

Probabilistic Safety Verification of Neural Policies via Predicate Abstraction

Technical Appendix

Marcel Vinzent¹, Holger Hermanns¹, Jörg Hoffmann^{1,2}

¹Saarland University, Saarland Informatics Campus, Saarbrücken, Germany

²German Research Center for Artificial Intelligence (DFKI), Saarbrücken, Germany
 {vinzent, hermanns, hoffmann}@cs.uni-saarland.de

Additional Experimental Results

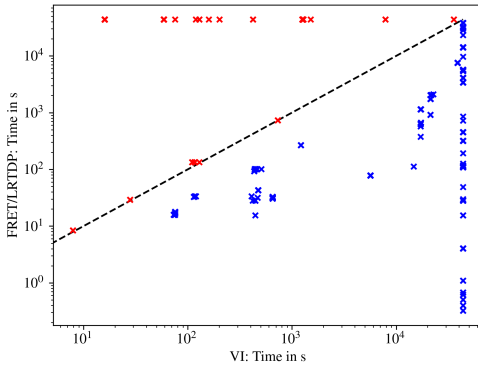


Figure 2: Abstract MaxProb with FRET-LRTDP vs. VI (value iteration). For each benchmark instance there is a data point for the verification time of each bound p_u in $\{0, 0.05, \dots, 0.95\}$ during prob-CEGAR-PPA-inc with $(p_\Delta, p_u^{init}) = (0.05, 0)$, indicating whether FRET-LRTDP or VI is faster.

Abstract MaxProb. Figure 2 compares abstract MaxProb search with FRET-LRTDP vs. VI (value iteration), specifically the verification time of individual p_u during prob-CEGAR-PPA-inc. FRET-LRTDP largely outperforms VI. Constructing the reachable fragment of $\Theta_{\mathcal{P}}^\pi$ is often too expensive. FRET-LRTDP mitigates this thanks to incremental exploration of $\Theta_{\mathcal{P}}^\pi$. VI dominates on some instances. We observe that, thanks to its complete exploration of $\Theta_{\mathcal{P}}^\pi$, VI may select abstract counterexamples with short and realizable paths in early CEGAR iterations.

Automated bound derivation. Table 3 shows results for prob-CEGAR-PPA-inc for different configurations of (p_Δ, p_u^{init}) . Overall, the runtime performance of all configurations is similar, especially for $p_u^{init} = 0$. This indicates (1) there is a sweet spot during CEGAR after which $[P_u^{lo}, P_u^{up}]$ can be tightened heavily (2) since linear search is conducted incrementally the impact of p_Δ is alleviated. That said, $p_\Delta = 0.01$ never drastically faster than all other configuration, but can be significantly slower (e.g., 4 Blocks,

CA, 32) showing that overly conservative p_Δ may still be counter-beneficial. $(p_\Delta = 0.05, p_u^{init} = 1)$ has a smaller coverage (12 instances) than $(p_\Delta = 0.05, p_u^{init} = 0)$ (16 instances). For large (initial) p_u , CeAna is more prone to enumerate an infeasible number of paths until the accumulated abstract probability $Pr(ups_{\mathcal{P}})$ exceeds p_u . In contrast, for small p_u , CeAna may proceed to the computation of $maxPr[ups_{\mathcal{P}}]$; establishing non-trivial P_u^{lo} . However, there are also instances where $p_u^{init} = 1$ dominates $p_u^{init} = 0$.

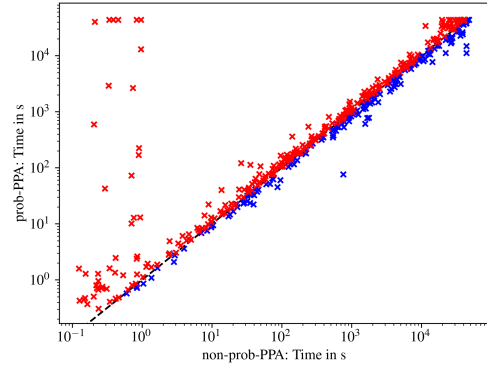


Figure 3: Abstract state space computation prob-PPA vs. non-prob-PPA, indicating whether prob-PPA or non-prob-PPA is faster. Predicate set scales as per Vea (2022).

Computing $\Theta_{\mathcal{P}}^\pi$. Figure 3 compares the time to compute the reachable fragment of $\Theta_{\mathcal{P}}^\pi$ for prob-PPA vs. non-prob-PPA over predicate sets of increasing size. It completes the results presented in the main text (Figure 1), showing results for all benchmarks. In line with our observations in the main text, prob-PPA overall loses little performance in abstract state space building compared to non-prob-PPA. There are outliers for very coarse abstractions, i.e., small predicate sets \mathcal{P} , where non-prob-PPA finishes in seconds while prob-PPA takes hours. For such \mathcal{P} , SMT under π is notoriously expensive (Vea 2022), and, thereby, more prone to outliers: Some SMT calls take drastically longer than others. The additional intricacy in prob-PPA leads to an increase of such outliers.

Benchmark	NN	App	$p_\Delta = 0.01$		$p_\Delta = 0.05$		$p_\Delta = 0.05, p_u^{init} = 1$		$p_\Delta = 0.1$	
			$[P_u^{lo}, P_u^{up}]$	Time	$[P_u^{lo}, P_u^{up}]$	Time	$[P_u^{lo}, P_u^{up}]$	Time	$[P_u^{lo}, P_u^{up}]$	Time
4 Blocks (CI)	16	×	[0.4, 0.41]	15	[0.37, 0.41]	18	[0.41, 1]	-	[0.37, 0.41]	16
	16	✓	[0.41, 0.64]	-	[0.41, 0.65]	-	[0.41, 0.66]	-	[0.41, 0.65]	-
	32	×	[0.27, 0.27]	29	[0.25, 0.27]	33	[0.25, 0.27]	30	[0.2, 0.27]	29
	32	✓	[0.19, 1]	-	[0.19, 1]	-	[0.19, 1]	-	[0.19, 1]	-
	64	×	[0.19, 0.19]	7200	[0.17, 0.19]	7553	[0.19, 0.19]	7157	[0.17, 0.26]	7560
6 Blocks (CI)	16	×	[0.99, 1]	66	[0.95, 1]	79	[0.95, 1]	67	[0.9, 1]	67
	16	✓	[0.99, 1]	313	[0.95, 1]	321	[0.95, 1]	310	[0.9, 1]	315
	32	×	[0.27, 0.27]	12397	[0.27, 0.3]	14227	[0.27, 0.9]	-	[0.27, 0.3]	11762
8 Blocks (CI)	16	×	[0.99, 1]	392	[0.95, 1]	454	[0.95, 1]	383	[0.9, 1]	386
	16	✓	[0.99, 1]	3493	[0.95, 1]	3396	[0, 1]	-	[0.9, 1]	3457
	32	×	[0.99, 1]	27260	[0.95, 1]	31376	[0.95, 1]	25743	[0.9, 1]	25873
8 Puzzle (CI)	16	×	[0.99, 1]	4938	[0.95, 1]	5434	[0.95, 1]	33433	[0.92, 1]	4908
	32	×	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-
	32	✓	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-
	64	✓	[0.95, 1]	-	[0.95, 1]	9630	[0.95, 1]	13135	[0.9, 1]	48
Transport	16	×	[0.1, 1]	-	[0.1, 1]	-	[0, 1]	-	[0.1, 1]	-
	16	✓	[0.1, 1]	-	[0.1, 1]	-	[0, 1]	-	[0.1, 1]	-
	32	×	[0.1, 1]	-	[0.1, 1]	-	[0, 1]	-	[0.1, 1]	-
	32	✓	[0.1, 1]	-	[0.1, 1]	-	[0, 1]	-	[0.1, 1]	-
4 Blocks (CA)	16	×	[0.34, 0.34]	44	[0.34, 0.36]	44	[0.02, 1]	-	[0.26, 0.36]	41
	16	✓	[0.34, 0.42]	-	[0.34, 0.42]	-	[0.03, 1]	-	[0.34, 0.47]	-
	32	×	[0.27, 0.27]	19690	[0.27, 0.28]	1145	[0.27, 0.3]	422	[0.27, 0.3]	785
	32	✓	[0.22, 0.44]	-	[0.17, 0.44]	-	[0.14, 1]	-	[0.14, 1]	-
6 Blocks (CA)	16	×	[0.86, 1]	-	[0.86, 1]	-	[0, 1]	-	[0.86, 1]	-
	32	×	[0.99, 1]	3589	[0.95, 1]	4121	[0.95, 1]	3591	[0.9, 1]	3637
	32	✓	[0.99, 1]	38160	[0.95, 1]	38289	[0.95, 1]	40772	[0.9, 1]	39409
8 Blocks (CA)	16	×	[0.99, 1]	5283	[0.99, 1]	5692	[0.99, 1]	5029	[0.91, 1]	4866
	32	×	[0.84, 1]	-	[0.84, 1]	-	[0, 1]	-	[0.8, 1]	-
8 Puzzle (CA)	16	×	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-
	16	✓	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-
	32	×	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-
	32	✓	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-	[0.1, 1]	-

