
Automatic Risk Adaptation in Distributional Reinforcement Learning

Frederik Schubert*

Leibniz University Hannover
schubert@tnt.uni-hannover.de

Theresa Eimer*

Leibniz University Hannover
eimer@tnt.uni-hannover.de

Bodo Rosenhahn

Leibniz University Hannover
rosenhahn@tnt.uni-hannover.de

Marius Lindauer

Leibniz University Hannover
lindauer@tnt.uni-hannover.de

Abstract

The use of Reinforcement Learning (RL) agents in practical applications requires the consideration of suboptimal outcomes, depending on the familiarity of the agent with its environment. This is especially important in safety-critical environments, where errors can lead to high costs or damage. In distributional RL, the risk-sensitivity can be controlled via different distortion measures of the estimated return distribution. However, these distortion functions require an estimate of the risk level, which is difficult to obtain and depends on the current state. In this work, we demonstrate the suboptimality of a static risk level estimation and propose a method to dynamically select risk levels at each environment step. Our method ARA (Automatic Risk Adaptation) estimates the appropriate risk level in both known and unknown environments using a Random Network Distillation error. We show reduced failure rates by up to a factor of 7 and improved generalization performance by up to 14% compared to both risk-aware and risk-agnostic agents in several locomotion environments.

1 Introduction

Many real-world applications of Reinforcement Learning (RL) are risk sensitive, including robotics, autonomous driving and healthcare. This means that some states such as a car crash are associated with a very high cost and must be avoided not only during deployment but already while training an RL agent. Furthermore, the risk of failure might change because of variations in the agent's task [Pinto et al., 2017, Zhao et al., 2020]. So in order to be applicable in practice, agents need to avoid fatal mistakes in particular in the face of a changing environment.

Prior work assumed that information about risky elements in an environment, such as failure states, is explicitly provided via a model of the environment [García and Fernández, 2015, Turchetta et al., 2020, Jansen et al., 2020]. However, knowledge about the target domain might not always be readily available. Therefore, there is a need for agents that can learn a risk-aware behaviour on their own. To reduce the need for external information, we propose that an agent can extract safety-relevant information from its own interactions with the environment.

One natural approach is to consider a risk-level of an action in a given state for making a decision. The problem faced in practice is how to choose an appropriate risk-level without having a lot of prior knowledge about the environment. Too conservative behavior could lead to suboptimal exploration

*Equal Contribution

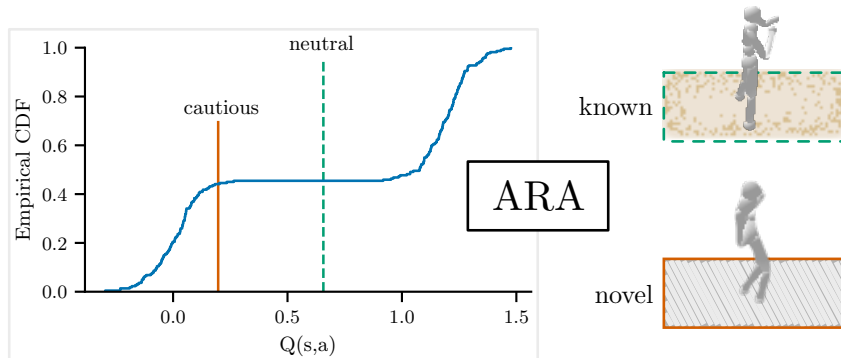


Figure 1: Illustration of an agent that is faced with changed dynamics during training and its effect on the value estimation through ARA. Shown is the estimated empirical CDF of Q-Values for an action a in state s . There are two equally likely outcomes. However, the estimated average return (dashed green) is less likely with smaller ground friction. As the current state in the novel environment (solid brown) has not been seen before, ARA causes the agent to act more cautiously.

and therefore reduced reward. In contrast, too risky behavior could lead to many failure cases during training. While it is possible to use intuition in order to select such risk-levels, we demonstrate both the difficulty of selecting an appropriate risk level and the suboptimality of a one-size-fits-all approach to risk estimation.

To relieve users from the burden of choosing a risk-level, we propose Automatic Risk Adaptation (ARA)². It adjusts the agent’s risk-aware behaviour under changing risk levels without additional information (see Figure 1). The idea of ARA is based on the parametric uncertainty of the agent’s predictions to adapt its risk-awareness, thus encouraging more cautious behaviour in states that have not been visited as often. As an additional advantage over prior work [Choi et al., 2021], it does not add significant overhead to training as it requires no additional episode evaluations.

Since risk-awareness is especially interesting for agents moving around in an open and potentially continuously changing environment, we focus on robot tasks to study the properties of ARA. In particular, we demonstrate the effectiveness of ARA on a variety of robot locomotion tasks.

Our **contributions** therefore are:

1. We show exemplary that static risk levels are suboptimal for generalization to variations of an environment on a windy grid-world.
2. We propose a risk sensitivity regulation mechanism, called ARA, to apply distributional RL algorithms in real world applications.
3. ARA reduces failure rates in risk-sensitive tasks by up to 7 times compared to static CVaR policies.
4. Our method improves generalization performance by up to 14% across varying risk levels.

