

1. PROBLEM SETTING

How should a CTRL agent schedule measurements?

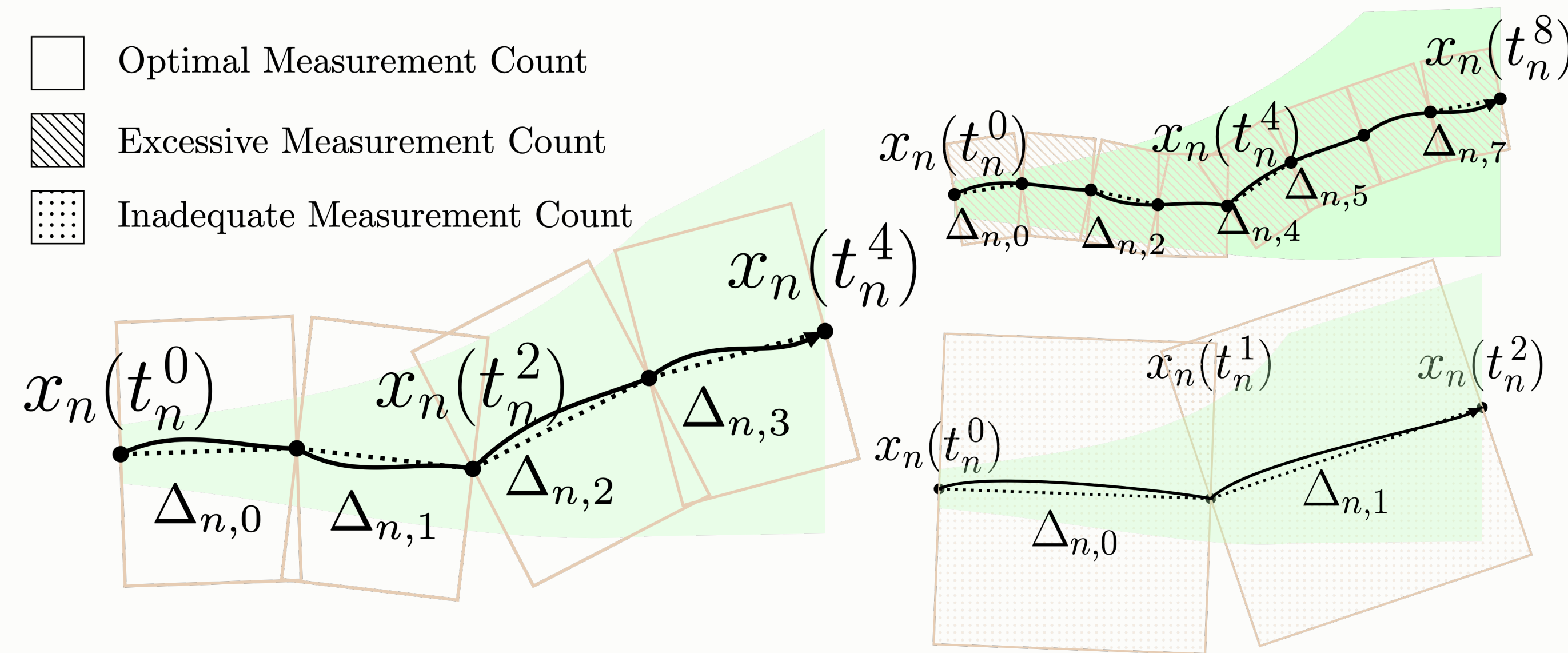
We learn continuous-time dynamics from finite measurements.

LOW TOTAL VARIANCE More frequent measurements per episode	HIGH TOTAL VARIANCE Larger measurement gaps, more episodes
CTRL vs. DTRL: measurement scheduling DTRL Fixed uniform grid: $0, \Delta, 2\Delta, \dots, T$. Timing is preset. Takeaway: DTRL fixes when to measure; CTRL chooses how long to wait before measuring again.	CTRL Flexible gaps: $\Delta_{n,k} = t_{k+1}^n - t_k^n$. Timing is a decision.

Why does timing matter?

Continuous dynamics are observed through finite measurements, so the measurement grid itself becomes part of the decision.

Measurement-grid intuition



Black dots are observations. Brown rectangles mark measurement-grid cost $\Delta_{n,k}^2$. Green shading marks total variance Var^{u_n} .

Key terms

Episode	one trajectory over $[0, T]$	Measurement	state observed at t_n^k
Gap Δ	time between measurements	Total variance	variance of integrated reward

Two ways to measure complexity

Episode complexity	Measurement complexity
Number of episodes N needed to learn an ϵ -optimal policy.	Total measurements $\sum_{n=1}^N m_n$, where m_n is the episode- n count.
More episodes give more independent trajectories.	Finer grids spend more measurements inside episodes.
Scheduling trades off more episodes against more measurements within each episode.	

Mathematical model

Continuous-time controlled dynamics:

$$dx(t) = f^*(x(t), u(x(t))) dt + g^*(x(t), u(x(t))) dw(t).$$

Instance-dependent variance:

$$\text{Var}_{|_0^T} [b(x(t), u(x(t))) dt].$$

Takeaway: balance within-episode resolution against more episodes.

