

APPROXIMATING FULL CONFORMAL PREDICTION FOR NEURAL NETWORK REGRESSION WITH GAUSS-NEWTON INFLUENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Uncertainty quantification is an important prerequisite for the deployment of deep learning models in safety-critical areas. Yet, this hinges on the uncertainty estimates being *useful* to the extent the predictive prediction intervals are well-calibrated and sharp. In the absence of inherent uncertainty estimates (*e.g.* pretrained models), popular approaches that operate post-hoc include Laplace’s method and split conformal prediction (split-CP). However, Laplace’s method can be miscalibrated when the model is misspecified and split-CP requires sample splitting, and thus comes at the expense of statistical efficiency. In this work, we construct prediction intervals for neural network regressors post-hoc without held-out data. This is achieved by approximating the *full* conformal prediction method (full-CP). Whilst full-CP nominally requires retraining the model for every test point and candidate label, we propose to train just once and locally perturb model parameters using Gauss-Newton influence to approximate the effect of retraining. Coupled with linearization of the network, we express the absolute residual non-conformity score as a piecewise linear function of the candidate label allowing for an efficient procedure that avoids the exhaustive search over the output space. On standard regression benchmarks and [bounding box localization](#), we show the resulting prediction intervals are locally-adaptive and often tighter than those of split-CP.

1 INTRODUCTION

Despite the impressive advancements in machine learning over recent years, most models, particularly neural networks, are still designed and trained to provide only point estimates, lacking the ability to rigorously quantify uncertainty in their predictions. This poses a significant challenge, as reliable decision-making depends on having a trustworthy representation of the uncertainty of each prediction. This need has drawn increased attention to uncertainty quantification in machine learning research, especially as the use of these models becomes more widespread and starts to permeate safety-sensitive fields such as healthcare and autonomous driving.

Conformal Prediction (CP) (Vovk et al., 2005; Angelopoulos & Bates, 2021) is a class of uncertainty quantification methods that represent uncertainty through prediction intervals. These intervals intuitively convey the degree of uncertainty—the larger the interval, the greater the uncertainty. What sets CP methods apart is their rigorous, *distribution-free coverage* guarantee: these intervals will encompass the true label with a probability of at least $1 - \alpha$, where α is a user-defined miscoverage rate. In recent years, a variant known as split or inductive conformal prediction (Papadopoulos et al., 2002) has gained traction within the machine learning community due to its ease of use and low computational cost. Split-CP, as the name suggests, divides the available data into training and calibration sets, using the former to fit a model and the latter to construct prediction intervals. However, this approach is statistically inefficient because it does not leverage all available data for model fitting, which would ideally ensure the best possible model fit while maintaining coverage. Full (or transductive) CP addresses this inefficiency, utilizing all data for both training and calibration, but it incurs a significant computational cost. It requires retraining the model multiple times—once for each test point and possible label—which is prohibitively expensive for deep learning applications.

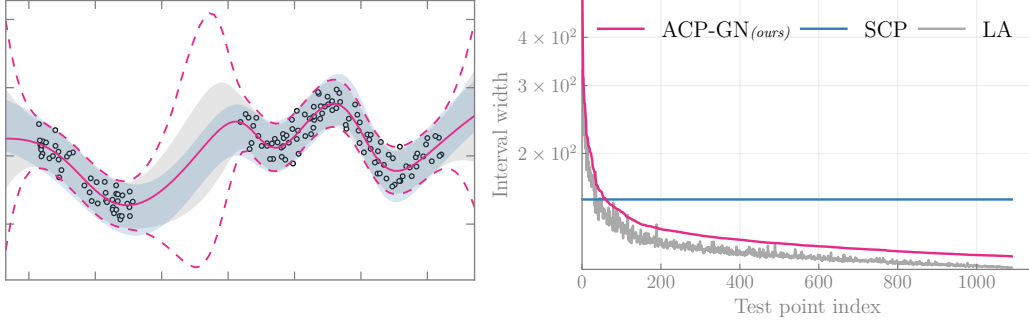


Figure 1: Our approx. full-CP via Gauss-Newton influence (ACP-GN) produces *adaptive* intervals (right figure)—similar to Bayes via Laplace approximation (LA)—while satisfying coverage as seen in the high-overlap with split-CP (SCP) close to the data (white dots with black edge in left figure).

In this work, we scale full-CP to neural network regression by (i) eliminating the need for retraining through a Gauss-Newton influence approximation and (ii) avoiding an exhaustive search over the label space via network linearization. This gives a new and scalable full-CP method for neural network regression we dub *approximate full-CP via Gauss-Newton influence* (ACP-GN). As a second contribution, we show the same tools can be applied to improve split-CP via normalization (Johansson et al., 2021), resulting in a new adaptable split-CP method we call (SCP-GN). We empirically validate ACP-GN and SCP-GN in multiple regression tasks [and bounding box localization](#), showing they not only satisfy coverage but also produce tight and adaptable prediction intervals (see Fig. 1 for results with ACP-GN).

2 BACKGROUND ON CONFORMAL PREDICTION

Neural network regression follows the canonical empirical risk minimization framework. Given a dataset $\mathcal{D}_N := \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ consisting of N examples with inputs $\mathbf{x}_i \in \mathbb{R}^D$ and targets $y_i \in \mathbb{R}$, we minimize the (regularized) empirical risk:

$$\theta_* = \arg \min_{\theta} \left(\sum_{i=1}^N \ell(y_i, f_i(\theta)) + \frac{1}{2} \delta \|\theta\|^2 \right) \quad (1)$$

where $f_i(\theta)$ is shorthand for $f(\mathbf{x}_i; \theta)$, the output of a deep neural network (DNN) at \mathbf{x}_i with parameters $\theta \in \mathbb{R}^D$, $\ell(y, f)$ is a loss function, and δ is an L_2 regularizer (*i.e.* weight decay). In this work, we restrict our attention to the squared-error loss: $\ell(y, f) = \frac{1}{2} (y - f)^2$. As is standard practice, the problem in Eq. (1) is approached using stochastic-gradient methods. For an unseen input \mathbf{x}_{N+1} , this gives us a point prediction $f_{N+1}(\theta_*)$. However, in an ideal scenario, we would like a prediction *interval* whose width reflects the uncertainty associated with that input. This is where conformal prediction (Vovk et al., 2005) comes into play. It provides a framework for constructing prediction intervals while satisfying the following frequentist coverage guarantee known as *marginal coverage*

$$\mathcal{P}(y_{N+1} \in C_{\alpha}(\mathbf{x}_{N+1})) \geq 1 - \alpha, \quad (2)$$

where y_{N+1} is the unseen target, $C_{\alpha}(\mathbf{x}_{N+1})$ is the prediction interval with desired miscoverage rate $\alpha \in (0, 1)$, and the probability is over all samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^{N+1}$, hence the name marginal coverage.

The prediction interval $C_{\alpha}(\mathbf{x}_{N+1})$ is constructed using *nonconformity* scores, which quantify how unusual a sample (\mathbf{x}_i, y_i) is compared to other samples in a set. In the context of regression, which is the focus of this paper, the absolute residual $R_i = R(\mathbf{x}_i, y_i) = |y_i - f_i(\theta)|$ is the most common score (Kato et al., 2023). Given the data \mathcal{D}_N and a score function, $C_{\alpha}(\mathbf{x}_{N+1})$ is defined as

$$C_{\alpha}(\mathbf{x}_{N+1}) = \{y \in \mathbb{R} : \pi(y) \leq \lceil (1 - \alpha)(N + 1) \rceil\}, \quad (3)$$

where $\pi(y) = \sum_{i=1}^{N+1} \mathbb{1}\{R_i \leq R(\mathbf{x}_{N+1}, y)\}$ is the rank of $R(\mathbf{x}_{N+1}, y)$ among the other N residuals. Remarkably, the only requirement for prediction intervals constructed as in Eq. (3) to satisfy marginal coverage is that the set of scores $\{R(\mathbf{x}_i, y_i)\}_{i=1}^{N+1}$ be exchangeable. This is equivalent to assuming the data is exchangeable and the score function (and consequently the regressor) preserves

this exchangeability by treating data points symmetrically. Thus, Eq. (2) is a distribution-free guarantee that remains valid even if the model is misspecified. This contrasts with Bayesian methods, the dominant approach for uncertainty quantification in deep learning, which are often poorly calibrated under model misspecification (Dawid, 1982; Fraser, 2011; Grünwald & van Ommen, 2017).