2 Related Work

Safety in risk-sensitive settings is a well recognized challenge in Reinforcement Learning [García and Fernández, 2015]. One way to address this challenge is the risk-aware distortion of quantile functions of the return distribution [Dabney et al., 2018a], as for example shown by Ma et al. [2020] and Keramati et al. [2020]. Such approaches often use Conditional Value at Risk (CVaR) [Rockafellar and Uryasev, 2000] as it is an effective method to optimize worst-case performance. However, these works keep the risk level α fixed which inhibits generalization to different risk levels. In contrast, Lyu and Amato [2020] use a simple schedule to adapt the risk level over time in a multi-agent setting in order to foster cooperation. While this is beneficial for overall performance, a hand-crafted schedule, similarly to a hand-crafted static value, can only ever be an approximation of the true risk and is therefore constrained in the domains it can be applied to. We aim to improve over both by providing

²Code can be found in the supplementary and will be made public upon acceptance.

a dynamic risk-adjustment method that is sensitive to change in both the policy and environment and thus broadly applicable. Choi et al. [2021] employ a dynamic approach by conditioning the policy on the risk level α . Their method requires policy evaluation and training across different uniformly distributed α levels, which dramatically increases the number of observations and therefore the computational complexity needed for training.³ In contrast, our method requires very little overhead and no additional episode evaluations for risk level estimation.

Where domain knowledge is available, it can be used to estimate the appropriate risk level. In model-based RL, the world model can function both as a simulation environment and a source of knowledge about the risk of the setting [Jansen et al., 2020, Yu and Rosendo, 2021]. Input from a teacher is another way to convey risk in RL. In active learning, agents can query human input more frequently in risky parts of the state space [Brown et al., 2019]. Similarly, a teacher can be used to teach the agent environment constraints like fatal actions without the agent needing to experience them [Turchetta et al., 2020]. We assume that such domain knowledge or teacher experience is not generally available and thus aim to find other sources for risk-adaption.

As more conservative action choices conflict with exploring the environment, safe exploration is a field of research that complements risk-aware policies, such as ARA, especially in environments that require thorough exploration of the state space. Just as for risk-aware action selection, there are many possible approaches to solve this problem, from optimism in the face of uncertainty [Keramati et al., 2020] to learnt exploration mechanisms using e.g. Gaussian Processes [Sui et al., 2015].

3 Risk-Aware Distributional Reinforcement Learning

Before describing the specifics of ARA, we provide the necessary background on quantile function distortion for risk-awareness in distributional RL agents.

Distributional Reinforcement Learning (DistRL) approximates not only the expected return with respect to a given state-action pair, but the whole return distribution [Bellemare et al., 2017]. This enables DistRL agents to integrate the uncertainty of the return into their decision process.

Distributional Deep Q-Learning via Quantile Regression [Dabney et al., 2018b] learns to estimate the quantile function of Q-values for a given state-action pair. This is done by modeling the Q-value distribution in state s for an action a as a random variable $Z(s, a)$. Let Z_τ denote the inverse cumulative distribution function (CDF) of Z . If we query Z_τ using $\tau \sim \mathcal{U}(0, 1)$, we obtain a Q-value sampled from Z . Each percentile rank τ corresponds to a possible Q-value for a pair (s, a) that is distributed according to Z , so in expectation:

$$Q(s, a) = \mathbb{E}_{\tau \sim \mathcal{U}(0,1)}[Z_\tau(s, a)] \tag{1}$$

Estimating the quantiles of the Q-value distribution lifts the constraint of setting fixed Q-value ranges of previous approaches such as the C51 DistRL agent [Bellemare et al., 2017].

While in Equation (1) all possible Q-values in a state s are still equally weighted, it is also possible to distort the distribution using a non-linear distortion function $\beta : [0, 1] \rightarrow [0, 1]$ [Balbás et al., 2009]. This results in a more pessimistic or more optimistic estimate of the Q-value distribution in a given state, depending on the chosen function.

$$Q_\beta(s, a) = \mathbb{E}_\tau[Z_{\beta(\tau)}(s, a)] \tag{2}$$

Distorting the quantile distribution can thus control the kind of policy that a DistRL agent aims to learn, e.g. one that optimizes the worst-case performance in a given environment.

Quantile Function Distortion [Chow and Ghavamzadeh, 2014] has frequently been used in distributional RL for risk-sensitive settings [Dabney et al., 2018a, Singh et al., 2020]. In contrast to other forms of risk assessment, e.g. worst outcome performance or uncertainty measures, which are often hard to measure or take a considerable amount of compute, distortion functions can be easily applied to distributional value functions. A common example is CVaR [Chow and Ghavamzadeh, 2014], also known as *expected shortcoming*.

³We were not able to compare against Choi et al. [2021] as their code is not publicly available.

CVaR describes the expectation that the value of a random variable, in this setting the Q-value, is lower than the value at risk (VaR) at a given confidence level α . The VaR in turn is an originally economic risk measure that describes the maximum financial loss with probability α . We can view it as a measure of regret in this application. The CVaR is used as a mapping β to distort the quantile function for action selection and compute the gradients for the value function updates.

$$\beta_{\text{CVaR}}(\tau; \alpha) = \tau \cdot \alpha \tag{3}$$

Therefore we take an estimation of the regret with a given probability into account and make more conservative choices overall, as indicated by the confidence level α . We can also view this confidence level as a risk-level since high confidence implies a low risk setting.

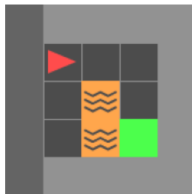
4 ARA: Automatic Risk Adaptation

While CVaR is an easy-to-use and effective risk measure, it has an additional hyperparameter in α . As α is the risk level of the distortion function, it has a significant impact on the success of training and needs to be set correctly for the problem at hand. For this reason, we propose a way to automatically adapt the risk sensitivity of the setting depending on the agent’s current situation, allowing for an adaptive behaviour that is cautious in risky areas of the state space and less so in safe zones.

4.1 Static Risk Levels are Suboptimal for Generalization

To show that statically chosen risk levels in quantile function distortion are generally not a good choice in varying environment conditions, we use the GridWorld LavaGap [Chevalier-Boisvert et al., 2018] environment with added wind that can turn the agent (see Figure 2). The agent can turn left or right and move forward. There is no time limit to reach the goal state, so the only way for the agent to fail is by falling into the lava.

