



GitHub Repository

# Reliable Uncertainties for Stellar Properties from Gaia XP Spectra

Joshua D. Ingram<sup>1</sup>, James Carzon<sup>1</sup>, Joshua S. Speagle<sup>2</sup>, and Ann B. Lee<sup>1</sup><sup>1</sup>Department of Statistics and Data Science, Carnegie Mellon University; <sup>2</sup>David A. Dunlap Department of Astronomy & Astrophysics, University of Toronto, CanadaCarnegie  
Mellon  
University

## Overview

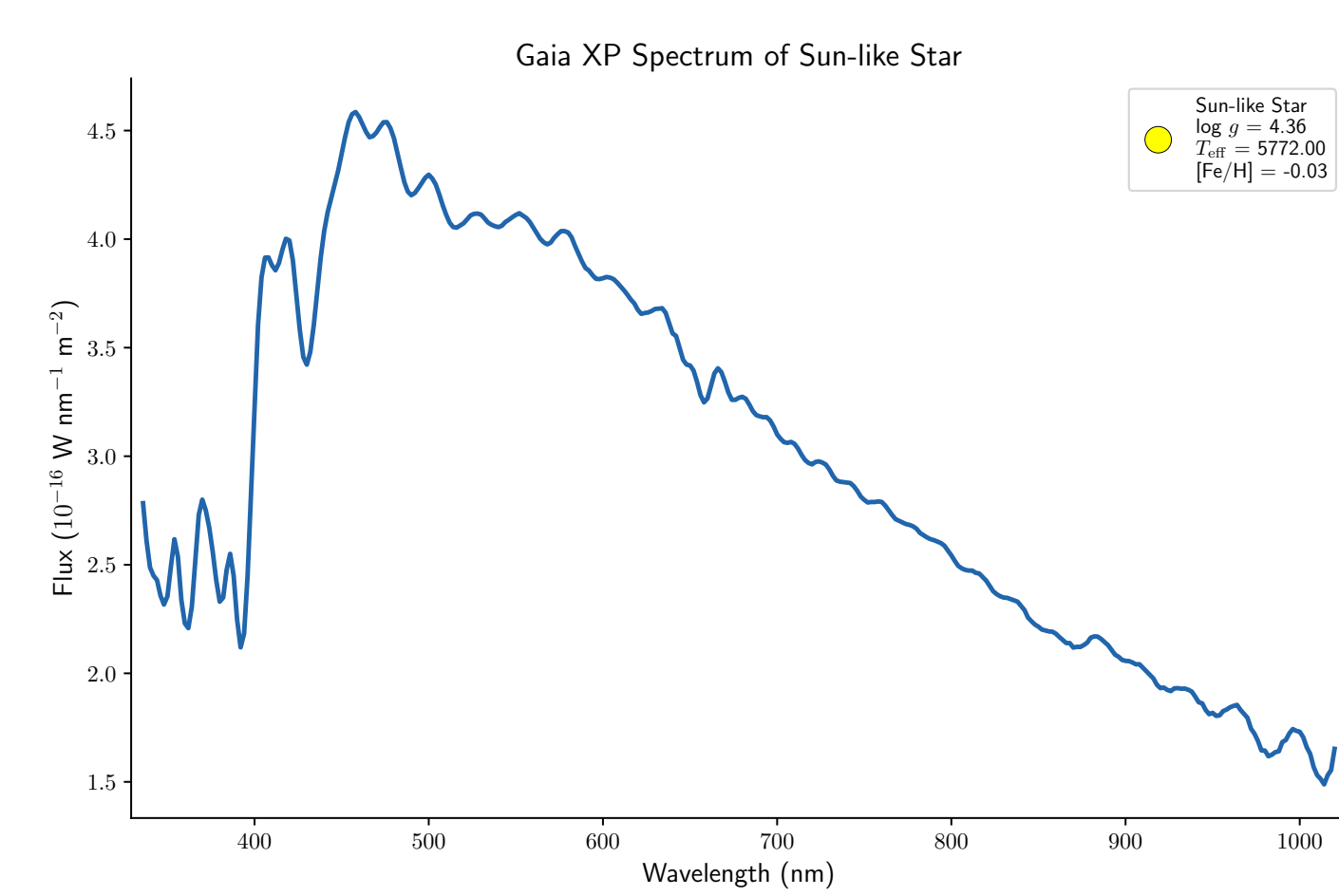
- ▶ **Gaia Data Release 3 (DR3)** has low-resolution XP spectra for over 220 million stars
- ▶ Cross-matching with high-quality measurements from **APOGEE Data Release 17 (DR17)** enables estimation of stellar properties ( $T_{\text{eff}}$ ,  $\log g$ ,  $[\text{Fe}/\text{H}]$ ) at scale
- ▶ Cross-matched data are often **biased due to systematics**: instrument limitations and mission priorities result in **distribution shift**, where training and target data distributions do not match
- ▶ Conventional approaches often provide point estimates with no uncertainties [1]

**Goal:** Given a Gaia XP spectrum  $\mathbf{x} \in \mathbb{R}^{110}$ , estimate stellar properties  $\theta = (\log g, T_{\text{eff}}, [\text{Fe}/\text{H}])$  with **reliable uncertainties having guaranteed local coverage** at any desired confidence level  $1 - \alpha$