2. ALGORITHM + THEORY

CT-MLE in one sentence: CT-MLE estimates marginal state transition densities by MLE, not path derivatives.

Four-step CT-MLE workflow

1. Construct confidence set based on observed data
2. Compute new models and policy based on optimism
3. Plan with new models and policy
4. Collect new measurements and add them to the observed data

New measurements are added to the observed data before the next confidence set is constructed.

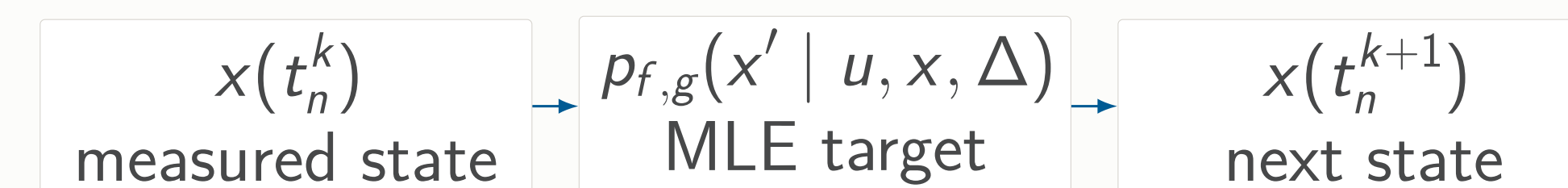
What is used for analysis?

Marginal density MLE	learn $p_{f,g}(x' u, x, \Delta)$, not dx/dt .
Optimistic planning	plan over statistically plausible models.
Randomized additional measurement	capture reward-integral behavior.
Variance-aware measurement rule	match measurement-grid cost to total variance.

What CT-MLE estimates

$$p_{f,g}(x' | u, x, \Delta)$$

marginal transition density over a finite gap



Theory takeaway

$$\text{Regret}(N) \lesssim \iota \left(d_{\beta, \beta} m_N \beta + \sqrt{d_{m_N} \beta \left(\sum_{n=1}^N \Delta_n^2 + \sum_{n=1}^N \text{Var}^{u_n} \right)} \right).$$

Note that

m_N	total measurement count up to episode N .
$d_{m_N}, d_{\beta, \beta} m_N$	complexity terms evaluated at the indicated measurement budgets.
β, ι	confidence radius and logarithmic factor.
$\Delta_n^2 = \sum_k \Delta_{n,k}^2$	measurement-grid cost for episode n .
Var^{u_n}	total variance of the integrated reward under policy u_n .

Match measurement-grid cost to total variance.

$$\text{SET } \Delta_n^2 := \sum_{k=0}^{m_n-1} \Delta_{n,k}^2 \approx \text{Var}^{u_n}.$$

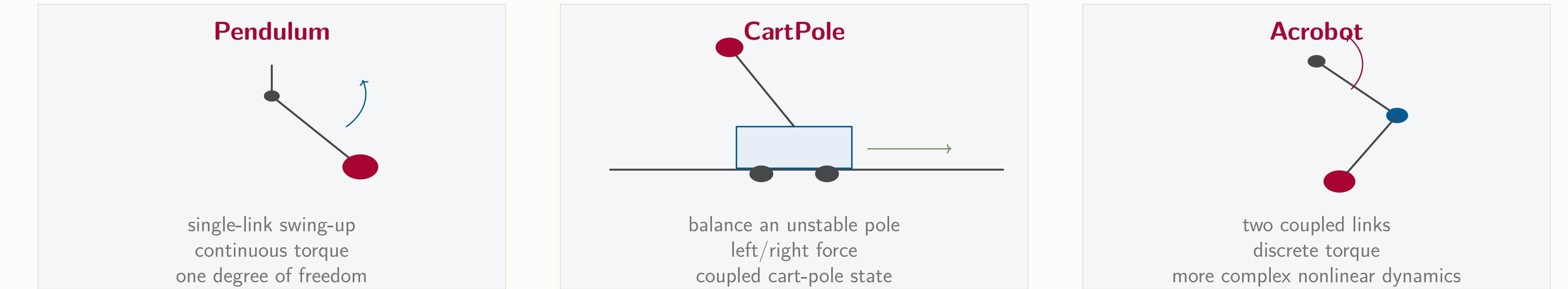
Low total variance
High total variance

smaller gaps, finer grid.
larger gaps, more episodes.