Table 3: Results for prob-CEGAR-PPA-inc with FRET-LRTDP for different configurations of (p_Δ, p_u^{init}) . Default $p_u^{init} = 0$.

Verification of p_u . Figure 2 compares the verification of a specific bound p_u with (prob-CEGAR)-PPA- p_u vs. the verification time of that p_u during (prob-CEGAR)-PPA-inc. It completes the results presented in the main text (Table 2), showing results for all benchmarks. The results are in line with the observations in the main text. PPA- p_u tends to dominate and can be significantly faster. That said, PPA-inc is overall competitive and can, in some cases, even make verification feasible in the first place (outliers).

Proofs – Probabilistic PPA

Wachter et al. (2007) show that probabilistic predicate abstraction preserves satisfaction of safe PCTL formulae (Bianco and de Alfaro 1995).

Theorem (Wachter et al. 2007). If $\Theta_{\mathcal{P}}$ satisfies a safe PCTL formula Φ , then Θ also satisfies Φ .

Proposition 2 (safety in $\Theta_{\mathcal{P}}$ implies safety in Θ) follows as a special case of Theorem . Proposition 4 (safety in $\Theta_{\mathcal{P}}^\pi$ implies safety in Θ^π) follows analogously.

Probabilistic Abstract State Expansion

Algorithm 3 shows pseudocode for probabilistic abstract state expansion. The algorithm is an adaption of state expansion

for non-probabilistic PPA (Vea 2022). Specifically, we deploy applicability (line 2) and \mathcal{T} -transition tests (line 30) to avoid costly SMT tests under π . Additionally, we lift the latter to the probabilistic case (line 15). EnumPaDist (line 7) shows the enumeration of abstract successor distribution (candidates). We first compute abstract successor states for each individual update $u \in \text{Supp}(\bar{u})$ (line 9) using Vea’s enumeration procedure (line 22). For the abstract transition problem under π (line 5), we apply Vea’s enhancements to efficiently handle π : continuous relaxation of integer state variable to use NN-tailored SMT solvers (Katz et al. 2019) embedded into a branch & bound framework. If the transition exists, it is “processed” for MaxProb search (line 6).

Proofs – prob-CEGAR-PPA

prob-CEGAR (Hermanns et al. 2008) maintains a set of abstract unsafe paths $ups_{\mathcal{P}} \subseteq \text{Path}(A_{\mathcal{P}}, s_{\mathcal{P}}^0, \phi_u)$. Han et al. (2009) show that for all probabilistic counterexamples (A, s^0) – concrete or abstract – there exist a finite subset of unsafe paths (with decreasing probability) that “witnesses” the violation of p_u . We use a slightly more general formulation for any (infinite) subset induced by (A, s^0) . The proof is adapted in a straight-forward manner.

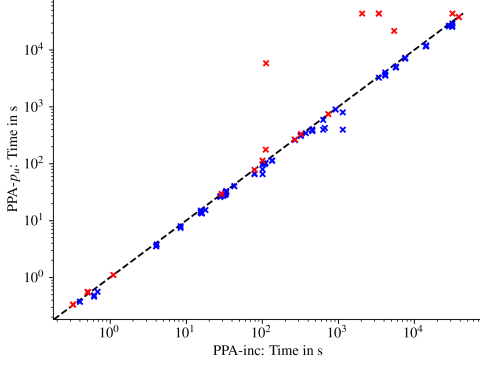


Figure 4: Verification time of $p_u \in 0, 0.05, 0.1, 0.2, 0.3, 0.5$ with (prob-CEGAR)-PPA- p_u vs. during (prob-CEGAR)-PPA-inc, indicating whether PPA- p_u or PPA-inc is faster.

Lemma 9. Let $ups \subseteq Path(A, s^0, \phi_u)$. $Pr(ups) > p_u$ iff there exists a finite set of unsafe paths $ups' \subseteq Path(A, s^0, \phi_u)$ such that $Pr(ups') > p_u$ and for all $\sigma' \in ups', \sigma \in ups \setminus ups'$ it holds $Pr(\sigma') \geq Pr(\sigma)$.

The original result then follows as a corollary, i.e., $Pr(A, s^0, \phi_u) > p_u$ iff there exists finite $ups \subseteq Path(A, s^0, \phi_u)$ (with decreasing probability) such that $Pr(ups) > p_u$.

Termination of CeAna. We first prove some lemma.

Lemma 10. Whenever Algorithm 1 enters loop (line 11 - 21) it holds $P_{max} \leq p_u$ and $can-add(ups_{\mathcal{P}})$.

Proof. The statement holds initially. Moreover, whenever \mathcal{T} -analysis iterates it holds $P_{max} \leq p_u$ (line 19) and $can-add(ups_{\mathcal{P}})$ (line 20). \square

Lemma 11. For any $p_\varepsilon > 0$. The set $ups_{\mathcal{P}} := \{\sigma_{\mathcal{P}} \in Path(A_{\mathcal{P}}, s_{\mathcal{P}}^0, \phi_u) \mid Pr(\sigma_{\mathcal{P}}) \geq p_\varepsilon\}$ is finite.

Proof. Proof by contradiction. For $ups_{\mathcal{P}}$ with infinitely many paths of probability at least p_ε , $Pr(ups_{\mathcal{P}})$ diverges. However, $Pr(ups) \leq Pr(A_{\mathcal{P}}, s_{\mathcal{P}}^0, \phi_u) \leq 1$. Contradiction. \square

Proof of Theorem 5. It suffices to show that \mathcal{T} -analysis (line 11 - 21) terminates. By Lemma 10, in each iteration loop (line 12 - 17) is entered. By Lemma 9 and Lemma 11, $Pr(ups_{\mathcal{P}}) > p_u$ and $\min_{\sigma_{\mathcal{P}} \in ups_{\mathcal{P}}} Pr(\sigma_{\mathcal{P}}) < p_\varepsilon$

for finite $ups_{\mathcal{P}}$ respectively. Hence loop (line 12 - 17) has only finitely many iterations over an invocation of CeAna. Hence, \mathcal{T} -analysis terminates. \square

Correctness of CeAna.

