# Extreme Risk Mitigation in Reinforcement Learning Using Extreme Value Theory

Karthik Somayaji NS   Yu Wang   Malachi Schram,

Jan Drgona   Mahantesh Halappanavar

Frank Liu   Peng Li
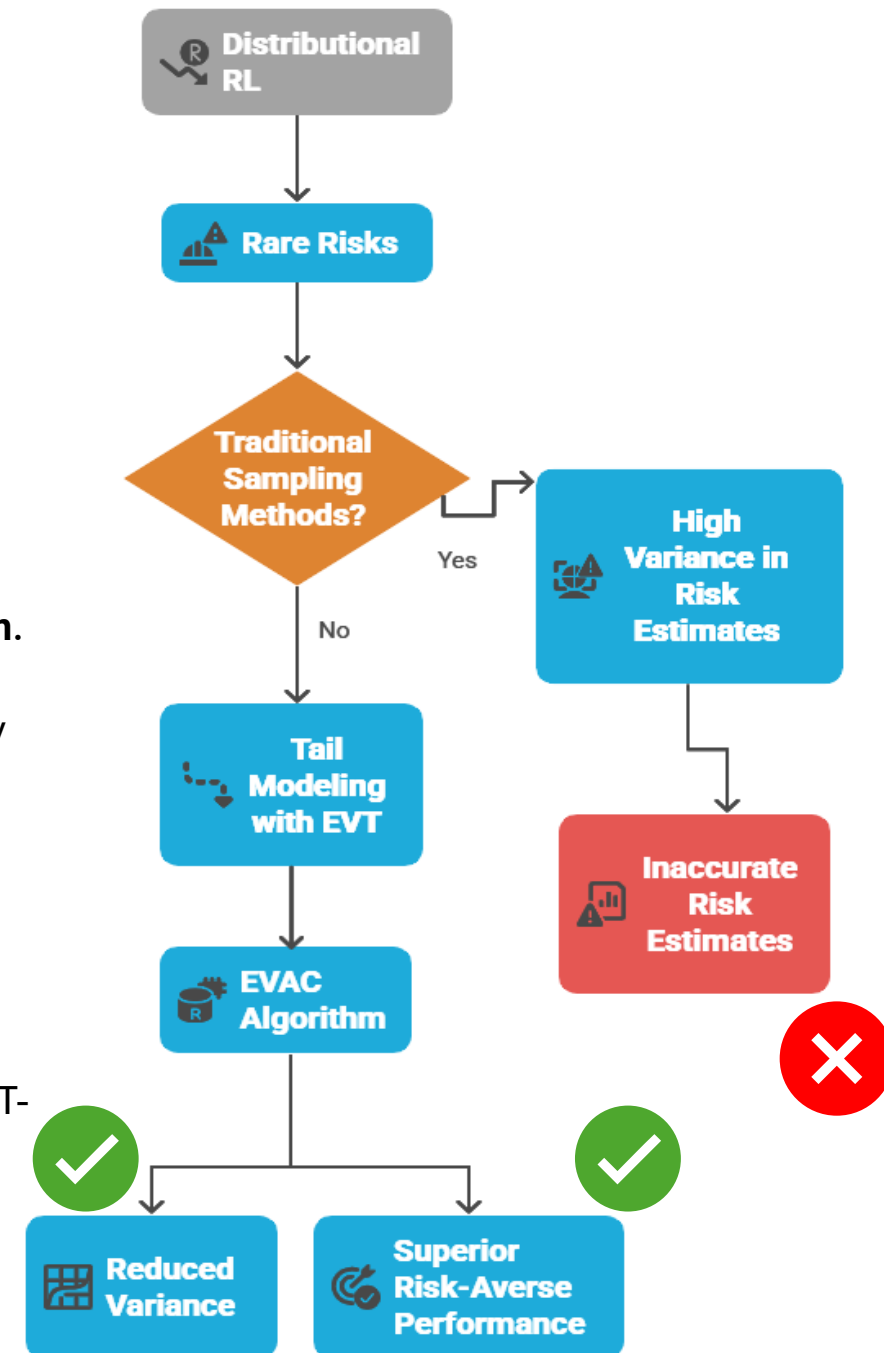
# Background

- **Distributional RL**
  - Instead of expected return - models the **entire distribution** of possible returns.

