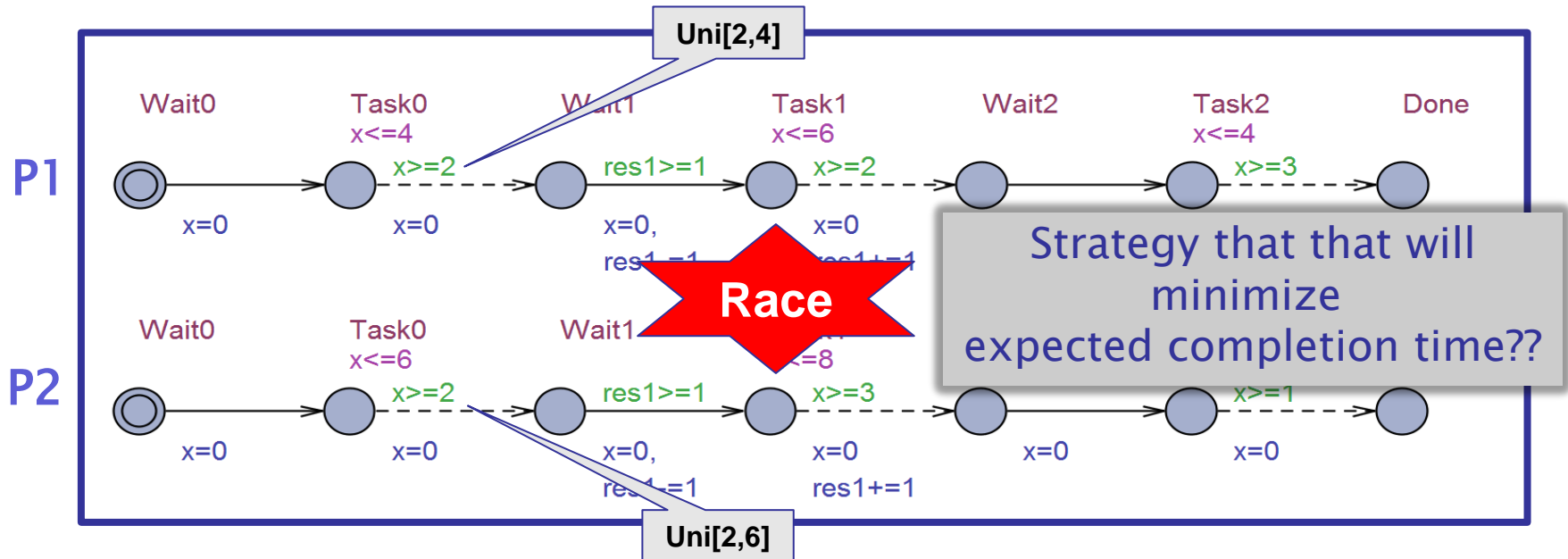Op                                                                    =60:
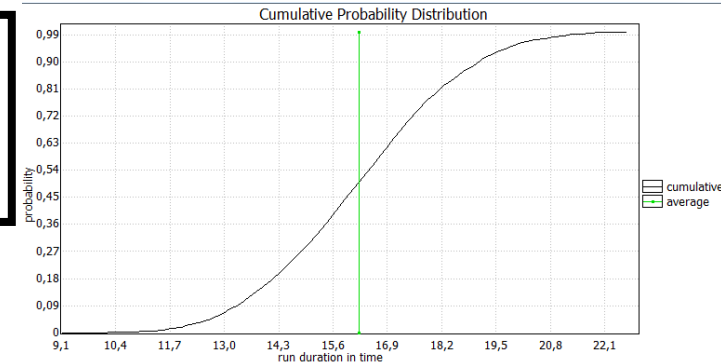
Bruyere, V., Filiot, E., Randour, M., Raskin, J.F.: Meet your expectations with guarantees: Beyond worst-case synthesis in quantitative games. STACS14
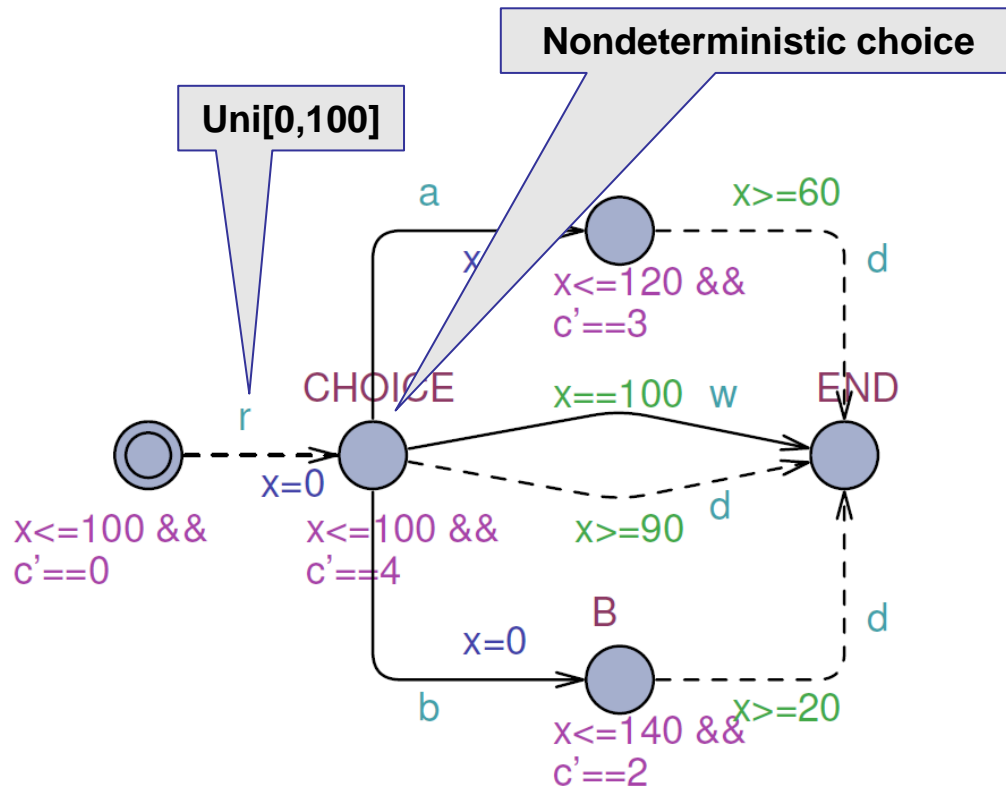
# Duration Probabilistic Automata



```
/* Resources */  res1:1;
/* Processes */  P1: [2,4].<res1:1>[2,6].[3,4];
                 P2: [2,6].<res1:1>[3,8].[1,5];
```