In this setting we can vary both wind direction and wind strength to adjust the risk level. We compute the tabular optimal value function for light south wind that turns the agent with a likelihood of 25% and construct a categorical distribution with 50 values for τ as a distributional value function. Additionally, we distort the value function with CVaR using values for α in increments of 0.1 between 0 and 1.



Setting	Best Choice	Mean Failures	Worst Choice
Light south wind	0.00 ($\alpha : 0.4$)	0.0	0.0
Strong south wind	0.04 ($\alpha : 0.5$)	0.05	0.07
Light north wind	0.22 ($\alpha : 0.1$)	0.27	0.33
Light east wind	0.34 ($\alpha : 0.9$)	0.44	0.49

Figure 2: **Left:** Windy GridWorld with Lava Crossing. **Right:** Differences in failure rates for several values of α over 100 evaluations of the same value function under CVaR risk distortion with different training (first row) and test settings. The columns show the best, average and worst failure rate that a CVaR policy reaches.

We evaluate 100 episodes for each policy in settings with greater wind strength (90% instead of 25%, 'Strong south wind') and different wind directions (90° turn in direction with 'light east wind', 180° with 'light north wind'). We can see the difference the choice of risk level makes in Figure 2.

We observe that different choices for α , even in this very simple setting, lead to different failure rates, with a widening disparity between the mean performance and the best one. While in the original 'light south wind' setting, there was no difference, it increases to around 2% for stronger wind to 5% and eventually 10% with winds from different directions the agent has not trained on. Furthermore, the mean performance of all choices for α moves closer to the performance of the worst choice with rising risk level. This means the performance gap between different values for α is widening significantly in riskier settings with unfamiliar wind directions, making the choice more important.

The values that performed best in each setting do not seem to follow an obvious pattern. While the best failure rate in the original 'light south wind' variation was achieved with $\alpha = 0.4$, it increased to 0.5 with stronger wind and changed drastically to 0.1 and 0.9 respectively for east and north wind.

This shows that users will not be able to easily identify the correct value for α even in such a simple setting, limiting the practical use of CVaR. At the same time our results emphasize that correctly choosing α is a major factor for performance and failure prevention, especially in high risk settings.

To effectively take advantage of CVaR policies across settings, α should therefore be dynamically adapted to the current environment conditions, like changes of wind direction or strength in the example above. As this complicates the selection process even further, however, we present an approach to find and adjust α automatically on the fly.

4.2 Dynamic Risk Level Estimation

In order to correctly adapt α in any given risk-distorted DistRL policy, we propose using Random Network Distillation to approximate the environment's current risk level.

Random Network Distillation (RND) [Burda et al., 2019] provides a measure of uncertainty that has previously been applied to the exploration problem in environments with sparse rewards. It uses a randomly initialized frozen target network f and a predictor network g . Given the same input, the predictor aims to match the random features of the target network. Because both networks receive the same input and the target network is in the same function class as the predictor, the prediction error u on the training data decreases as the amount of samples increases:

$$u(s) = \|f(s) - g(s)\|^2 \quad (4)$$

However, changes in the training distribution manifest themselves in an increase in the prediction error. This yields an uncertainty estimate for any state based on the current policy and environment dynamics. As we expect to see failure states less often in training due to penalization in the reward function, transitions around those states will be seen less often and produce a higher uncertainty.

From Uncertainty to Risk-Level As discussed above, quantile function distortion is an effective way to integrate risk-awareness into any distributional RL agent predicting state-action values at the cost of manually tuning its risk level. With ARA, however, we use RND to estimate the appropriate α . RND provides an estimate of how often a state has been visited in training. If there are many unknown states in a given environment, we prefer the agent to act conservatively in order to avoid fail states. Additionally, the agent is incentivized to avoid fail states during training because of the associated negative reward. Therefore the likelihood of a state with higher associated RND error being a fail state is comparatively large and we want the agent to increase its risk level in its proximity.

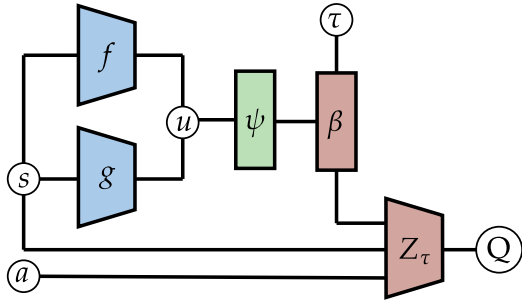


Figure 3: Components of ARA. The target f and the predictor g of the RND are shaded in blue, the uncertainty to risk mapping ψ of ARA in green and the DistRL components in red.

To reflect this in the DistRL algorithm, we use the RND error for the parameter α in the distortion function β in Equation (2):

$$\beta_{\text{ARA}}(\tau) = \beta_{\text{CVaR}}(\tau; \psi(u)) \quad (5)$$

with u being the RND error and $\psi(u) = e^{-u}$. In order to keep u on a consistent scale, we normalize it by its running average estimate.

The RND architecture is an additional hyperparameter of our algorithm, though far less subject to changes in the current setting. We have found the architecture as originally proposed by Burda et al. [2019] to work well. We explore the influence of ψ on the performance of our method in Appendix C.

Using RND for risk estimation in this way is independent of the action and state space size. Furthermore, the risk level estimation is appropriate for the policy even in a test setting as we are not

introducing a second learnt element but base the estimation on the policy itself. Thus, ARA can provide a reliable risk estimation in many settings without significant training overhead.

