Policy Validation | Score Learning & Calibration | Policy Deployment

Learning + Calibration

Generative Policy $\pi(\boldsymbol{A}|\boldsymbol{O})$

Observation-Based Score

Successful Rollouts

Calibration

Action-Based Score

$F_A(\boldsymbol{\tau}_{:t}) \wedge F_O(\boldsymbol{\tau}_{:t})$

**FIPER**

<u>Fai</u>lure <u>Pr</u>ediction at Runtim<u>e</u>

Score

Threshold

Failure Warning

Success