



Towards Systematically Engineering Adaptive Systems using Machine-Learning Techniques

Martin Wirsing** * Ludwig-Maximilians-Universität München

**In cooperation withLenz Belzner and Thomas GaborMaiborn & Wolff LMU München

*Thanks to Rolf Hennicker, Alexander Knapp, Matthias Hölzl and all other former members of the ASCENS project

Seminar Dynamische und Adaptive Systeme, WS 2019/20, LMU München





- What is an adaptive system?
- A development lifecycle for adaptive systems and ensembles
- Adapting behaviours by online planning
- Safe learning: Policy synthesis from requirements
- Concluding remarks





1. What is an adaptive system?

Martin Wirsing



Adaptive system







- An adaptive system adjusts its behaviour to multiple situations:
 - change of
 - environment
 - human/other agents behavior
 - network infrastructure
 - goals/requirements
 - the system itself
 - or any combination thereof
- We distinguish between
 - black-box adaptation and
 - white-box adaptation





- Goal satisfaction in an environment For any environment η , system S, goal γ : η , S |= γ
- An adaptive system adjusts to a change of
 - goals/requirements: η , S |= γ '
 - environment: η' , S |= γ
 - the system itself: η , S' |= γ
- Adaptation space
 - Adaptation domain $\mathcal{A} \subseteq Env \times G$
 - S can adapt to \mathcal{A} , written $S \parallel \mathcal{A}$:

 $S \parallel - A$ iff $\forall (\eta, \gamma) \in A$. $\eta, S \mid = \gamma$





Adaptation Modelling: White-Box Adaptation



- White-box adaptation
 - Modify control mechanism at runtime so that the (adaptation) requirements are satisfied
- Approaches
 - Classical adaptation techniques

E.g. programming using modes, reconfiguration, policies

• Al adaptation techniques

E.g. using planning, learning, reasoning, swarm algorithms





[Bruni, Corradini, Gaducci, Lluch Lafuente, Vandin : A Conceptual Framework for Adaptation. FASE 2012; Bruni, Corradini, Gaducci, Hölzl, Lluch Lafuente, Vandin, MW: Reconciling White-box and Black-Box Adaptation, ASCENS book 2015]





2. Recap: DevAdapt A development lifecycle for adaptive systems

Martin Wirsing



ASCENS Approach: "DevAdapt" Lifecycle





Agile lifecycle: DevAdapt

- Iterations at development time and runtime connected by deployment and feedback
- 3 feedback loops
- Cf. DevOps life cycle for integrating development and operations

Martin Wirsing

[N. Koch et al.: Life Cycle for the Development of Autonomic Systems: The e-Mobility Showcase, 3rd Workshop Self-Awareness in Autonomic Systems, 2013; Compuware Mainframe DevOps, https://www.compuware.com/lifecycle-overview/]





Develop "SOTA/Gem" requirements specification consisting of

- environment specification
- goal-oriented specification of the adaptive system with
 - Maintain goals of form $G_{maintain} \phi$
 - Achieve goals of form $F_{post} \phi$
- adaptation/awareness
 requirements (adaptation space)



[M. Hölzl, MW: Towards a System Model for Ensembles. Formal Modeling: Actors, Open Systems, Biological Systems 2011; D. Abeywickrama, M, Mamei, F. Zambonelli: Engineering Collectives of Self-driving Vehicles: The SOTA Approach. ISoLA (3) 2018: 79-93]







- Adaptation patterns such as
 - Reactive component
 - MAPE-K architecture of autonomic manager
- Modelling and programming/training: Whitebox adaptation
 - Classical adaptation, e.g.
 - Component-based programming with SCEL
 - Machine learning adaptation, e.g.
 - Reinforcement learning
 - Deep Neural Networks





MAPE-K Autonomic manager



Autonomic component







- Show that the system satisfies the (adaptation) requirements specification by
 - Testing
 - Simulation
 - (Statistical) Modelchecking using compositional techniques, fluid flow analysis, constraint programming, ...