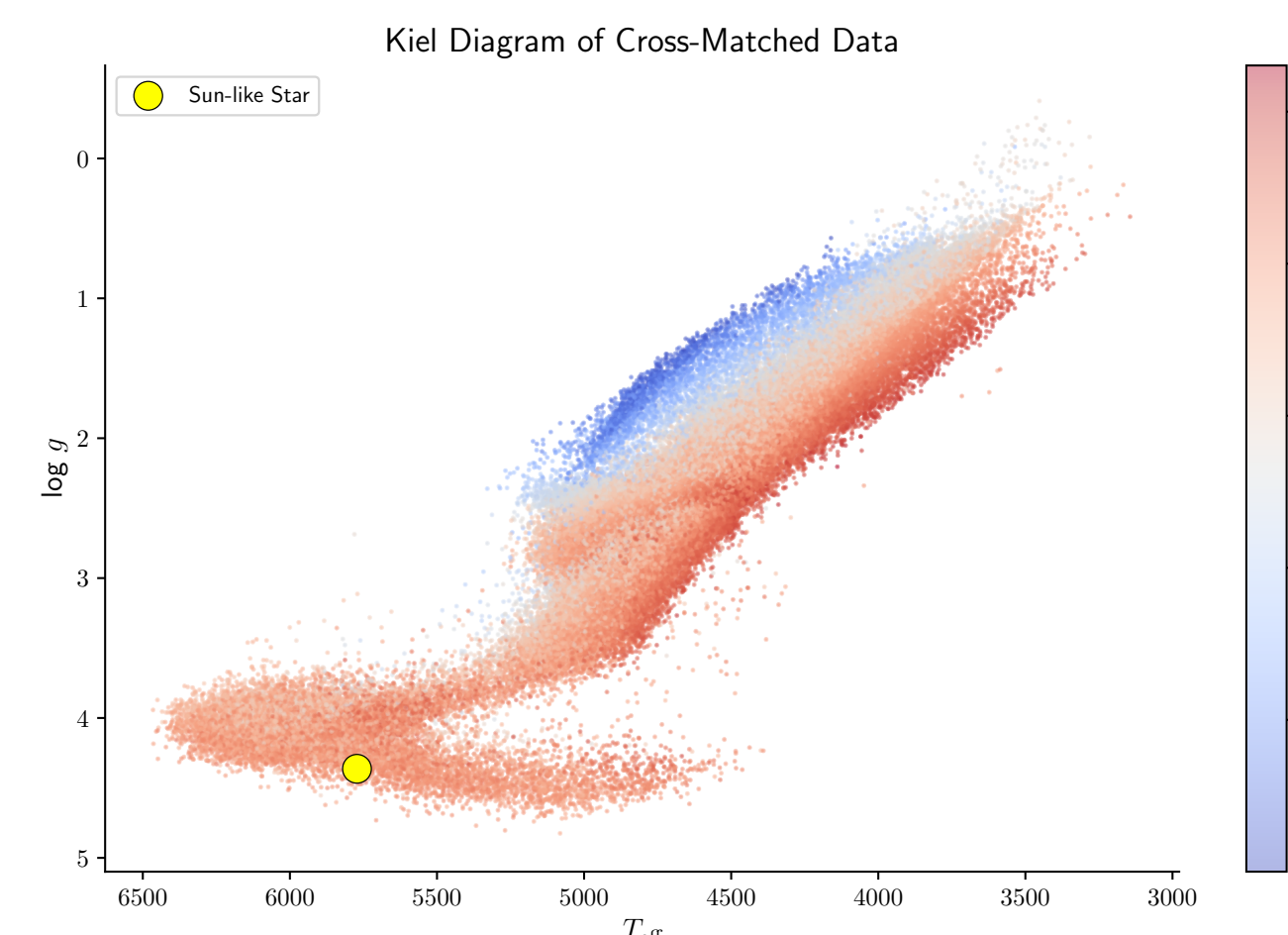
## Gaia-APOGEE Cross-Matched Catalog

**Data:** 202,970 Gaia XP spectra cross-matched with APOGEE DR17 stellar properties [2]

- ▶ **XP Spectra ( $\mathbf{X}$ ):** 110-dimensional Hermite polynomial coefficients for constructing BP and RP optical spectra
- ▶ **Properties ( $\theta$ ):**  $\log g$ ,  $T_{\text{eff}}$ ,  $[\text{Fe}/\text{H}]$  derived from high-resolution spectra



XP spectrum of a Sun-like star.