  - Captures the full uncertainty and variability in outcomes - essential for **safety-critical decisions**.

- **Rare Risks**
  - Catastrophic events that **occur infrequently** but can have severe consequences.

  - Represent infrequent but catastrophic events lying in the **tail of the return distribution**.

  - Traditional sampling methods struggle to accurately model these low-probability events, leading to **high variance** in risk estimates.

- **Our Contributions**
  - **Tail Modeling with EVT:** Leverage Extreme Value Theory to model the tail of the return distribution using a Generalized Pareto Distribution (GPD).

  - **EVAC Algorithm:** Develop a novel Extreme Valued Actor-Critic that integrates EVT-based tail modeling into distributional RL for improved risk aversion .

  - **Empirical Validation:** Demonstrate through experiments on benchmark environments that our method theoretically and practically reduces variance in risk metrics and achieves superior risk-averse performance.
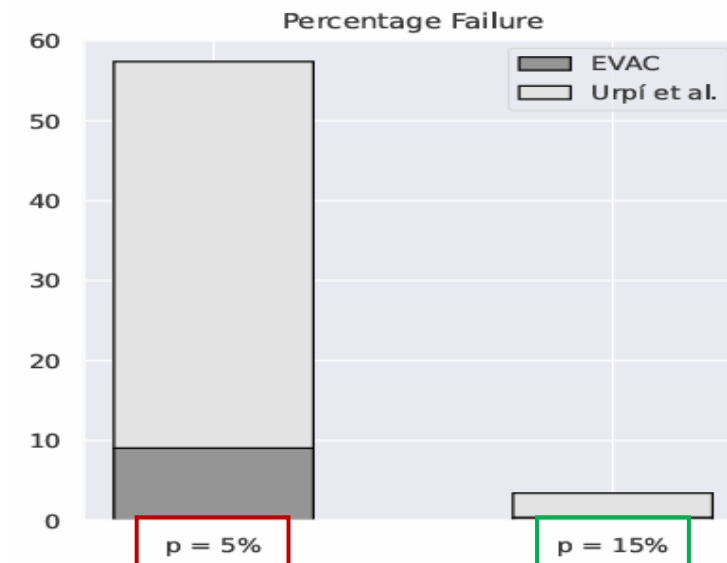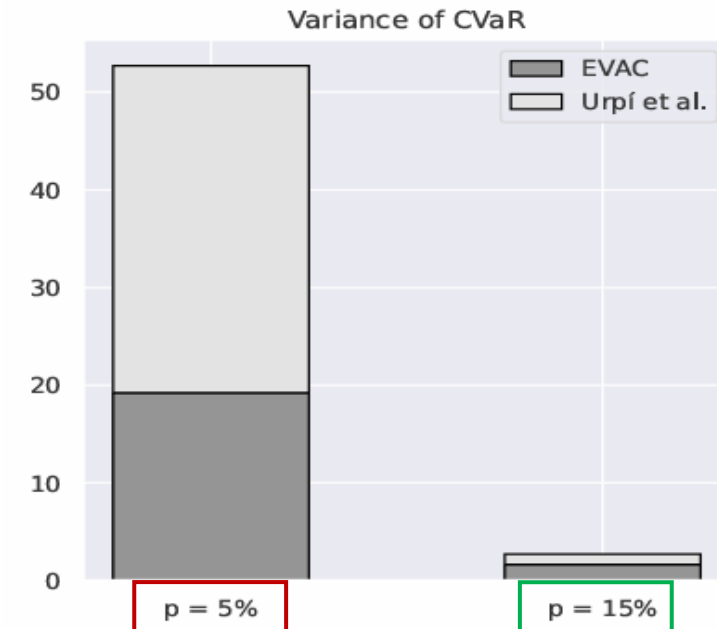
# The Challenge of Modeling Extreme Risks

- **Setup**
  - For example, in Mujoco - HalfCheetah environments, we simulate rare risk by applying a penalty via a Bernoulli variable
  - Two risk levels when velocity of agent exceeds threshold:
    - **15% Penalty:** Moderate rare risk
    - **5% Penalty:** Extreme rare risk
  - Evaluation via :

    - **Empirical risk aversion** (Percentage Failure – how many times does agent go to unsafe states)

    - Emprical **CVaR (conditional value at risk) variance**

- **Conclusion**
  – Accurate tail modeling is critical; EVAC reduces variance and enhances risk-sensitive policy convergence





* Núria Armengol Urpí, Sebastian Curi, and Andreas Krause. Risk-averse offline reinforcement learning. In Proc. International Conference on Learning Representations (ICLR), May 2021.

