

# Position: Don't Use the CLT in LLM Evals With Fewer Than a Few Hundred Datapoints (...it's really easy to do a lot better!)

Sam Bowyer<sup>1</sup> Laurence Aitchison<sup>1,\*</sup> Desi R. Ivanova<sup>2,\*</sup> <sup>1</sup>University of Bristol <sup>2</sup>University of Oxford \*Equal Contribution

## 1. Failures of the Central Limit Theorem

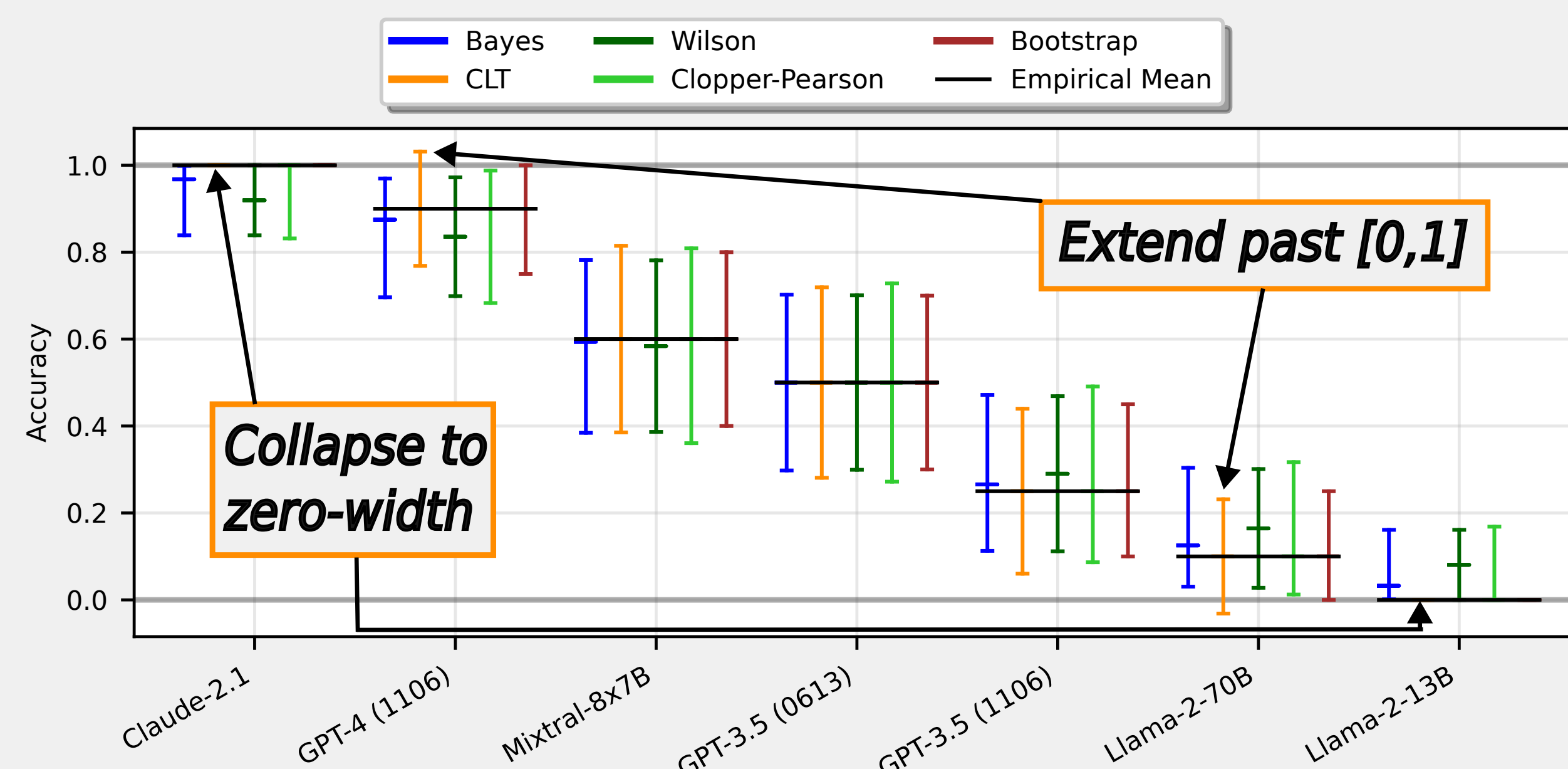
**Central Limit Theorem:** If  $X_1, \dots, X_N$  are IID r.v.s with mean  $\mu \in \mathbb{R}$  and finite variance  $\sigma^2$ , then with sample mean  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$ ,

$$\sqrt{N}(\hat{\mu} - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2) \text{ as } N \rightarrow \infty.$$

**CLT-based confidence interval** on binary data  $y_i \in \{0, 1\}$  (incorrect/correct) for  $i = 1, \dots, N$ :

$$CI_{1-\alpha} = \bar{y} \pm z_{\alpha/2} \sqrt{\bar{y}(1-\bar{y})/N},$$

where  $z_{\alpha/2}$  is the  $(1 - \alpha/2)$ -th quantile of  $\mathcal{N}(0, 1)$ .