Pr[<=1000](<> P1.Done && P2.Done)



Kempf, J.F., Bozga, M., Maler, O.: As soon as probable: Optimal scheduling under stochastic uncertainty. In: TACAS. pp. 385{400 (2013)

# Motivation



Minimize **expected cost** subject to **guaranteed time-bound**

$\times$ Priced Timed Game

$\checkmark$ Timed Game    TIGA

$\checkmark$ Timed Automata    UPPAAL

$\checkmark$ Priced Timed Automata    CORA

$\checkmark$ Stochastic (P)TA    SMC

$\times$ **Priced Timed MDP**    TIGA/SMC
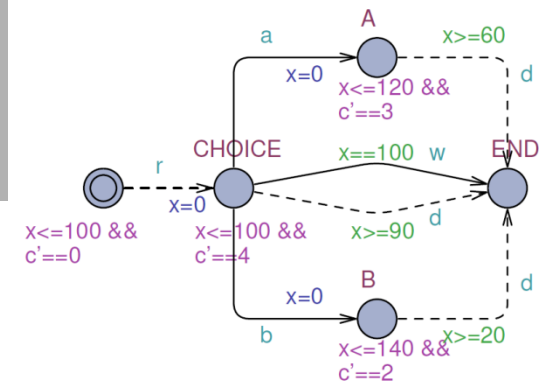~ **Decision Stochastic Priced Timed Automata**

# Overview

- ## Priced Timed Games
  - Time bounded reachability strategies
- ## Priced Timed Markov Decision Processes
  - Minimal expected cost reachability strategy

- ## Optimal Strategy Synthesis Using Reinforcement Learning
- ## Representation of Stochastic Strategies

- ## Experimental Results

# PTA and PTG



$$\mathcal{A} = (L, \ell_0, X, \Sigma, E, P, Inv)$$

is a tuple where

- $L$ is a finite set of locations,
- $\ell_0 \in L$ is the initial location,
- $X$ is a finite set of non-negative real-valued clocks,
- $\Sigma$ is a finite set of actions,
- $E \subseteq L \times \mathcal{B}(X) \times \Sigma \times 2^X \times L$ is a finite set of edges,
- $P : L \to \mathbb{N}$ assigns a price-rate to each location, and
- $Inv : L \to \mathcal{B}(X)$ sets an invariant for each location.

Priced Timed Game $\Sigma = \Sigma_c \uplus \Sigma_u$

# PTA Semantics

$$S_{\mathcal{A}} = (Q, q_0, \Sigma, \rightarrow)$$

- $(\ell, v) \in Q$ for $\ell \in L$ and $v \in \mathbb{R}_{\geq 0}^X$ st $v \models Inv(\ell)\}$,
- $q_0 = (\ell_0, 0)$ is the in

and[1]

- $(\ell, v) \xrightarrow{a}_0 (\ell', v')$ if
  $(\ell \xrightarrow{g,a,r} \ell') \in E$ st.

- $(\ell, v) \xrightarrow{d}_p (\ell, v + d$
  $p = P(\ell) \cdot d$, $v \models Inv(\ell)$ and $v + d \models Inv(\ell)$.

> - Set of runs of: $Exec_{\mathcal{A}}$.
> - Set of finite (maximal) runs: $Exec_{\mathcal{A}}^f$ ($Exec_{\mathcal{A}}^m$).
> - $\pi[i]$ the state $q_i$,
> - $\pi|_i$ ($\pi|^i$) the prefix (suffix) of $\pi$ ending (starting) at $q_i$.
> - $C(\pi)$ ($T(\pi)$) denotes total accumulated cost (time).

### Run $\pi$:

$$q_0 \xrightarrow{d_0}_{p_0} q_0' \xrightarrow{a_0}_0 q_1 \xrightarrow{d_1}_{p_1} q_1' \xrightarrow{a_1}_0 \cdots \xrightarrow{d_{n-1}}_{p_{n-1}} q_{n-1}' \xrightarrow{a_{n-1}}_0 q_n \cdots$$

$a_i \in \Sigma$, $d_i, p_i \in \mathbb{R}_{\geq 0}$, and $q_i$ is a state $(\ell_{q_i}, v_{q_i})$.

# Strategies & Outcome

Priced Timed Game $\Sigma = \Sigma_c \uplus \Sigma_u$

$$\sigma : \mathit{Exec}^f_{\mathcal{G}} \rightharpoonup \mathcal{P}\left(\Sigma_c \cup \{\lambda\}\right) \setminus \{\emptyset\}$$

such that for any finite run $\pi$, if $q = \mathit{last}(\pi)$ and $a \in \sigma(\pi) \cap \Sigma_c$, then $q \xrightarrow{a} q'$ fs $q'$.

## $\mathit{Out}(\sigma) \subseteq \mathit{Exec}_{\mathcal{G}}$

- $q_0 \in \mathit{Out}(\sigma)$
- If $\pi \in \mathit{Out}(\sigma)$ then $\pi' = \pi \xrightarrow{e} q' \in \mathit{Out}(\sigma)$ if $\pi' = \mathit{Exec}_{\mathcal{G}}$ and either one of the following three conditions hold:
  1. $e \in \Sigma_u$, or
  2. $e \in \Sigma_c$ and $e \in \sigma(\pi)$, or
  3. $e \in \mathbb{R}_{>0}$ and for all $e' < e$, $\mathit{last}(\pi) \xrightarrow{e'} q'$ for some $q'$ st $\sigma(\pi \xrightarrow{e'} q') \ni \lambda$.

# Cost Bounded Reachability Strategies

For $G \subseteq L$, $B \in \mathbb{R}_{\geq 0}$: $(G, B)$ is a cost-bounded reachability objective.

