



General, Learned, and Verified: What the Future of Learning Agents in Power Grids Could Be

Eric MSP Veith <eric.veith@uol.de> , 2025-03-10



% whoami



- Eric MSP Veith <eric.veith@uo1.de>
- Currently head of a junior research group at University of Oldenburg, Germany
- Computer scientist by heart: First ICT, then distributed heuristics, then Multi-Agent Systems, now advanced Deep Reinforcement Learning
- PhD in 2017: “Universal Smart Grid Agent for Distributed Power Generation Management.”
- Creator of the Adversarial Resilience Learning methodology (advanced DRL in CNIs)



% whereami



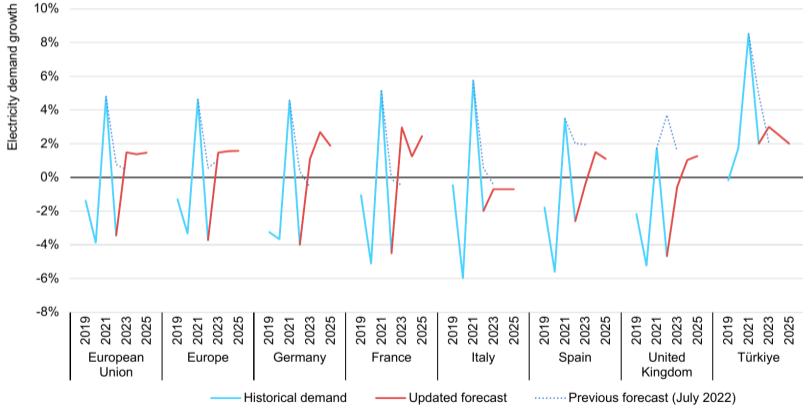
General, Learned, and Verified: What the Future of Learning Agents in Power Grids Could Be —
Eric MSP (with @eric.msp@tu-berlin.de) — Adversarial Resilience Learning



Electricity Demand Rising

After significant decline in 2022, European electricity demand is set to recover

Year-on-year relative change in electricity demand, Europe, 2019-2025

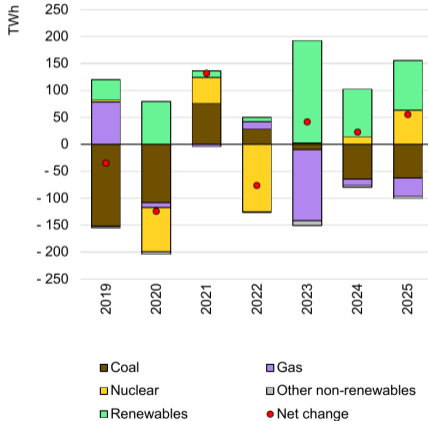




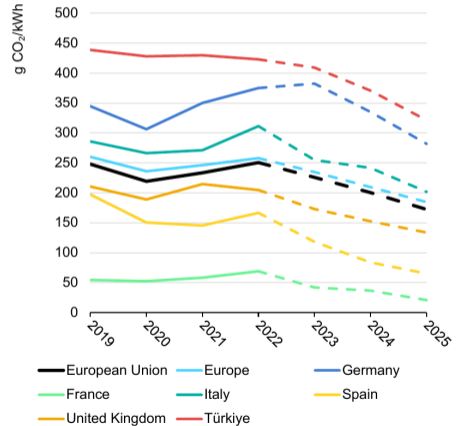
Renewables Are Replacing Fossil Fuels

Following two years of increases, CO₂ intensity starts to decline again from 2023 onward

Year-on-year change in electricity generation, European Union, 2019-2025



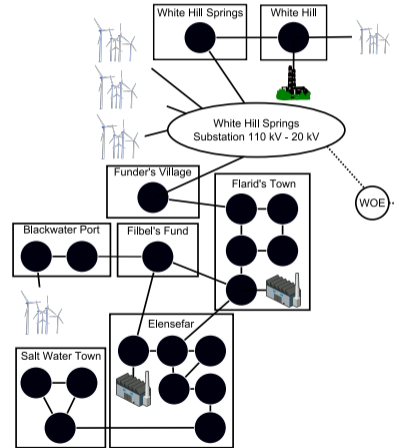
Development of average CO₂ intensity, Europe, 2019-2025





AI as Promise of an Alternative

- Multi-Agent Systems promise local, more more efficient grid operation
- Each node (subgrid, ...) an agent
- Nodes (agents) forecast local power generation/consumption
- On disequilibrium, match forecasts to achieve equilibrium
- An example, based on the literature [5, 3]





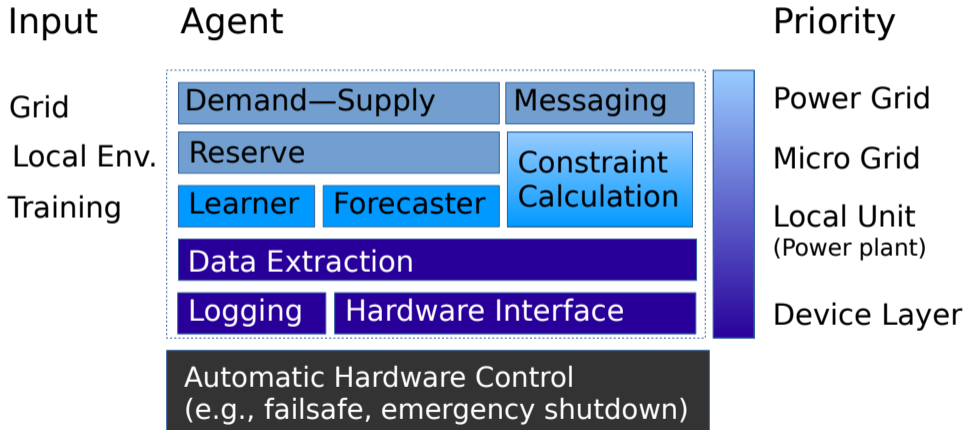
Pillars of CPES-MAS

An approach for a Multi-Agent System (MAS) that manages high shares of volatile generators and consumers in an energy system is based on three pillars:

1. Forecasting of local generation or demand
2. Communicating demand and generation (distributed snapshotting with the power grid in mind)
3. Solving the combinatorial problem of demand and supply