# Modeling the Tail Distribution

- **Pickands-Balkema-de Haan Theorem:**
  - Approximates the conditional excess distribution above a high threshold using the Generalized Pareto Distribution ($H_{\xi,\sigma}(x)$)

> **(Pickands-Balkema-de Haan Theorem )** *Pickands III (1975) Let $X_1 \cdots X_n$ be a sequence of IID random variables with distribution function (CDF) given by $F$ whose limiting behavior approaches the GEV distribution. Let $F_u(x) = P(X - u \leq x | X > u)$ be the conditional excess distribution. Then,*
>
> $$\lim_{u \to \infty} F_u(x) \xrightarrow{D} H_{\xi,\sigma}(x),$$

- **Key Idea:**

  - For values above a high threshold $u$, the conditional excess (i.e., the amount by which values exceed $u$) converges in distribution to a Generalized Pareto Distribution (GPD).

- **Tail Construction Process:**
  - **Threshold Selection:** Choose $u$ (or a quantile level $\eta$) that partitions the distribution into a non-tail region and a tail region.

  - **Parameter Estimation:** Represent excess values ($X - u \mid X > u$) using the GPD with parameters $\xi$ (shape) and $\sigma$ (scale). Fit the GPD to the observed excesses (via maximum likelihood or other methods) to obtain robust estimates of the tail behavior.

- **Benefits:**
  - Provides an efficient, parameterized model of extreme events that reduces variance in risk measures (e.g., CVaR) and supports more robust risk-averse policy learning.
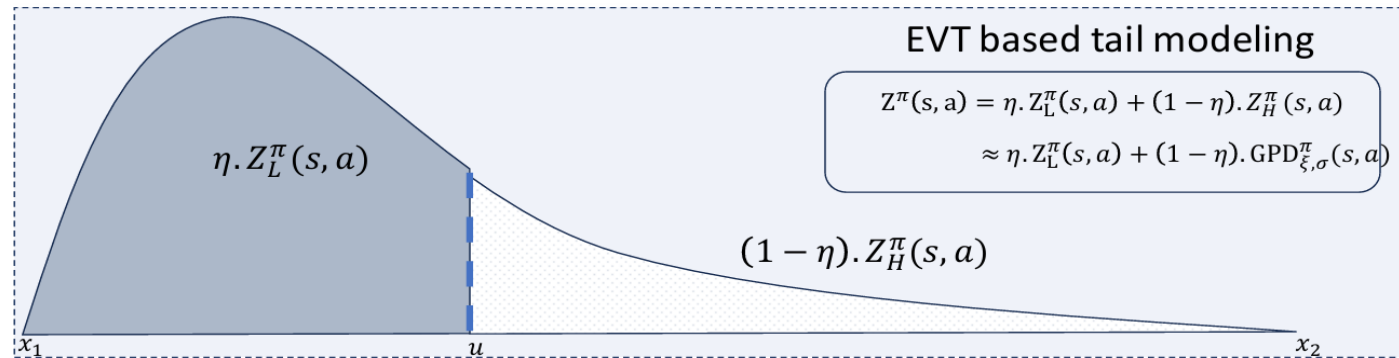
# Extreme Valued Actor Critic

- **Decomposed Critic Architecture**
  - **Non-Tail Region:**
    - Quantile-based modeling (e.g., standard quantile regression) for values below a chosen threshold ($\eta$).
  - **Tail Region:**
    - Parameterized by a Generalized Pareto Distribution (GPD) for excess values above the threshold.
    - GPD parameters ($\xi, \sigma$) updated via maximum likelihood estimation to capture rare, extreme events accurately.