Lemma 12. During an invocation of CeAna it holds $maxPr[ups_{\mathcal{P}}](\Theta) \leq P_{max} \leq Pr(ups)$.

Proof. The statement holds at initialization and is preserved as an invariant: (line 16, 17): Let $ups'_{\mathcal{P}} = ups_{\mathcal{P}} \cup \{\sigma_{\mathcal{P}}\}$. Let $P'_{max} = P_{max} + Pr(\sigma_{\mathcal{P}})$. $maxPr[ups'_{\mathcal{P}}] \leq maxPr[ups_{\mathcal{P}}] + Pr(\sigma_{\mathcal{P}}) \leq P_{max} + Pr(\sigma_{\mathcal{P}}) = P'_{max}$. $P'_{max} = P_{max} + Pr(\sigma_{\mathcal{P}}) \leq Pr(ups_{\mathcal{P}}) + Pr(\sigma_{\mathcal{P}}) = Pr(ups'_{\mathcal{P}})$.

(line 18): Trivial. \square

Lemma 12 shows that P_{max} is an upper bound on the maximal realizable probability as intended.

Lemma 13. While CeAna iterates loop (line 12 - 17) it holds $P_{max} + P_{unused} > p_u$.

Proof. The statement holds initially and is required whenever \mathcal{T} -analysis iterates (line 20). Moreover, while loop (line 12 - 17) continues the sum $P_{max} + P_{unused}$, and thereby the statement, are preserved (line 16, 17). \square

Lemma 13 guarantees that CeAna behaves well in that whenever path computation is invoked (line 13) the $|ups_{\mathcal{P}}| + 1$ -th path actually exists.

Proof of Theorem 6. By π -all realization of $ups_{\mathcal{P}}^*$ (line 25), there exist $s^0 \in [s_{\mathcal{P}}^0] \cap [\phi_0]$ such that for each $\sigma_{\mathcal{P}} \in ups_{\mathcal{P}}^*$, there exists $\sigma \in [\sigma_{\mathcal{P}} \wedge \phi_u](\Theta^\pi)$ with $s^0(\sigma) = s^0$, and thereby $\sigma \in Path([A_{\mathcal{P}}], s^0, \phi_u)$. Let $ups \subseteq Path([A_{\mathcal{P}}], s^0, \phi_u)$ denote the corresponding concrete path set. By construction $Pr(ups) = Pr(ups_{\mathcal{P}}^*) > p_u$ and hence π is unsafe. \square

Termination of prob-CEGAR-PPA.

Lemma 14. Let $S^0 \subseteq \mathcal{S}$ such that $|S^0| \geq 2$. Let

$$\bigcup_{v \in \mathcal{V}} \bigcup_{i \in 1, \dots, |S^0(v)|-1} \{v \leq S_i^0(v)\} \subseteq \mathcal{P},$$

where $S^0(v) = \{s(v) \mid s \in S^0\}$ and $S_i^0(v) \in S^0(v)$ denotes the i -th smallest value. For all distinct $s, t \in S^0$ it holds $s|_{\mathcal{P}} \neq t|_{\mathcal{P}}$.

Proof. $s \neq t$ and thereby $s(v) \neq t(v)$ for some $v \in \mathcal{V}$. Let, w.l.o.g., $s(v) < t(v)$. There exists $(v \leq c) \in \mathcal{P}$ such that $s \models v \leq c$ while $\neg(t \models v \leq c)$. Hence, $s|_{\mathcal{P}} \neq t|_{\mathcal{P}}$. \square

Lemma 15. Whenever CeAna invokes refinement for (ii-a) (line 21) (or for (ii-b) (line 26)) there exist distinct $\sigma_{\mathcal{P}}, \sigma'_{\mathcal{P}} \in ups_{\mathcal{P}}$ such that there exist unsafe concretizations $\sigma \in [\sigma_{\mathcal{P}} \wedge \phi_u](\Theta)$ and $\sigma' \in [\sigma'_{\mathcal{P}} \wedge \phi_u](\Theta)$ (or $\sigma \in [\sigma_{\mathcal{P}} \wedge \phi_u](\Theta^\pi)$ and $\sigma' \in [\sigma'_{\mathcal{P}} \wedge \phi_u](\Theta^\pi)$) with $s^0(\sigma) \neq s^0(\sigma')$.

Proof. (ii-a): Consider maximal \mathcal{T} -realizable subset $ups_{\mathcal{P}}^* \subseteq ups_{\mathcal{P}}$ with maximizing start state $s_{\mathcal{P}}^0$. By inner loop condition (line 12), Lemma 12, and definition of $can-add(ups_{\mathcal{P}})$, refinement for (ii-a) is invoked only once $Pr(ups_{\mathcal{P}}) > p_u$ while $maxPr[ups_{\mathcal{P}}](\Theta) \leq p_u$ (line 19). Hence, $ups_{\mathcal{P}} \setminus ups_{\mathcal{P}}^*$ is non-empty. Let $\sigma_{\mathcal{P}} \in ups_{\mathcal{P}}^*$ and $\sigma'_{\mathcal{P}} \in ups_{\mathcal{P}} \setminus ups_{\mathcal{P}}^*$. We know $\sigma_{\mathcal{P}}$ and $\sigma'_{\mathcal{P}}$ are \mathcal{T} -path realizable. Hence the statement holds for $s_{\mathcal{P}}^0$ and start state $s^0(\sigma')$ of any unsafe concretization $\sigma' \in [\sigma'_{\mathcal{P}} \wedge \phi_u]$.

(ii-b): The result for follows analogously. Since π -all realization fails (line 25), there exist $\sigma_P, \sigma'_P \in \text{ups}_P^*$, which are π -path realizable (line 23), such that $s^0(\sigma) \neq s^0(\sigma')$ for any two unsafe concretization paths $\sigma \in [\sigma_P \wedge \phi_u](\Theta^\pi)$ and $\sigma' \in [\sigma'_P \wedge \phi_u](\Theta^\pi)$. \square

Lemma 16. *In each iteration, CEGAR either terminates or strictly refines the abstraction $\mathcal{P} \subsetneq \mathcal{P}'$ in that there exist $s, t \in \mathcal{S}$ such that $s|_{\mathcal{P}} = t|_{\mathcal{P}}$ in the original abstraction while $s|_{\mathcal{P}'} \neq t|_{\mathcal{P}'}$ in the refined abstraction.*

Proof. The statement holds for refinement of spuriousness (i-a) and (i-b) as per Vea (2023) and by Lemma 14 and Lemma 15 for refinement of spuriousness (ii-a) and (ii-b). \square

Proof of Theorem 7. Let $s, t \in \mathcal{S}$. Since \mathcal{S} is finite (bounded-integer state variables) and by Lemma 16, $s|_{\mathcal{P}} = t|_{\mathcal{P}}$ iff $s = t$ within finitely many iterations ($\Theta_P^\pi \equiv \Theta^\pi$). Then also each abstract path is trivially realizable. Hence, either CEGAR returns SAFE or CeAna returns REAL, and hence CEGAR returns UNSAFE. \square

CEGAR for Automated Bound Derivation

Algorithm 2: prob-CEGAR-PPA-inc.