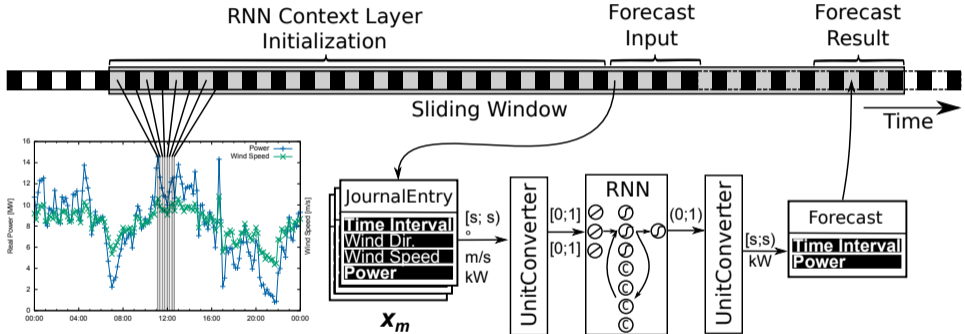


Agent Design





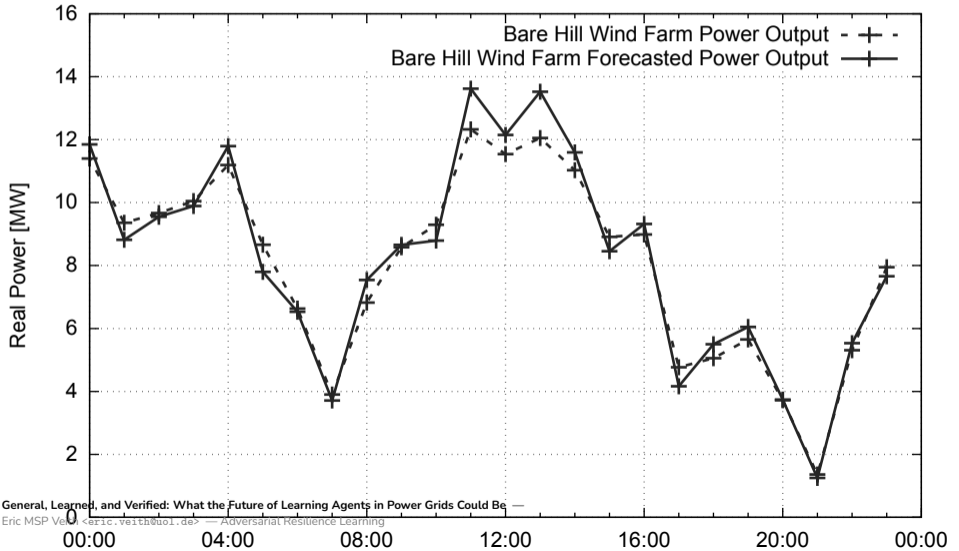
Forecasting



Forecasting is industry state of the art now.



Forecasting





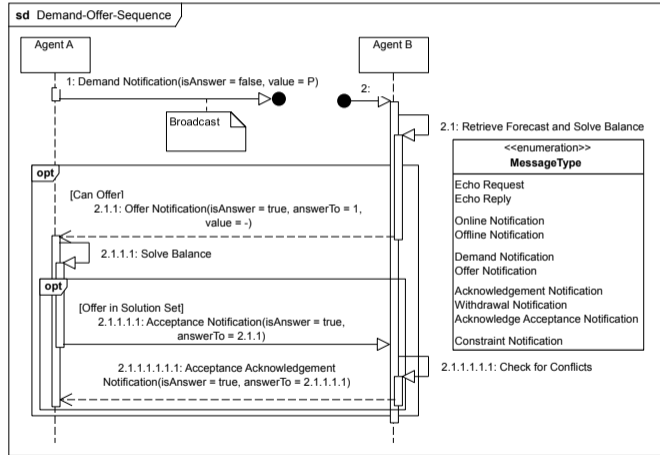
Communication

After forecasting works (i. e., imbalances are known in advance), each agent (= node) must ask its neighbors to help attain the equilibrium.

- Approach: Use overlay network – (virtual) communication lines between agents based on actual grid lines
- Requests (demand for power, or increase in feed-in that should be consumed) travel via selective broadcast
- Each node along the way must try to contribute!



Four-Way Handshake





LPEP Forwarding

- L_i : Links of the i -th agent
- $l_{i,k}$: k -th link of i -th agent
- $\text{distance}(l_{i,k})$: Distance metric
- m_j : j -th message
- M_i : Message Journal of the i -th agent

$$M_i = \{m_1 \mapsto \{(l_{i,1}, m_{1,\text{distance}(l_{i,1})}), \dots, (l_{i,n}, m'_{1,\text{distance}(l_{i,n})})\},$$

$\dots,$

$$m_n \mapsto \{(l_{i,1}, m_{n,\text{distance}(l_{i,1})}), \dots, (l_{i,n}, m'_{n,\text{distance}(l_{i,n})})\}$$

$$l_{i,1}(t) \leq l_{i,2}(t) \Leftrightarrow l_{i,1,\text{distance}}(t) \leq l_{i,2,\text{distance}}(t)$$

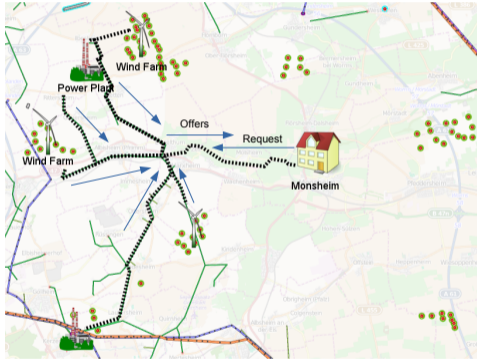


Forwarding

1. Respect *Constraint Notifications*:
 - 1.1 No answer if $\min(M(m))$ a constraint notification to m , additionally
 - 1.2 send *Withdrawal Notification* iff already answered
2. $m_{isAnswer}$: forward on best connect ($\min(M(m_{answerTo}))$)
3. *Selective Broadcast* for requests:
 - 3.1 Replace request with *Constraint Notification*, if necessary
 - 3.2 $M(m) = \emptyset$: forward on $|L| - 1$ links
 - 3.3 $m' = \min(M(m'))$: Update by forwarding
 - 3.4 Otherwise: no forwarding



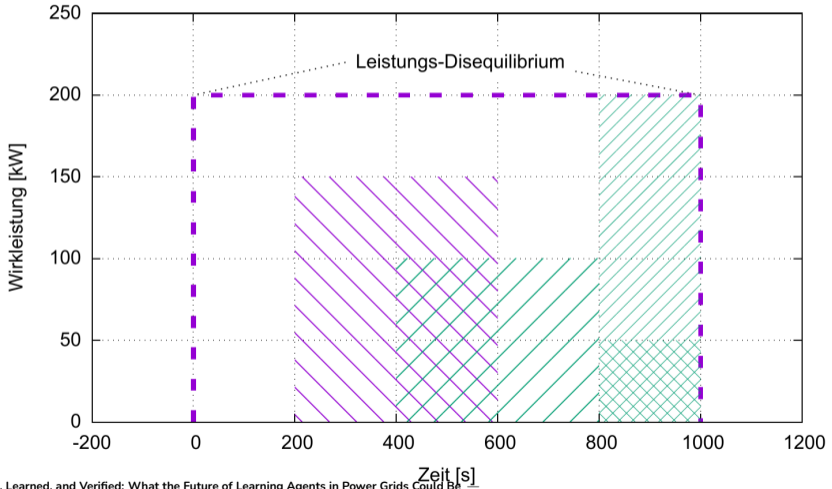
How to Decide...?



1. Local forecasting shows demand or oversupply of energy
2. Requests are sent
3. Other nodes make offers
4. Offers reach requestor
5. **Decision about offers?**



Power Balance Concept





Problem Statement

'Power Balance Algebra':

$$\{[t_1; t_3] \mapsto P_1\} \cup \{[t_2; t_4] \mapsto P_2\} = \{[t_1; t_2] \mapsto P_1, [t_2; t_3] \mapsto P_1 + P_2, [t_3; t_4] \mapsto P_2\}, \quad (1)$$

$$[t_1; t_2] \mapsto P_1 \subseteq [t_3; t_4] \mapsto P_2 \Leftrightarrow t_1 \geq t_3 \wedge t_2 \leq t_4 \wedge P_1 \leq P_2; \quad (2)$$

Distance Function: $d(r_i) : r_i \mapsto \mathbb{R} \quad (3)$

Problem Statement:

$$\sum_i b_i r_i \subseteq r_0, \quad i \neq 0, b_i \in \{0, 1\}, \quad (4)$$

$$\text{Subject to: } \min \sum_i b_i d(r_i), \quad i \neq 0, b_i \in \{0, 1\}. \quad (5)$$



Atomization

$$\mathbf{P} = (|P_0|, |P_1|, \dots, |P_i|, |P_C|),$$

$$\mathbf{t} = (t_{2,0} - t_{1,0}, t_{2,1} - t_{1,1}, \dots, t_{2,i} - t_{1,i}),$$