5 Experiments

We show improved generalization performance and lower failure rates in high risk settings compared to conventional static risk level estimations using ARA in combination with DSAC [Ma et al., 2020] on several 3D locomotion tasks.

5.1 Robustness to Changing Dynamics

Given enough training data, we expect an RL agent to benefit from training on a diverse distribution of its task [Cobbe et al., 2019]. When faced with a variation of the problem during training, however, risk-aware agents have to adjust their behaviour, i.e. their risk level α , to take these different dynamics into account, as we showed on a small example in Section 4.1. This is of course crucial for real-world applications, where we usually cannot guarantee a perfectly stable training environment. Here we demonstrate that this principle also applies to more complex problems and that the risk adaption through ARA therefore can accelerate training and robustness.

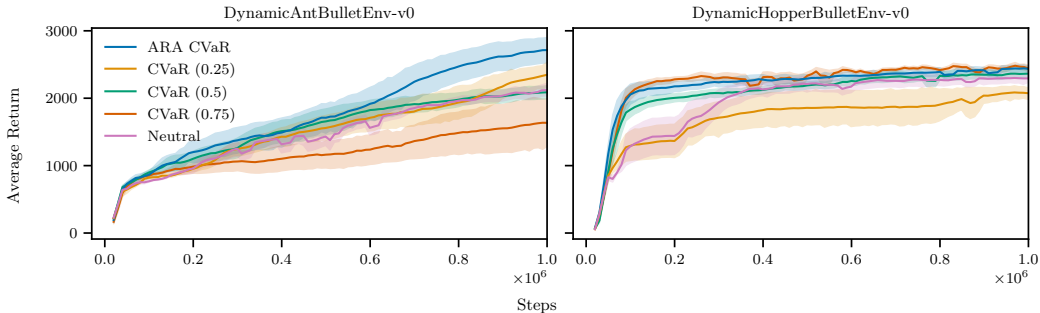


Figure 4: Average return (\pm standard error) over 50 episodes (evaluated at each 10 000 step interval) during training with changing dynamics.

Experimental Setup We evaluate ARA on variations of the *AntBulletEnv-v0* and *HopperBulletEnv-v0* environments that are based on the *PyBullet* physics engine [Ellenberger, 2018–2019] with changing environment dynamics. We vary the friction parameters of the feet and the mass⁴ of the torso randomly by up to 20% of their default values in each episode. For a more thorough exploration of different variation levels, see Appendix D. Since the environment dynamics are not part of the agent’s observation, it needs to learn to distinguish between different levels of friction and mass purely through its interactions with the environment. As defined by the environments, the reward is comprised of a movement cost, cost of collision and a reward for moving forward.

Falling down is an explicit failure case in these environments, but this is not the only way overly risky or cautious behaviour can lead to performance reduction. It is just as important to judge movement speed and direction appropriately for the current friction and mass levels in order to move as quickly and safely possible. Therefore, risk-awareness is closely connected to the reward in this setting.

We compare to a non-risk-aware DSAC agent and a DSAC agent with added static CVaR distortion with an α of 0.25 (as is common in the field [Ma et al., 2020, Keramati et al., 2020]), 0.5 and 0.75. This allows us to judge ARA against agents of varying static risk levels. All results show averages and standard error over 5 runs with different random seeds. For hardware specifications and algorithm hyperparameters, please refer to Appendix A.

Empirical Results are shown in Figure 4. On *DynamicAntBullet*, ARA is able to achieve a better final performance than both the statically risk-aware CVaR agents and the risk neutral agent by at

⁴The mass and friction parameters are called `mass`, `lateralFriction`, `spinningFriction` and `rollingFriction` in *PyBullet* and can be changed via the `changeDynamics` method.

least 14%. Clearly a higher risk level is beneficial for this task, as an α of 0.25 does improve final performance compared to the higher 0.5 or the neutral agent slightly. α of 0.75 on the other hand is not a good choice in this setting, as the agent here can only reach about 60% of ARA’s performance.

While ARA’s final performance in the *DynamicHopperBullet* is only slightly better than neutral and CVaR with α of 0.5 and around the same as an α of 0.75, it reaches this performance as fast as the best performing static CVaR with $\alpha = 0.75$. In addition, ARA is three times faster than the neutral agent. In contrast to the *DynamicAntBullet*, the more conservative CVaR with 0.25 falls short on this task.

Overall we observe that the performance ranks of the static CVaR agents are reversed on these environments, with the *DynamicAntBullet* requiring a more cautious approach compared to the *DynamicHopperBullet*. Even so, the risk-awareness is beneficial on both as demonstrated by the neutral agent performing only averagely. We can conclude that the Hopper robot is either easier to control or less likely to fail, giving the tasks a different risk distribution. This is reflected in the average risk parameter which is 10% higher (i.e. less risky) at the end of training for *DynamicHopperBullet*.

Our agents with static CVaR distortion can not outperform ARA on either environments or match its overall performance. As they lack ARA’s ability to detect the changing dynamics and automatically adapt the risk level, they are limited to performing well on only a subset of environments and environment conditions. We can see this effect on CVaR with α of 0.25 and 0.75, that each performs significantly better in one environment than the other. The fact that they also cannot adapt to the changing conditions furthers the gap between ARA and the best static value we tested, by not only matching the best static CVaR agent in each environment but beating it considerably in the case of the *DynamicAntBullet* environment. Overall, our approach is more robust and able to generalize over dynamic changes, regardless of the inherent risk level in the environment.

Further Insights To confirm whether the gains of ARA stem from an improved estimation of the returns under the changing dynamics, we analyze the average Q-value estimation error of the algorithms during training. For this, we plot the average difference between the maximum empirical discounted return and the agent’s estimated Q-value.