Input: $\langle \mathcal{V}, \mathcal{L}, \mathcal{O} \rangle, (\phi_0, \phi_u, p_u), p_\Delta \in (0, 1]$.

```

1  $[P_u^{lo}, P_u^{up}] \leftarrow 0, 1$  //  $\maxPr(\Theta^\pi, \phi_0, \phi_u)$  interval.
2  $\mathcal{P} \leftarrow \{\phi_u\}$ 
3 while 1 do
4    $(A_P^\pi, s_P^0) \leftarrow \max\text{-prob}(\langle \mathcal{V}, \mathcal{L}, \mathcal{O} \rangle, (\phi_0, \phi_u, p_u), \mathcal{P}, \pi)$ 
5    $P_u^{up} \leftarrow \min(P_u^{up}, Pr(A_P^\pi, s_P^0, \phi_u))$ 
6   if  $P_u^{up} - P_u^{lo} \leq p_\Delta$  then return
7   while  $P_u^{up} \leq p_u$  do  $p_u \leftarrow p_u - p_\Delta$  // Safe  $p_u$ .
8
9    $\text{ups}_P, P_{max} \leftarrow \emptyset, 0$ 
10  while 1 do
11     $\text{result} \leftarrow \text{CeAna}(A_P^\pi, s_P^0, \text{ups}_P, P_{max})$ 
12    if  $\text{result} = \text{SPURIOUS}$  then break
13     $P_u^{lo} \leftarrow \max(P_u^{lo}, P_{max})$ 
14    if  $P_u^{up} - P_u^{lo} \leq p_\Delta$  then return
    // Unsafe  $p_u$ .
15    while  $P_u^{lo} > p_u$  do  $p_u \leftarrow p_u + p_\Delta$ 
```

Algorithm 2 shows the adapted CEGAR loop with linear search on p_u for automated bound derivation (prob-CEGAR-PPA-inc). CEGAR continues until the interval $[P_u^{lo}, P_u^{up}]$ is sufficiently tightened (line 6, 14). P_u^{up} is updated to the (maximal) abstract unsafety probability (line 5). While p_u is safe, it is decreased (line 7). ups_P and P_{max} (line 9) are maintained over multiple invocations of CeAna within a CEGAR iteration. The linear search pauses and CEGAR iterates whenever CeAna returns SPURIOUS (line 12). If not, P_u^{lo} is updated according to the maximal π -realizable probability (line 13) computed by CeAna. p_u is increased (line 15).

By termination of prob-CEGAR-PPA- p_u , each p_u is proven safe or unsafe within finitely many iterations. When this case occurs, p_u is decreased or increased until $p_u < P_u^{up}$ or $P_u^{lo} \leq p_u$ respectively. Moreover, due to the termination condition (line 6, 14), it also holds $P_u^{lo} < p_u$ and $p_u < P_u^{up}$ respectively. In other words, p_u is updated until $p_u \in [P_u^{lo}, P_u^{up}]$ and hence safety is no longer deducible. Since P_u^{lo} and P_u^{up} are monotonically increasing and decreasing respectively, each p_u is iterated at most once during linear search. Hence, prob-CEGAR-PPA-inc terminates.

Neural Action Policy

A (ReLU) feed-forward *neural network* for Θ is a (real-valued) function

$$f_\pi: \mathcal{S} \rightarrow \mathbb{R}^{d_d}, s \mapsto f_d(\dots f_2(f_1(s))),$$

where d denotes the number of layers in the NN, d_i for $i \in \{1, \dots, d\}$ denotes the size of layer i , and

- $f_1: \mathcal{S} \rightarrow \mathbb{R}^{d_1}, s \mapsto (s(v_\pi^1), \dots, s(v_\pi^{d_1}))$ is the *input layer* function, where $v_\pi^j \in \mathcal{V}$ for $j \in \{1, \dots, d_1\}$ denotes the state variable associated with input neuron j .
- $f_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}, V \mapsto \text{ReLU}(W_i \cdot V + B_i)$, for $i \in \{2, \dots, d-1\}$, is the function of *hidden layer* i . W_i is the weight matrix of layer i , i.e., $(W_i)_{j,k}$ denotes the weight of the output of neuron k in layer $i-1$ to the input of neuron j in layer i . B_i is the bias vector, i.e., $(B_i)_j$ denotes the bias of neuron j in layer i .
- $f_d: \mathbb{R}^{d_{d-1}} \rightarrow \mathbb{R}^{d_d}, V \mapsto W_d \cdot V + B_d$ is the function of output layer d . Here, no ReLU activation is applied.

We distinguish two variants of *neural action policies* implemented by f_π . The *no-filter* neural action policy implemented by f_π is the function $\pi_{napp}: \mathcal{S} \rightarrow \mathcal{L}, s \mapsto \arg\max_{l \in \mathcal{L} \setminus \{\tau\}} f_\pi^l(s)$, where f_π^l denotes the output of f_π associated with l . The *app-filter* neural action policy implemented by f_π is the function $\pi_{app}: \mathcal{S} \rightarrow \mathcal{L}$,

$$s \mapsto \begin{cases} \arg\max_{\{l \in \mathcal{L} \setminus \{\tau\} \mid \exists o \in \mathcal{O}_l: s \models g(o)\}} f_\pi^l(s) & \exists o \in \mathcal{O}: l(o) \neq \tau \\ \tau & \wedge s \models g(o) \\ & \text{otherwise} \end{cases}$$

where $\mathcal{O}_l = \{o \in \mathcal{O} \mid l(o) = l\}$ and τ is a special label, defined to be selected (by π_{app}) iff there does not exist a label $l \in \mathcal{L} \setminus \{\tau\}$ that is applicable in s .

SMT Encodings

We provide a specification of the various SMT checks deployed during PPA computation as well as CeAna. The encodings are largely compositional, and hence presented in a modular manner.

Encodings without π . Given a copy V of state variables \mathcal{V} , we denote the encoding of a constraint $\phi \in C(\mathcal{V})$ over V by $\Phi(\phi, V)$. The *state variable bound encoding* $\Phi(V)$ is the conjunction $\bigcup_{v \in V} lo(v) \leq v \wedge v \leq up(v)$, where $lo(v)$

and $up(v)$ denote the lower and upper bound of the corresponding state variable respectively. The *abstract state encoding* of $s_p \in \mathcal{S}_p$, denoted $\Phi(s_p, V)$, is the conjunction

$$\bigcup_{p \in s_p} \begin{cases} \Phi(p, V) & s_p(p) = 1 \\ \neg \Phi(p, V) & s_p(p) = 0. \end{cases}$$

For update $u \in U$ let V^u denote a u -indexed copy of V . The *update encoding* of u , denoted $\Phi(u, V, V^u)$, is the conjunction $\bigwedge_{v \in V} v^u = \Phi(u(v), V)$, where $v^u \in V^u$ denotes the u -copy of v . The *probabilistic update encoding* of \bar{u} , denoted $\Phi(\bar{u}, V)$, is the conjunction

$$\bigwedge_{u \in \text{Supp}(\bar{u})} \Phi(u, V, V^u). \text{ The abstract distribution encoding of } \mu_p \in \text{Dist}(U \times \mathcal{S}_p), \text{ denoted } \Phi(\mu_p, V), \text{ is the conjunction}$$

$$\bigwedge_{(u, s'_p) \in \text{Supp}(\mu_p)} \Phi(s'_p, V^u).$$

The (no- π) **abstract transition problem encoding** of $(s_p, o, \mu_p) \in \mathcal{S}_p \times \mathcal{O} \times \text{Dist}(U \times \mathcal{S}_p)$, denoted $\Phi((s_p, o, \mu_p), \mathcal{T}_p, V)$ conjoins: $\Phi(V)$, $\bigwedge_{u \in \text{Supp}(\bar{u}(o))} \Phi(V^u)$, $\Phi(s_p, V)$, $\Phi(g(o), V)$, $\Phi(\bar{u}(o), V)$, $\Phi(\mu_p, V)$. We assume, w.l.o.g., that $\bar{u}(o)$ and μ_p agree on U , i.e., $(u, s'_p) \in \text{Supp}(\mu_p)$ only if $u \in \text{Supp}(\bar{u}(o))$, and for each $u \in \text{Supp}(\bar{u}(o))$ there exists exactly one $s'_p \in \mathcal{S}_p$ such that $(u, s'_p) \in \text{Supp}(\mu_p)$ with $\mu_p(u, s'_p) = \bar{u}(o)(u)$.