Time steps (t)

Probability of rescuing the victim within a given time

 Note: Environment assumptions and adaptation space requirements must also be validated

Martin Wirsing

[J. Combaz,S. Bensalem, F.Tiezzi ,A. Margheri, R. Pugliese, J. Kofron: Correctness of Service Components and Service Component Ensembles. ASCENS Book 2015:105-158; R. De Nicola , M. Loreti, R. Pugliese, F. Tiezzi: SCEL: a Language for Autonomic Computing, ACM TAAS, 2014; L. Bulej, T. Bures, I. Gerostathopoulos, V. Horký, J. Keznikl, L. Marek, M. Tschaikowski, M. Tribastone, P. Tuma: Supporting Performance Awareness in Autonomous Ensembles. ASCENS Book 2015: 291-322]





3. Adapting Behaviours by Simulation-based Online Planning

Martin Wirsing

[Belzner, R. Hennicker, MW: Onplan: A framework for simulation-based online planning. In C. Braga and P. C. Ölveczky, eds., FACS 2015 LNCS 9539, Springer, 2015. 1-30.]



Search and Rescue Case Study Requirements



Environment

- Victims, fires and ambulances
- Unknown topology
- Unknown initial situation

Agent actions

- Noop, Move
- Load or drop a victim
- Extinguish fire if adjacent

Agent Requirements

- Find victims and bring them to an ambulance
 - $F(AND_{i=1,...n} victim_i at ambulance)$
- Extinguish fire if adjacent
 - − G (robot.position.fire →

action.robot = extinguishFire)







Adaptation Space

Environment and system may change

- in a nondeterministic way
 - Fires ignite and cease
 - Actions of agent may fail

by unexpected events

New fires break out and agents drop victims

Goals of agent may change

- Goal "Save victims from fire"
 - *F* (*AND*_{*i*=1,...n} ¬ *victim*_{*i*} *at fire*)
 changes to
 "Save victims and bring them to ambulance"







Online Planning

• Perform planning and execution of actions iteratively at runtime

Architecture

• MAPE-K loop with parallel planning

Online Planning Pseudo-Code

- while true do
 - observe state
 - execute || plan
- end while







Simulation-based Online Planning

simulates future episodes at runtime

- Architecture
 - Iterated MAPE-K loop
- Simulation-based planning
 - Simulate by generating *n* episodes and their rewards
 - After each episode update the current policy according to strategy
- Strategies
 - Monte Carlo Tree Search
 - Cross Entropy Planning







Monte Carlo Tree Search for Discrete Domains



Simulation with Monte Carlo Tree Search

- Policy as tree
 - Nodes represent states and action choices
 - Add a node per simulation
 - Aggregate simulation data in nodes
 - Reward and frequency
 - Sample actions w.r.t. aggregated data





[Cameron B Browne et al.: A survey of Monte Carlo tree search methods. Computational Intelligence and AI in Games, IEEE Transactions on, 4(1):1 - 43, 2012.]



[Kocsis, Levente, and Csaba Szepesvári. Bandit-based Monte Carlo planning. Machine Learning: ECML 2006. Springer: Berlin Heidelberg, 2006, 282-293].



- X_j : Average reward of child node j
- *n*: Nr. of episodes from current node
- n_j : Nr. of episodes from child node j
- C: UCT exploration constant

[Kocsis, Levente, and Csaba Szepesvári. Bandit-based Monte Carlo planning. Machine Learning: ECML 2006. Springer: Berlin Heidelberg, 2006, 282-293].



Expand the Tree



Add a new node

• When an episode leaves the tree











Search and Rescue Case Study: Experiment



Experiment

- Topology randomly generated with 20 nodes, 6-7 connections/nodes
- Action failure probability = 5%
- Planning depth = 20
- Reward for victims at safe positions:

 $R(s) = 100 \cdot \Sigma_{v \in victims(s)}$ #(v.position.safe)









Provided reward

• Victim at ambulance: +100

System synthesizes sensible behaviour

Results in 0.95 confidence interval

Checked with MultiVeStA



Measured (in %): Victims at ambulance (blue), in a fire (red) Positions on fire (green)

[Stefano Sebastio and Andrea Vandin. MultiVeStA: statistical model checking for discrete event simulators. ValueTools '13. 2013. 310-315.]