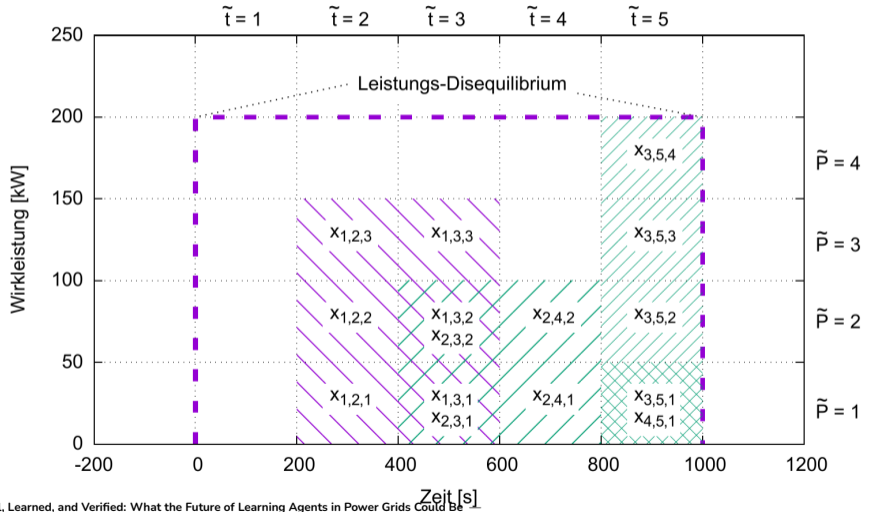
$$\Delta P = \text{ggT}(\mathbf{P}),$$

$$\Delta t = \text{ggT}(\mathbf{t}),$$

$$x_{i, \tilde{t}, \tilde{P}} = \begin{cases} 1 & \text{if agent } i \text{ influences the grid in time-subinterval } \tilde{t} \text{ with} \\ & \text{power from the power-subinterval } \tilde{P}, \\ 0 & \text{else.} \end{cases}$$



Atomization Illustrated





Model of the Disequilibrium

A symmetric function for each time-subinterval:

$$S_k^n(x_{i,\tilde{t}=k}, \tilde{p}) = \begin{cases} 1 & \text{if } n \text{ variables in } x_{i,\tilde{t}=k}, \tilde{p} \text{ equal } 1, \\ 0 & \text{else;} \end{cases}$$

Full Disequilibrium:

$$S = \bigcap_{k=1}^m S_k^n(x_{i,\tilde{t}=k}, \tilde{p})$$

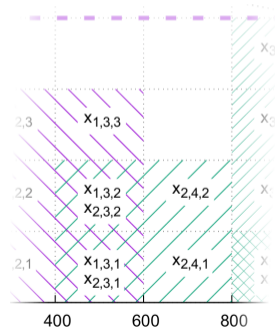


Modelling Responses

Acceptance Function:

$$r_i(x_{i,\tilde{t},\tilde{p}}) = \begin{cases} 1 & \text{if } x_{i,\tilde{t},\tilde{p}} \text{ describes a valid interval for accepting the re-} \\ & \text{ponse of } i, \\ 0 & \text{else.} \end{cases}$$

$$\begin{aligned} r_2(x_{i,\tilde{t},\tilde{p}}) &= \bar{x}_{2,3,1} \wedge \bar{x}_{2,3,2} \wedge \bar{x}_{2,4,1} \wedge \bar{x}_{2,4,2} \\ &\quad \vee x_{2,3,1} \wedge x_{2,3,2} \wedge \bar{x}_{2,4,1} \wedge \bar{x}_{2,4,2} \\ &\quad \vee x_{2,3,1} \wedge x_{2,3,2} \wedge x_{2,4,1} \wedge x_{2,4,2} \end{aligned}$$





Equilibrium

$$S = \bigcap_{k=1}^m S_k^n(\mathbf{x}_{i, \tilde{t}=k}, \tilde{\mathbf{p}})$$

$$R = \bigcap_{i \in I', \tilde{t}, \tilde{\mathbf{p}}} r_i(\mathbf{x}_{i, \tilde{t}}, \tilde{\mathbf{p}}),$$

$$C = S \cap R.$$



Equilibrium

$$S = \bigcap_{k=1}^m S_k^n(\mathbf{x}_i, \tilde{\mathbf{t}}=k, \tilde{\mathbf{p}})$$

$$R = \bigcap_{i \in I', \tilde{\mathbf{t}}, \tilde{\mathbf{p}}} r_i(\mathbf{x}_i, \tilde{\mathbf{t}}, \tilde{\mathbf{p}}),$$

$$C = S \cap R.$$

- Best solution through ordering: $r_i \leq r_{i'}$ \Leftrightarrow $d(r_i) \leq d(r_{i'})$
- Generating next vector in S through permutation
- Exploiting the commutative property of the intersection operator:
 $R_n \cap (\dots \cap (R_2 \cap (R_1 \cap S)))$



Efficiency

Data Effect

$$\kappa = \frac{W}{D} \left[\frac{\text{kWh}}{\text{kB}} \right]$$

Data Efficiency

$$\xi = \frac{\Delta P}{D} \left[\frac{\text{kW}}{\text{kB}} \right]$$



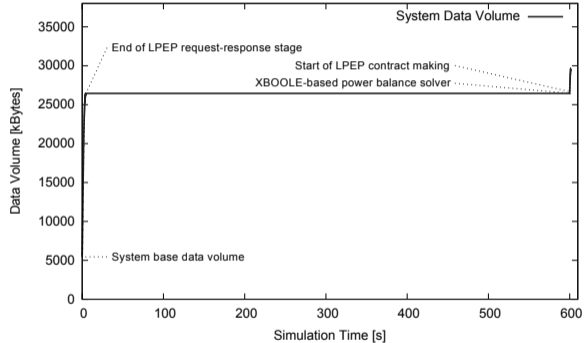
Comparison

Comparison with BDD approach by Inoue *et al.* (2014):

	BDD	Universal Agent
Loss Avoided (ΔP)	17 208 kW	17 208 kW
Runtime	> 16 min	< 11 min (simulated)
D	100 MB	28.9 MB
ξ	0.168 kW/kB	0.581 kW/kB



Universal Agent Efficiency



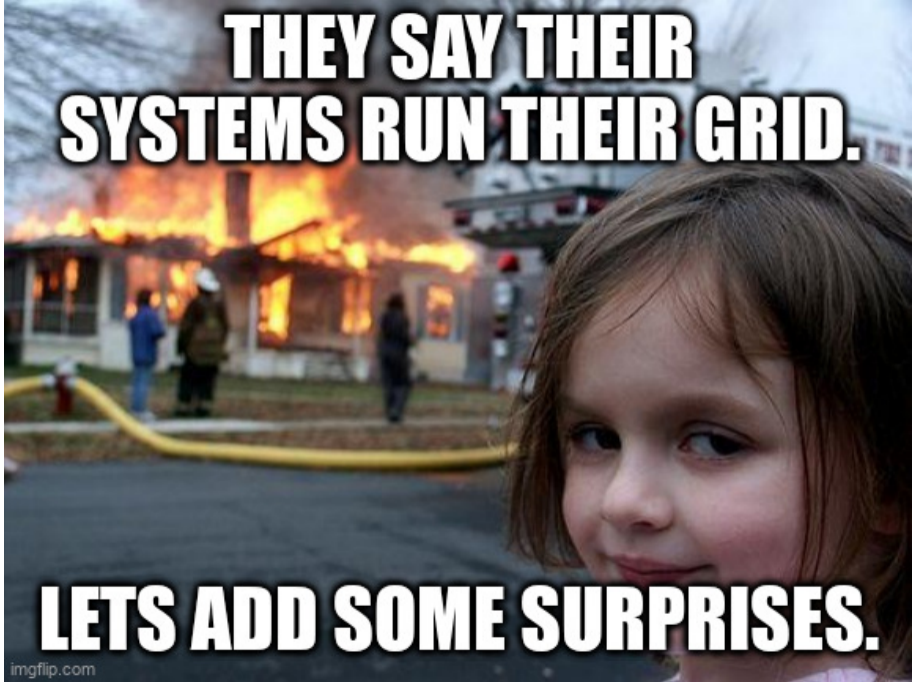
- BDD approach in low-load situation: 100 kB
- *Universal Agent* concept especially useful in complex load situations

AND THIS, GENTLEMEN

IS HOW YOU RUN YOUR GRID.