$\pi$ is winning w.r.t. $(G, B)$, if $last(\pi) \in G \times \mathbb{R}_{\geq 0}^X$ and $C(\pi) \leq B$. A strategy $\sigma$ over $\mathcal{G}$ is a winning strategy if all runs in $Out(\sigma)$ are winning.
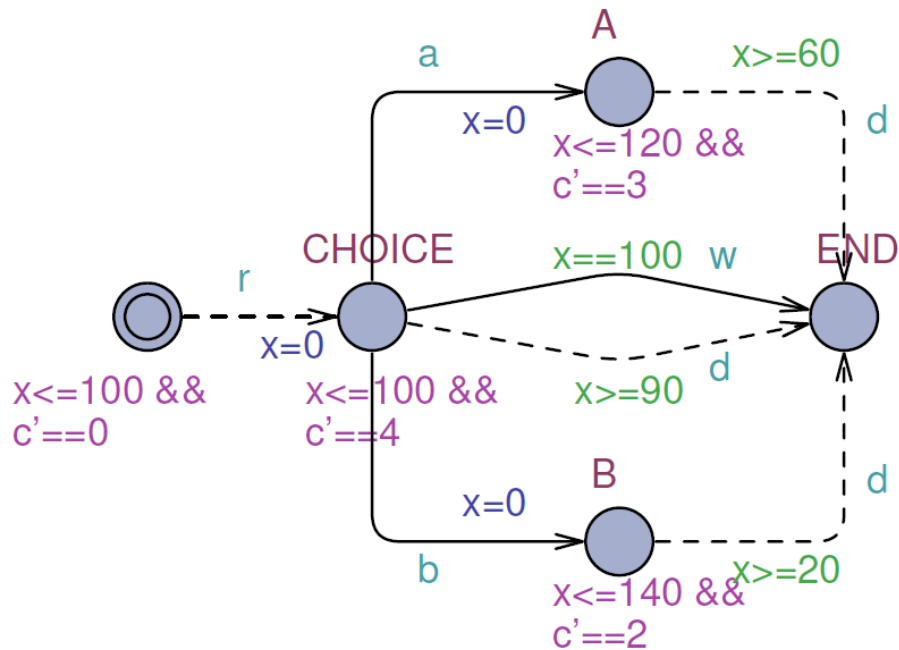
## Theorem (Memoryless, Most Permissive Strategies)

Let $\mathcal{G}$ be a non-Zeno, clocked TG. If a time-bounded reachability objective $(G, T)$ has a winning strategy, then it has

1. a deterministic, memoryless winning strategies, and
2. a (unique) most permissive, memoryless winning strategy $\sigma_{\mathcal{G}}^p(G, T)$.

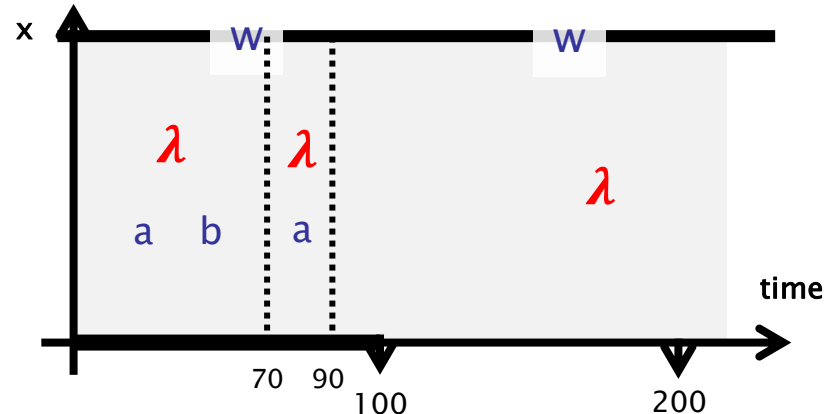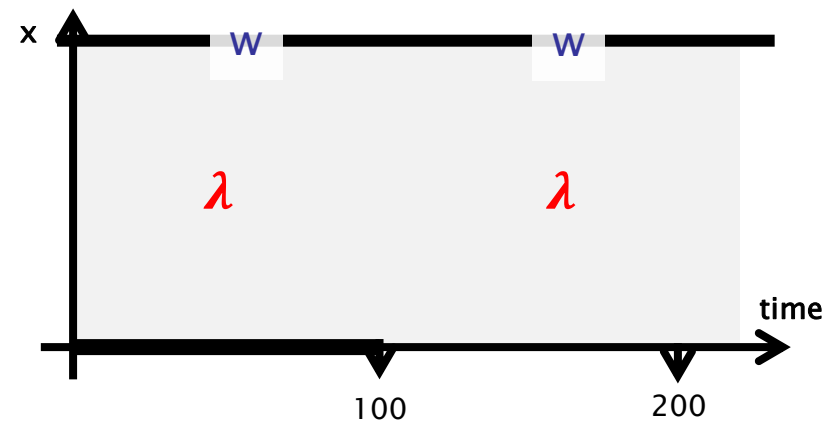# Motivation



**Objective**: $A\langle\rangle(\textbf{END} \wedge \textbf{time} \leq 210)$

**Deterministic, memoryless strategy:**

**Most permissive, memoryless strategy**

# Priced Timed MDPs

$$\mathcal{M} = \langle \mathcal{G}, \mu^u \rangle$$

where

- $\mathcal{G} = (L, \ell_0, X, \Sigma_c, \Sigma_u, E, P, Inv)$ is a PTG, and
- $\mu^u$ is a family of density-functions, $\{\mu_q^u : \exists \ell \exists v . q = (\ell, v)\}$, with $\mu_q^u(d, u) \in \mathbb{R}_{\geq 0}$ assigning the *density* of the environment aiming at taking the uncontrollable action $u \in \Sigma_u$ after a delay of $d$ from state $q$.

Assumptions:

1. $\mu_q^u(d, u) > 0$ only if $q \xrightarrow{d,u}$ in $\mathcal{G}$.
2. $\sum_u (\int_{t \geq 0} \mu_q^u(t, u) dt) = 1$

# Stochastic Strategies

$$\mu^c \text{ for a PTMDP } \mathcal{M} = \langle \mathcal{G}, \mu^u \rangle$$

is a family of density-functions, $\mu^c = \{\mu^c_q : \exists \ell \exists v.q = (\ell, v)\}$, where $\mu^c_q(d, c) \in \mathbb{R}_{\geq 0}$ assigns the *density* of the controller aiming at taking the controllable action $c \in \Sigma_c$ after a delay of $d$ from state $q$.

- Repeated races between $\mu^u$ and $\mu^c$,
- Induced probability measure $\mathbb{P}_{\langle \mathcal{G}, \mu^u \rangle, \mu^c}$ on (certain) sets of runs.

# Induced Probability Measure

Cylinder set $\mathcal{C}(q, I_0\ell_0 I_1 \cdots I_n\ell_n)$ with $\ell_i \in L$ and $I_i = [l_i, u_i]$ with $l_i, u_i \in \mathbb{Q}$, $i = 0..n$, consists of all maximal runs having a prefix of the form:

$$q \xrightarrow{d_0}\xrightarrow{a_0} (\ell_0, v_0) \xrightarrow{d_1}\xrightarrow{a_1} \cdots \xrightarrow{d_n}\xrightarrow{a_n} (\ell_n, v_n)$$

where $d_i \in I_i$ for all $i < n$.