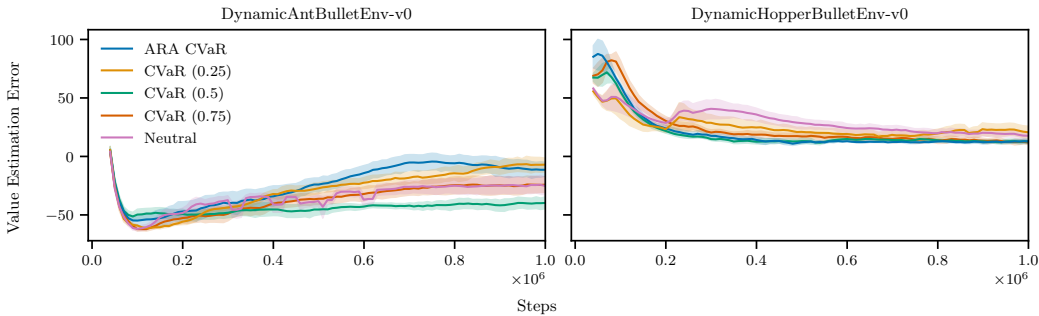


Figure 5: Average Q-value estimation error during training (\pm standard error). Error close to 0 is better.

From Figure 5 it is apparent that ARA leads the RL agent to more accurately estimate the Q-values. In the *DynamicAntBullet*, ARA and the static CVaR with $\alpha = 0.5$ show a less pronounced overestimation of the Q-values in the beginning than the other agents. However, only ARA is able to further reduce the value estimation error later in training at around 600,000 steps, resulting in a better performance.

Interestingly, the error in Q-value estimation does not directly correlate to the static CVaR agents’ performance in this environment. The CVaR with α of 0.25 should perform as well as ARA and α of 0.75 should perform much better than 0.5 according to how well they approximate the Q-values. We suspect this is due to the CVaR distortion preventing the policy from acting optimally in situations the agent is confident in whereas ARA will trust the policy in these states.

The results for the *DynamicHopperBullet* are not as clear, but again, ARA and the CVaR with α of 0.5 and 0.75 agents have the smallest estimation errors. It should be noted that the neutral DSAC has more variance in its Q-value estimate in both environments.

Our hypothesis for the low Q-value estimation error of ARA is that the agent explores a part of the state space until it has seen enough samples to act less cautiously and discover novel states. So even if ARA does not directly affect the Q-networks through any gradient updates (as it is only used in the update of the policy network), it enables the agent to safely experience the environment.

5.2 Effects of Adaptive Risk Levels during Training

Dynamic risk adaption has an additional benefit besides acting according to the current setting. It is able to prioritize risky or cautious behaviour according to how much the agent has learnt already. We want to be more cautious in the beginning of training, but trust the agent more as it performs better.

Experimental Setup We show this effect on the *BipedalWalkerHardcore-v3* locomotion environment [Brockman et al., 2016]. In this environment the agent has to learn to walk over a terrain that it observes via several sensors, the reward corresponding to the distance traveled (at most 300). The terrain varies between episodes, so the agent has to learn to generalize over possible obstacles on the ground. If the agent falls, the episode is counted as a failure and the agent receives a reward of -100 .

Unlike in the scenarios in Section 5.1, reward and failure rate are not as closely connected here. The reward only takes distance into account, so therefore failing very early on a few runs but finishing the rest may result in the same performance as failing on most runs near the end.

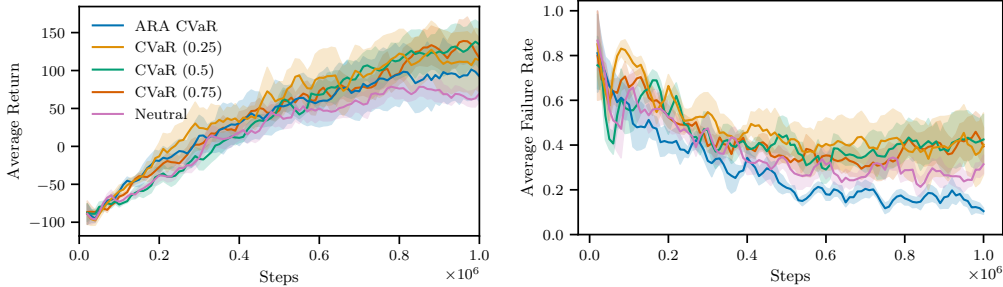


Figure 6: Average return (left) and failure rate (right) on *BipedalWalkerHardcore-v3* over 50 episodes (evaluated at each 10,000 step interval) during training of Neutral DSAC, DSAC with fixed CVaR α (0.25, 0.5 and 0.75) and DSAC with ARA. The shaded area indicates the standard error over 5 runs.

Empirical Results In Figure 6 the average return and failure rate of ARA are shown during training. We see that there is indeed a trade-off between performance and safety in this scenario. It is apparent that a higher average return performance does not always lead to lower failure rates as here, all approaches are fairly close in terms of reward, with the CVaR using α of 0.5 and 0.75 outperforming ARA, and $\alpha = 0.25$ and the neutral agent performing worst.

ARA is able to adapt the risk level to the current state and achieves consistently lower failure rates than the baselines, however. It has a lower failure rate than every other method at every point during training and eventually improving by a factor of 4 and 7 over the neutral and static risk-aware agents respectively. As ARA uses the parametric uncertainty to adapt the risk level, it can learn to significantly reduce failures even in this setting where reward function and failure rate are not as closely connected.

The risk parameter ($\psi(u)$ in Equation (5)) during training shows in Figure 7 how risk-level and policy performance interact. After a drop in the beginning that is caused by the robot reaching novel states, it slowly rises as the agent learns to control the robot. The stagnation in the end is caused by new obstacles that the agent has to overcome to continue.