For index i let V^i denote a i -indexed copy of V . The **\mathcal{T} -path realization encoding (i-a)** of abstract unsafe path $\sigma_p \in \text{Path}(\Theta_p, \phi_u)$, denoted $\Phi(\sigma_p, \Theta, V)$ conjoins:

$$\bigwedge_{0 \leq i \leq |\sigma_p|} \Phi(V^i), \quad \Phi(\phi_0, V^0), \quad \bigwedge_{0 \leq i \leq |\sigma_p|} \Phi(s_p^i(\sigma_p), V^i), \\ \bigwedge_{0 \leq i < |\sigma_p|} \Phi(g^i(\sigma_p), V^i) \quad \wedge \quad \bigwedge_{0 \leq i < |\sigma_p|} \Phi(u^i(\sigma_p), V^i, V^{i+1}), \\ \Phi(\phi_u, V^{|\sigma_p|}).$$

Note that (1) $\Phi(\neg \phi_u, V^i)$ for $0 \leq i < |\sigma_p|$ is entailed by over-approximation of σ_p and hence *not* explicitly encoded; (2) $\Phi(\sigma_p, \Theta, V)$ – and thereby all compositional encodings thereof – is *distribution-insensitive* in that it constrains s_p^i but not μ_p^i . A solution to this encoding is a concrete path σ such that $|\sigma| = |\sigma_p|$ and $o^i(\sigma) = o^i(\sigma_p)$, $u^i(\sigma) = u^i(\sigma_p)$, $s^i(\sigma) \in [s_p^i(\sigma_p)]$ – but not necessarily $\mu^i(\sigma) \in [\mu_p^i(\sigma_p)]$ – for $i \in \{0, \dots, |\sigma|\}$. This is a commonplace enhancement (Hermanns et al. 2008), which allows off-the-shelf checking and refining of path spuriousness without adaptations for spuriousness induced by μ_p^i . It preserves correctness in that when counterexample analysis derives *realizable*, there exists $ups \subseteq \text{Path}(A, s^0, \phi_u)$ such that $Pr(ups) > p_u$ for some adversary A but not necessarily concretization $[A_p]$.

For abstract path σ_p let b_{σ_p} denote a fresh binary variable. The **\mathcal{T} -probabilistic maximal realization encoding (ii-a)** of abstract unsafe path set $ups_p \subseteq \text{Path}(\Theta_p, \phi_u)$, denoted $\Phi_{max}(ups_p, \Theta, V)$ has **maximization objective** $\sum_{\sigma_p \in ups_p} Pr(\sigma_p) \cdot b_{\sigma_p}$, and conjoins

$$\bigwedge_{\sigma_p \in ups_p} \left(\neg b_{\sigma_p} \vee (\Phi(\sigma_p, \Theta, V^{\sigma_p}) \wedge \bigwedge_{v \in V} v = v^{\sigma_p, 0}) \right).$$

$\Phi(ups_p, \Theta, V)$ maximizes the probability mass of paths that are realizable from some common concrete start state s_*^0 encoded over V . Intuitively, b_{σ_p} constrains that, only if σ_p

is realizable from s_*^0 , it contributes to the maximization.

Encodings with π . Given a copy V of state variables V and a neural network f_π , x_j^i and z_j^i denote *real-valued auxiliary variables* of neuron (i, j) , which are implicitly V -indexed. The *NN encoding* of neural network f_π , denoted $\Phi(f_\pi, V)$, conjoins: $\bigwedge_{1 \leq j \leq d_1} x_j^1 = v_\pi^j$, $\bigwedge_{2 \leq i \leq d-1} \bigwedge_{1 \leq j \leq d_i} z_j^i =$

$$\sum_{k=1}^{d_{i-1}} (W_i)_{j,k} \cdot x_k^{i-1} + (B_i)_j, \quad \bigwedge_{2 \leq i \leq d-1} \bigwedge_{1 \leq j \leq d_i} x_j^i = \text{ReLU}(z_j^i),$$

$$\text{and } \bigwedge_{1 \leq j \leq d_d} x_j^d = \sum_{k=1}^{d_{d-1}} (W_d)_{j,k} \cdot x_k^{d-1} + (B_d)_j. \text{ Note that this}$$

encoding is solver specific in that it assumes a special construct for ReLU constraints (Katz et al. 2019).

Let $l \in \mathcal{L} \setminus \{\tau\}$. Let x_l^d denote the auxiliary variable of the output neuron associated with label l . The *policy selection encoding* of π and l , denoted $\Phi(\pi, l, V)$, conjoins $\Phi(f_\pi, V)$ and $\bigwedge_{l' \in \mathcal{L} \setminus \{l, \tau\}} x_l^d > x_{l'}^d$ (if π is a no-filter policy) or

$$\bigwedge_{l' \in \mathcal{L} \setminus \{l, \tau\}} \left(x_l^d > x_{l'}^d \vee \neg \bigvee_{o \in \mathcal{O}_l} \Phi(g(o), V) \right) \text{ (if } \pi \text{ is an app-filter policy).}$$

The **π -abstract transition problem encoding** of $(s_p, o, \mu_p) \in \mathcal{S}_p \times \mathcal{O} \times \text{Dist}(U \times \mathcal{S}_p)$, denoted $\Phi((s_p, o, \mu_p), \mathcal{T}_p^\pi, V)$, is the conjunction $\Phi((s_p, o, \mu_p), \mathcal{T}_p, V) \wedge \Phi(\pi, l(o), V)$.

The **π -path realization encoding (i-b)** of abstract unsafe path $\sigma_p \in \text{Path}(\Theta_p^\pi, \phi_u)$, denoted $\Phi(\sigma_p, \Theta^\pi, V)$ is the conjunction $\Phi(\sigma_p, \Theta, V) \wedge \bigwedge_{0 \leq i < |\sigma_p|} \Phi(\pi, l(o^i(\sigma_p)), V^i)$.

The **π -all realization encoding (ii-b)** of abstract unsafe path set $ups_p \subseteq \text{Path}(\Theta_p, \phi_u)$, denoted $\Phi(ups_p, \Theta^\pi, V)$ conjoins $\bigwedge_{\sigma_p \in ups_p} \Phi(\sigma_p, \Theta^\pi, V^{\sigma_p}) \wedge \bigwedge_{v \in V} v = v^{\sigma_p, 0}$.

Automata Network

Veas code base supports PTS encoded in the automata language JANI (Budde et al. 2017). We now show the automata network structure underlying the generic PTS description $\langle \mathcal{V}, \mathcal{L}, \mathcal{O} \rangle$ in the main text.

A *network of automata* is a tuple $\langle \mathcal{V}, \mathcal{L}, \mathcal{A}, \text{Sync} \rangle$, where

- \mathcal{V} is a finite set of variables, each with a bounded integer domain.
- \mathcal{L} is a finite set of labels, excluding the *silent label* $\tau \notin \mathcal{L}$.
- \mathcal{A} is a finite set of automata.
- $\text{Sync} \subseteq (\mathcal{A} \rightarrow \mathcal{L}) \times (\mathcal{L} \cup \{\tau\})$ is a finite set of *synchronization constraints*.

An *automaton* a is a tuple $\langle \text{Loc}, E \rangle$, where Loc is a non-empty finite set of *locations*, and E is a finite set of *edges* of a . An *edge* e of a is a tuple $(l, \text{loc}_s, g, \bar{u}) \in E$ with

- label $l \in \mathcal{L}$ for labeled edges or $l = \tau$ for silent edges.
- source location $\text{loc}_s \in \text{Loc}$.
- guard $g \in C(\mathcal{V})$.