EVT based tail modeling

$$Z^\pi(s, a) = \eta \cdot Z_L^\pi(s,a) + (1 - \eta) \cdot Z_H^\pi(s, a)$$
$$\approx \eta \cdot Z_L^\pi(s,a) + (1 - \eta) \cdot \text{GPD}_{\xi,\sigma}^\pi(s, a)$$

$\eta \cdot Z_L^\pi(s, a)$

$(1 - \eta) \cdot Z_H^\pi(s, a)$

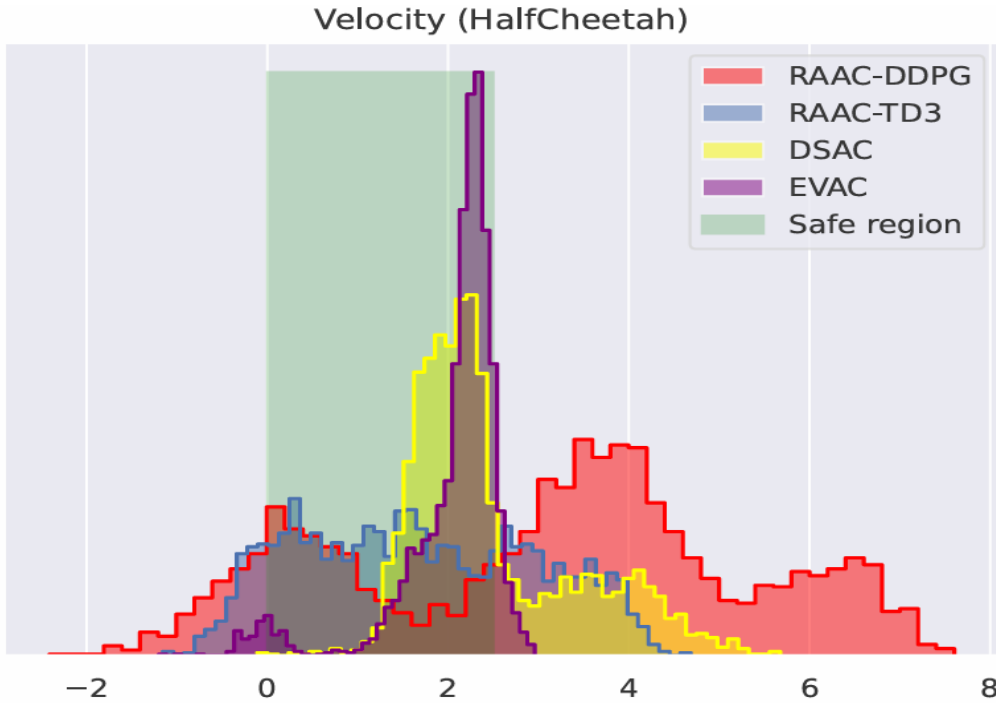$x_1$  $u$  $x_2$

- **Convergence Guarantee**
  We prove that the critic modeling above under a fixed policy is a contraction

- **Integration with Actor-Critic**
  We optimize the actor using a risk measure (e.g., CVaR) derived from the refined (tail + non-tail) distribution.

- **Variance Reduction in quantile regression estimate of the tail region**

# Experiments



Velocity (HalfCheetah)

**Rare Risk Setup (HalfCheetah):**
- Penalize the agent whenever its velocity exceeds a threshold (e.g., 2.5).

- Penalty is triggered with a small probability ($p=5\%$), simulating extreme rare events.

- Safe region is velocity < 2.5

**Rare Risk Setup (Navigation with rare obstacles):**
- Penalize the agent whenever it enters the red circles.

- Penalty is triggered with a small probability ($p=5\%$), simulating extreme rare events

| Algorithm | Percentage Failure | Cumulative Reward | CVaR |
|---|---|---|---|
| RAAC-DDPG | $16.55 \pm 4.43$ | $637.81 \pm 319.78$ | $135.18 \pm 13.02$ |
| RAAC-TD3 | $41.3 \pm 16.6$ | $836.8 \pm 195.85$ | $129.81 \pm 38.75$ |
| D-SAC | $30.04 \pm 22.26$ | $1356.05 \pm 269.63$ | $149.76 \pm 27.38$ |
| EVAC | $\mathbf{2.87 \pm 1.3}$ | $\mathbf{1502.46 \pm 94.25}$ | $\mathbf{156.71 \pm 11.07}$ |