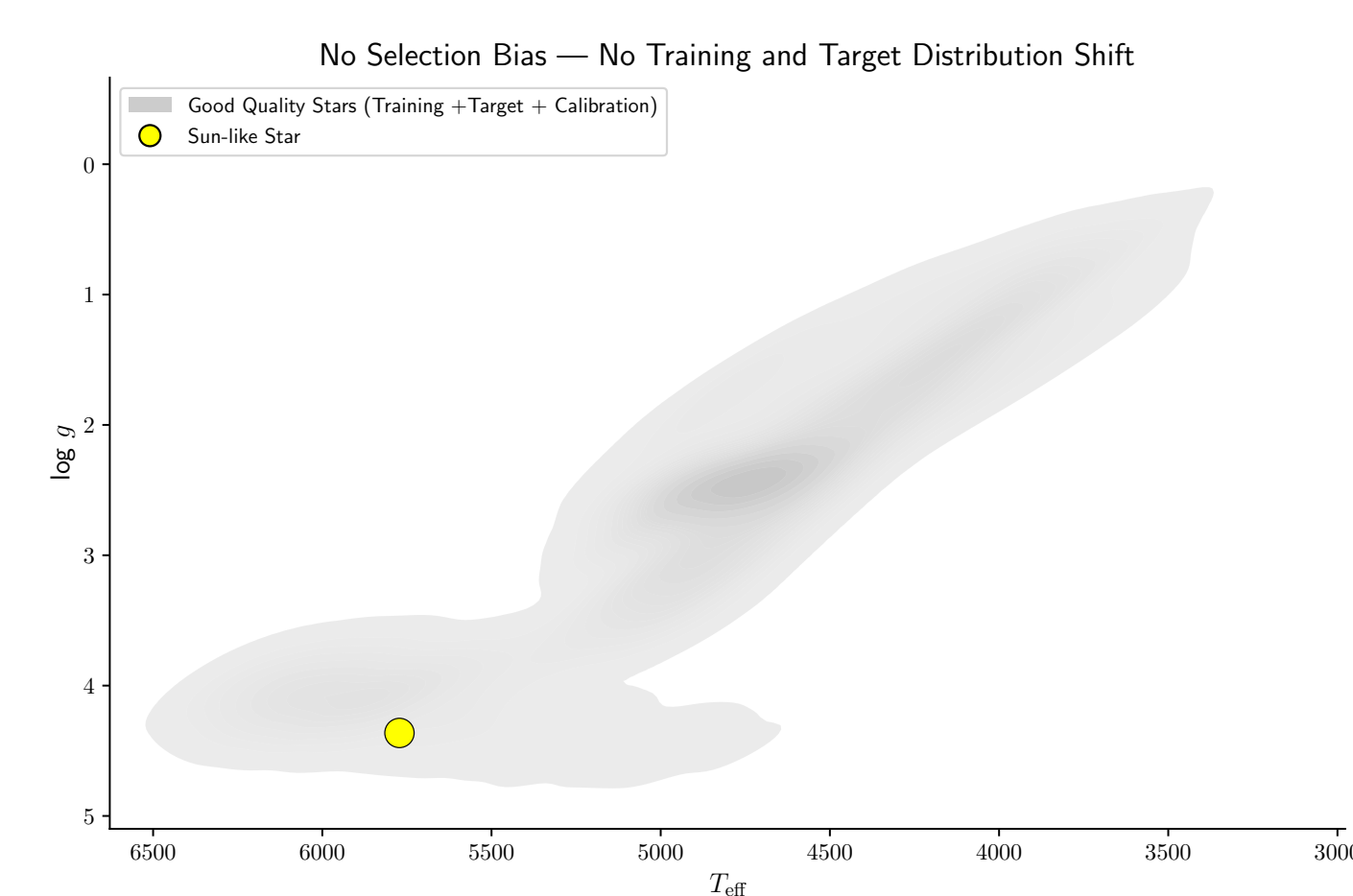
Kiel diagram of cross-matched data, colored by metallicity.

## Inference Setting

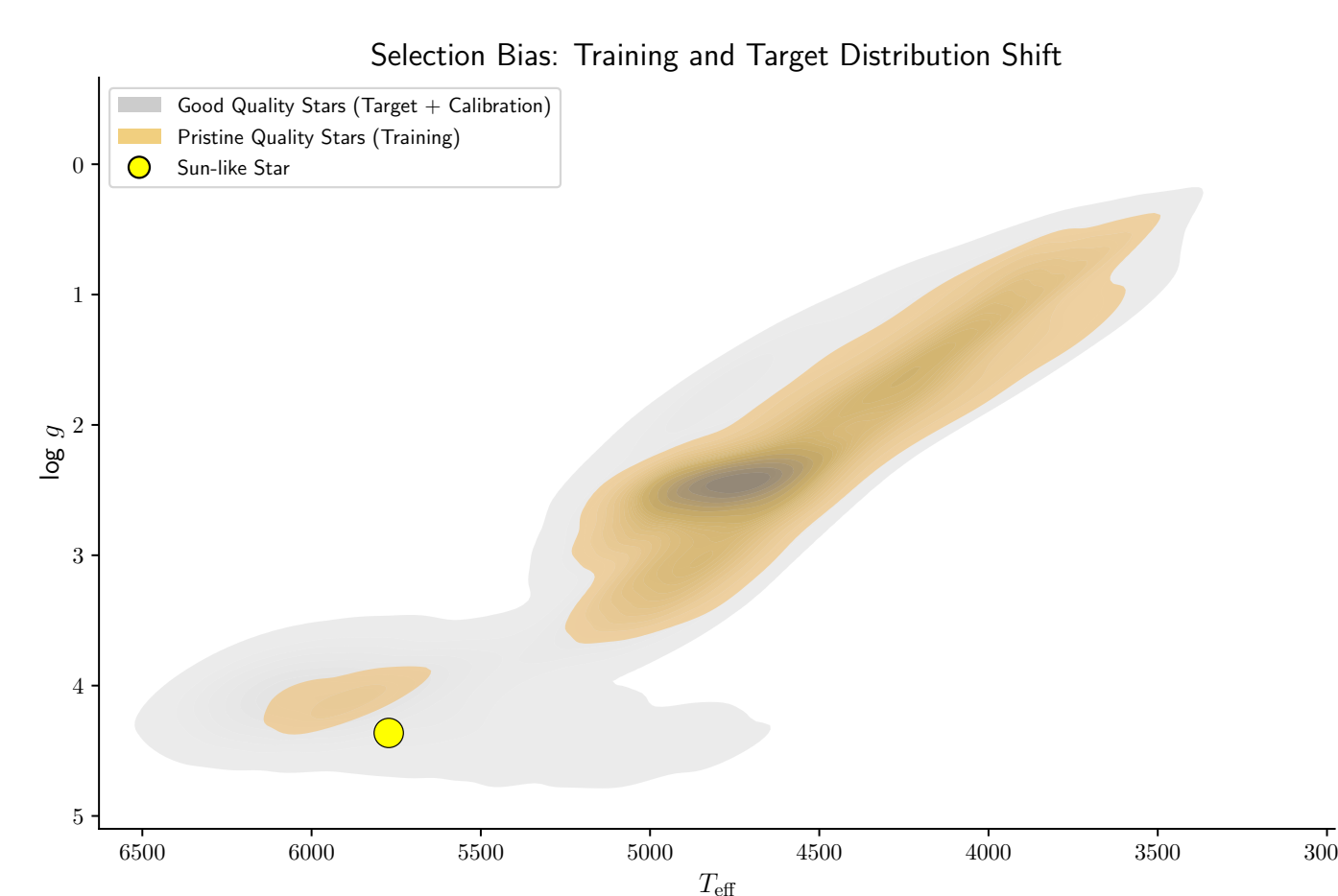
**Challenge:** Construct uncertainties with desired level of local coverage for stellar properties that are reliable in the presence of selection bias