**THEY SAY THEIR
SYSTEMS RUN THEIR GRID.**



LETS ADD SOME SURPRISES.



Knowns vs. Unknowns

- The well-defined way of traditional MAS can guarantee a (theoretical) optimal solution
- They are robust: Cases known at design time can be handled
- Unknown unknowns and even some known problems can not be handled.

... We need a system that can act universally.

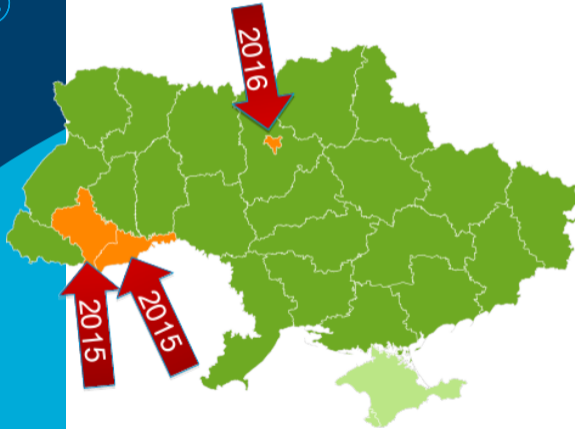


**“There are only two types of companies:
those who have been hacked,
and those who don’t yet know
they have been hacked.”**

— John T. Chambers

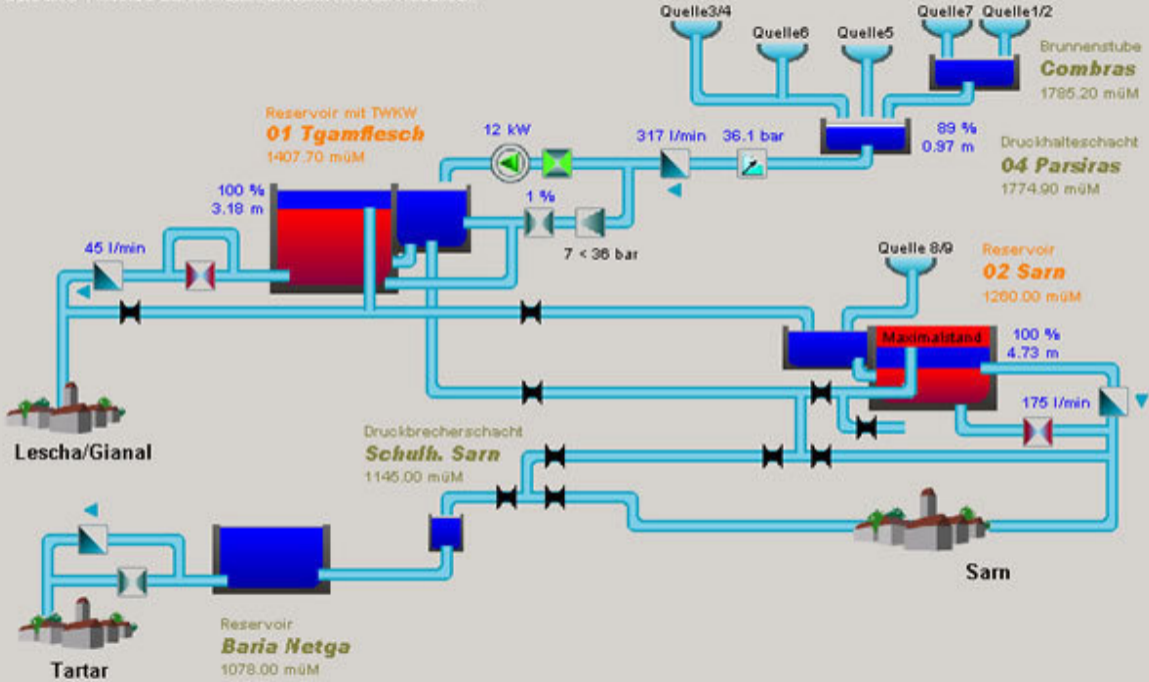


Energy Systems Fit The Bill Just As Well

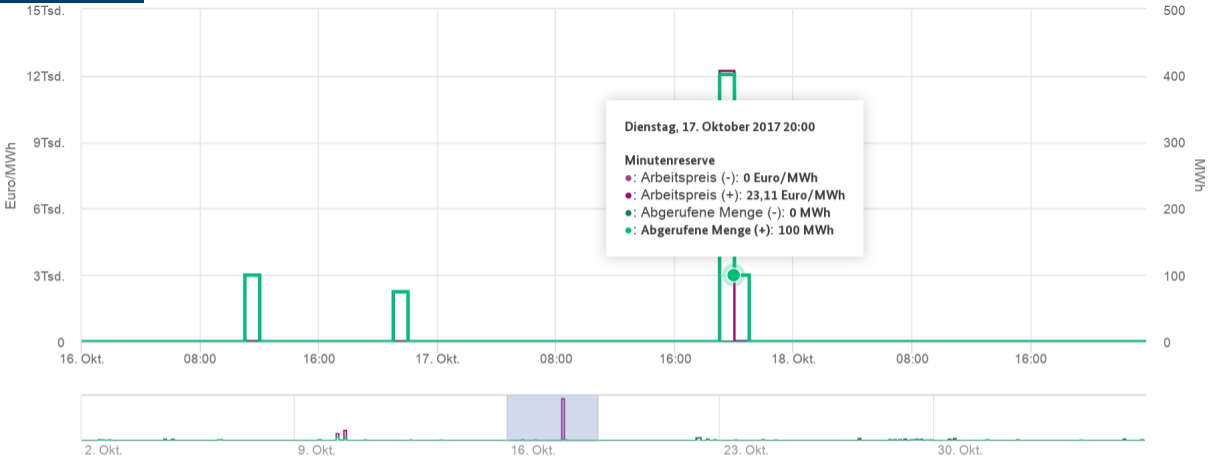


Dec 23rd, 2005

- **Cyber attack** causes blackout in the Ukraine
- **3 DSOs** targeted
- **High level of automation** helps attackers
- Operative intrusion in **OT**; disconnection of **several substations**
- Several months in preparation







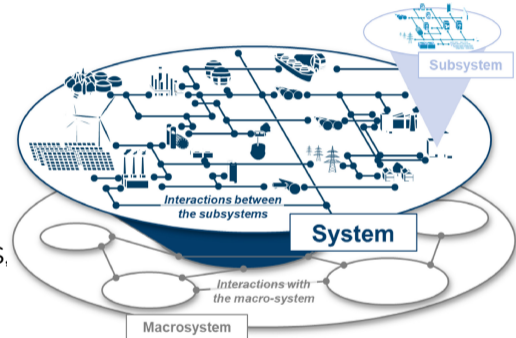
Systemstabilität - Minutenreserve

- Abgerufene Menge (+)
- Abgerufene Menge (-)
- Arbeitspreis (+)
- Arbeitspreis (-)
- Vorgehaltene Menge (+)
- Vorgehaltene Menge (-)
- Leistungspreis (+)
- Leistungspreis (-)