## Probability Measure

$$\mathbb{P}_{\langle \mathcal{G}, \mu^u \rangle, \mu^c} \left( \mathcal{C}(q, I_0\ell_0 I_1\ell_1 \cdots I_{n-1}\ell_n) \right) =$$

$$\sum_{\substack{p \in \{u,c\} \\ }} \sum_{\substack{a \in \Sigma_p \\ \ell_q \xrightarrow{a} \ell_1}} \int_{t \in I_0} \mu_q^p(t, a) \cdot \left( \int_{\tau > t} \mu_q^{\overline{p}}(\tau) d\tau \right) \cdot$$

$$\mathbb{P}_{\langle \mathcal{G}, \mu^u \rangle, \mu^c} \left( \mathcal{C}((q^t)^a, \mathcal{C}(I_1 \cdots I_{n-1}\ell_n)) \right) dt$$

where $\mu_q^p(\tau) = \sum_{a \in \Sigma_p} \mu_q^p(\tau, a)$.

# Minimum Expected Cost

Let $\pi \in Exec^m$ and let $G$ be as set of goal locations.

$$C_G(\pi) = min\{C(\pi|_i) : \pi[i] \in G\}$$

denotes the accumulated cost before $\pi$ reaches $G$.

**Expected Value of $C_G$ given $\mu^c$:**

$$\mathbb{E}_{\mu^c}^{\langle \mathcal{G}, \mu^u \rangle}(C_G) = \int_{\pi \in Exec^m} C_G(\pi) \mathbb{P}_{\langle \mathcal{G}, \mu^u \rangle, \mu^c}(d\pi)$$
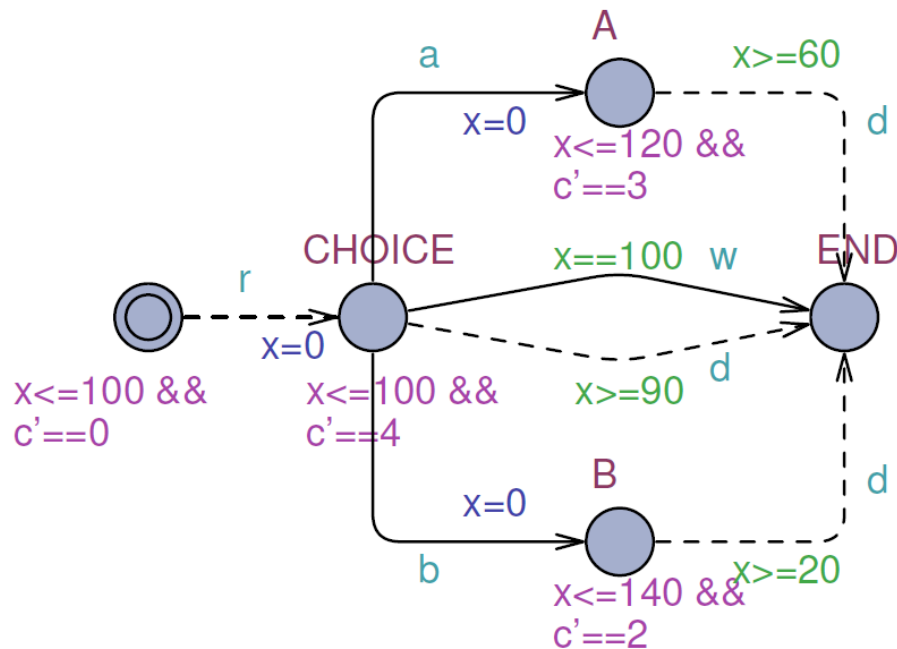
**Optimal strategy $\mu^o$**

$$\mathbb{E}_{\mu^o}^{\langle \mathcal{G}, \mu^u \rangle}(C_G) = \inf\left\{ \mathbb{E}_{\mu^c}^{\langle \mathcal{G}, \mu^u \rangle}(C_G) \mid \mu^c \prec \sigma^P(G, T) \right\}$$

where $\sigma^P(G, T)$ is the most permissive $T$ time-bounded reachability strategy.

# Motivation



Minimal Expected Cost
Strategy (0,b)  **2*80**=160

Expected Cost for TIGA
Strategy (100,w)  **4*95**=380

Minimal Expected Cost while guaranteeing END is reached within time 210:

Strat.:    t>90→    (100,w)
           t>70→    (0,b)
           ow  →    (0,a)
           =
           204

# Reinforcement Learning



Time Bounded Reachability $(G,T)$

# Strategies

## Nondeterministic Strategies   (UPPAAL TIGA)

$R_\ell = \{(Z_1, a_1), \ldots, (Z_k, a_k)\}$, where $a_i \in \Sigma_c \cup \{\lambda\}$. Now $R$ represents the strategy $\sigma_R$ where $\sigma_R((\ell, v)) \ni a$ iff $(Z, a) \in R_\ell$ for some $Z$ with $v \in Z$.

## Stochastic Strategies  (non-lazy *)

- Urgent: $\mu^c_{(\ell,v)}(d, a) = 0$ if $d > 0$, or
- Wait: $\mu^c_{(\ell,v)}(d, a) = 0$ whenever $\sigma^P(\ell, v + d) \ni \lambda$.

$$\mu^c_{(\ell,v)} : (\Sigma_c \cup \{w\}) \to [0, 1].$$

Classes allowing for efficient representation and learning

* Non-lazy strategies suffices for DPAs

# Learning

Given a set of runs $\Pi$ the relevant information for the sub-strategy $\mu_\ell^c$ is given as $In_\ell$:

$$In_\ell = \{(s_n, v) \in (\Sigma_c \cup \mathbb{R}) \times \mathbb{R}_{\geq 0}^X \mid (q_0 \xrightarrow{s_0}_{p_0} \dots \xrightarrow{s_{n-1}}_{p_{n-1}} (\ell, v) \xrightarrow{s_n}_{p_n} \dots) \in \Pi\}$$