Overall, ARA successfully matches the risk parameter to the policy development, making the training much safer while obtaining similar performance to the other agents. The adaptive risk parameter supports the safe training of the agent, regulating trust in the learnt policy according to its quality.

6 Limitations

While we expect ARA to be an improvement over existing approaches in risk-aware DistRL, it is no silver bullet. As previously noted by Keramati et al. [2020], exploration heavy environments are hard to learn for risk-aware policies and we do not expect dynamic risk levels to solve this issue completely. In fact, in Figure 8 we can see the same exploration issue in ARA as in the conservative CVaR with $\alpha = 0.25$ if we look at a single mass and friction setting of the Ant environment. Runs that perform worse show a more conservative risk setting and are therefore impeding exploration. Existing measures to mitigate the exploration problem should thus be considered in the future in combination with ARA if exploration is a limiting factor in a given environment.

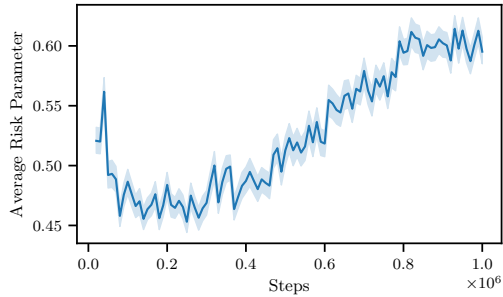


Figure 7: Average risk parameter of ARA over 5 seeds on *BipedalWalkerHardcore-v3* (\pm standard error).

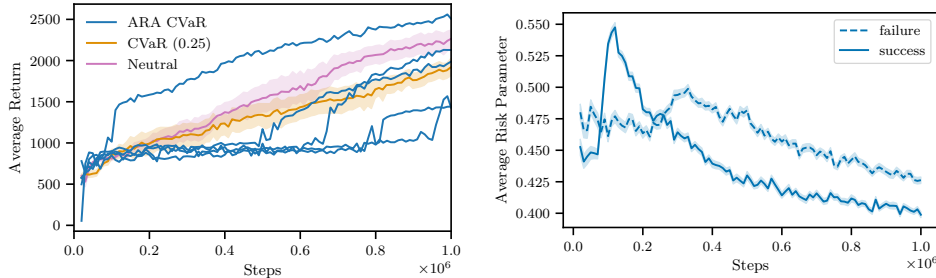


Figure 8: Average risk parameter of 5 seeds of ARA during training in the *AntBulletEnv-v0* environment **without** changing dynamics. There is a large performance gap between runs of ARA due to a lack of exploration. This is expressed by the difference in the risk level around step 150,000 where the successful agent learns to walk, leading to new areas of the state space.

Furthermore, although we remove the need to tune α itself, we also introduce a new hyperparameter by the RND architecture. It only depends on the overall distribution of observations, however, and is thus robust to the individual tasks that the agents learn to solve. While we have found the original RND architecture to work well without any tuning for the settings presented in this paper, this may of course not be the case when applying ARA in new domains.

7 Conclusion

We showed that risk-aware quantile distortion policies with static α are suboptimal both in training and test settings across changing risk levels. To solve this problem, we propose Automatic Risk Adaptation (ARA), which deploys Random Network Distillation (RND) as an uncertainty measure in order to automatically find the appropriate risk level for a given task. We demonstrate that ARA can improve generalization and lower the failure rate of agents considerably compared to both risk-aware and non-risk-aware agents without adding a significant training overhead. Our method makes no assumption about the specific distributional RL algorithm and is thus generally applicable to other methods such as Implicit Quantile Networks [Dabney et al., 2018a]. Tying risk-estimation to uncertainty is at the core of our approach and could be developed further in future work by e.g. developing more accurate representations of uncertainty or exploring more complex mappings from uncertainty to risk-level that could make policies risk-aware only under certain conditions. As risk-aware RL is important to broad and reliable application of RL in the real-world, we believe these steps will contribute to improving the robustness of RL in the face of dynamic environments.

References

- L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta. Robust adversarial reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML, 2017*.
- W. Zhao, J. P. Queralta, and T. Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE Symposium Series on Computational Intelligence, SSCI, 2020*.
- J. García and F. Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16, 2015.
- M. Turchetta, A. Kolobov, S. Shah, A. Krause, and A. Agarwal. Safe reinforcement learning via curriculum induction. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS, 2020*.
- N. Jansen, B. Könighofer, S.n Junges, A. Serban, and R. Bloem. Safe reinforcement learning using probabilistic shields (invited paper). In *31st International Conference on Concurrency Theory, CONCUR, 2020*.
- J. Choi, C. R. Dance, J.-E. Kim, S. Hwang, and K. Park. Risk-conditioned distributional soft actor-critic for risk-sensitive navigation. *CoRR*, 2021.
- W. Dabney, G. Ostrovski, D. Silver, and R. Munos. Implicit quantile networks for distributional reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning, ICML, 2018a*.
- X. Ma, L. Xia, Z. Zhou, J. Yang, and Q. Zhao. DSAC: Distributional soft actor critic for risk-sensitive reinforcement learning, 2020.
- R. Keramati, C. Dann, A. Tamkin, and E. Brunskill. Being optimistic to be conservative: Quickly learning a CVaR policy. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI, 2020*.
- R. T. Rockafellar and S. Uryasev. Optimization of conditional value-at-risk. *The Journal of Risk*, 2, 2000.
- X. Lyu and C. Amato. Likelihood quantile networks for coordinating multi-agent reinforcement learning. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems, AAMAS '20, 2020*.
- C. Yu and A. Rosendo. Risk-aware model-based control. *Frontiers Robotics AI*, 8, 2021.
- D. S. Brown, Y. Cui, and S. Niekum. Risk-aware active inverse reinforcement learning. *CoRR*, 2019.
- Y. Sui, A. Gotovos, J. W. Burdick, and A. Krause. Safe exploration for optimization with gaussian processes. In *Proceedings of the 32nd International Conference on Machine Learning, ICML, 2015*.
- M. G. Bellemare, W. Dabney, and R. Munos. A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML, 2017*.
- W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, AAAI, 2018b*.
- A. Balbás, J. Garrido, and S. Mayoral. Properties of distortion risk measures. *Methodology and Computing in Applied Probability*, 11, 2009.
- Y. Chow and M. Ghavamzadeh. Algorithms for CVaR optimization in MDPs. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems NeurIPS, 2014*.
- R. Singh, Q. Zhang, and Y. Chen. Improving robustness via risk averse distributional reinforcement learning. In *Proceedings of the 2nd Annual Conference on Learning for Dynamics and Control, LADC, 2020*.