We study performance under two data settings:

1. **No selection bias:** Training data match target distribution
2. **Selection bias:** Train on only “pristine” quality data, predominantly giant-branch (GB) stars, with targets including main sequence (MS) and metal-poor stars; held-out “good” quality data used for calibration



1. **No Selection bias:** training distribution spans the target distribution of stars (no distribution shift).



2. **Selection bias:** “pristine” quality cuts concentrate training data on the giant branch (distribution shift).

**Example Sun-like star:** We illustrate performance with a held-out Sun-like star:

Gaia DR3 Source ID	Distance [pc]	$\log g$ [dex]	$T_{\text{eff}}$ [K]	$[\text{Fe}/\text{H}]$ [dex]
4660210013529490176	334.15	4.36	5772	-0.03

## Reliable Uncertainties Via Confidence Distributions

**Key Idea:** Learn confidence distribution that encodes the confidence sets for all levels  $\alpha$  simultaneously

Let  $F(t | \theta) = P_{X|\theta}(\lambda(\theta, X) \leq t)$  be the CDF of the test statistic under  $\theta$ . Then the  $1 - \alpha$  confidence set

$$C_{1-\alpha}(\mathbf{x}) = \{\theta : F(\lambda(\theta, \mathbf{x}) | \theta) \geq \alpha\}$$

satisfies  $P_{\theta^*}(\theta^* \in C_{1-\alpha}(\mathbf{x})) \geq 1 - \alpha$  for all  $\theta^*$  and all  $\alpha \in (0, 1)$ .

**Framework:** [3] [4]

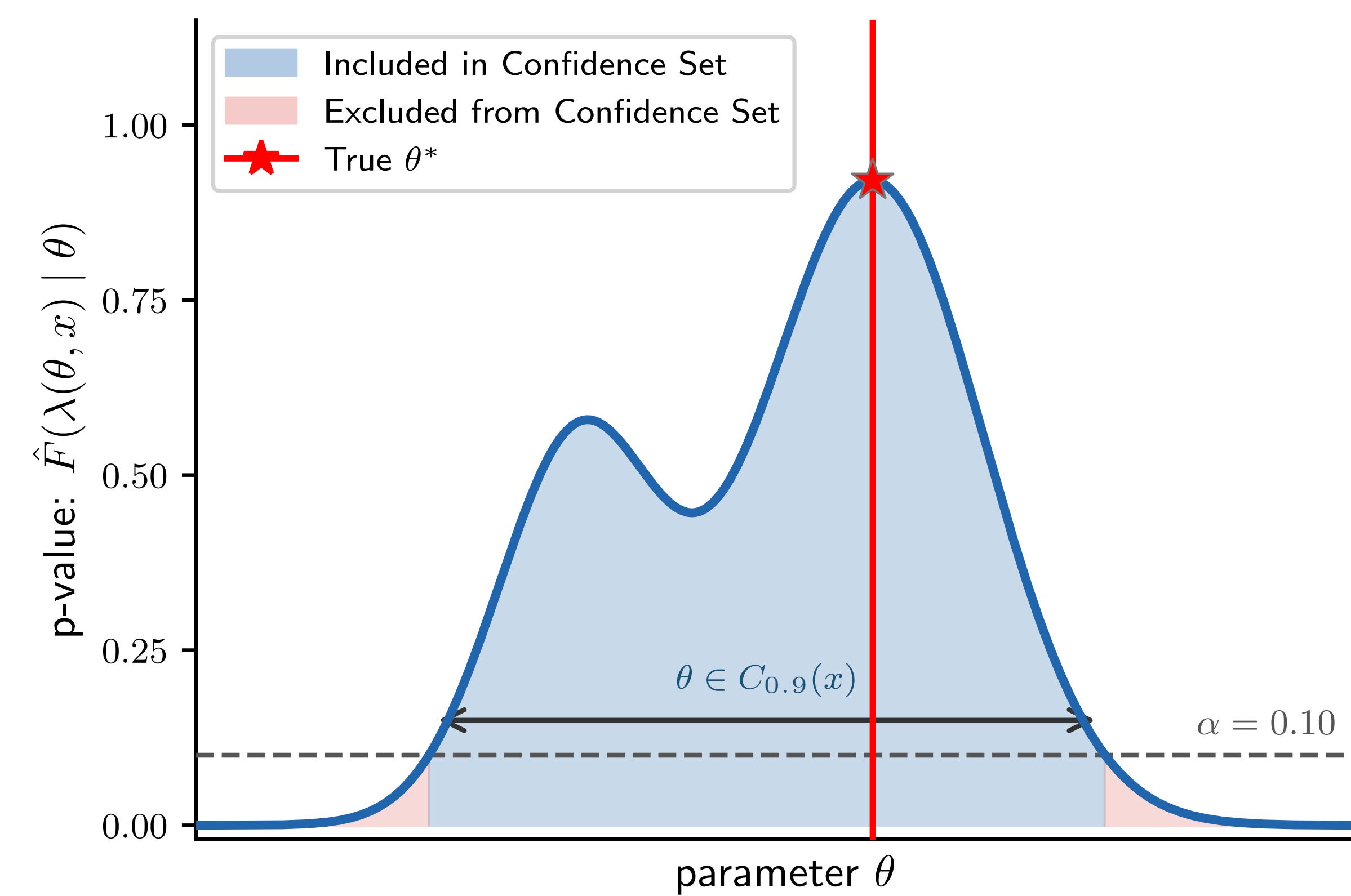
1. Train a **neural posterior estimator** (normalizing flow)  $\hat{p}(\theta | \mathbf{x})$  on labeled spectra
2. Define test statistic  $\lambda(\theta, \mathbf{x}) = \log \hat{p}(\theta | \mathbf{x})$
3. For each  $\theta$ , evaluate the **learned CDF** of  $\lambda$  on a held-out calibration set to obtain the p-value:

$$\hat{F}(\lambda(\theta, \mathbf{x}) | \theta) = P_{X|\theta}(\lambda(\theta, X) \leq \lambda(\theta, \mathbf{x}))$$

4. Confidence sets at any level  $1 - \alpha$  follow by thresholding:

$$C_{1-\alpha}(\mathbf{x}) = \{\theta : \hat{F}(\lambda(\theta, \mathbf{x}) | \theta) \geq \alpha\}$$

**Key Advantage:** Unlike posterior credible sets, confidence distributions achieve **guaranteed frequentist coverage** by construction, even when the posterior is trained under distribution shift.



1-dimensional confidence distribution sliced at level  $\alpha$ ; the set of  $\theta$  with  $\hat{F}(\log \hat{p}(\theta | \mathbf{x}_{\text{obs}}) | \theta) > \alpha$  are included in the confidence set.

## Learning the CDF of the Test Statistic

We impose a **parametric model on the conditional CDF of the test statistic**  $\hat{F}(\lambda | \theta)$ , learned using calibration data  $\{(\theta_i, \lambda_i)\}_{i=1}^n$

**Sigmoid CDF Model:** We model  $F(\lambda | \theta)$  as a logistic function with  $\theta$ -dependent location and scale parameters:

$$\hat{F}(\lambda | \theta) = \sigma(\kappa(\theta)(\lambda - \mu(\theta))), \quad \sigma(t) = \frac{1}{1 + e^{-t}}$$

where  $\mu(\theta)$  is the location (median) and  $\kappa(\theta) > 0$  controls the spread/sharpness.

**Shallow MLP:** A shallow multi-layer perceptron (MLP) maps  $\theta$  to the CDF parameters:

$$\text{MLP}(\theta) = (\mu(\theta), \log \kappa(\theta)), \quad \theta \in \mathbb{R}^4 \rightarrow \mathbb{R}^2$$

- ▶ Architecture: 2 hidden layers  $\times$  64 units,  $\tanh$  activations;  $\theta$  inputs normalized