Simulation of $Uni(\sigma^p)$ for $A\langle\rangle($**END** $\wedge$ **time** $\leq 210)$

wait

a

b

`time(`$\pi$`)@Choice`

$C(\pi)$

- Covariance Matrices

- Logistic Regression



$$f(v) = \frac{1}{1 + e^{-(-1.131 + 0.647v(x))}}$$

- Splitting



Using Learning Determinization

$$\mu^c_{(\ell,v)} : (\Sigma_c \cup \{w\}) \to [0,1].$$

# Experiments

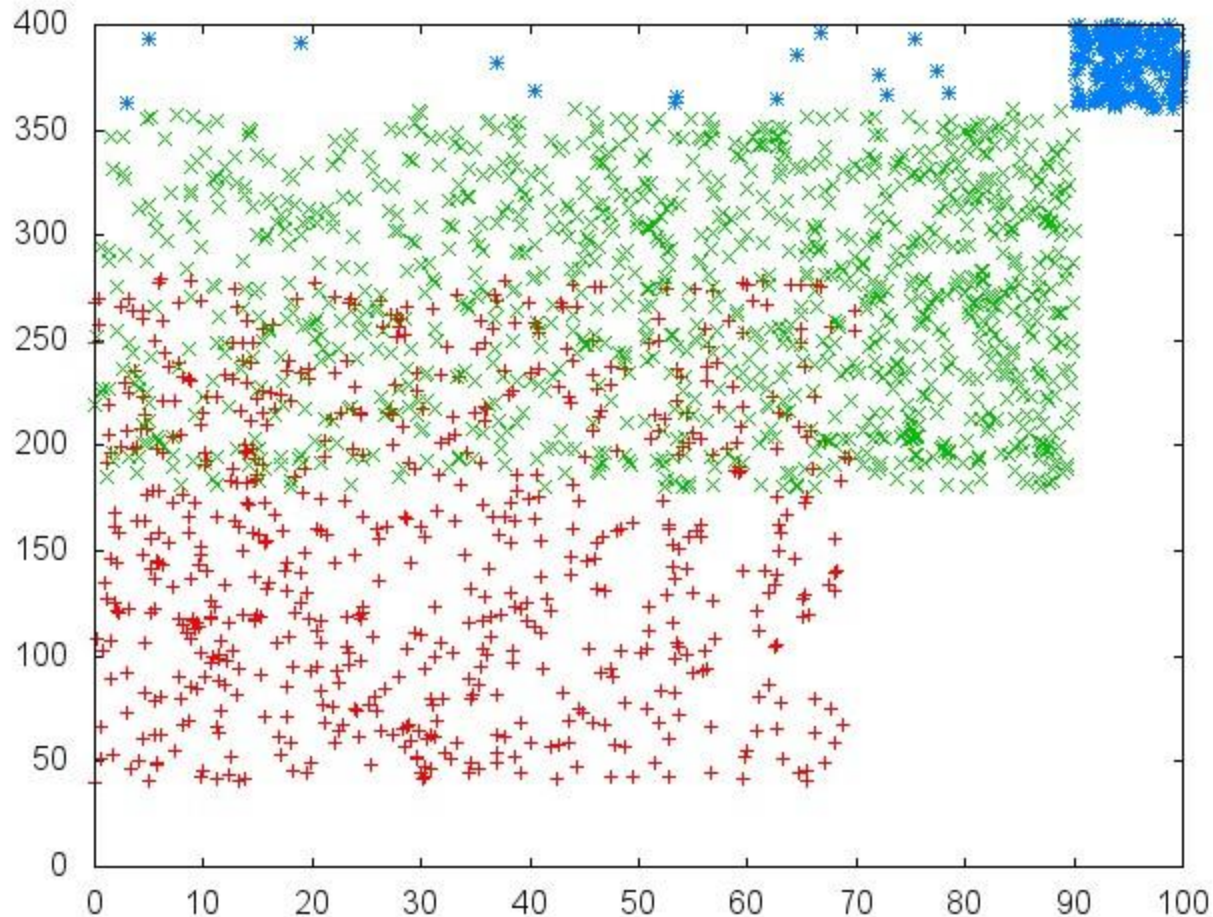| Model | Uniform | Co-variance | Splitting | Regression | Exact [?] |
|---|---|---|---|---|---|
| Motivational example | 410.60 | 200.54 10.57s 6.09MB 0/27 | 204.21 13.16s 6.23MB 0/50 | 200.65 15.27s 6.34MB 0/10 | |
| GoWork | 38.62 | 37.83 16.89s 6.47MB 0/32 | 37.80 12.99s 6.43MB 0/29 | 37.90 19.41 6.56MB 0/9 | |