- probabilistic update $\bar{u} \in \text{Dist}(\text{Loc} \times (\mathcal{V} \rightarrow \text{Exp}(\mathcal{V})))$, where each $(\text{loc}_d, u) \in \text{Supp}(\bar{u})$ is composed of a destination location loc_d and a partial variable update u with $\text{dom}(u) \subseteq \mathcal{V}$.

An automaton consists of a set of *locations* connected by *edges*. Each edge links from a *source* location to finitely many *destination* locations, each weighted with some non-zero probability. An edge can be taken, i.e., the automaton can transit from the source to some destination of the edge, only if its *guard* evaluates to true over the current variable assignment. If an edge is taken, the variables are updated according to the *update* associated with the destination location. While *silent* edges can be taken independently, *labeled* edges can only be taken as part of a *synchronization*. Here, a synchronization constraint specifies for each automaton – possibly from a subset of participating automata – an action label. Additionally, it specifies the label of the synchronization. Under this label, the participating automata may synchronize taking edges whose label combination agrees with the synchronization constraint.

The state space description $\langle \mathcal{V}, \mathcal{L}, \mathcal{O} \rangle$ of an automata network $\langle \mathcal{V}, \mathcal{L}, \mathcal{A}, \text{Sync} \rangle$ is obtained as follows:

- $\mathcal{V} = \mathcal{V} \cup \{v_a \mid a \in \mathcal{A}\}$, where v_a is the location variable of automaton a . $D(v_a) = \text{Loc}(a)$ is interpreted as a bounded-integer interval.
- $\mathcal{L} = \mathcal{L} \cup \{\tau\}$.
- \mathcal{O} contains an operator (l, g, \bar{u})
 - for each silent edge $(\tau, \text{loc}_s, g, u, \text{loc}_d) \in E(a)$ in each automaton $a \in \mathcal{A}$ where $g := v_a = \text{loc}_s \wedge g$, $l := \tau$, and $u \in \text{Supp}(\bar{u})$ for each $(\text{loc}_d, u) \in \text{Supp}(\bar{u})$ with $u = u \cup \{v_a \mapsto \text{loc}_d\} \cup \{v \mapsto v \mid v \in \mathcal{V} \setminus (\text{dom}(u) \cup \{a\})\}$ and $\bar{u}(u) = \bar{u}(u)$,
 - for each synchronization constraint $(\lambda, l) \in \text{Sync}$ with $\text{dom}(\lambda) = \{a^1, \dots, a^n\}$ and each combination of edges $e^1 \in E(a^1), \dots, e^n \in E(a^n)$ such that $e^i = (\lambda(a^i), \text{loc}_s^i, g^i, \bar{u}^i)$ for $i \in \{1, \dots, n\}$, where $l = l$, $g := \bigwedge_{i=1}^n (v_{a^i} = \text{loc}_s^i \wedge g^i)$, and $u \in \text{Supp}(\bar{u})$ for each combination $u^1 \in \text{Supp}(\bar{u}^1), \dots, u^n \in \text{Supp}(\bar{u}^n)$ with $u = (\bigcup_{i=1}^n u^i \cup \{v_{a^i} \mapsto \text{loc}_d^i\}) \cup \{v \mapsto v \mid v \in \mathcal{V} \setminus \bigcup_{i=1}^n (\text{dom}(u^i) \cup \{v_{a^i}\})\}$ and $\bar{u}(u) = \prod_{i=1}^n \bar{u}^i(u)$.

The silent label τ models edges and synchronizations that can be taken independent of a control policy π . That is, $\mathcal{T}^\pi = \{(s, o, \mu) \in \mathcal{T} \mid \pi(s) = l(o) \vee l(o) = \tau\}$. We silently ignore this special case in our approach. Its implementation, e.g., for abstract state expansion or realization checks, is straight forward.

Encoding the Policy Into the Automata Network

In our approach, the neural policy π is a distinctive component of the safety problem. Alternatively, the composition of π and π -controlled PTS (modeled as an automata network), can again be encoded as an automata network. This enables us to feed the neural policy safety analysis into off-the-shelf probabilistic model checkers.

Many SOA probabilistic model checkers, e.g., STORM, focus on discrete state systems in that they do not fully support real-valued state variables. We thus provide two distinct encodings of the neural policy, one with *real-valued variables* and another that *discretizes* the NN structure with finite precision n based on *rational arithmetic*.

Encoding with real-valued variables. Let π be a neural action policy, and let $n = \langle \mathcal{V}, \mathcal{L}, \mathcal{A}, \text{Sync} \rangle$ be the π -controlled automata network. The $(n \mid \pi)$ -composition is the automata network $\langle \mathcal{V} \cup \mathcal{V}_\pi, \mathcal{L}, \mathcal{A} \cup \{a_\pi\}, \text{Sync}_\pi \rangle$ where

- \mathcal{V}_π is the set of *neural policy variables*. $x_j \in \mathcal{V}_\pi$ for each neuron index $j \in 1, \dots, d_i$ over the **even-index** layers $i \in 2, \dots, d$, where $D(x_j)$ is the **least tight** interval as per interval arithmetic (Moore et al. 2009) over all corresponding neurons (i, j) . $y_j \in \mathcal{V}_\pi$ for each neuron index $j \in 1, \dots, d_i$ over the **odd-index** layers $i \in 3, \dots, d$, with $D(y_j)$ analogously to even-index variables. If π is **app-filter**, then $v_l \in \mathcal{V}_\pi$ for each $l \in \mathcal{L}$, where $D(v_l) = \{0, 1\}$.
- $a_\pi = \langle \text{Loc}_\pi, E_\pi \rangle$ is the *policy automaton*.
- $\text{Sync}_\pi = \{(\lambda \cup \{a_\pi \mapsto l\}, l) \mid (\lambda, l) \in \text{Sync} \wedge l \neq \tau\} \cup \{(\lambda, \tau) \in \text{Sync}\}$ is the *policy-annotated synchronization set*.

The policy automaton $a_\pi = \langle \text{Loc}_\pi, E_\pi \rangle$ is composed of locations $\text{Loc}_\pi = \{\text{loc}_i \mid i \in \{1, \dots, d\}\}$ and edges E_π ,

- for each layer $1 < i \leq d-1$ a deterministic *to-hidden-layer* edge $(\tau, \text{loc}_{i-1}, 1, \bar{u}) \in E_\pi$ where $\text{Supp}(\bar{u}) = \{(\text{loc}_i, u)\}$ and $u(z_j^i) = \text{ReLU}(\sum_k^{d_i-1} (W_i)_{j,k} \cdot z_k^{i-1} + (B_i)_j)$ for each j in $1, \dots, d_i$. z_j^i denotes v_π^j for the input layer $i = 1$, x_j for even layers i and y_j for odd layers $i > 1$.
- a deterministic *to-output-layer* edge $(\tau, \text{loc}_{d-1}, 1, \bar{u}) \in E_\pi$ where $\text{Supp}(\bar{u}) = \{(\text{loc}_d, u)\}$ and $u(z_j^d) = \sum_k^{d_d-1} (W_d)_{j,k} \cdot z_k^{d-1} + (B_d)_j$ for each j in $1, \dots, d_d$, and if π is **app-filter**, then $u(v_l) = \bigvee_{o \in \mathcal{O}_l} g(o)$ for each $l \in \mathcal{L}$ with \mathcal{O}_l as induced by n .
- for each $l \in \mathcal{L}$ a *synchronization* edge $(l, \text{loc}_d, g, \bar{u}) \in E_\pi$ with $\text{Supp}(\bar{u}) = \{(\text{loc}_1, u)\}$, $u = \emptyset$ and $g = \bigwedge_{l' \in \mathcal{L} \setminus \{l\}} z_{l'} > z_l \vee \neg v_{l'}$ for **no-filter** and alternatively $g = \bigwedge_{l' \in \mathcal{L} \setminus \{l\}} z_{l'} > z_{l'} \vee \neg v_{l'}$ for **app-filter**, where z_l denotes the encoding variable of the output neuron associated with l .