2.1 SPLIT CONFORMAL PREDICTION

As stated in the introduction, CP methods can generally be categorized into split-CP and full-CP variants. The primary distinction between these two families of methods lies in how they ensure the exchangeability of the set of scores $\{R(\mathbf{x}_i, y_i)\}_{i=1}^N \cup \{R(\mathbf{x}_{N+1}, y)\}$. Essentially, this means that all samples, including the test sample $(\mathbf{x}_{N+1}, y_{N+1})$, must exert the same influence over the prediction model $f(\cdot; \theta)$. Split-CP offers the simplest and most computationally efficient solution. By keeping the model fixed when computing the scores of each of the $N + 1$ samples, it ensures the scores are exchangeable, as the score function is applied elementwise. However, this approach is data inefficient because it prevents using the first N samples to fit the model, requiring a separate training dataset. This statistical inefficiency can negatively impact the quality of the final prediction intervals in two key ways. Since any conformal prediction method is designed to achieve coverage as described in Eq. (2), different approaches are compared in terms of their *efficiency*—how small the prediction intervals are—and *adaptability*—how much the size of the prediction intervals varies across samples. To be informative about both the true label and predictive uncertainty, prediction intervals must be both efficient and adaptable. However, in split-CP, efficiency is compromised because not all data is used for calibration, while fixing the model parameters reduces adaptability.

2.2 FULL CONFORMAL PREDICTION

Full-CP uses all available data for both training the model and computing prediction intervals. The main challenge in full-CP is that exchangeability of scores requires treating all data points symmetrically, meaning the model should be fit on $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ as well as $(\mathbf{x}_{N+1}, y_{N+1})$. Since y_{N+1} is unknown a priori, this necessitates retraining the model for all possible values y_{N+1} can take. To be precise, full-CP requires the following modification to the optimization problem in Eq. (1)

$$\theta_*^+(y) = \arg \min_{\theta} \left(\sum_{i=1}^N \ell(y_i, f_i(\theta)) + \ell(y, f_{N+1}(\theta)) + \frac{1}{2} \delta \|\theta\|^2 \right), \quad (4)$$

where $\theta_*^+(y)$ is the optimal model parameters for the augmented training set $\mathcal{D}_{N+1}(y) := \mathcal{D}_N \cup \{(\mathbf{x}_{N+1}, y)\}$ that includes the test point \mathbf{x}_{N+1} plus a candidate label y for y_{N+1} . With a slight abuse of notation, we use $R_i(y)$ to denote the residual with model parameters $\theta_*^+(y)$, that is,

$$R_i(y) = |y_i - f_i(\theta_*^+(y))| \quad \forall i = 1, \dots, N \quad \text{and} \quad R_{N+1}(y) = |y - f_{N+1}(\theta_*^+(y))|. \quad (5)$$

From here we can construct prediction intervals as in Eq. (3), but now the residuals vary for each test point \mathbf{x}_{N+1} and candidate label $y \in \mathbb{R}$. As mentioned before, this has two major limitations. Firstly it requires retraining the model for every candidate label y which is infeasible for DNNs. Secondly in theory the method asks to consider an uncountable set (*i.e.* all real numbers). Therefore, in practice a finite grid of possible labels for y is used, typically delimited by the training targets. Evidently, the grid imposes computation-precision trade-off and has implications for the coverage if a valid candidate happens to lie between two grid points (Chen et al., 2018). In a few cases, the prediction set can be computed efficiently and exactly without the need to trial candidate labels of y . In addition, this procedure only depends on a single fit on the original (unaugmented) dataset which can be efficiently updated not only for variations in y but also for different inputs \mathbf{x}_{N+1} . These include the Lasso (Lei, 2019), k-Nearest Neighbours Regression (Papadopoulos et al., 2011), and ridge regression (Noureddin et al., 2001; Burnaev & Vovk, 2014). In the following section, we review conformalized ridge regression which is the basis of our approximate full-CP procedure.

2.3 CONFORMALIZED RIDGE REGRESSION (CRR)

Ridge regression is a special case of Eq. (1) where we have a linear model $f_i(\theta) := \mathbf{x}_i^\top \theta$. In this case, the absolute residual can be written as a piecewise linear function of the candidate y with $R_i(y) = |a_i + b_i y|$, where a_i and b_i are coefficients capturing information from the training data and test point, resp. It is then convenient to express the rank as $\pi(y) = \sum_{i=1}^{N+1} \mathbb{1}\{y \in S_i\}$, with

$$S_i = \{y : R_i(y) \leq R_{N+1}(y)\} = \{y : |a_i + b_i y| \leq |a_{N+1} + b_{N+1} y|\}. \quad (6)$$

Each set S_i can be an interval, a point, a ray, a union of two rays, the real line, or the empty set. Eq. (6) suggests that the rank for a given y can only change at points where $R_i(y) = R_{N+1}(y)$ to which we refer as “changepoints”. For the absolute residual, these changepoints fall into one of the following two cases, where we assume $b_i \geq 0$ (if needed, multiplying a_i and b_i by -1):

- If $b_i \neq b_{N+1}$, then S_i is an interval or a union of two rays, and we have two changepoints, $-(a_i - a_{N+1})/(b_i - b_{N+1})$ and $-(a_i + a_{N+1})/(b_i + b_{N+1})$.
- If $b_i = b_{N+1} \neq 0$, then S_i is a ray, unless $a_i = a_{N+1}$, in which case $S_i = \mathbb{R}$. Here we have single changepoint $-(a_i + a_{N+1})/2b_i$.

This leads to an exact form of the prediction set by taking the union of finitely many intervals and rays whose endpoints are given by the changepoints sorted in increasing order. Given the changepoints, different implementations have been proposed with varying time complexity. For the smaller datasets, we use the conformal ridge regression confidence machine algorithm (Noureddinov et al., 2001; Vovk et al., 2005), which uses the absolute residual and changepoints described above.

There is also a simpler, asymmetric version of CRR that uses the signed residuals (Burnaev & Vovk, 2014). In this case, we can compute the lower and upper bounds of the prediction interval separately, using residuals $f_i(\theta) - y_i$ for the lower bound, and $y_i - f_i(\theta)$ for the upper bound. This is the version we use for the larger datasets and that appears in the depiction of our method in Alg. 2, where we use l_i and u_i to denote the changepoints for the lower and upper bounds, resp. When using signed residuals, S_i is either a ray with changepoint $l_i = u_i = (a_i - a_{N+1})/(b_{N+1} - b_i)$ if $b_{N+1} - b_i > 0$ or otherwise, $S_i = \mathbb{R}$ with $l_i = -\infty$ and $u_i = \infty$.

Finally, we need to write down the expression for coefficients a_i and b_i . Using the Sherman-Morrison formula, a widely-used tool of the regressions diagnostics literature (Cook, 1977), Burnaev & Vovk (2014) showed that the required coefficients a_1, \dots, a_{N+1} and b_1, \dots, b_{N+1} for the CRR procedure can be efficiently computed for different \mathbf{x}_{N+1} by a simple rank-1 update or *perturbation* to the ridge solution on \mathcal{D}_N (see App. A for derivation):

$$y_i - \mathbf{x}_i^T \theta_*^+(y) = y_i - \underbrace{\mathbf{x}_i^T \theta_*}_{a_i} + \underbrace{\frac{h_{i,N+1}}{1 + h_{N+1}} \mathbf{x}_{N+1}^T \theta_*}_{b_i} - \frac{h_{i,N+1}}{1 + h_{N+1}} y \quad (7)$$

$$y - \mathbf{x}_{N+1}^T \theta_*^+(y) = -\underbrace{\frac{1}{1 + h_{N+1}} \mathbf{x}_{N+1}^T \theta_*}_{a_{N+1}} + \underbrace{\frac{1}{1 + h_{N+1}} y}_{b_{N+1}} \quad (8)$$

where $h_{N+1} = \mathbf{x}_{N+1}^T \mathbf{H}_*^{-1} \mathbf{x}_{N+1}$ with Hessian matrix $\mathbf{H}_* = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \delta \mathbf{I}$.