# Learned Strategies

Covariance

# Learned Strategies

# Learned Strategies

# Learned Strategies
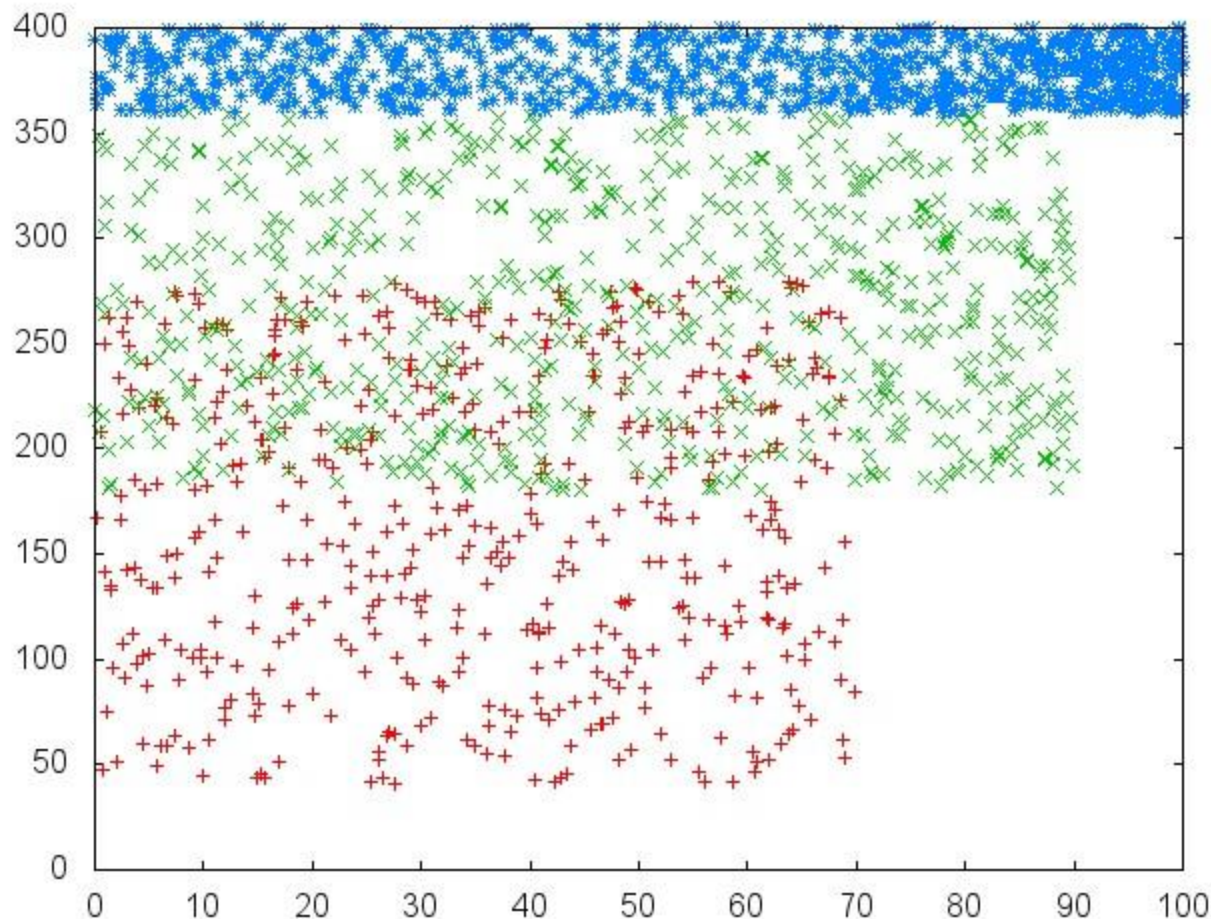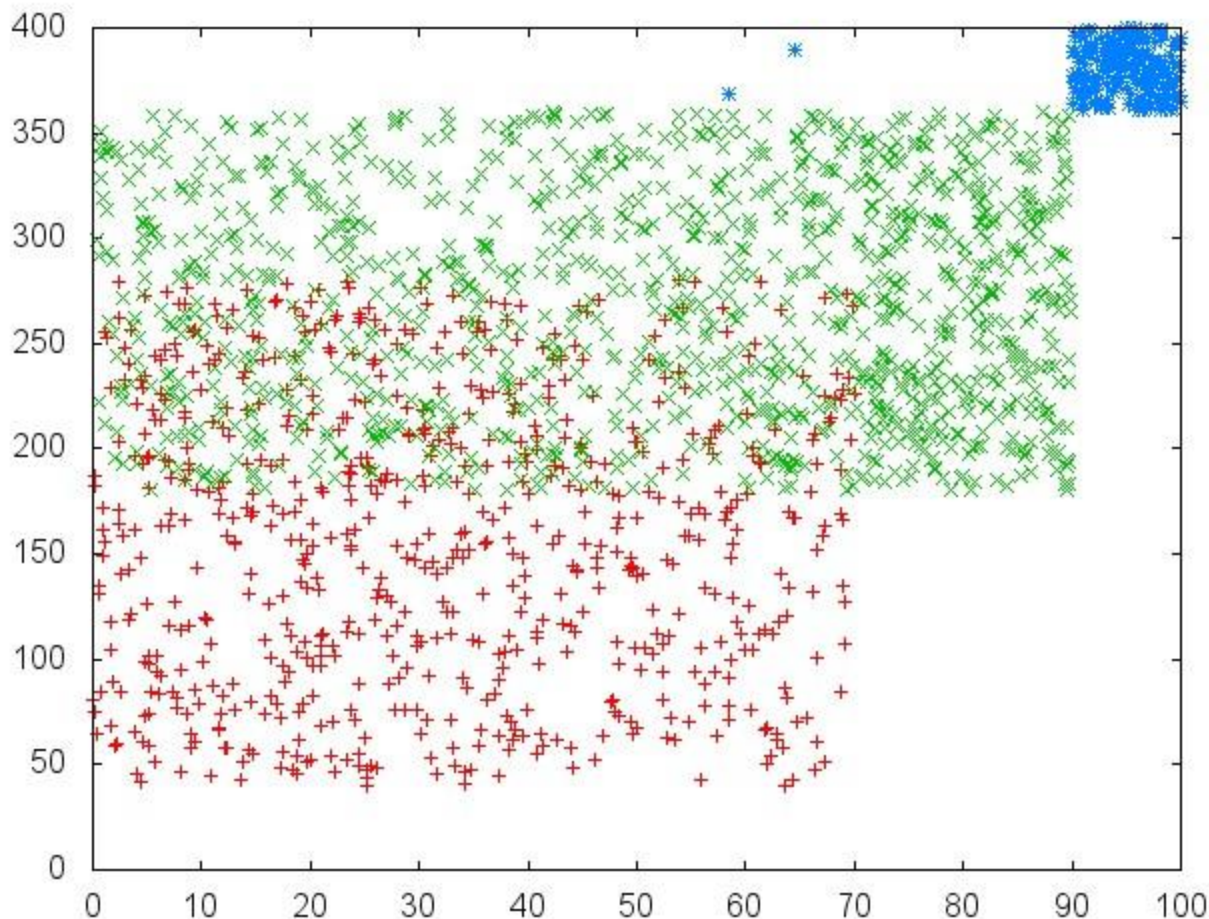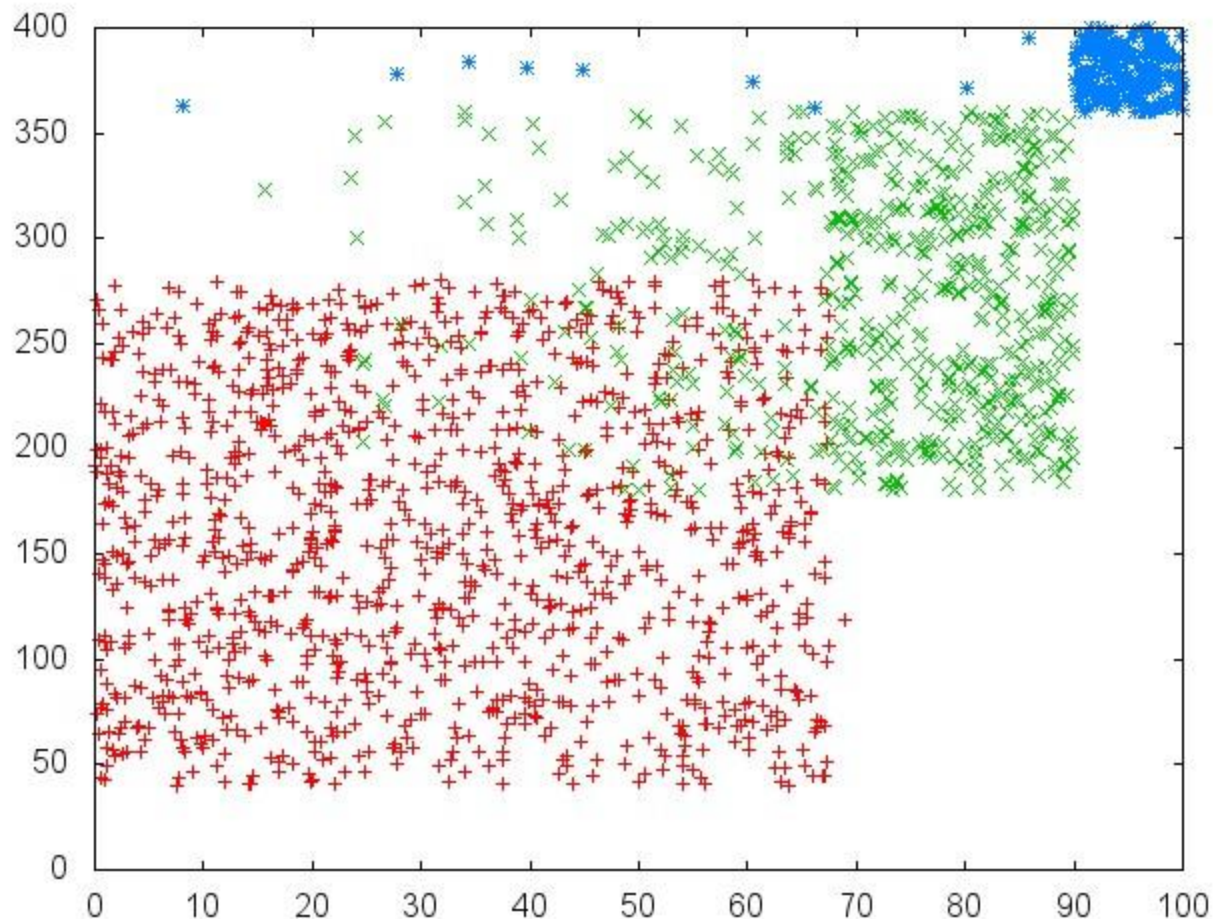
# Learned Strategies