Learning Resilient Control

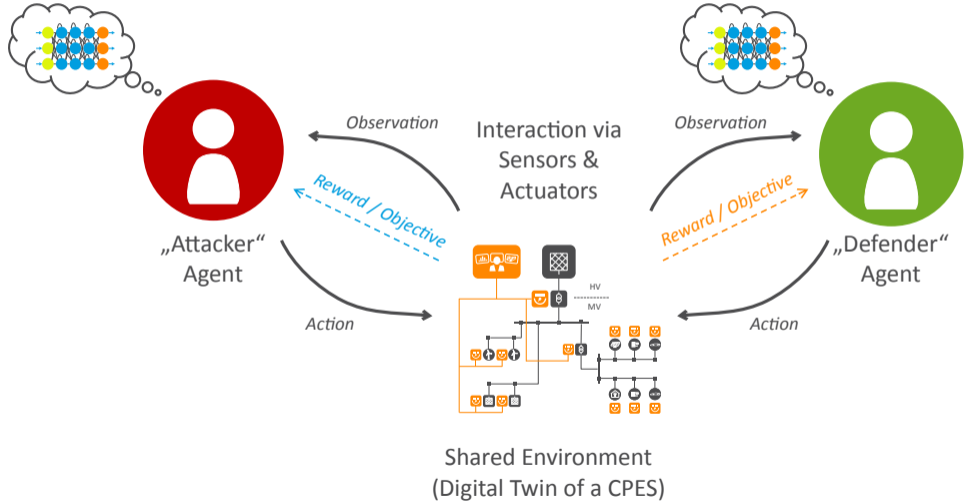
- **Interconnected CPS have always attack surface due to their inherent complexity**
- Low latency of ICT and OT
- High interdependence
- Complexity in breadth and depth
- Critical Services as SPOF (DNS, BGP, SCADA, SDL)
- **Learning Strategies for automatic issue mangement**
- “Adversarial Resilience Learning”



Kotzur, Leander, et al. “A modeler’s guide to handle complexity in energy systems optimization.” *Advances in Applied Energy* 4 (2021): 100063.

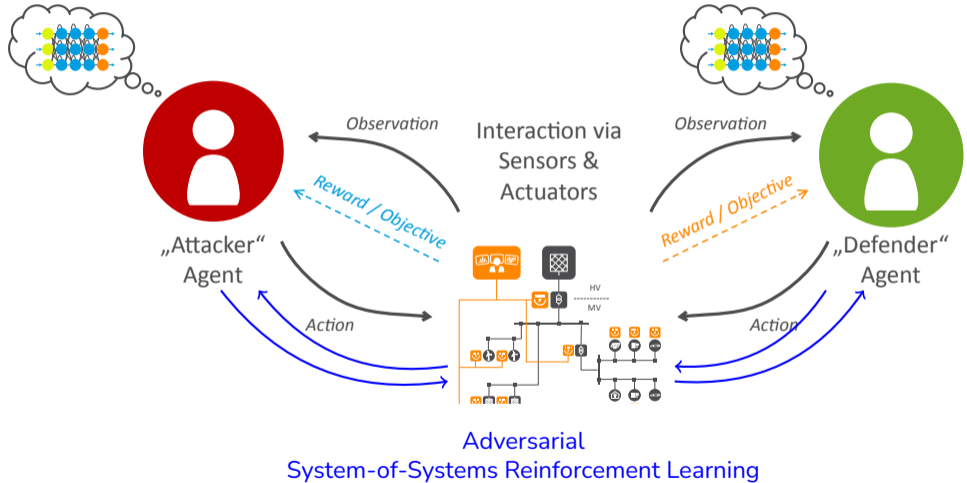


Adversarial Resilience Learning



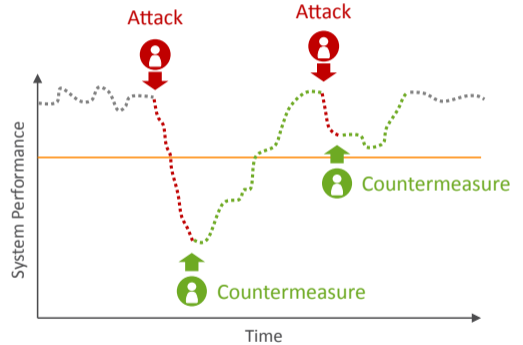
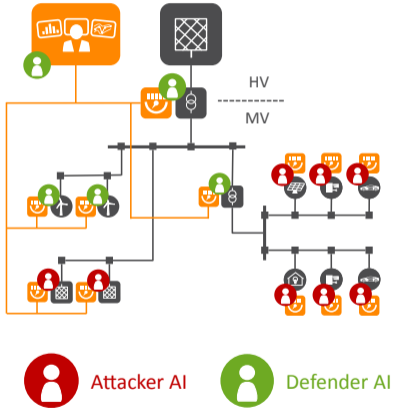


Adversarial Resilience Learning





ARL Agent Interaction



Defender Points

2656

Loads Connected



Generators Connected



Buses Connected



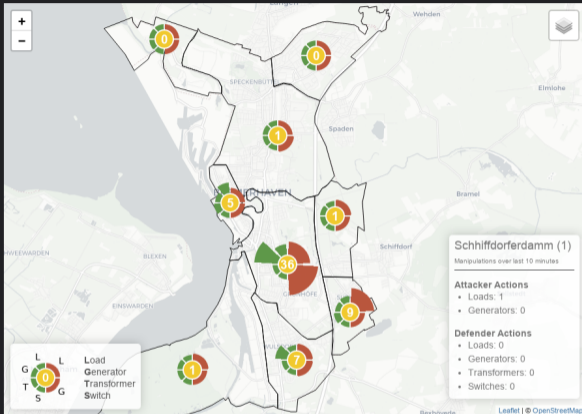
Transformers Connected...



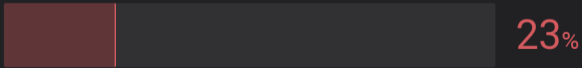
Most Valuable Actions (Defender)

Info	Time	Points
Changed Scaling from Lehe Households - 4 to 0.6667	2018-01-01 02:00:10	-0.25
Changed Scaling from Leherheide Industrielast to 0.8889	2018-01-01 02:26:10	-0.07
Changed Tap_pos from trafo to 1.0000	2018-01-01 01:47:50	-0.07
Changed Scaling from PV Fischereihafen to 0.0000	2018-01-01 02:14:30	-0.07
Changed Tap_pos from trafo to 1.0000	2018-01-01 01:45:30	-0.06
Changed Scaling from MFCG-Kindergarten	2018-01-01	

Map



Time Left (Coins Left in %)



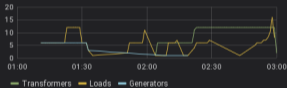
Attacker Points

7344

Constraint Violations



Malfunctions



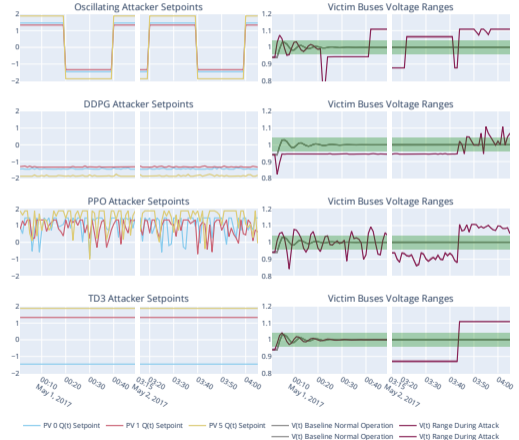
Most Valuable Actions (Attacker)

Info	Time	Points
Changed Scaling from Geestemünde Households - 0 to 0.5000	2018-01-01 01:00:00	0.96
Changed Scaling from Geestemünde Households - 0 to 0.5000	2018-01-01 01:00:00	0.93
Changed Scaling from Geestemünde Households - 0 to 0.5000	2018-01-01 01:00:00	0.91