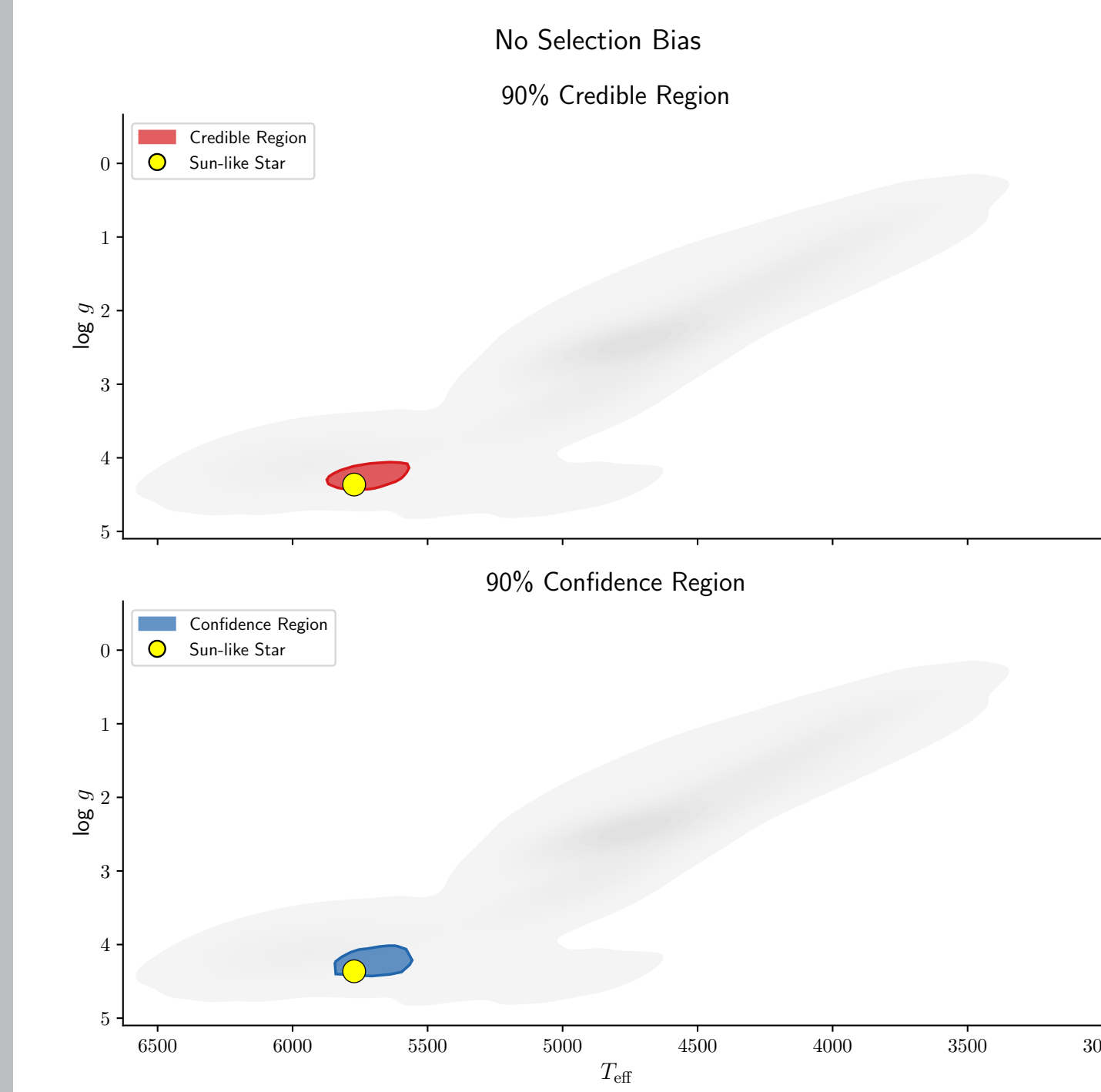
**Inference:** For a new observation  $\mathbf{x}_{\text{obs}}$ , evaluate the learned CDF over a grid of  $\theta$ , and construct  $1 - \alpha$  confidence sets by thresholding at level  $\alpha$ :

$$C(\theta | \mathbf{x}_{\text{obs}}) = \hat{F}(\log \hat{p}(\theta | \mathbf{x}_{\text{obs}}) | \theta) > \alpha$$