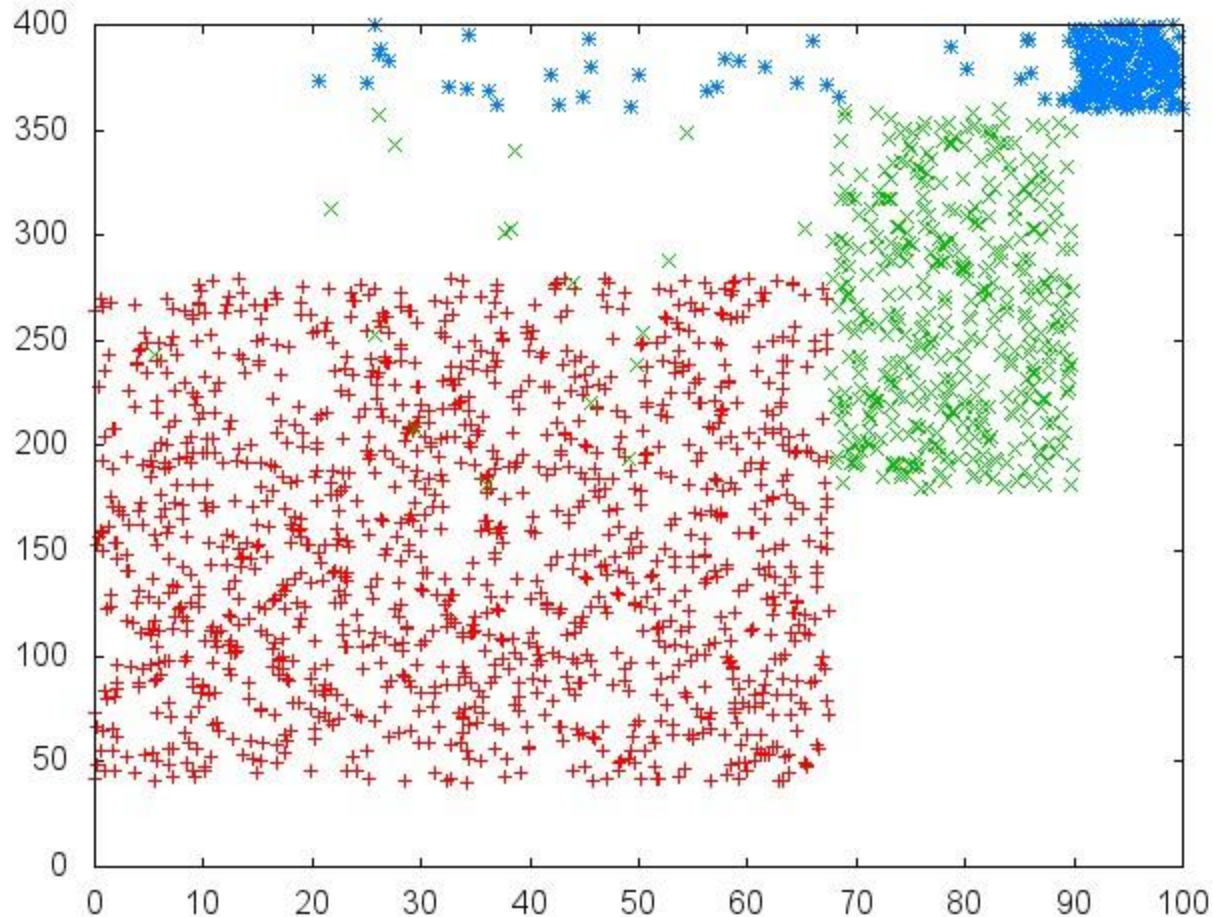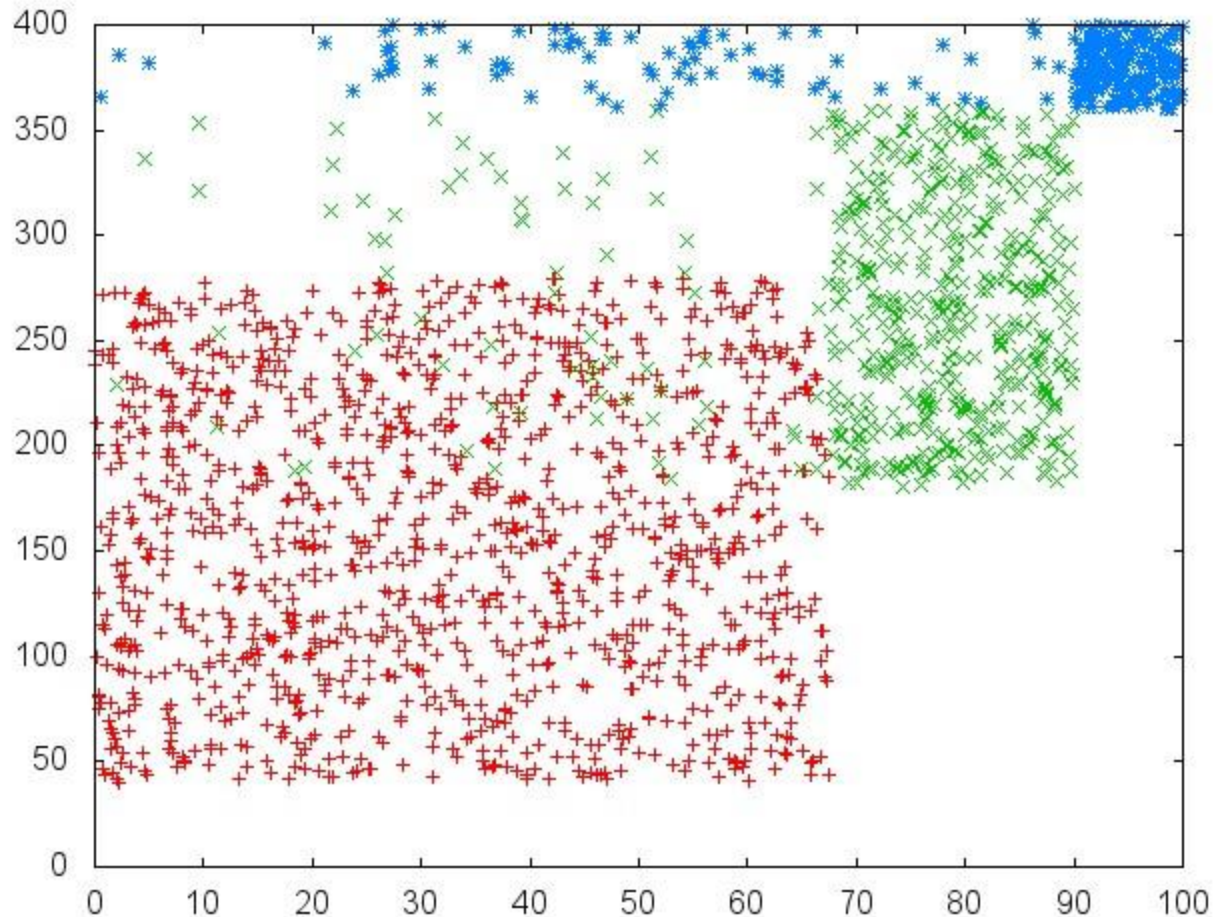# Learned Strategies          Splitting