The encoding distinguishes between neural policies *with* and *without applicability filter*. π is encoded as an additional automaton in the network. Labeled transitions are controlled by π via synchronization constraints. We introduce additional state variables to encode the NN output computation.¹ The layer-wise structure allows to introduce two sets of real-valued variables only, one for even-index layers and one for odd-index layers, rather than one variable per neuron in the network. Given the domain bounds of the state

¹A compact symbolic encoding of the NN output, i.e., recursively inlining the layer structure, results in an infeasible blow up in encoding size due to the piecewise-linear ReLU activations.

variable inputs, interval arithmetic (Moore et al. 2009) can be applied to derive bounds on each neuron in the NN. The domain of z_j^i is then the least tight interval, i.e., the smallest lower bound and the largest upper bound over all corresponding neurons. Alternatively, if supported by the model checker, e.g., ePMC and MODEST, the domain can be unbounded. The NN output computation is encoded by the silent and deterministic *to-hidden-layer* and *to-output-layer* edges. $ReLU(z)$ is syntactic sugar for $\max(z, 0)$. Alternatively, one could also use an *if-then-else* construct. If strictly limited to linear expressions, the piecewise-linear case distinction ($z > 0$) could be encoded via two distinct edges.

For applicability-filtering, we introduce an additional set of binary applicability flags. These are set by the *to-output-layer* edge. In the start constraint ϕ_0 of a safety property, each variable $v \in V_\pi$ is fixed to an arbitrary value. The initial location is loc_0 .

Discretized encoding with rational arithmetic. The discretized encoding has finite precision $n \in \mathbb{N}$ and uses rational arithmetic. It maintains a unified denominator $q = 10^n$. NN weights and biases, in floating representation, are approximated with precision n by a rational number with denominator q . For instance, let $n = 4$, then $2.12341\dots$ is approximated by $\frac{21234}{10000}$. The denominator is maintained implicitly.²

The neuron policy variables $z \in V_\pi$ encode the numerators of the underlying rational values. $D(z)$ is an integer interval. The lower and upper bound is the numerator of the *floor* and *ceil* value of the real-valued interval bounds transformed to denominator q respectively.

The encoding of *to-hidden-layer* edges for layer $i < d - 1$ and neuron j is adapted from the real-valued encoding as follows

$$u(z_j^i) = ReLU\left(\sum_k^{d_{i-1}} \left\lfloor \frac{\text{num}((W_i)_{j,k}) \cdot z_k^{i-1}}{q} \right\rfloor + \text{num}((B_i)_j)\right)$$

where $\text{num}(w)$ denotes the numerator of the rational number approximation of w with precision n . The encoding of the *to-output-layer* edge is adapted accordingly. The division by q preserves the unified denominator in that $ReLU\left(\frac{p_1}{q} \cdot$

$\frac{p_2}{q} + \frac{p_0}{q}\right) \equiv ReLU\left(\frac{p_1 \cdot p_2}{q} + \frac{p_0}{q}\right) \equiv ReLU\left(\frac{p_1 \cdot p_2 + p_0}{q}\right) \equiv \frac{ReLU\left(\frac{p_1 \cdot p_2 + p_0}{q}\right)}{q}$. Note that, due to the division, the update encoding is non-linear. That said, the SOA model checkers we deploy all support this language fragment. The encoding of the synchronization edge remains syntactically unmodified, with respect to the real-valued encoding, since $\frac{z}{q} > \frac{z'}{q} \equiv z > z'$.

Due to the finite precision approximation on the encoding level, the discretization does not guarantee to mimic the true policy-restricted system faithfully. Imprecision may accumulate over systems executions. Given the checkers' limitations, this is the best possible basis for a performance com-

parison. In particular, the size of the NN-controlled system is equal across encodings.

In our evaluation, we experiment with $n = 3, \dots, 7$. We obtain similar results independent of n : All configuration exceed either time or memory limits on all problem instances. The only exception is STORM's sparse engine which successfully terminates on 6 Blocks (CI) NN 16 (without app-filter) after 150 to 180 seconds for $n \leq 5$. In the main text, we report the runtime for $n = 5$.

JANI-to-PRISM translation. The PRISM and the PET support input models in PRISM's own input language only. We translate automata networks encoded in JANI to PRISM in a straight-forward manner. Each (JANI) automaton becomes a (PRISM) *module*, with the automaton location encoded as an additional bounded-integer module variable. In PRISM all modules that share a common label must participate in the synchronization under this label. While this is a strict subset of the synchronization constraints supported in JANI, it is sufficient for our purposes (on our benchmarks).

Start state enumeration. The probabilistic model checkers that we experiment with support sets of start states compactly represented by a constraint ϕ_0 . The only exception is MODEST. Here, we encode *rejection-based* enumeration of ϕ_0 directly into the automata network. In a initialization step, the automata network non-deterministically assigns a value to each variable $v \in V$. If the resulting state s satisfies ϕ_0 , the system execution proceeds from s . Otherwise s is terminal.

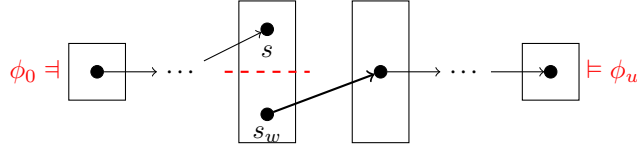
References

- Bianco, A.; and de Alfaro, L. 1995. Model Checking of Probabilistic and Nondeterministic Systems. In *Foundations of Software Technology and Theoretical Computer Science, 15th Conference, Bangalore, India, December 18-20, 1995, Proceedings*, volume 1026 of *LNCS*, 499–513. Springer.
- Budde, C. E.; Dehnert, C.; Hahn, E. M.; Hartmanns, A.; Junges, S.; and Turrini, A. 2017. JANI: Quantitative Model and Tool Interaction. In *TACAS*, volume 10206 of *LNCS*, 151–168.
- Han, T.; Katoen, J.; and Damman, B. 2009. Counterexample Generation in Probabilistic Model Checking. *IEEE Trans. Software Eng.*, 35(2): 241–257.
- Hermanns, H.; Wachter, B.; and Zhang, L. 2008. Probabilistic CEGAR. In *CAV*, volume 5123 of *LNCS*, 162–175. Springer.
- Katz, G.; Huang, D. A.; Ibeling, D.; Julian, K.; Lazarus, C.; Lim, R.; Shah, P.; Thakoor, S.; Wu, H.; Zeljic, A.; Dill, D. L.; Kochenderfer, M.; and Barrett, C. 2019. The Marabou Framework for Verification and Analysis of Deep Neural Networks. In *CAV*, volume 11561 of *LNCS*, 443–452. Springer.
- Moore, R. E.; Kearfott, R. B.; and Cloud, M. J. 2009. *Introduction to Interval Analysis*. SIAM.
- Vincent, M.; Sharma, S.; and Hoffmann, J. 2023. Neural Policy Safety Verification via Predicate Abstraction: CEGAR. In *AAAI*, 15188–15196. AAAI Press.