- M. Chevalier-Boisvert, L. Willems, and S. Pal. Minimalistic gridworld environment for openai gym. <https://github.com/maximecb/gym-minigrid>, 2018.
- Y. Burda, H. Edwards, A. J. Storkey, and O. Klimov. Exploration by random network distillation. In *7th International Conference on Learning Representations, ICLR'19*, 2019.
- K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman. Quantifying generalization in reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning, ICML*, 2019.
- B. Ellenberger. Pybullet gymperium. <https://github.com/benelot/pybullet-gym>, 2018–2019.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section 6
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We include the code in the supplemental material
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Training details can be found in the Appendix A as well as the code in the supplemental material
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We include standard deviation over 10 random seeds for our experiments
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We use existing RL environments as well as extend an existing code base. We cite both,
 - (b) Did you mention the license of the assets? [Yes]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Experiment Hardware & Hyperparameters

Hardware All experiments with were conducted on a slurm GPU cluster consisting of 6 nodes with eight Nvidia RTX 2080 Ti each. The maximum memory was 10GB.

Hyperparameters We used the same hyperparameters for all experiments. Following Ma et al. [2020], we did not tune the policy temperature α .

Hyperparameter	Value
Policy Temperature	0.2
Learning Rate	0.0003
Batch Size	256
Discount Factor	0.99
Target Smoothing	0.005
Replay Buffer Size	10^6
Minimum Steps Before Training	10^4
Quantile Fractions	32
Quantile Fraction Embedding Size	64
Huber Regression Threshold	1

Table 1: Hyperparameters of DSAC

The RND target architecture is the same as in [Burda et al., 2019], with 3 convolutional layers followed by a dense layer with 512 units, and the learning rate is set at 0.0003. The predictor has the same architecture with two additional dense layers with 512 units each.

B Extended Experiments on Non-Risky Environments

A task not mentioned so far with a risk of falling over is the *Walker2DBulletEnv-v0* environment, which is based on the PyBullet physics engine [Ellenberger, 2018–2019]. The result of ARA compared again with the two baselines is shown in Figure 9.

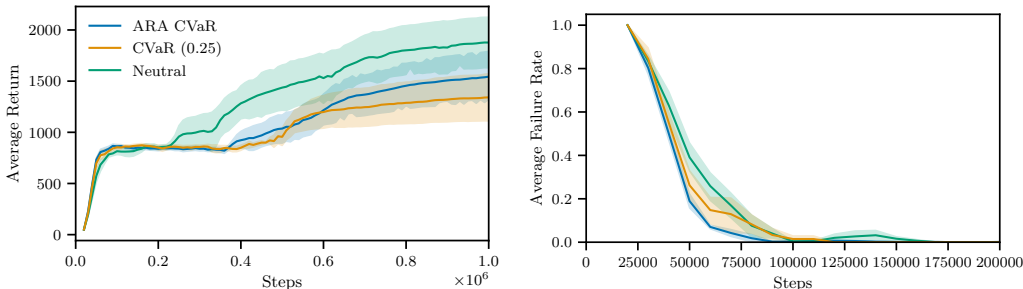


Figure 9: Average return and failure rate on *Walker2DBulletEnv-v0* over 50 episodes (evaluated at each 10,000 step interval) during training of Neutral DSAC, DSAC with fixed CVaR (0.25) and DSAC with ARA. Results are averaged over 5 seeds with the standard error indicated by the shaded area. We only show the failure rate over the first 200,000 steps to better visualise the effect of ARA.

As this environment is easier than *BipedalWalkerHardcore-v3*, the agents learn to control the robot much faster and keep it from falling early on in training. Here both risk-sensitive policies achieve low average failure rates faster than the neutral baseline, with ARA being ahead of the CVaR policy by about 1000 steps. Because *Walker2DBulletEnv-v0* is more difficult in terms of exploration, their average return performance increases more slowly in turn. However, as ARA adapts its risk level during training, it obtains higher rewards than the CVaR agent with a fixed risk level by about 23%, closing the gap to the neutral agent’s performance by almost a third.

C Effect of the RND Mapping

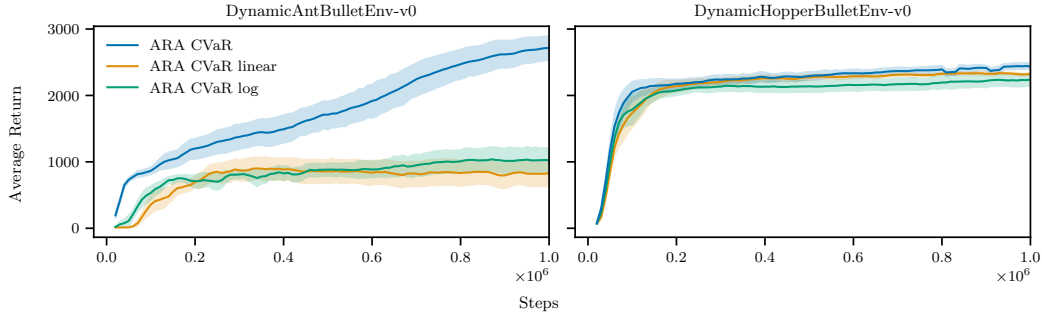


Figure 10: Average return over 50 episodes during training of ARA on the dynamic environments using different RND error mappings (\pm standard error).