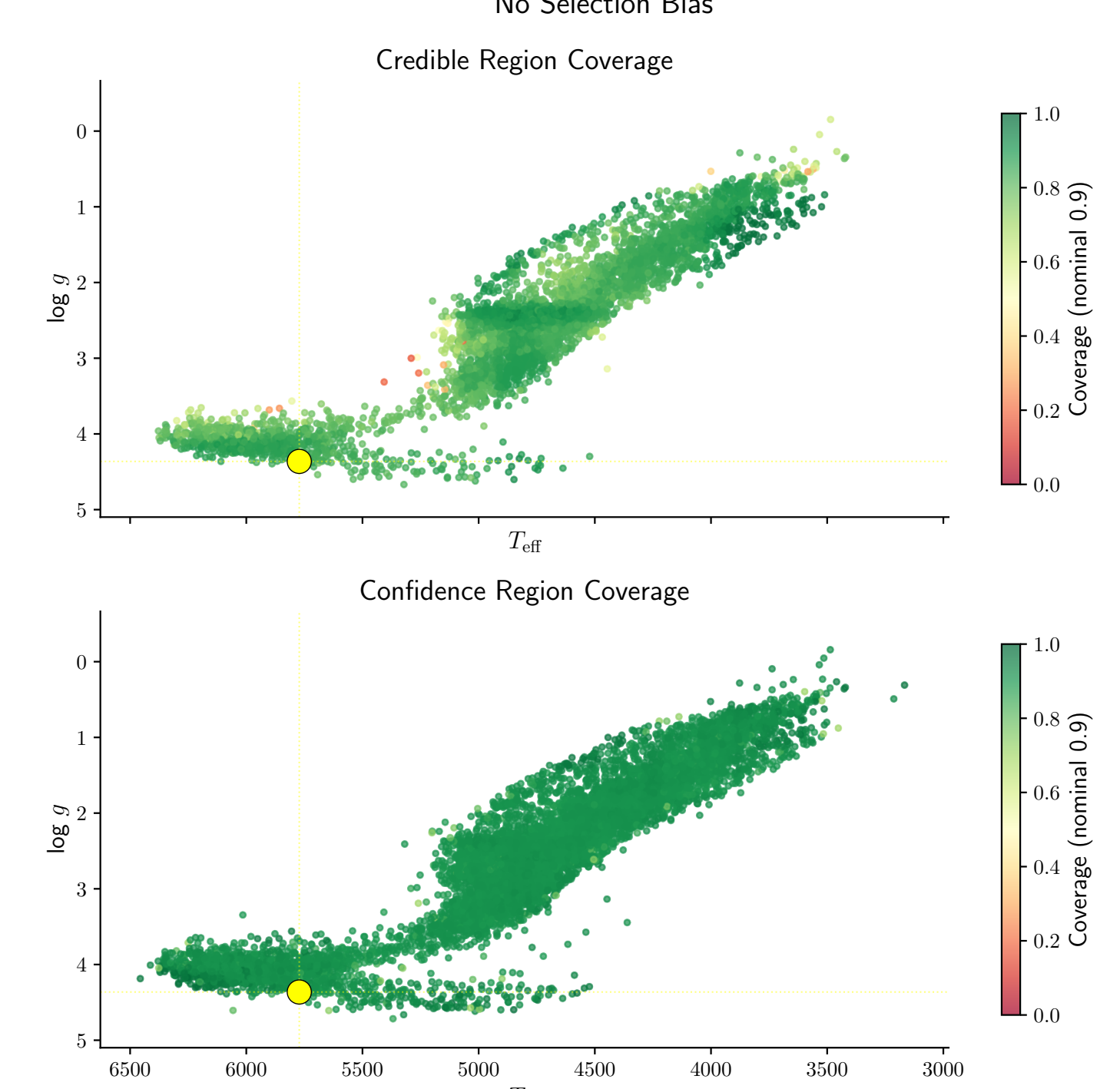
## Implementation

**GitHub:** Available in the open-source lf2i package: [github.com/lee-group-cmu/lf2i](https://github.com/lee-group-cmu/lf2i)

## Results: No Selection Bias

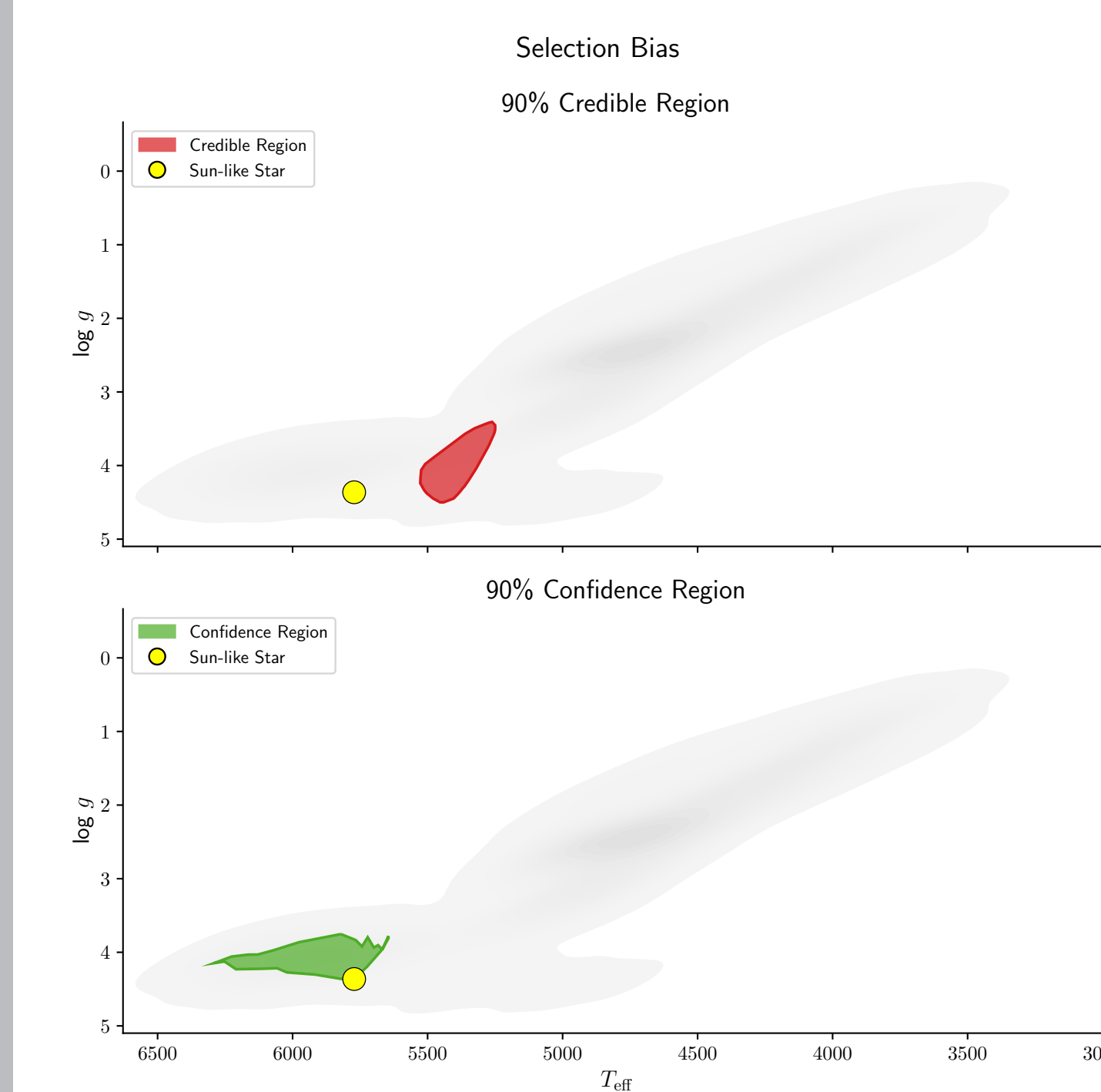


90% confidence region (blue) and HPD credible region (red) for the Sun-like star under selection bias.