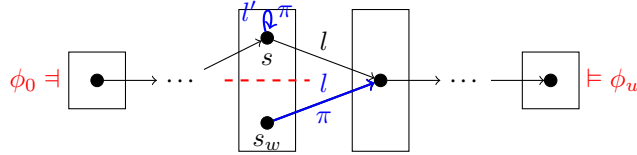
²Explicitly encoding the denominator via an additional variable for each neuron value variable is bound to overflow issues due to denominator unification for addition.

Vinzent, M.; Steinmetz, M.; and Hoffmann, J. 2022. Neural Network Action Policy Verification via Predicate Abstraction. In *ICAPS*. AAAI Press.

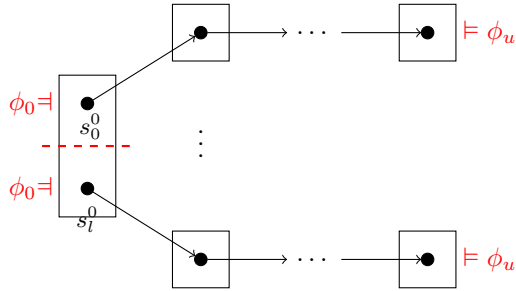
Wachter, B.; Zhang, L.; and Hermanns, H. 2007. Probabilistic Model Checking Modulo Theories. In *Fourth International Conference on the Quantitative Evaluation of Systems (QEST 2007)*, 17-19 September 2007, Edinburgh, Scotland UK, 129–140. IEEE Computer Society.



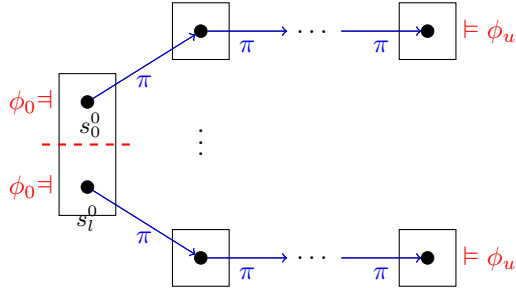
(a) (i-a) \mathcal{T} -path spuriousness: Path is not realizable in Θ .



(b) (i-b) π -path spuriousness: Path is realizable in Θ but **not** in Θ^π .



(c) (ii-a) \mathcal{T} -probabilistic spuriousness: Individual paths of $ups_{\mathcal{P}}$ are realizable in Θ but **not** from a common start state $s^0 \models \phi_0$.



(d) (ii-b) π -probabilistic spuriousness: $ups_{\mathcal{P}}$ is realizable in Θ and individual paths of $ups_{\mathcal{P}}$ are realizable in Θ^π , but they are **not** realizable from a common start state $s^0 \models \phi_0$ in Θ^π .

Figure 5: Sources of spuriousness in prob-CEGAR-PPA.

Algorithm 3: Probabilistic abstract state expansion.

Input: $s_{\mathcal{P}} \in \mathcal{S}_{\mathcal{P}}$.

```

1 for each  $o \in \mathcal{O}$  do
  // Applicability test:
2   if  $\neg(s_{\mathcal{P}} \models g(o))$  then continue
3   for each  $\mu_{\mathcal{P}} \in \text{EnumPaDist}(s_{\mathcal{P}}, o)$  do
4      $g, \bar{u}, l \leftarrow g(o), \bar{u}(o), l(o)$ 
5     //  $\pi$ -transition test:
6     if  $\exists s \in [s_{\mathcal{P}}]: s \models g \wedge s[\bar{u}] \in [\mu_{\mathcal{P}}] \wedge \pi(s) = l$  then
7       process  $(s_{\mathcal{P}}, o, \mu_{\mathcal{P}})$ 

7 Procedure EnumPaDist  $(s_{\mathcal{P}}, o)$ :
8   for each  $u \in \text{Supp}(\bar{u}(o))$  do
9      $S'_{\mathcal{P}}(u) \leftarrow \text{EnumPaSuc}(s_{\mathcal{P}}, o, u)$ 
10   $\text{Dist}_{\mathcal{P}}, \mu_{\mathcal{P}} \leftarrow \emptyset, \{(u, s'_{\mathcal{P}}) \mapsto 0 \mid (u, s_{\mathcal{P}}) \in U \times \mathcal{S}_{\mathcal{P}}\}$ 
11  EnumPaDist  $(\text{Dist}_{\mathcal{P}}, \mu_{\mathcal{P}})$ 
12  return  $\text{Dist}_{\mathcal{P}}$ 

13 Procedure EnumPaDist  $(\text{Dist}_{\mathcal{P}}, \mu_{\mathcal{P}})$ :
14   if  $\text{Supp}(\mu_{\mathcal{P}}) = \text{Supp}(\bar{u}(o))$  then
15     //  $\mathcal{T}$ -transition test:
16     if  $\exists s \in [s_{\mathcal{P}}]: s \models g(o) \wedge s[\bar{u}(o)] \in [\mu_{\mathcal{P}}]$  then
17        $\text{Dist}_{\mathcal{P}} \leftarrow \text{Dist}_{\mathcal{P}} \cup \{\mu_{\mathcal{P}}\}$ 
18     return
19   Select  $u \in \text{Supp}(\bar{u}(o)) \setminus \text{Supp}(\mu_{\mathcal{P}})$ 
20   for each  $s'_{\mathcal{P}} \in S'_{\mathcal{P}}(u)$  do
21      $\mu'_{\mathcal{P}} \leftarrow \mu_{\mathcal{P}} \cup \{(u, s'_{\mathcal{P}}) \mapsto \bar{u}(o)(u)\}$ 
22     EnumPaDist  $(\text{Dist}_{\mathcal{P}}, \mu'_{\mathcal{P}})$ 

// Adopted from Vea (2022):
23 Procedure EnumPaSuc  $(s_{\mathcal{P}}, o, u)$ :
24    $S'_{\mathcal{P}}, s'_{\mathcal{P}} \leftarrow \emptyset, \emptyset$ 
25   // Operator entailment:
26   for each  $(p, b) \in \mathcal{P} \times \{0, 1\}$  do
27     if  $\forall s \in [s_{\mathcal{P}}] \cap [g(o)]: p(s[\bar{u}]) = b$  then
28        $s'_{\mathcal{P}}(p) \leftarrow b$ 
29   EnumPaSuc  $(S'_{\mathcal{P}}, s'_{\mathcal{P}})$ 
30   return  $S'_{\mathcal{P}}$ 

28 Procedure EnumPaSuc  $(S'_{\mathcal{P}}, s'_{\mathcal{P}})$ :
29   if  $\text{dom}(s'_{\mathcal{P}}) = \mathcal{P}$  then
30     // Non-prob.  $\mathcal{T}$ -transition test:
31     if  $\exists s \in [s_{\mathcal{P}}]: s \models g(o) \wedge s[\bar{u}] \in [s'_{\mathcal{P}}]$  then
32        $S'_{\mathcal{P}} \leftarrow S'_{\mathcal{P}} \cup \{s'_{\mathcal{P}}\}$ 
33     return
34   Select  $p \in \mathcal{P} \setminus \text{dom}(s'_{\mathcal{P}})$ 
35   for each  $b \in \{0, 1\}$  do
36      $s''_{\mathcal{P}} \leftarrow s'_{\mathcal{P}} \cup \{p \mapsto b\}$ 
37     // Predicate entailment:
38     for each  $(p', b') \in (\mathcal{P} \setminus \text{dom}(s''_{\mathcal{P}})) \times \{0, 1\}$  do
39       if  $\forall s \in \mathcal{S}: p(s) = b \rightarrow p'(s) = b'$  then
40          $s''_{\mathcal{P}}(p') \leftarrow b'$ 
41     EnumPaSuc  $(S'_{\mathcal{P}}, s''_{\mathcal{P}})$ 

```