2.4 NORMALIZED NONCONFORMITY SCORES

In the previous section, we derived CRR using the nonconformity score given in Eq. (5), which is often referred to as add-one-in (AOI). In this section, we also consider the leave-one-out (LOO) and studentized scores, which lead to the variants CRR-deleted and CRR-studentized and their corresponding extensions to neural network regression given by our approximate full-CP method. All those variants are valid choices (Vovk et al., 2005; Shafer & Vovk, 2008), and the literature is not conclusive regarding which one is to be preferred (Fong & Holmes, 2021; Fontana et al., 2023). Yet, in our neural network regression experiments, the studentized variant outperformed the other two.

In the leave-one-out (LOO) variety (also known as jackknife), scores R_i^{LOO} are computed by excluding the i^{th} data point from the augmented data $\mathcal{D}_{N+1}(y)$ before retraining the model. The only exception being R_{N+1}^{LOO} , which requires no retraining, as θ_* already ignores the $(N+1)^{\text{th}}$ data point. In the case of ridge regression with absolute residuals, we can derive the jackknife score from the standard one as (Vovk et al., 2005)

$$R_i^{\text{LOO}} = R_i / (1 - \bar{h}_i), \quad (9)$$

where we have introduced the leverage score defined with respect to the augmented problem $\bar{h}_i = \mathbf{x}_i^T \bar{\mathbf{H}}_*^{-1} \mathbf{x}_i$, where $\bar{\mathbf{H}}_* = \sum_{i=1}^{N+1} \mathbf{x}_i \mathbf{x}_i^T + \delta \mathbf{I}$. The leverage score can be viewed as a kind of diagnostics measure for measuring the sensitivity of the prediction to changes in the target. The formula is a consequence of the exact rank-1 updates available in ridge regression. This leads to the

deleted-CRR method which proceeds in the same way as standard CRR except it uses the normalized coefficients: $a_i \leftarrow a_i/(1 - \bar{h}_i)$ and $b_i \leftarrow b_i/(1 - \bar{h}_i)$ for $i = 1, \dots, N + 1$. Such relations can also be used to recover the standard score starting from the jackknife one:

$$R_i = R_i^{\text{LOO}} / (1 + \bar{h}_i^{\setminus i}), \quad (10)$$

where $\bar{h}_i^{\setminus i} = \mathbf{x}_i^\top \bar{\mathbf{H}}_*^{\setminus i} \mathbf{x}_i$ with $\bar{\mathbf{H}}_*^{\setminus i} = \sum_{j=1, j \neq i}^{N+1} \mathbf{x}_j \mathbf{x}_j^\top + \delta \mathbf{I}$. Papadopoulos (2024) gave an interpretation of this relation as locally-weighted conformal prediction (Papadopoulos et al., 2008) where the expression in the denominator can be seen as the leave-one-out predictive variance from a Bayesian perspective. This can be seen as a measure of the difficulty of the i^{th} example. Similar normalizations can also be applied to nonconformity scores in a split-CP framework (Vovk et al., 2005). Finally, the studentized-CRR can be interpreted as a compromise between standard and deleted-CRR. It uses a similar normalization and defines nonconformity scores as

$$R_i^{\text{student}} = R_i / \sqrt{1 - \bar{h}_i}. \quad (11)$$

This transformation is applied analogously to the CRR coefficients with a_i and b_i .

3 APPROXIMATE FULL-CP FOR NEURAL NETWORK REGRESSION

We propose to leverage the conformalized ridge regression framework (Noureddin et al., 2001) along with a carefully constructed perturbation to the NN solution θ_* that approximates the solution to the augmented problem. The former avoids the need to specify a grid of candidate labels y and the latter allows us to express the absolute residual nonconformity score as a piecewise linear function of the candidate label akin to that in ridge regression. Most importantly, these choices mean that a neural network is trained only once on the original data \mathcal{D}_N .

We propose to use *Gauss-Newton influence* (St Laurent & Cook, 1992) to approximate the solution to the augmented problem $\theta_*^+(y)$ in Eq. (4). Whilst originally proposed for leave-one-out (LOO) estimation to evaluate the influence of individual (or groups of) training examples on the model, we adapt this for add-one-in (AOI) estimation. This approximation $\hat{\theta}_*^+(y) \approx \theta_*^+(y)$ is given by

$$\hat{\theta}_*^+(y) = \theta_* + \frac{\hat{e}_{N+1}(y)}{1 + \hat{h}_{N+1}} \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1}, \quad (12)$$

where $\phi_i := \nabla_{\theta} f_i(\theta_*)^\top$ is the Jacobian of the neural network at θ_* , $\mathbf{H}_{\text{GN}} = \sum_{i=1}^N \phi_i \phi_i^\top + \delta \mathbf{I}$ is the Gauss-Newton approximation to the Hessian, $\hat{e}_{N+1}(y) = y - f_{N+1}(\theta_*)$ is the residual for the $(N + 1)^{\text{th}}$ example and $\hat{h}_{N+1} = \phi_{N+1}^\top \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1}$ can be interpreted as a leverage score. Eq. (12) can be understood as taking a single step of the Gauss-Newton algorithm at θ_* on the augmented objective in Eq. (4). St Laurent & Cook (1992) showed that this can be expressed as a least-squares solution to a surrogate “linear” model allowing it to be reformulated as a deviation from θ_* .

Next we show that using Eq. (12) along with linearization of the neural network, we can approximate the residual as a linear function of the postulated label y recovering an identical form to that of conformalized ridge regression in Eqs. (7) and (8),

$$y_i - f_i(\theta_*^+(y)) \approx \underbrace{y_i - f_i(\theta_*) + \frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} f_{N+1}(\theta_*)}_{a_i} - \underbrace{\frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} y}_{b_i} \quad (13)$$

$$y - f_{N+1}(\theta_*^+(y)) \approx \underbrace{-\frac{1}{1 + \hat{h}_{N+1}} f_{N+1}(\theta_*)}_{a_{N+1}} + \underbrace{\frac{1}{1 + \hat{h}_{N+1}} y}_{b_{N+1}} \quad (14)$$

where $\hat{h}_{i,N+1} = \phi_i^\top \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1}$ (see App. B for the derivation and App. C for the extension to the multi-output setting). The coefficients can be readily adapted for the case of normalized nonconformity scores in the same way as outlined in Sec. 2.4.

In Alg. 2, we provide a simplified but complete algorithmic depiction of our method, which we dub *approximate full-CP via Gauss-Newton influence* (ACP-GN). We contrast it against standard

full-CP in Alg. 1, highlighting the expensive grid search and retraining steps. Note that in ACP-GN there is no grid search and the model is fit only once using all available \mathcal{D}_N , as highlighted in blue. At test time, when a new point \mathbf{x}_{N+1} is received, instead of retraining the model, ACP-GN only computes the CRR coefficients $\{(a_i, b_i)\}_{i=1}^{N+1}$ using Gauss-Newton influence as in Eq. (13) and Eq. (14), which are then used to derive the prediction set in closed form either using the method of (Noureddinov et al., 2001) or the asymmetric one of (Burnaev & Vovk, 2014) that we show in Alg. 2.

Algorithm 1: Standard Full-CP.

```

for each test point  $\mathbf{x}_{N+1}$  do
  for each  $y$  in a given grid do
    optimize  $\theta_*^+(y)$  as in Eq. (4)
     $R_{N+1}(y) = |y - f_{N+1}(\theta_*^+(y))|$ 
    for  $i \in \{1, \dots, N\}$  do
       $|R_i(y) = |y_i - f_i(\theta_*^+(y))|$ 
     $\pi(y) = \sum_{i=1}^{N+1} \mathbb{1}\{R_i(y) \leq R_{N+1}(y)\}$ 
    if  $\pi(y) \leq \lceil (1 - \alpha)(N + 1) \rceil$  then
      include  $y$  in  $C_\alpha(\mathbf{x}_{N+1})$ 
```

Algorithm 2: ACP-GN (ours).

```

optimize  $\theta_*$  as in Eq. (1)
for each test point  $\mathbf{x}_{N+1}$  do
  compute  $a_{N+1}, b_{N+1}$  as in Eq. (14)
  for  $i \in \{1, \dots, N\}$  do
    compute  $a_i, b_i$  as in Eq. (13)
    if  $b_{N+1} - b_i > 0$  then
       $| l_i = u_i = (a_i - a_{N+1}) / (b_{N+1} - b_i)$ 
    else
       $| l_i = -\infty$  and  $u_i = \infty$ 
  sort  $\{l_i\}_{i=1}^N$  and  $\{u_i\}_{i=1}^N$  in ascending order
   $C_\alpha(\mathbf{x}_{N+1}) =$ 
     $| [l_{\lfloor (N+1)(\alpha/2) \rfloor}, u_{\lceil (N+1)(1-\alpha/2) \rceil}]$ 
```

3.1 CONFORMALIZING LINEARIZED LAPLACE

We can show that our approximate full conformal regression by Gauss-Newton influence can be interpreted as conformalizing Linearized Laplace (MacKay, 1992; Khan et al., 2019; Immer et al., 2021b) for regression. This relates to the result of Burnaev & Vovk (2014) who showed the CRR procedure can be viewed as conformalizing or “de-Bayesing” Bayesian Linear Regression (with Gaussian assumptions) (Burnaev & Vovk, 2014). They provide asymptotic results indicating that in well-specified settings, the conformal prediction intervals and Bayesian credible intervals closely align. This places our method within the wider context of *Conformal Bayes* (Melluish et al., 2001; Wasserman, 2011) for recalibrating Bayesian intervals in case of model misspecification.