ARL Agent Can Discover Attacks

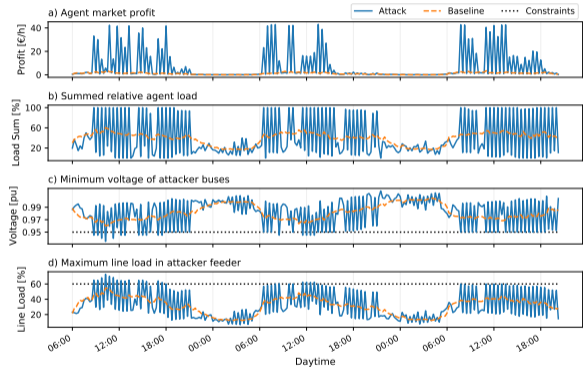
- Attack on voltage level
- Attacker controls Q feed-in
- Known attack: Oscillating behavior
- ARL agent independently discovers attack, but also finds variant





Transactive Energy Can Be Gamed

- Economic and control techniques, based on market standard values
- There is no “sound” market design yet than cannot be gamed
- Worse yet: Agents can find weaknesses & gain market dominance without system knowledge

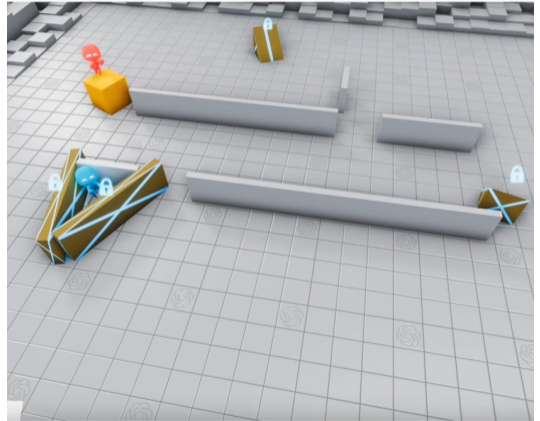


Agents learn to “game” local energy markets
Wolgast, Veith, and Nieße [6]



Multi-Agent Autocurricula

- ARL is an autocurriculum setup
- Independently known & verified to work
- Example Setup: Two groups of agents play hide and seek
- No domain information; agents learn strategies and tool use independently
- Result: Agents learn to exploit bugs in the underlying game engine
 - Holes in walls
 - Sliding boxes
 - Edge/corner jumps





Autocurricula Helpful in Theory

- DRL agents collect initial samples from random actions
- However, random actions over correlated actuators lead to convolution problem, i. e., if $X, Y \sim \mathcal{U}$, then

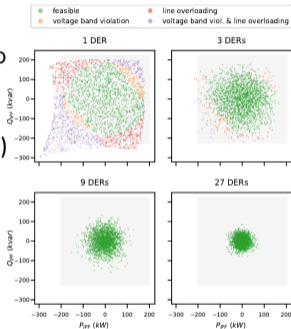
$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx, \quad (6)$$

which is a triangle distribution

- Equally, consider SAC's entropy maximization,

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(R(s_t, a_t, s_{t+1}) + \alpha \underbrace{H(\pi(\cdot | s_t))}_{\text{Entropy term}} \right) \right] \quad (7)$$

- ... obviously, a “push” is required





Autocurricula Helpful in Theory II

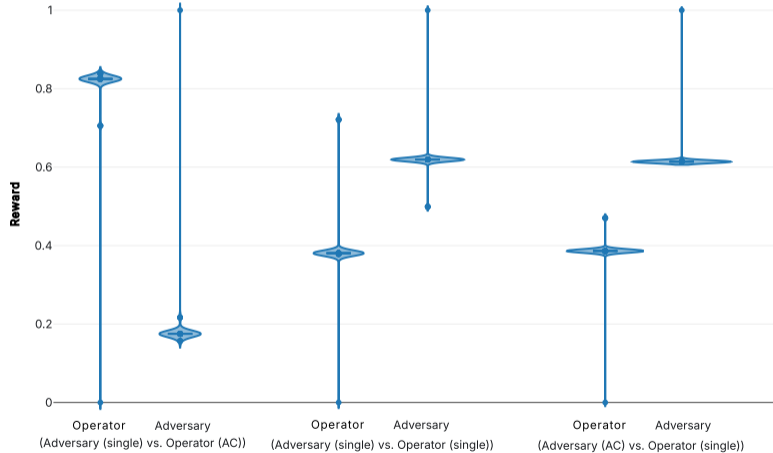
1. Formally, DRL approximates the unknown environment distribution p with q , i. e.,

$$\begin{array}{ll} \text{minimize} & KL(p, q) \\ \text{subject to} & q \end{array} \quad (8)$$

2. Learn a policy to exploit q , π_Ω
3. (Single agent: get stuck in local optimum because p is mostly unknown because of missing sample data)
4. Adversary agent: Observe p as influenced by π_Ω
5. $R_A(\mathbf{s}_t \sim p) = -R_\Omega(\mathbf{s}_t \sim p)$, therefore $\pi_A \hat{=} -\pi_\Omega$
6. Result: agents observe adversarial samples from the “other end” of p 's spectrum
7. Agents try to counter adverse effects: efficient state/action space exploration



... and in Practice





ARL Works

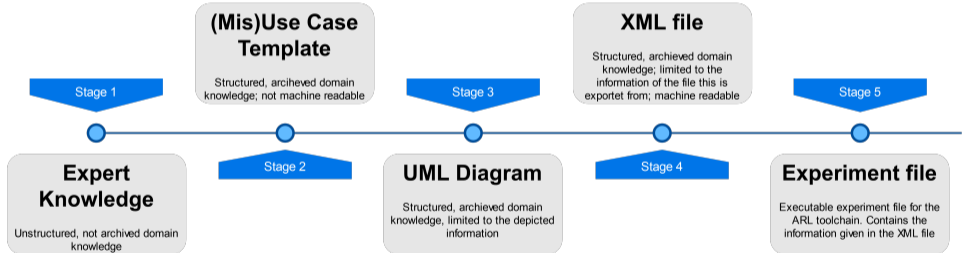
To summarize...

- ARL works for finding attack vectors (“easy”)
- ARL defender learn resilient control (“not quite so easy, but still...”)
- ARL agents learn faster & more robust strategies through the autocurriculum setup (“prove me, I’m only circumstantial evidence!”)
- ARL defender agents can control modern power grids (“ha-ha, as if that would be acceptable...”)
- **There is still a lot missing:**
 - Behavior guarantees
 - Adhere to constraints (rulesets)
 - Learn from existing domain knowledge
 - Adapt during production use (not just retraining)
 - ...



Learning from Domain Knowledge

Example: Misuse Cases





Trajectories from (Mis-) Use Cases

- Annotate UML diagrams to allow sampling; construct:
 - Experiment file
 - State machine from transitions

$$M_{tg} = (Q, \Sigma, \delta, q_0, F)$$

with

$$(q, (\{c_q\} \in \text{ActuatorSetpoints},$$

$$\{i_q\} \in \text{TimeStepIntervals})) \in Q$$

$$(i \in Q, n \in Q, \{sc\} \in \text{StepConstraints}) \in \delta$$

(9)

Relevant properties:

- Non-determinism
- State/actuator constraints c_q (think Gymnasium spaces)
- Time step intervals (sync to simulation semantics)
- Constrained steps (e. g., grid codes)



Combined AWAC and State Machine Sampling

Initialize Simulation S

State Machine $M_{tg} = (Q, \Sigma, \delta, s_0, F)$

$maximum_steps \leftarrow x$

for $j \leq x$ **do**

$s \leftarrow S.state$

$a \leftarrow \{c\} \in (M_{tg}.state, cM_{tg}.state) \in Q_{M_{tg}}$

$r \leftarrow R(a)$

$s' \leftarrow S.step(a)$

$db \leftarrow db \cup \{(s, a, s', r)\}$

$advance(M_{tg})$



Combined AWAC and State Machine Sampling II

Dataset $D = \{(s, a, s', r)_j\} \sim db$

Initialize buffer $\beta = D$

Initialize π_θ, Q_ϕ

for iteration $i = 1, 2, \dots, n$ **do**

 Sample batch $(s, a, s', r) \sim \beta$

$y = R(s, a) + \gamma \mathbb{E}_{s', a'} [Q_{\phi_{k-1}}(s', a')]$

$\phi \leftarrow \arg \min_{\phi} \mathbb{E}_D [Q_\phi(s, a) - y]^2$

$\theta \leftarrow \arg \max_{\theta} \mathbb{E}_{s, a \sim \beta} [\log \pi_\theta(a|s) \exp(\frac{1}{\lambda} A^{\pi_k}(s, a))]$