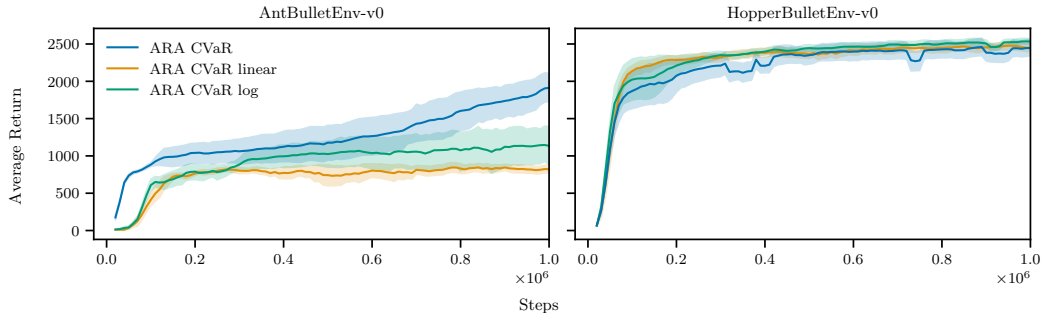


Figure 11: Average return over 50 episodes during training of ARA on the environments **without** changing dynamics using different RND error mappings (\pm standard error).

We designed ARA with an exponential mapping from RND error to risk parameter as we believe it is important not to overlook possible failure states in safety-critical settings, even if we only see little uncertainty from the RND error. Comparing this original mapping with two alternatives, we can validate this intuition.

To contrast the original $\psi = e^{-u}$, we also evaluate a linear mapping in $\psi = -u + 1$ and a logarithmic one in $\psi = -u^2 + 1$. While ARA is very sensitive even to small RND error values, the linear mapping behaves neutrally and the logarithmic mapping is less sensitive to small values for u , instead reacting more strongly for large ones.

We compare the three mappings on the *AntBullet* and *HopperBullet* environments, both with our 20% friction and mass variations and without. As is the main paper, all results are averaged over 5 random seeds.

In Figure 10 and Figure 11 we can see that high sensitivity to uncertainty is beneficial both in the case of static and varying environment dynamics. While the performance on the *HopperBullet* environment is very similar, the difference on the *AntBullet* environment is considerable.

Both the linear and logarithmic versions hardly improve after the initial 200,000 steps when training on changing environment dynamics, although logarithmic mapping is slightly better in the case of static dynamics.

In the case of the logarithmic mapping, this behaviour is to be expected. It leads to small uncertainties being largely ignored and generally acts in a more risk seeking way. We believe this is the reason for its slight performance increase on the static version of the *AntBullet* environment as here the overall risk of falling is lower.

Likely it is also the reason why the linear and logarithmic mappings perform so similar. The linear mapping places no particular importance on any value range of the RND error, which apparently

results in a poor risk adaption for most cases. The logarithmic mapping also cannot react appropriately to smaller uncertainties, but it seems better suited for states the agent knows well.

While an alternative logarithmic mapping might make sense in an explicitly risk seeking context, for our settings, the exponential version is by far the best choice. As the linear mapping places no particular importance on either known or unknown states, it does not provide meaningful risk estimation for any situation and should therefore be avoided.

D Ablation

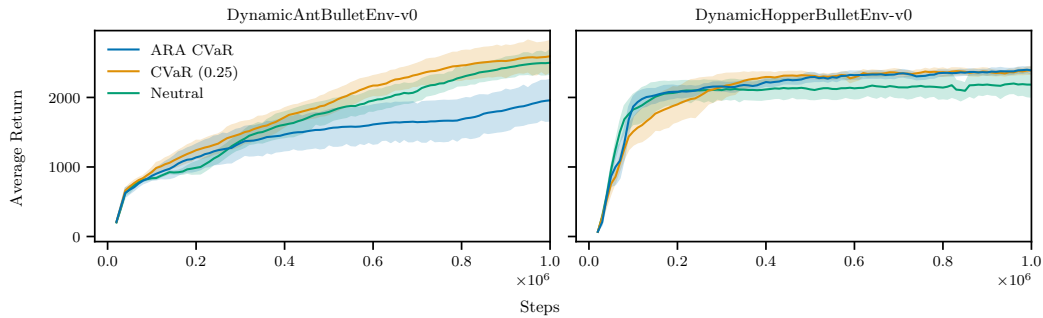


Figure 12: Average return over 50 episodes during training of ARA on dynamic environments with less variation (\pm standard error).

Finally, we assessed how our choice of the strength of variation in the environment dynamics affects ARA. In Figure 12, the friction and mass parameters are varied by only up to $\pm 10\%$. In this setting, the inhibited exploration of ARA on the *AntBullet* task is similar to the static environment which is expressed in the delayed increase in performance.

On *HopperBullet*, ARA is able to achieve the same performance as CVaR with α of 0.25 and is reaching an average return of 2000 as fast as the neutral version.

These experiments indicate that ARA is more applicable in environments with a higher risk level but its automatic risk adaptation also makes it a viable solution in environments with less uncertainty.