The Laplace approximation constructs a Gaussian posterior approximation centered around the point estimate θ_* and covariance given by the inverse of the Hessian (local curvature) of the empirical risk evaluated at θ_* . When the Hessian is approximated by the generalized Gauss-Newton matrix, we refer to the resulting Laplace approximation as the Laplace-GGN posterior, $q_*(\theta) = N(\theta | \theta_*, \Sigma_*)$ where $\Sigma_* = \mathbf{H}_{\text{GN}}^{-1}$. This is often accompanied by linearizing the output of the neural network about θ_* (Foong et al., 2019; Immer et al., 2021b):

$$f_i(\theta) \approx f_i^{\text{lin}}(\theta) = f_i(\theta_*) + \nabla_{\theta} f_i(\theta_*)^{\top} (\theta - \theta_*). \quad (15)$$

The overall method is referred to as Linearized Laplace and retains the original NN point prediction as the mean of the posterior predictive. Analogous to the add-one-in estimate in Eq. (12), we can derive an approximation to the add-one-in posterior $\hat{q}_*^+(\theta)$ by perturbing the Laplace-GGN posterior (see App. D for derivation) whose mean is equal to $\hat{\theta}_*^+(y)$:

$$\hat{q}_*^+(\theta) = \mathbb{N}(\theta | \hat{\theta}_*^+(y), \hat{\Sigma}_*^+) \quad \text{where } \hat{\Sigma}_*^+ = (\mathbf{H}_{\text{GN}} + \phi_{N+1} \phi_{N+1}^{\top})^{-1} \quad (16)$$

This is a simple extension of the leave-one-out results in (Nickl et al., 2023) to the add-one-in case. Using this perturbed posterior in combination with the linearized predictor in Eq. (15), we recover Eqs. (13) and (14). From the perspective of Linearized Laplace, the Gauss-Newton influence gives the exact AOI solution with respect to the linearized network and in particular its feature expansion given the Jacobian of the network. However, it is often the case in practice that θ_* is not a minimum of the empirical risk (*i.e.* neural network not trained to convergence). Consequently, θ_* is not a minima of the linearized network’s objective. This can be corrected by solving for the following objective,

$$\tilde{\theta} = \arg \min_{\theta} \left(\sum_{i=1}^N \frac{1}{2} (\tilde{y}_i - \phi_i^{\top} \theta)^2 + \frac{1}{2} \delta \|\theta\|^2 \right) \quad (17)$$

where $\tilde{y}_i := \phi_i^{\top} \theta_* + e_i$ with residual $e_i = y_i - f_i(\theta_*)$. This is a linear-Gaussian system and hence can be solved for in a single step. This process is often referred to as “refinement” in the Linearized

Laplace literature and has been shown to improve predictions (Immer et al., 2021b). The exact AOI solution with respect to Eq. (17) is then given by,

$$\tilde{\theta}_*^+(y) = \tilde{\theta} + \frac{y - f_{N+1}^{\text{lin}}(\tilde{\theta})}{1 + \hat{h}_{N+1}} \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1} \quad (18)$$

which reduces to Eq. (12) when $\tilde{\theta} = \theta_*$. We can derive analogous expressions to Eqs. (13) and (14) which are exact for the linearized neural network in Eq. (15),

$$y_i - f_i^{\text{lin}}(\tilde{\theta}_*^+(y)) = y_i - f_i^{\text{lin}}(\tilde{\theta}) + \frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} f_{N+1}^{\text{lin}}(\tilde{\theta}) - \frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} y \quad (19)$$

$$y - f_{N+1}^{\text{lin}}(\tilde{\theta}_*^+(y)) = -\frac{1}{1 + \hat{h}_{N+1}} f_{N+1}^{\text{lin}}(\tilde{\theta}) + \frac{1}{1 + \hat{h}_{N+1}} y \quad (20)$$

where the main distinction is the use of linearized versions of the original neural network prediction.

Normalized Split-CP with Gauss Newton Influence It turns out the tools we used to derive our approximate full-CP method are also effective to improve the adaptability and efficiency of split-CP. In vanilla split conformal regression, all prediction intervals have the same width because the model remains fixed. One way to alleviate this issue is to *normalize* the scores as $R_i = R_i/\sigma_i$, where σ_i estimates the difficulty in predicting the i^{th} data point correctly (Papadopoulos et al., 2008; Johansson et al., 2021). As observed in (Papadopoulos, 2024), the normalized scores described in Sec. 2.4 are scaled by the predictive variance, which is closely related to how difficult y_i is to predict. This motivates our Gauss-Newton split-CP variant, with scores,

$$R_i = |y_i - f_i(\theta_*)|/\sqrt{1 + \hat{h}_i}, \quad (21)$$

where $\hat{h}_i = \phi_i^\top \mathbf{H}_{\text{GN}}^{-1} \phi_i$ is marginal variance given by Linearized Laplace. This is similar to Eq. 4.10 in (Vovk et al., 2005) and is analogous to the studentized scores in the full-CP case.

4 RELATED WORK

Ours is not the first work to alleviate the need for model refitting in full conformal prediction. In fact, split-CP (Papadopoulos et al., 2002; Lei et al., 2018) and cross-conformal predictors (Vovk, 2015) were also explicitly designed to avoid retraining the model, albeit at the cost of statistical efficiency. Some other notable examples that also try to preserve statistical efficiency include the use of homotopy continuation technique (Ndiaye & Takeuchi, 2019) and algorithmic stability Ndiaye (2022). Recently, Guha et al. (2024) showed one can bypass some of the challenges in conformal regression by framing it as a classification problem (Guha et al., 2024), but that requires binning the output space, incurring an accuracy/precision trade-off.

The work of (Martinez et al., 2023) is particularly relevant, since it leverages influence function (Hampel, 1974; Cook & Weisberg, 1980) to approximate the full-CP algorithm. However, they only considered classification problems and their solution still requires an exhaustive search over possible labels, which is critical for the regression setting we consider in this paper. Moreover, in this work, we apply the closely-related Gauss-Newton influence rather than influence function. Firstly, the Hessian is often approximated by the Gauss-Newton matrix when used in influence functions for deep learning (Bae et al., 2022) since it is guaranteed to be positive semi-definite (Martens, 2010). Secondly, it is possible to recover the typical inverse Hessian vector product form of influence function by taking an *infinitesimal* variant of the Gauss-Newton influence.

5 EXPERIMENTS & RESULTS

We evaluate the performance of our method, approximate full-CP via Gauss-Newton influence (ACP-GN), against Linearized Laplace (LA) (MacKay, 1992; Immer et al., 2021b), a recently popular Bayesian method for post-hoc uncertainty quantification, split conformal prediction (SCP) (Papadopoulos et al., 2002), [conformalized residual fitting \(CRF\)](#) (Papadopoulos et al., 2002) and [conformalized quantile regression \(CQR\)](#) (Romano et al., 2019). We also evaluate two further proposals, “ACP-GN (split + refine)” and “SCP-GN” which we proceed to describe along with the aforementioned methods:

Table 1: Our proposed approximate full-CP via Gauss-Newton influence (ACP-GN) almost always gives the tightest intervals in limited data regimes whilst satisfying the desired coverage (cf. *yacht*, *boston*, *energy*). On larger datasets, ACP-GN remains competitive on efficiency compared with other conformal methods but can sometimes miscover. As a remedy, we propose two variants inspired by ACP-GN that generally fix the miscoverage issue. Average prediction interval width and coverage for our proposed approaches (shaded gray) against baselines (non-shaded) for three different settings of the confidence level. The best average widths over well-calibrated approaches (indicated by ✓/✗) appear in bold. Reported metrics are accompanied by standard error from repeated runs.