# Learned Strategies

# Experiments /DPA

| Model | Uniform | Co-variance | Splitting | Regression | Exact [?] |
|---|---|---|---|---|---|
| p0s3p1s4_4 | 18.07 | 17.61 19.31s 6.15MB 2/40 | 17.54 18.28s 6.20MB 0/7 | 17.56 20.87s 6.30MB 2/33 | 1062.77s 145.47MB |
| p0s3p1s4_16 | 18.41 | 17.63 12.13s 6.06MB 1/11 | 17.88 13.21s 6.23MB 2/27 | 17.73 24.27s 6.36MB 1/18 | 176.15s 35.60MB |
| p0s4p1s4_5 | 19.80 | 19.25 20.67s 6.43MB 1/21 | 19.22 21.38s 6.64MB 0/11 | 19.23 29.02s 6.62MB 1/23 | 8547.52s 486.92MB |

Kempf, J.F., Bozga, M., Maler, O.: As soon as probable: Optimal scheduling under stochastic uncertainty. In: TACAS. pp. 385{400 (2013)
http://www-verimag.imag.fr/PROJECTS/TEMPO/DATA/201304_dpa/

# Experiments /DPA Random

| Model | Uniform | Co-variance | Splitting | Regression | Exact [?] |
|---|---|---|---|---|---|
| ran-4-3 | 3944.58 | 2379.90<br>62.63s<br>12.01MB<br>0/10 | 2370.75<br>41.34s<br>13.12MB<br>2/32 | 2346.28<br>74.13s<br>12.23MB<br>1/24 | |
| ran-4-4 | 8092.31 | 5035.81<br>56.97s<br>21.99MB<br>2/33 | 5050.73<br>52.17s<br>22.33MB<br>2/30 | 5029.37<br>112.34s<br>16.60MB<br>2/55 | |
| tiga-ran-4-3 | 3168.30 | 2789.67<br>64.07s<br>13.44MB<br>3/32 | 2778.92<br>71.25s<br>14.64MB<br>2/25 | 2774.52<br>71.48s<br>13.60MB<br>3/31 | |
| tiga-ran-4-4 | 6978.53 | 6358.83<br>124.68s<br>21.31MB<br>1/40 | 6291.49<br>118.67s<br>22.43MB<br>2/43 | 6330.04<br>88.43s<br>18.04MB<br>0/2 | |

# Experiments /DPA Random

| Model | Uniform | Co-variance | Splitting | Regression | Exact [?] |
|---|---|---|---|---|---|
| ran-5-10 | 22030.00 | 15010.20<br>220.93s<br>931.96MB | 13603.70<br>347.84s<br>480.51MB | 14162.10<br>412.31s<br>265.81MB | |
| ran-5-15 | 39569.70 | 29642.20<br>332.06s<br>2042.07MB | 30890.90<br>387.16s<br>804.52MB | 24121.90<br>965.80s<br>1231.08MB | |
| ran-5-3 | 11538.70 | 6109.22<br>52.37s<br>29.45MB | 6305.93<br>72.01s<br>28.69MB | 6118.35<br>116.03s<br>18.34MB | |
| ran-5-4 | 9175.81 | 3888.85<br>97.34s<br>90.61MB | 3796.84<br>92.88s<br>43.18MB | 3697.70<br>135.72s<br>31.00MB | |
| ran-5-5 | 6693.26 | 3766.95<br>122.72s<br>145.17MB | 3515.98<br>151.66s<br>108.07MB | 3570.11<br>207.10s<br>62.79MB | |

# Conclusion & Future Work

- Efficient synthesis of strategies for PTMDP ensuring time-bounds and minimizing expected cost.

- If not time-bound needed we can omit the UPPAAL TIGA synthesis

- Extension to Hybrid MDPs utilizing UPPAAL SMCs support for SHAs.

- Make TIGA/SMC available to you!

- Datastructures supporting general stochastic strategies – not just non-lazy ones.

- More clever filtrations of runs.