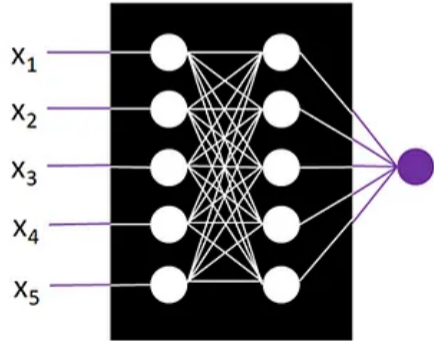
if $i > \text{num_offline_steps}$ **then**

$\tau_1, \dots, \tau_K \sim p_{\pi_\theta}(\tau)$

$\beta \leftarrow \beta \cup \{\tau_1, \dots, \tau_K\}$



“The Black Box”





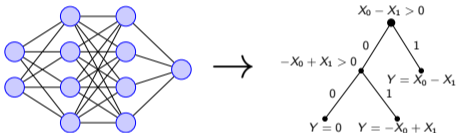
Explanation goals

- Motivation: No trust without explanation of learned strategies of agents
- Idea: Use Decision Trees (DTs) with extraction of rulesets for explanation
 - DTs are transparent and somewhat interpretable
 - They can be trained directly (no need for black-box Deep Neural Network (DNN) models)
 - But DNNs are better regularized, which increases trainability [2]
- Conflicting goals:
 - Construction of powerful (Deep Reinforcement Learning (RL) (RL)) learning system
 - (Post-hoc) Explainability with comprehensible model (e. g. DTs)



Learned Policy Explanation

- Equivalent transformation of efficient-learnable Feed-Forward DNNs (DNNs) into compressed DTs



- NN2EQCDT algorithm heavily relies on equivalence description of DNNs and DTs [1], but still addressed research gaps to better use it for explainability:
 - Transformation algorithm and actual implementation proposed for PyTorch models
 - Exponential growth is addressed by lossless pruning
 - Dynamic compression reduces computation time significantly and may reduce inference time
 - Option to directly include global constraints for further pruning



NN₂EQCDT algorithm

```

1:  $\hat{\mathbf{W}} = \mathbf{W}_0$ 
2:  $\hat{\mathbf{B}} = \mathbf{B}_0^\top$ 
3:  $rules = \text{calc\_rule\_terms}(\hat{\mathbf{W}}, \hat{\mathbf{B}})$ 
4:  $T, new\_SAT\_leaves = \text{create\_initial\_subtree}(rules)$ 
5:  $\text{set\_hat\_on\_SAT\_nodes}(T, new\_SAT\_leaves, \hat{\mathbf{W}}, \hat{\mathbf{B}})$ 
6: for  $i = 1, \dots, n - 1$  do
7:    $SAT\_paths = \text{get\_SAT\_paths}(T)$ 
8:   for  $SAT\_path$  in  $SAT\_paths$  do
9:      $\mathbf{a} = \text{compute\_a\_along}(SAT\_path)$ 
10:     $SAT\_leave = SAT\_path[-1]$ 
11:     $\hat{\mathbf{W}}, \hat{\mathbf{B}} = \text{get\_last\_hat\_of\_leave}(T, SAT\_leave)$ 
12:     $\hat{\mathbf{W}} = (\mathbf{W}_i \odot [(\mathbf{a}^\top)_{\times k}])\hat{\mathbf{W}}$ 
13:     $\hat{\mathbf{B}} = (\mathbf{W}_i \odot [(\mathbf{a}^\top)_{\times k}])\hat{\mathbf{B}} + \mathbf{B}_i^\top$ 
14:     $rules = \text{calc\_rule\_terms}(\hat{\mathbf{W}}, \hat{\mathbf{B}})$ 
15:     $new\_SAT\_leaves =$ 
       $\text{add\_subtree}(T, SAT\_leave, rules, invariants)$ 
16:     $\text{set\_hat\_on\_SAT\_nodes}(T, new\_SAT\_leaves,$ 
       $\hat{\mathbf{W}}, \hat{\mathbf{B}})$ 

```

Finally:

- Converting final rules to expressions
- Pruning the (temporary) UNSAT nodes



Effective weight matrix calculation

```

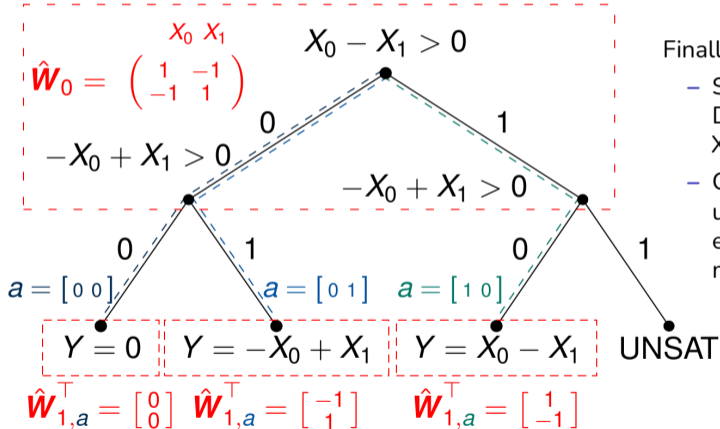
1:  $\hat{W} = W_0$ 
2:  $\hat{B} = B_0^\top$ 
3: for  $i = 0, \dots, n - 2$  do
4:    $\mathbf{a} = [ ]$ 
5:   for  $j = 0, \dots, m_i - 1$  do
6:     if  $(\hat{W}_j \mathbf{x}_0^\top + B_j^\top)^\top > 0$  then
7:        $\mathbf{a}.append(1)$ 
8:     else
9:        $\mathbf{a}.append(0)$ 
10:     $W_{i+1} \in \mathbb{R}^{m_i \times k}, \mathbf{a} \in \mathbb{Z}_2^{m_i}$ 
11:     $\hat{W} = (W_{i+1} \odot [(\mathbf{a}^\top)_{\times k}]) \hat{W}$ 
12:     $\hat{B} = (W_{i+1} \odot [(\mathbf{a}^\top)_{\times k}]) \hat{B} + B_{i+1}^\top$ 
13: return  $(\hat{W} \mathbf{x}_0^\top + \hat{B})^\top$ 

```

- Using right-handed linear transformation with bias
- Tailored to ReLU(-like) activation functions (e.g. ReLU, PReLU, LeakyReLU)



XOR model: DT Construction



Finally:

- Simple example of an DT representing an XOR function
- Construction of DT using calculated effective weight matrices



XOR model: DT Pruning

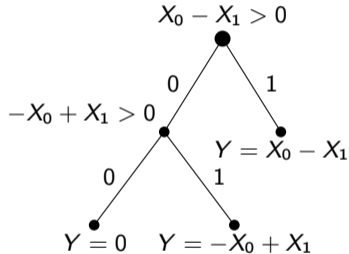


Figure: Simple pruning example

Pruning UNSAT node by

- remove parent and
- connecting sibling subgraph to parent of parent



Comparison of construction methods



Figure: Boxplot ($n = 30$) for the computation time of the NN2EQCDT algorithm for the simple model