		Avg. Width		Avg. Coverage			
		90%	95%	90%	95%	90%	99%
yacht N=308 I=6	LA	1.690±0.017	2.014±0.020	2.647±0.027	88.73±0.61 (✓)	90.78±0.59 (✗)	93.89±0.60 (✗)
	SCP	2.553±0.093	4.001±0.115	10.018±0.361	89.56±0.66 (✓)	94.07±0.39 (✓)	99.32±0.08 (✓)
	CRF	2.526±0.092	3.947±0.115	9.674±0.294	89.53±0.64 (✓)	94.10±0.38 (✓)	99.29±0.10 (✓)
	CQR	4.090±0.105	5.845±0.187	18.650±0.484	89.94±0.42 (✓)	94.42±0.32 (✓)	99.02±0.17 (✓)
	ACP-GN	1.594±0.016	2.385±0.029	6.915±0.067	87.36±0.58 (✓)	92.56±0.68 (✓)	99.03±0.11 (✓)
	SCP-GN	2.270±0.086	3.349±0.098	7.216±0.254	89.85±0.51 (✓)	94.91±0.32 (✓)	99.19±0.15 (✓)
boston N=506 I=13	ACP-GN (split + refine)	1.993±0.020	2.954±0.037	7.307±0.178	89.35±0.62 (✓)	94.90±0.51 (✓)	99.45±0.08 (✓)
	LA	9.398±0.046	11.199±0.055	14.718±0.072	91.24±0.31 (✓)	94.34±0.22 (✓)	97.53±0.11 (✗)
	SCP	10.635±0.123	14.509±0.171	36.272±1.847	89.56±0.42 (✓)	94.64±0.32 (✓)	99.11±0.13 (✓)
	CRF	11.932±0.605	16.073±0.862	40.690±3.333	90.01±0.33 (✓)	94.77±0.22 (✓)	99.30±0.08 (✓)
	CQR	11.692±0.129	15.115±0.213	31.628±1.822	90.10±0.33 (✓)	95.12±0.24 (✓)	99.07±0.14 (✓)
	ACP-GN	9.182±0.046	12.111±0.038	20.512±0.057	90.64±0.26 (✓)	95.49±0.16 (✓)	99.11±0.08 (✓)
energy N=768 I=8	SCP-GN	10.301±0.089	13.418±0.151	24.714±0.865	89.52±0.50 (✓)	94.82±0.32 (✓)	99.05±0.12 (✓)
	ACP-GN (split + refine)	13.103±0.072	16.729±0.134	27.561±0.445	90.12±0.26 (✓)	95.41±0.20 (✓)	99.27±0.10 (✓)
	LA	1.502±0.006	1.790±0.007	2.353±0.009	88.96±0.35 (✓)	92.92±0.33 (✗)	96.95±0.23 (✗)
	SCP	1.942±0.032	2.486±0.046	3.772±0.093	89.44±0.28 (✓)	94.80±0.20 (✓)	99.18±0.08 (✓)
	CRF	1.923±0.031	2.454±0.046	3.728±0.092	89.39±0.28 (✓)	94.78±0.22 (✓)	99.14±0.08 (✓)
	CQR	4.670±0.030	5.139±0.029	6.438±0.120	90.08±0.26 (✓)	95.24±0.21 (✓)	98.96±0.09 (✓)
bike N=10,886 I=18	ACP-GN	1.462±0.006	1.884±0.008	3.076±0.015	88.28±0.33 (✓)	93.69±0.33 (✓)	98.88±0.11 (✓)
	SCP-GN	1.911±0.029	2.449±0.044	3.609±0.071	89.69±0.29 (✓)	94.79±0.18 (✓)	99.21±0.09 (✓)
	ACP-GN (split + refine)	1.745±0.016	2.174±0.021	3.300±0.045	90.54±0.25 (✓)	94.96±0.22 (✓)	99.18±0.10 (✓)
	LA	100.451±2.394	119.694±2.853	157.305±3.749	89.82±0.39 (✓)	93.29±0.33 (✗)	96.83±0.16 (✗)
	SCP	131.138±0.812	180.477±1.244	324.756±4.635	90.33±0.21 (✓)	95.17±0.15 (✓)	99.00±0.07 (✓)
	CRF	127.836±0.894	174.362±1.376	311.580±5.077	90.39±0.19 (✓)	95.26±0.14 (✓)	99.01±0.07 (✓)
protein N=45,730 I=9	CQR	141.329±5.943	167.682±5.835	244.863±4.952	89.83±0.23 (✓)	94.80±0.14 (✓)	98.89±0.07 (✓)
	ACP-GN	98.813±2.485	130.893±3.231	213.131±5.630	89.36±0.43 (✓)	94.41±0.27 (✗)	98.67±0.09 (✗)
	SCP-GN	122.245±1.073	160.505±1.761	254.409±3.767	90.34±0.24 (✓)	95.26±0.15 (✓)	99.02±0.08 (✓)
	ACP-GN (split + refine)	128.336±4.336	170.782±5.859	281.632±10.176	89.98±0.22 (✓)	94.94±0.16 (✓)	99.01±0.06 (✓)
	LA	9.385±0.022	11.183±0.027	14.697±0.035	85.43±0.18 (✗)	89.69±0.15 (✗)	94.81±0.10 (✗)
	SCP	13.041±0.088	17.161±0.098	26.181±0.119	89.78±0.08 (✓)	94.83±0.06 (✓)	98.94±0.04 (✓)
facebook_2 N=81,311 I=53	CRF	12.645±0.127	16.931±0.146	26.973±0.202	89.86±0.09 (✓)	94.84±0.07 (✓)	98.93±0.04 (✓)
	CQR	13.541±0.144	14.798±0.129	18.239±0.041	90.07±0.10 (✓)	95.07±0.09 (✓)	98.96±0.04 (✓)
	ACP-GN	10.243±0.019	13.294±0.027	20.101±0.053	87.54±0.15 (✗)	93.04±0.11 (✗)	98.24±0.05 (✗)
	SCP-GN	12.426±0.085	16.102±0.096	24.032±0.138	89.78±0.10 (✓)	94.86±0.08 (✓)	98.94±0.03 (✓)
	ACP-GN (split + refine)	12.660±0.028	16.073±0.031	23.445±0.057	89.83±0.09 (✓)	94.90±0.09 (✓)	98.97±0.05 (✓)
	LA	66.088±2.760	78.749±3.289	103.493±4.322	97.47±0.12 (✗)	98.01±0.09 (✗)	98.65±0.06 (✗)
facebook_2 N=81,311 I=53	SCP	16.387±0.208	35.387±0.462	152.706±1.591	89.97±0.07 (✓)	95.00±0.06 (✓)	99.06±0.03 (✓)
	CRF	15.088±0.188	29.679±0.396	102.326±2.552	89.94±0.07 (✓)	94.98±0.06 (✓)	99.03±0.02 (✓)
	CQR	17.605±0.645	21.571±0.660	30.852±1.303	90.16±0.28 (✓)	95.13±0.12 (✓)	99.01±0.03 (✓)
	ACP-GN	18.396±0.546	40.088±1.091	166.792±5.632	90.47±0.11 (✗)	95.45±0.08 (✗)	99.35±0.04 (✗)
	SCP-GN	16.287±0.202	33.655±0.563	118.489±4.305	89.99±0.07 (✓)	94.98±0.06 (✓)	99.00±0.03 (✓)
	ACP-GN (split + refine)	21.469±0.906	42.184±0.788	152.460±2.499	90.14±0.08 (✓)	95.10±0.06 (✓)	99.13±0.02 (✗)

- **LA**: The Laplace approximation with the Hessian approximated by the generalized Gauss-Newton matrix. The linearized predictive in Eq. (15) is used for inference. We use the implementation provided in the Laplace PyTorch library (Daxberger et al., 2021a) as well as to implement our ACP-GN method.
- **SCP**: Uses absolute residual nonconformity score in the procedure outlined in Sec. 2.1.
- **CRF**: Trains an additional network to predict the absolute residuals of the original network which is then used to normalize the absolute residual nonconformity score.
- **CQR**: Trains a quantile regression network using the pinball loss function. The predicted lower and upper quantile functions are then used in the split conformal quantile regression algorithm.
- **ACP-GN**: Uses the *studentized* nonconformity score of Eq. (11). We found this score performed best, but we also report results for the standard and deleted versions in App. H.
- **SCP-GN**: Normalizes the absolute residual nonconformity score by the posterior predictive standard deviation given by the LA method (trained only on the same split as SCP) as shown in Eq. (21).
- **ACP-GN (split + refine)**: Uses a train-calibration split like in SCP, where we pretrain the model on the training set before running ACP-GN on the calibration set. More precisely,

it solves the linearized network’s objective in Eq. (17) but defined on the calibration set. It then uses the linearized network prediction in Eq. (15) in lieu of the original network to evaluate the coefficient in Eqs. (19) and (20), again on the calibration set.