As LLM capabilities improve, it's becoming more common to run small- $N$  benchmarks, such as the Langchain Typewriter tool-use benchmark shown above ( $N = 20$ ).

## 2. Simulation Setup

- Synthetic datasets:**  $N$  samples  $y_i \sim \text{Ber}(\theta)$ ,  $\theta \sim \text{Uniform}[0, 1]$ .
- Construct intervals with various confidence levels**  $1 - \alpha \in [0.8, 0.995]$ .
- Repeat the above**, and compare different interval-construction methods via **coverage** (proportion of intervals that contain true value of  $\theta$ ; should equal  $1 - \alpha$ ) and **interval width**.

## 3. Alternative Interval Construction

**Bayesian Beta-Bernoulli credible interval** – uniform prior on  $\theta$ :

$$\theta \sim \text{Beta}(1, 1) = \text{Uniform}[0, 1]$$

$$y_i \sim \text{Bernoulli}(\theta) \text{ for } i = 1, \dots, N$$

Use quantiles of closed-form posterior to construct  $1 - \alpha$  CIs:

$$\theta \mid y_{1:N} \sim \text{Beta} \left( 1 + \sum_{i=1}^N y_i, 1 + N - \sum_{i=1}^N y_i \right)$$

Beta-Bernoulli Bayesian Credible Interval

```
1posterior = scipy.stats.beta(1 + sum(y), 1 + N - sum(y))
2bayes_ci = posterior.interval(confidence=0.95)
```

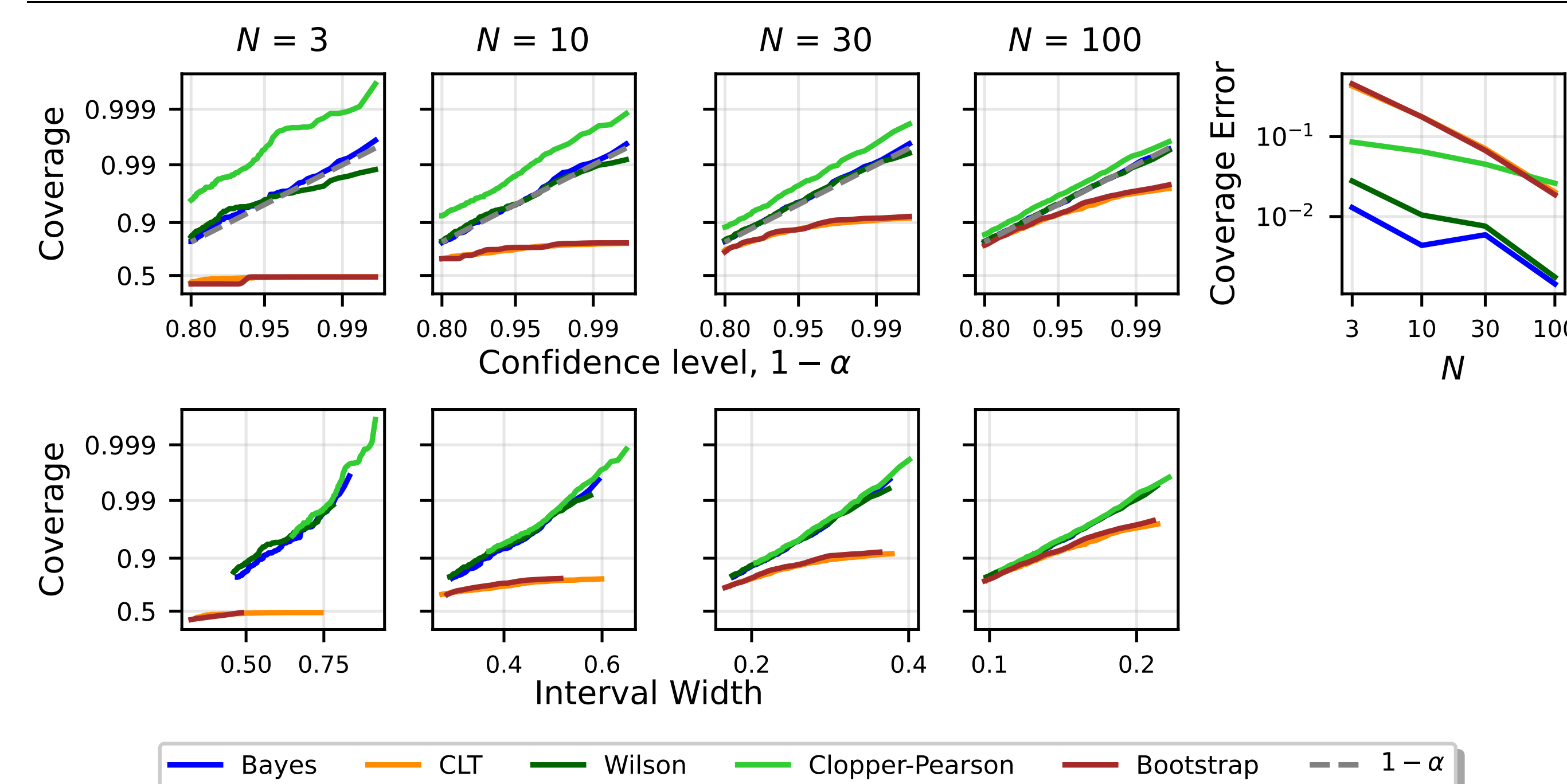
**Wilson-Score confidence interval** – based on binomial distribution:

$$CI_{1-\alpha, \text{Wilson}} = \frac{\hat{\theta} + \frac{z_{\alpha/2}^2}{2N}}{1 + \frac{z_{\alpha/2}^2}{N}} \pm \frac{\frac{z_{\alpha/2}}{2N}}{1 + \frac{z_{\alpha/2}^2}{N}} \sqrt{4N\hat{\theta}(1-\hat{\theta}) + z_{\alpha/2}^2}.$$

Wilson-Score Confidence Interval

```
1result = scipy.stats.binomtest(k=sum(y), n=N)
2wilson_ci = result.proportion_ci("wilson", 0.95)
```

## 4. Results

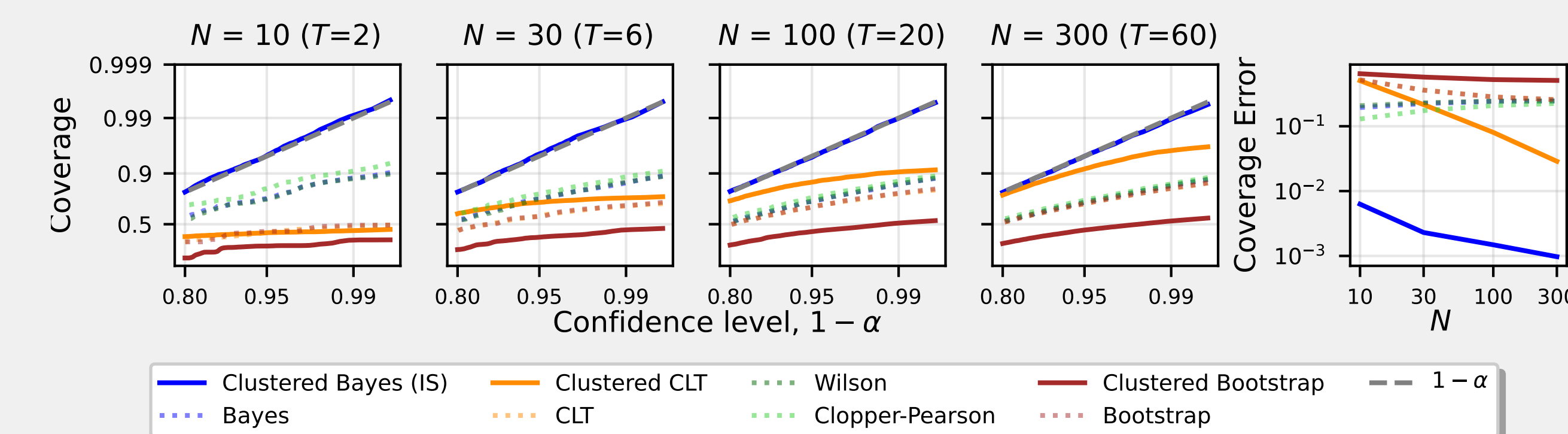


## 5. Other Settings

- Clustered Questions** – instead of  $N$  IID questions, we have  $T$  tasks, each with  $N_t$  IID questions. Bayesian model:

$$d \sim \text{Gamma}(1, 1), \quad \theta \sim \text{Beta}(1, 1),$$

$$\theta_t \sim \text{Beta}(d\theta, d(1-\theta)), \quad y_{i,t} \sim \text{Bernoulli}(\theta_t)$$



- Independent Comparisons** – Compare  $\theta_A$  and  $\theta_B$  for two different models, with access *only* to  $N_A, N_B, \hat{\theta}_A$ , and  $\hat{\theta}_B$ .
- Paired Comparisons** – Compare  $\theta_A$  and  $\theta_B$  for two different models, each with the same  $N$  IID questions and access to question-level successes  $\{y_{A,i}\}_{i=1}^N$  and  $\{y_{B,i}\}_{i=1}^N$ .
- Metrics that aren't simple averages of binary results** – e.g. F1 score (harmonic mean of precision and recall).

We construct Bayesian credible intervals that outperform CLT-based intervals in all settings; implemented in **bayes\_evals**.

## 6. Recommendations

**IID setting:** use Bayesian Beta-Bernoulli credible intervals or Wilson-score confidence intervals via **scipy** or equivalent.

**Other settings:** use Bayesian credible intervals as implemented in our simple package **bayes\_evals**.