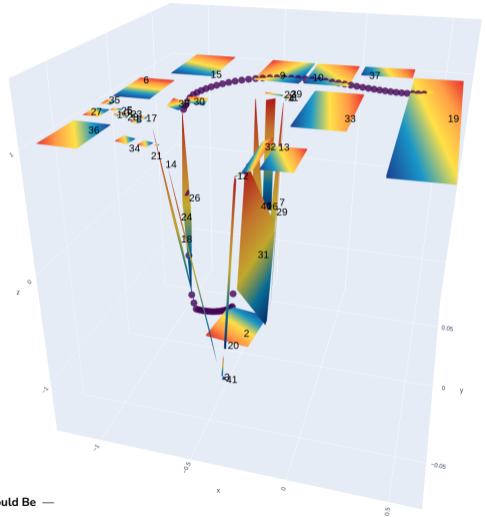
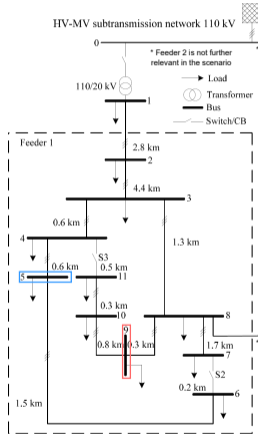
Table: Comparison of results or calculations for the construction of a DT from the simple model without and with compression of the NN2EQCDT algorithm

Pruning	#nodes	Computation time
<input type="checkbox"/>	262143	> 1.5h
<input checked="" type="checkbox"/>	83	9.75s

- Pruning ratio (amount of nodes) of 99.97%



Applications in Practice





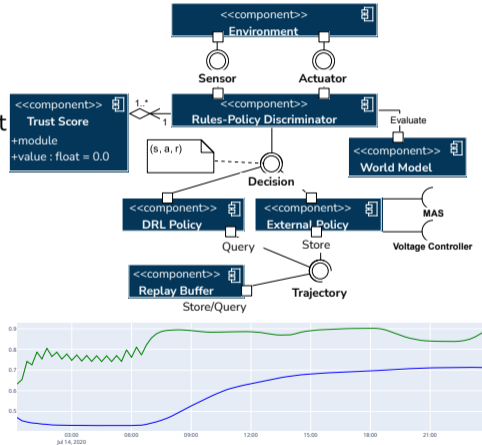
Co-Existence of MAS and DRL

- Hybrid systems out of focus, mostly either DRL, or MARL, or MAS.
- However, any agent isn't alone in its environment!
 - Game-theoretical models focus on a form of interaction (cooperation, competition, conflict, ...), but not on co-existence
 - Underrated in literature: controller conflicts
- Many possible hybrid architectures, e. g.,
 - Hierarchies
 - Imitation Learning
 - Safeguarding (research gap!)



“Cover me:” A Practical Example for Safeguarding MAS

1. Observe MAS, **imitation learn** nominal behavior
2. For every t , internally propose actions
3. Check: MAS action proposal, ARL agent proposal against world model, note projected future states & rewards
4. Update trust by averaging reward over an LTI function
5. Apply actions from proposal with highest trust value
6. Observe state, learn from all three transitions





Deep Reinforcement Learning is not the Only Answer

The state of the art has many nice features:

- Offline learning (learning from domain knowledge)
- Imitation learning (learn existing control strategies by example)
- Model-based and model-free DRL
- eXplainable Reinforcement Learning to explain each action with low computational overhead

... however, this agent is still far from being safe.



“Good” Agents Fail to Apply Learned Strategies





Catastrophic Forgetting on Topology Changes

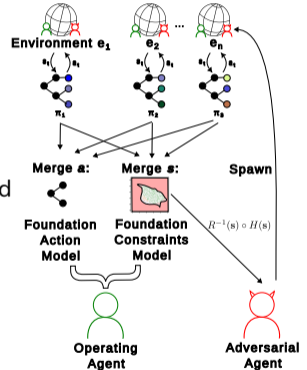
A simple topology change screws the agent completely. Countermeasures:

1. Train the agent on as many scenarios as possible.
2. Verify the DRL agent.
3. Create a *Foundation Model* for actions
4. ... Combine all of the above!



Training Strategy

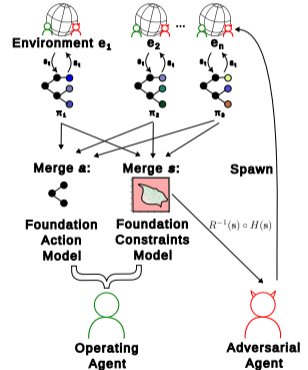
- Take the aut curriculum approach to spawning environments
- Two adversaries: One “spawner” and n “workers”
- Operator agent trains on all of them (traditional multi-worker)
- Adversary spawns environments based on inverted reward and entropy $R^{-1}(s) \circ H(s)$





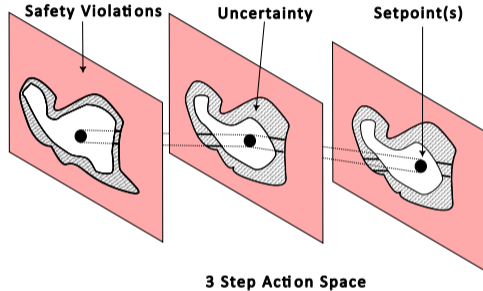
Latent Spaces

- Core idea: Use Graph Neural Networks to learn representations of underlying grids
- Graph space is our feature space: $\mathcal{X} = G(\mathbf{A}, \mathbf{K})$
- Train encoder for latent space representation of all G_i , where i is an environment instance we encountered:
 $\gamma : \mathcal{X} \mapsto \mathcal{L}$
- Use transformer to work directly on latent space
- Result: A foundation model for actions





And Verification...?



- Alongside the Foundation Action Model, train a foundation model for constraints: *Foundation Constraints Model*
- Use the Foundation Constraints Model for N-step verification of trajectories to provide safety guarantees



A Lookout

- The journey towards highly automated grid operation & extension has just begun.
- AI can help testing future grids, be part of certification processes
- AI itself needs safeguards: Rulesets, explainability, and eventually certification, too. (Insurance...?)
- We will see sophisticated agent architectures in the near future.
- If you want to see interesting code, head over to <http://palaestr.ai> or shout out to eric.veith@uol.de!



Bibliography I

- [1] Çağlar Aytekin. “Neural Networks are Decision Trees”. In: *CoRR abs/2210.05189* (2022). [retrieved: 05, 2023], pp. 1–8. arXiv: 2210.05189. URL: <https://arxiv.org/abs/2210.05189>.
- [2] Jimmy Ba and Rich Caruana. “Do deep nets really need to be deep?” In: *Advances in Neural Information Processing Systems 27* (2014), pp. 2654–2662.
- [3] Emilie Frost, Eric Veith, and Lars Fischer. “Robust and Deterministic Scheduling of Power Grid Actors”. In: *2020 7th International Conference on Control, Decision and Information Technologies (CoDIT)*. Vol. 1. 2020, pp. 100–105. DOI: [10.1109/CoDIT49905.2020.9263948](https://doi.org/10.1109/CoDIT49905.2020.9263948).
- [4] Eric Veith, Arlena Wellßow, and Mathias Uslar. “Learning new attack vectors from misuse cases with deep reinforcement learning”. In: *Frontiers in Energy Research* (2023).
- [5] Eric MSP Veith. *Universal Smart Grid Agent for Distributed Power Generation Management*. Logos Verlag Berlin GmbH, 2017.



Bibliography II

- [6] Thomas Wolgast, Eric MSP Veith, and Astrid Nieße. “Towards Reinforcement Learning for Vulnerability Analysis in Power-Economic Systems”. In: *DACH+ Energy Informatics 2021: The 10th DACH+ Conference on Energy Informatics*. Freiburg, Germany, Sept. 2021, pp. 1–20.