3. EXPERIMENT

Control environments and setup

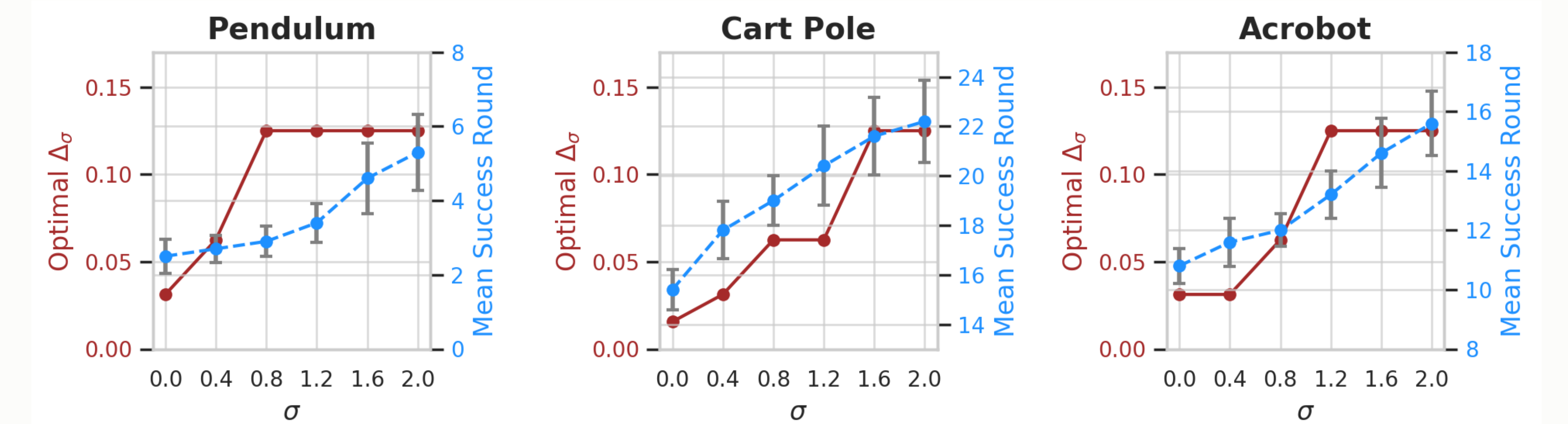
Reward curves use $\sigma = 2.0$. The gap sweep uses fixed uniform gaps $\Delta = 2^{-i}, i = 0, \dots, 7$.



Tasks differ in control structure; reward curves use the same injected noise level.

Evidence for theory: optimal gap increases with stochasticity

Higher stochasticity \rightarrow larger selected gap Δ_σ .

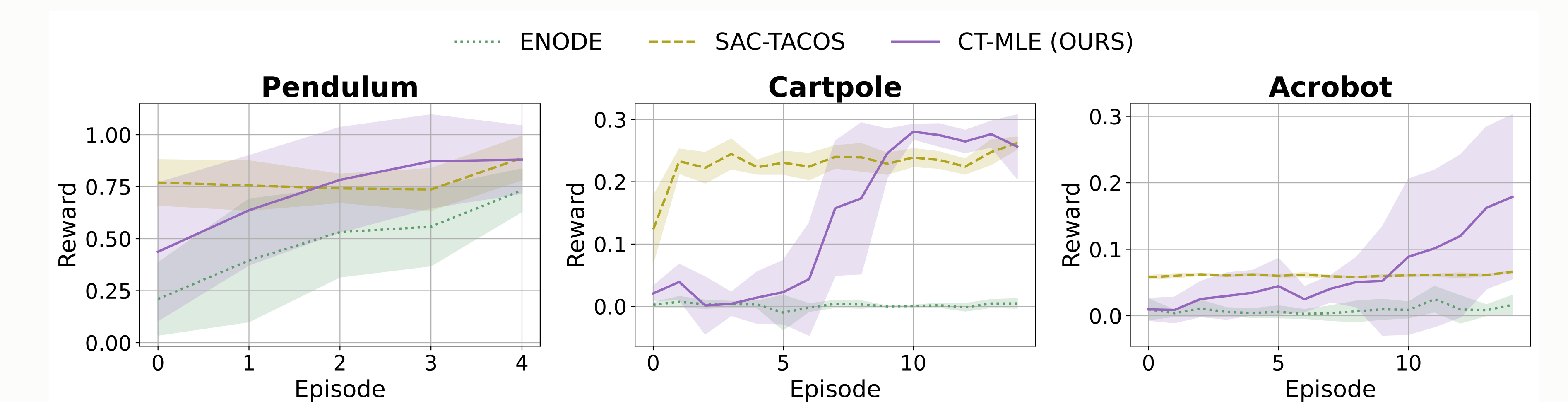


Red series: selected optimal gap Δ_σ ; the trend is consistent with the variance-aware rule.

Illustration.

- ▶ Within each task, increase σ ; red gives Δ_σ , blue gives the mean success round.
- ▶ **Theorem link:** higher $\epsilon_n \text{Var}^{u_n}$ supports larger $\epsilon_n \Delta_n^2$, hence larger measurement gaps.
- ▶ **Intuition:** more stochastic tasks need more episodes and coarser within-episode observations.

Reward curves confirm competitive policy quality



At $\sigma = 2.0$, CT-MLE remains competitive; separation is largest on Acrobot.

Reward curves check policy quality; the gap sweep checks the measurement rule.

Experiment takeaway

Adaptive measurement schedules allocate computation according to instance difficulty: harder stochastic environments benefit from more training episodes rather than denser measurements within each episode.

Measure finely when variance is low; collect more episodes when variance is high.

Marginal-density MLE supports a variance-aware measurement schedule.