5.1 UCI REGRESSION

We conduct experiments on popular benchmark datasets for regression taken from the UCI Machine Learning repository (Nottingham et al., 2024). These are varying in size and consequently lead to slight variations in experimental setup. Hence we place them into 3 groups for easy referencing: *small* (boston, concrete, energy, wine, yacht); *medium* (kin8nm, power); *large* (bike, community, protein, facebook_1, facebook_2). A subset of these are shown in Table 1 and the remainder are reported in App. G.1. For the small datasets, the reported metrics result from 10 repeats of a 10-fold cross-validation process. For the other datasets, we perform 20 different train-test splits. In both cases, 90% of the examples are used for the training set and the remaining 10% for the test set. When the method requires a calibration set, the training set is divided into two chunks of equal size. We show results for a desired miscoverage rate of $\alpha \in \{0.1, 0.05, 0.01\}$.

All methods use neural networks which are trained to convergence using the Adam optimizer (Kingma, 2014). Throughout we use fully-connected networks with 50 units/layer and GeLU activation function. The small and medium datasets have a single hidden layer whereas the large datasets use 3 hidden layers. For these architectures, it is feasible to evaluate the Gauss-Newton matrix without any approximations. However, we expect this to be prohibitive for larger architectures – inversion of the Gauss-Newton matrix scales cubically in the number of parameters. For this reason in App. H.1, we repeat the experiments using two scalable approximations: Kronecker-factored approximate curvature (KFAC) (Martens & Grosse, 2015) and last-layer approximation (Daxberger et al., 2021b) (*i.e.* neural linear model approach (Ober & Rasmussen, 2019)). As described in Sec. 2.3, to construct the predictive intervals given the coefficients in ACP-GN we use the ridge regression confidence machine algorithm (Noureddin et al., 2001) on the small datasets and the asymmetric version (Burnaev & Vovk, 2014) for the medium and large datasets. See App. G.1 for further details on the experimental setup.

To assess the efficiency (tightness) and well-calibratedness of our proposed methods for obtaining prediction intervals, we report their average prediction interval width and coverage against the baselines in Table 1. A method is reported as satisfying validity if its empirical coverage lies within the 1% and 99% quantiles of the exact marginal coverage distribution as given by the train/calibration set size (depending on the method) (Angelopoulos & Bates, 2021; Vovk, 2012). In the case of limited data regimes (yacht, boston, energy), ACP-GN gives the tightest intervals, with the exception of boston at 95% target coverage where LA is the most efficient (being one of the few cases when LA does not miscover). On the larger datasets, ACP-GN remains competitive on efficiency with the exception of facebook_2 at the higher values of target coverage, but we find it miscovers. Our proposed variant, “ACP-GN (split+refine)”, generally gives the correct coverage although incurring a trade-off in efficiency due to the sample splitting. We also observe that our novel normalization strategy inspired by ACP-GN, SCP-GN, improves over CRF for most datasets and settings of the target coverage.

5.2 BOUNDING BOX LOCALIZATION

We consider single-object localization and in particular adapt the task from Phan et al. (2018), which predicts bounding boxes localizing the face of different breeds of cats and dogs in varying conditions not limited to scale, pose and lighting. Conformal methods have recently been adapted for this task (De Grancey et al., 2022; Timans et al., 2024). We construct two-sided intervals similar to Timans et al. (2024) but without considering uncertainty in the classifier, following De Grancey et al. (2022).

All images with ground-truth bounding box annotations in the Oxford-IIIT Pet dataset (Parkhi et al., 2012) are extracted resulting in 3 686 images overall. The experiment is repeated with 20 different train-test splits where 20% of the data is used for testing. When calibration data is needed, a 25% split is partitioned from the train set. The VGG-19 architecture (Simonyan & Zisserman, 2015) pretrained on ImageNet is used as the object detection backbone but with the original output layer removed. The network is trained jointly with two heads, a regression head predicting 2D bounding box coordinates (4 outputs) and a binary classification head predicting between cat or dog. We only

Table 2: On a bounding box regression task using a deep convolutional neural network, our split and refine variant of ACP-GN results in the most efficient confidence regions whilst achieving comparable coverage to the split-CP baselines. We target coverage rates of {85%, 90%, 95%} and reported metrics are accompanied by standard error from repeated runs.

	Avg. Volume ($\times 10^{-2}$)			Avg. Coverage		
	85%	90%	95%	85%	90%	95%
LA	0.710 \pm 0.009	0.945 \pm 0.013	1.405 \pm 0.019	92.11 \pm 0.18 (✗)	94.49 \pm 0.18 (✗)	96.65 \pm 0.16 (✗)
SCP	0.809 \pm 0.025	1.166 \pm 0.035	2.261 \pm 0.106	87.80 \pm 0.42 (✓)	91.82 \pm 0.41 (✓)	95.85 \pm 0.32 (✓)
CRF	0.713 \pm 0.021	1.098 \pm 0.044	2.130 \pm 0.087	87.47 \pm 0.44 (✓)	91.61 \pm 0.34 (✓)	95.85 \pm 0.23 (✓)
CQR	0.947 \pm 0.028	1.451 \pm 0.046	3.680 \pm 0.113	87.45 \pm 0.51 (✓)	91.29 \pm 0.38 (✓)	96.06 \pm 0.26 (✓)
ACP-GN	0.384 \pm 0.004	0.605 \pm 0.007	1.225 \pm 0.017	90.97 \pm 0.20 (✗)	94.09 \pm 0.20 (✗)	97.41 \pm 0.16 (✗)
SCP-GN	0.768 \pm 0.020	1.071 \pm 0.030	1.854 \pm 0.066	87.17 \pm 0.45 (✓)	91.43 \pm 0.40 (✓)	95.58 \pm 0.28 (✓)
ACP-GN (split + refine)	0.311\pm0.006	0.472\pm0.014	1.031\pm0.046	87.34 \pm 0.41 (✓)	91.31 \pm 0.35 (✓)	96.12 \pm 0.27 (✓)

consider the regression head for constructing predictive intervals. Without calibration split, the network achieves 99.6% classification accuracy and 20.1% localization error with 0.5 IoU (Intersection over Union) threshold. With calibration split, the model achieves 99.5% classification accuracy and 27.5% localization error.

The training setup is adapted from Girshick (2015) using SGD optimizer and simple data augmentation involving random horizontal flips of probability 0.5. The robust L1 loss is used for the bounding box regression head and logistic loss for the classification head. With L1 loss, Eqs. (12) to (14) no longer hold but we demonstrate the efficacy of our procedure when the objectives for training and predictive interval construction are different. We use the last-layer approximation to the Gauss-Newton matrix throughout for computational reasons. The asymmetric implementation of CRR is used due to its efficiency and we target miscoverage rates $\alpha \in \{0.15, 0.1, 0.05\}$. See App. G.2 for further details on the experimental setup. Since bounding box localization is a type of multi-output regression, we obtain confidence *regions* given by hyperrectangles except for Laplace approximation (LA) that results in a hyperellipsoid. All conformal prediction methods are run for each output dimension independently and we evaluate the confidence region volumes (as shown in Table 2) by simply taking the product over interval widths per output dimension. We apply a multiple testing correction, the Bonferroni correction, to mitigate the miscoverage that arises when conformalizing the outputs separately. Although we find that all methods still consistently overcover suggesting further improvements are still possible.