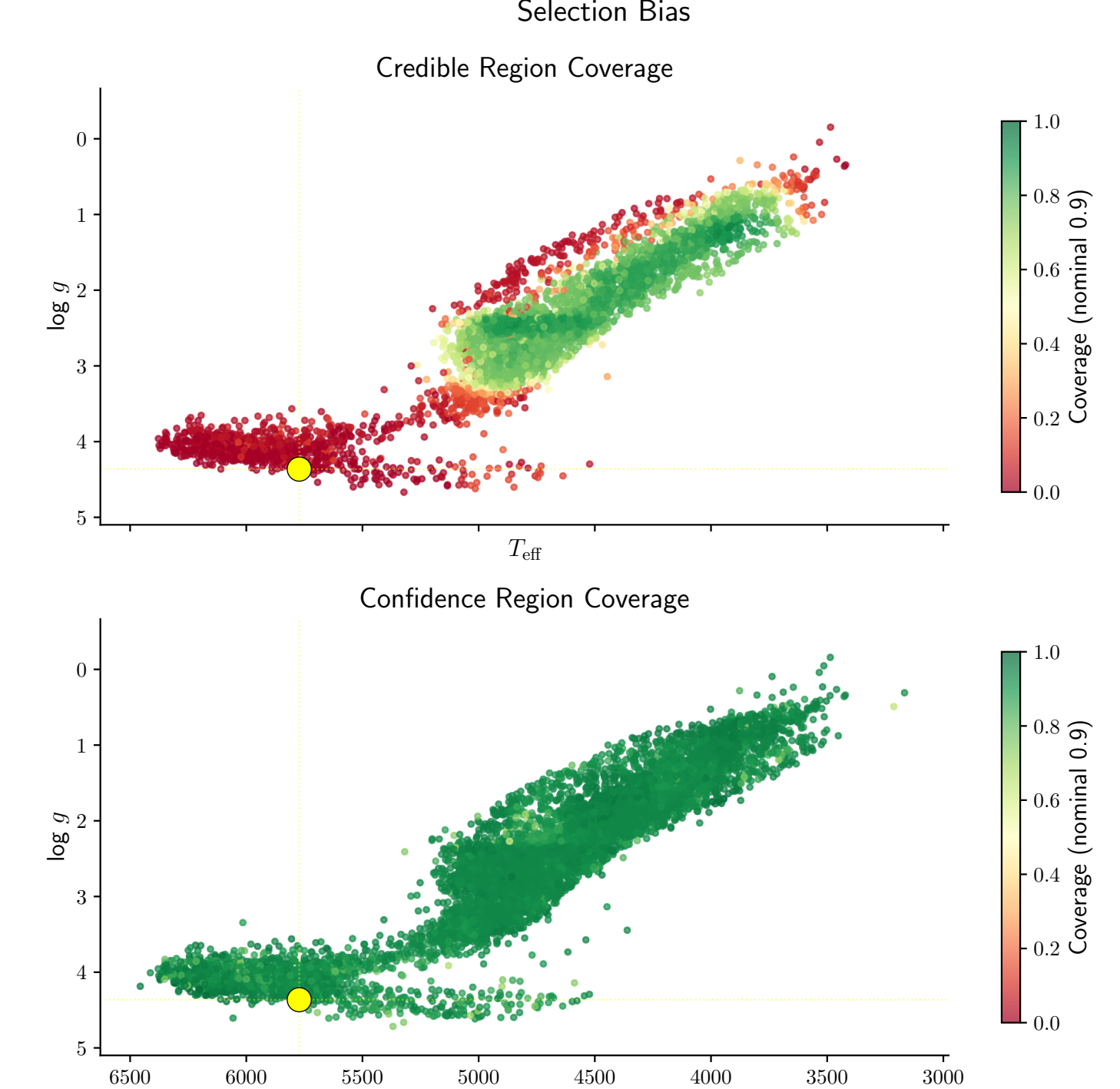


Local coverage at 90% level. Confidence regions achieves desired coverage while credible regions may under-cover.

## Results: Selection Bias



90% confidence region (green) and credible region (red) for the Sun-like star under selection bias.



Local coverage at 90% level. Confidence regions achieve desired coverage while credible regions severely under-cover.

## Takeaways

- ▶ **Confidence sets achieve nominal local coverage** in both settings
- ▶ **Credible regions under-cover**, especially under selection bias where the posterior fails to cover main-sequence and metal-poor stars
- ▶ **Calibration is modular and amortized:** any pre-trained posterior estimator can serve as the test statistic; only a held-out calibration set is required, no retraining needed

## Next Steps

- ▶ Extend to estimate additional stellar properties, such as  $[\alpha/M]$
- ▶ Explore richer parametric CDF families beyond the logistic model and alternative regressors
- ▶ Apply to larger set of the Gaia DR3 catalog (up to 220M+ stars) for large-scale inference

## References

[1] Andrae et al. (2023) [2] Laroche & Speagle (2025) [3] Carzon et al. (2026) [4] Dalmaso et al. (2021)