Autonomy



Validation (II)



Expose system to **<u>unexpected events</u>**

- At steps 20, 40, 60, 80
- All carried victims are dropped
- New fires break out
- Events NOT simulated by planner
- New situation incorporated by planner

System shows sensible reactions

Results in 0.95 confidence interval

Robustness



Positions on fire (green)



Validation (III)



Change system goals while operating

- Change of reward function
 - Steps 0-40: Reward for victims not in a fire
 - Steps 40-80: Reward for victims at ambulance
- Change NOT simulated by planner
- But planner incorporates new situation

System adapts behaviour wrt. goals

Results in 0.95 confidence interval

Flexibility



Measured (in %): Victims at ambulance (blue), in a fire (red) Positions on fire (green)





4. Safe Learning

PSyCo: Policy SYnthesis with safety COnstraints

Martin Wirsing





Grid world with two passages

 broad long passage and short narrow passage

Robot agent

can move into 4 directions but slips with prob
 0.05

Requirements

- Achieve goal
 - Go to bottom right corner:

F atBottomCorner

- Safety ("maintain") goal
 - Don't collide with walls:

 $p(G \neg collision) > p_{sat}$ iff $p(F collision) <= p_{sat}$







Idea: Generate reinforcement learning algorithm from requirements so that the safety requirements are (mostly) respected

- Achieve goals
 - Optimise goal satisfaction by Q-learning as long as safety requirements hold
- Safety goals
 - Control goal satisfaction
 in case of violation return into the "safety corridor"







Negated (!) safety property

Reinforcement learning

- Learn in parallel Q-function of achieve goal and safety distribution
- Then the combined policy
 - maximises expected future reward when being safe and
 - minimises future expected constraint satisfaction probability when violating constraints
- Formally:

Let $A_{ ext{safe}}(s,\pi) = \{a | p(F\phi|s,a,\pi) \leq p_{ ext{sat}}\}.$

$$\pi(s) = egin{cases} rg\max_{a \in A_{ ext{safe}}(s,\pi)} Q(s,a) & ext{if } A_{ ext{safe}}(s,\pi)
eq \emptyset \ rg\min_{a} p(F\phi|s,a,\pi) & ext{otherwise} \end{cases}$$





Learning safety-aware policy

• Learn in parallel

Q-function of achieve goal and safety distribution

• **Q-value iteration*:**

$$Q(s,a) \leftarrow (1-lpha)Q(s,a) + lpha(r+\gamma \max_{a' \in A(s')} Q(s',a'))$$

where *r* reward, α learning rate, γ discount factor

• Safety distribution:

$$\begin{array}{ll} \mbox{Initialize }\forall s,a,\pi:p(F\phi|s,a,\pi)=0. \\ p(F\phi|s,a,\pi) \leftarrow \begin{cases} (1-\alpha)p(F\phi|s,a,\pi)+\alpha 1 & \mbox{if }s'\models\phi\\ (1-\alpha)p(F\phi|s,a,\pi)+\alpha p(F\phi|s',\pi(s'),\pi) & \mbox{otherwise} \end{cases} \end{array}$$

[*C.Watkins: Learning from Delayed Rewards. Ph.D. thesis, Cambridge University, 1989]





Experiment

- Vary p_{sat} from 0.1 to 0.5
- Learning rate $\alpha = 0.1$, discount factor $\gamma = 0.9$,
- Reward r = 1.0 if robot at bottom right corner





Training: 100.000 episodes (each up to 100 steps) State visitation count



 $p_{sat} = 0.3$ Mostly the long passage is chosen to avoid too many collisions



 $p_{sat} = 0.5$ Mostly the short passage is chosen due to lax collision requirements





Validation: 1000 episodes using learned policy State visitation count



 $p_{sat} = 0.3$ Learned policy chooses long passage



 $p_{sat} = 0.5$ Learned policy chooses short passage





Validation: 1000 episodes with learned policy # collision free episodes vs. # episodes with collisions:



 $p_{sat} = 0.3$ 80% no collision on long passage





Validation: 1000 episodes using learned policy

Validation fails for strong collision freeness ($p_{sat} = 0.1, 0.2$)!





 $p_{sat} = 0.1$ State visitation count long passage preferred

 $p_{sat} = 0.1$ but learned policy follows short passage and has >30% collisions

Further design cycle is necessary!







| ,, | Safe Learning" PSyCo | C |
|-------------|-----------------------|---|
| Adaptation: | "slow" at design time | |
| Execution: | FAST | |

Online simulation planning

"fast" at runtime SLOW

- Many challenges
 - How to engineer complex safe adaptive multi-agent systems?
 - In adaptive environments?

Martin Wirsing