In Table 2, we observe that ACP-GN gives the tightest intervals as compared with the baseline conformal methods despite having a greater empirical coverage than those methods. Surprisingly our split and refine variant of ACP-GN gives even tighter intervals whilst matching the coverage of the conformal baselines. We observed a similar effect in our ablation of the previous UCI experiments with the last-layer approximation in App. H.1.

6 CONCLUSION

In this work, we show how to efficiently construct predictive prediction intervals with full conformal prediction (CP) for neural networks in regression tasks. While full-CP requires retraining the model from scratch on all of the training data and for each postulated label for the test point, we show that we can efficiently construct approximate full-CP predictive intervals with the Gauss-Newton influence (St Laurent & Cook, 1992) and the methods described at Nourtdinov et al. (2001); Burnaev & Vovk (2014) without retraining the model. In doing so, we also avoid creating a grid over possible regression targets for the test point, which incurs an undesirable accuracy/precision trade-off due to the uncountable set of real numbers. We demonstrate how such an approach corresponds to exact full-CP on a linearized version of the neural network and further show how it corresponds to “conformalizing” the linearized Laplace method (MacKay, 1992; Khan et al., 2019; Immer et al., 2021b), a popular method for getting uncertainty estimates in a Bayesian deep learning setting. Finally, we consider several alternative nonconformity scores, which lead to different variants of our method, and, through the conformal linearized Laplace lens, we also introduce a novel adaptive split-CP method. Experimentally, we see that our approximate full-CP methods typically provide tighter prediction intervals in limited data regimes across the well-calibrated approaches. For future work, we would like to extend our method to other real-world tasks such as pose estimation (Yang & Pavone, 2023) and tracking (Lindemann et al., 2023).

REFERENCES

- Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- Javier Antorán, David Janz, James U Allingham, Erik Daxberger, Riccardo Rb Barbano, Eric Nalisnick, and José Miguel Hernández-Lobato. Adapting the linearised laplace model evidence for modern deep learning. In *International Conference on Machine Learning*, pp. 796–821. PMLR, 2022.
- Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence functions are the answer, then what is the question? *Advances in Neural Information Processing Systems*, 35:17953–17967, 2022.
- Evgeny Burnaev and Vladimir Vovk. Efficiency of conformalized ridge regression. In *Conference on Learning Theory*, pp. 605–622. PMLR, 2014.
- Wenyu Chen, Kelli-Jean Chun, and Rina Foygel Barber. Discretized conformal prediction for efficient distribution-free inference. *Stat*, 7(1):e173, 2018.
- R Dennis Cook. Detection of influential observation in linear regression. *Technometrics*, 19(1): 15–18, 1977.
- R Dennis Cook and Sanford Weisberg. Characterizations of an empirical influence function for detecting influential cases in regression. *Technometrics*, 22(4):495–508, 1980.
- A Philip Dawid. The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610, 1982.
- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in Neural Information Processing Systems*, 34:20089–20103, 2021a.
- Erik Daxberger, Eric Nalisnick, James U Allingham, Javier Antorán, and José Miguel Hernández-Lobato. Bayesian deep learning via subnetwork inference. In *International Conference on Machine Learning*, pp. 2510–2521. PMLR, 2021b.
- Florence De Grancey, Jean-Luc Adam, Lucian Alecu, Sébastien Gerchinovitz, Franck Mamalet, and David Vigouroux. Object detection with probabilistic guarantees: A conformal prediction approach. In *International Conference on Computer Safety, Reliability, and Security*, pp. 316–329. Springer, 2022.
- Edwin Fong and Chris C Holmes. Conformal bayesian computation. *Advances in Neural Information Processing Systems*, 34:18268–18279, 2021.
- Matteo Fontana, Gianluca Zeni, and Simone Vantini. Conformal prediction: a unified review of theory and new challenges. *Bernoulli*, 29(1):1–23, 2023.
- Andrew YK Foong, Yingzhen Li, José Miguel Hernández-Lobato, and Richard E Turner. ‘in-between’ uncertainty in bayesian neural networks. *arXiv preprint arXiv:1906.11537*, 2019.
- DAS Fraser. Is bayes posterior just quick and dirty confidence? *Statistical Science*, 26(3):299–316, 2011.
- R Girshick. Fast r-cnn. *arXiv preprint arXiv:1504.08083*, 2015.
- Peter Grünwald and Thijs van Ommen. Inconsistency of bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4):1069–1103, 2017.
- Etash Kumar Guha, Shlok Natarajan, Thomas Möllenhoff, Mohammad Emtiyaz Khan, and Eugene Ndiaye. Conformal prediction via regression-as-classification. In *The Twelfth International Conference on Learning Representations*, 2024.
- Frank R Hampel. The influence curve and its role in robust estimation. *Journal of the american statistical association*, 69(346):383–393, 1974.

- Alexander Immer, Matthias Bauer, Vincent Fortuin, Gunnar Rätsch, and Mohammad Emtiyaz Khan. Scalable marginal likelihood estimation for model selection in deep learning. In *International Conference on Machine Learning*, pp. 4563–4573. PMLR, 2021a.
- Alexander Immer, Maciej Korzepa, and Matthias Bauer. Improving predictions of bayesian neural nets via local linearization. In *International conference on artificial intelligence and statistics*, pp. 703–711. PMLR, 2021b.
- Ulf Johansson, Henrik Boström, and Tuwe Löfström. Investigating normalized conformal regressors. In *2021 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 01–08. IEEE, 2021.
- Yuko Kato, David MJ Tax, and Marco Loog. A review of nonconformity measures for conformal prediction in regression. *Conformal and Probabilistic Prediction with Applications*, pp. 369–383, 2023.
- Mohammad Emtiyaz Khan and Håvard Rue. The bayesian learning rule. *Journal of Machine Learning Research*, 24(281):1–46, 2023.
- Mohammad Emtiyaz Khan, Alexander Immer, Ehsan Abedi, and Maciej Korzepa. Approximate inference turns deep networks into gaussian processes. *Advances in neural information processing systems*, 32, 2019.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jing Lei. Fast exact conformalization of the lasso using piecewise linear homotopy. *Biometrika*, 106(4):749–764, 2019.
- Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- Lars Lindemann, Matthew Cleaveland, Gihyun Shim, and George J Pappas. Safe planning in dynamic environments using conformal prediction. *IEEE Robotics and Automation Letters*, 2023.
- David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- James Martens. Deep learning via hessian-free optimization. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 735–742, 2010.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Javier Abad Martinez, Umang Bhatt, Adrian Weller, and Giovanni Cherubin. Approximating full conformal prediction at scale via influence functions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 6631–6639, 2023.
- Thomas Melliush, Craig Saunders, Ilia Nouretdinov, and Volodya Vovk. Comparing the bayes and typicalness frameworks. In *Machine Learning: ECML 2001: 12th European Conference on Machine Learning Freiburg, Germany, September 5–7, 2001 Proceedings 12*, pp. 360–371. Springer, 2001.
- Eugene Ndiaye. Stable conformal prediction sets. In *International Conference on Machine Learning*, pp. 16462–16479. PMLR, 2022.
- Eugene Ndiaye and Ichiro Takeuchi. Computing full conformal prediction set with approximate homotopy. *Advances in Neural Information Processing Systems*, 32, 2019.
- Peter Nickl, Lu Xu, Dharmesh Tailor, Thomas Möllenhoff, and Mohammad Emtiyaz Khan. The memory-perturbation equation: Understanding model’s sensitivity to data. *Advances in Neural Information Processing Systems*, 36, 2023.

- Kolby Nottingham, Rachel Longjohn, and Markelle Kelly. UCI Machine Learning Repository, 2024. URL <https://archive.ics.uci.edu/datasets>. Accessed: September, 2024.
- Ilia Nouretdinov, Thomas Melliush, and Volodya Vovk. Ridge regression confidence machine. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pp. 385–392, 2001.
- Roman Novak, Jascha Sohl-Dickstein, and Samuel S Schoenholz. Fast finite width neural tangent kernel. In *International Conference on Machine Learning*, pp. 17018–17044. PMLR, 2022.
- Sebastian W Ober and Carl E Rasmussen. Benchmarking the neural linear model for regression. In *Second Symposium on Advances in Approximate Bayesian Inference*, 2019.
- Harris Papadopoulos. Guaranteed coverage prediction intervals with gaussian process regression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Harris Papadopoulos and Haris Haralambous. Reliable prediction intervals with regression neural networks. *Neural Networks*, 24(8):842–851, 2011.
- Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *Machine learning: ECML 2002: 13th European conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13*, pp. 345–356. Springer, 2002.
- Harris Papadopoulos, Alex Gammerman, and Volodya Vovk. Normalized nonconformity measures for regression conformal prediction. In *Proceedings of the IASTED International Conference on Artificial Intelligence and Applications (AIA 2008)*, pp. 64–69, 2008.
- Harris Papadopoulos, Vladimir Vovk, and Alex Gammerman. Regression conformal prediction with nearest neighbours. *Journal of Artificial Intelligence Research*, 40:815–840, 2011.
- Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- Buu Phan, Rick Salay, Krzysztof Czarnecki, Vahdat Abdelzad, Taylor Denouden, and Sachin Vernekar. Calibrating uncertainties in object localization task. *arXiv preprint arXiv:1811.11210*, 2018.
- Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. In *6th international conference on learning representations, ICLR 2018-conference track proceedings*, volume 6. International Conference on Representation Learning, 2018.
- Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. *Advances in neural information processing systems*, 32, 2019.
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Roy T St Laurent and R Dennis Cook. Leverage and superleverage in nonlinear regression. *Journal of the American Statistical Association*, 87(420):985–990, 1992.
- Alexander Timans, Christoph-Nikolas Straehle, Kaspar Sakmann, and Eric Nalisnick. Adaptive bounding box uncertainties via two-step conformal prediction. In *Proceedings of the European Conference on Computer Vision*, 2024.
- Vladimir Vovk. Conditional validity of inductive conformal predictors. In *Asian conference on machine learning*, pp. 475–490. PMLR, 2012.
- Vladimir Vovk. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74: 9–28, 2015.

Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*, volume 29. Springer, 2005.

Larry Wasserman. Frasian inference. *Statistical Science*, 26(3):322–325, 2011.

Heng Yang and Marco Pavone. Object pose estimation with statistical guarantees: Conformal key-point detection and geometric uncertainty propagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8947–8958, 2023.

A DERIVATION OF CRR COEFFICIENTS

The ridge solution on \mathcal{D}_N is given by $\theta_* = \mathbf{H}_*^{-1} \mathbf{X}^\top \mathbf{y}$ where \mathbf{X} is the $(N \times I)$ feature matrix with \mathbf{x}_i^\top as rows and \mathbf{y} is the N -dim vector of targets. We can express the ridge solution on $\mathcal{D}_{N+1}(y)$, referred to as the add-one-in (AOI) solution, as a deviation from θ_* :

$$\theta_*^+(y) = (\mathbf{H}_* + \mathbf{x}_{N+1} \mathbf{x}_{N+1}^\top)^{-1} (\mathbf{X}^\top \mathbf{y} + \mathbf{x}_{N+1} y) \quad (22)$$

▷ use Sherman-Morrison formula

$$= \left(\mathbf{H}_*^{-1} - \frac{\mathbf{H}_*^{-1} \mathbf{x}_{N+1} \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1}}{1 + \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}} \right) (\mathbf{X}^\top \mathbf{y} + \mathbf{x}_{N+1} y) \quad (23)$$

$$= \mathbf{H}_*^{-1} \mathbf{X}^\top \mathbf{y} + \mathbf{H}_*^{-1} \mathbf{x}_{N+1} \left(y - \frac{\mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{X}^\top \mathbf{y}}{1 + \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}} - \frac{y \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}}{1 + \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}} \right) \quad (24)$$

▷ substitute $\theta_* = \mathbf{H}_*^{-1} \mathbf{X}^\top \mathbf{y}$ and $h_{N+1} = \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}$

$$= \theta_* + \mathbf{H}_*^{-1} \mathbf{x}_{N+1} \left(y - \frac{\mathbf{x}_{N+1}^\top \theta_*}{1 + h_{N+1}} - \frac{y h_{N+1}}{1 + h_{N+1}} \right) \quad (25)$$

$$= \theta_* + \mathbf{H}_*^{-1} \mathbf{x}_{N+1} \left(\frac{y - \mathbf{x}_{N+1}^\top \theta_*}{1 + h_{N+1}} \right) \quad (26)$$

Using this, it is easy to show the residuals can be expressed in terms of the postulated label y :

$$y_i - \mathbf{x}_i^\top \theta_*^+(y) = y_i - \mathbf{x}_i^\top \theta_* - \mathbf{x}_i^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1} \left(\frac{y - \mathbf{x}_{N+1}^\top \theta_*}{1 + h_{N+1}} \right) \quad (27)$$

▷ substitute $h_{i,N+1} = \mathbf{x}_i^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}$

$$= y_i - \mathbf{x}_i^\top \theta_* + \frac{h_{i,N+1}}{1 + h_{N+1}} \mathbf{x}_{N+1}^\top \theta_* - \frac{h_{i,N+1}}{1 + h_{N+1}} y \quad (28)$$

$$y - \mathbf{x}_{N+1}^\top \theta_*^+(y) = y - \mathbf{x}_{N+1}^\top \theta_* - \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1} \left(\frac{y - \mathbf{x}_{N+1}^\top \theta_*}{1 + h_{N+1}} \right) \quad (29)$$

▷ substitute $h_{N+1} = \mathbf{x}_{N+1}^\top \mathbf{H}_*^{-1} \mathbf{x}_{N+1}$

$$= -\frac{1}{1 + h_{N+1}} \mathbf{x}_{N+1}^\top \theta_* + \frac{1}{1 + h_{N+1}} y \quad (30)$$

B DERIVATION OF ACP-GN COEFFICIENTS

We present the derivation for expressing the residual in the neural network regression setting as a linear function of the postulated label y ,

$$y_i - f_i(\theta_*^+(y)) \quad (31)$$

▷ linearize neural network about θ_*

$$\approx y_i - [f_i(\theta_*) + \nabla_{\theta} f_i(\theta_*) (\theta_*^+(y) - \theta_*)] \quad (32)$$

▷ approximate $\theta_*^+(y) \approx \hat{\theta}_*^+(y)$ and substitute Gauss-Newton influence

$$= y_i - f_i(\theta_*) - \phi_i^\top \left(\frac{\hat{e}_{N+1}(y)}{1 + \hat{h}_{N+1}} \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1} \right) \quad (33)$$

▷ definition of residual $\hat{e}_{N+1}(y)$ to reveal postulated label y and substitute $\hat{h}_{i,N+1} = \phi_i^\top \mathbf{H}_{\text{GN}}^{-1} \phi_{N+1}$

$$= y_i - f_i(\theta_*) + \underbrace{\frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} f_{N+1}(\theta_*)}_{a_i} - \underbrace{\frac{\hat{h}_{i,N+1}}{1 + \hat{h}_{N+1}} y}_{b_i} \quad (34)$$

The residual of the postulated point proceeds in an almost identical manner.

C EXTENSION OF ACP-GN FOR MULTI-OUTPUT REGRESSION

In the case of multi-output regression with vector-valued targets $\mathbf{y}_i \in \mathbb{R}^O$ and outputs of a DNN $\mathbf{f}_i(\boldsymbol{\theta}) \in \mathbb{R}^O$, the coefficients needed in the CRR procedure $\{\mathbf{a}_i, \mathbf{B}_i\}$ now correspond to a O -dim vector and a $O \times O$ matrix respectively. Since CRR was only proposed for single-output regression, we also need to adapt the procedure. Considering the asymmetric version of CRR (Burnaev & Vovk, 2014), we simply need to solve for the O -dim changepoints which are given by $(\mathbf{B}_{N+1} - \mathbf{B}_i)^{-1}(\mathbf{a}_i - \mathbf{a}_{N+1})$ as long as $\mathbf{B}_{N+1} - \mathbf{B}_i$ is positive-definite (generalizing the positivity constraint in the single-output case). Then we propose to sort the set of changepoints followed by taking the quantile component-wise analogous to the existing algorithm. Now we proceed to derive $\{\mathbf{a}_i, \mathbf{B}_i\}$. Firstly the multi-output analogue of the AOI estimator in Eq. (12) is given by,

$$\hat{\boldsymbol{\theta}}_*^+(\mathbf{y}) = \boldsymbol{\theta}_* + \mathbf{H}_{\text{GN}}^{-1} \boldsymbol{\Phi}_{N+1}^\top (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \hat{\mathbf{e}}_{N+1}(\mathbf{y}) \quad (35)$$

where $\boldsymbol{\Phi}_i := \nabla_{\boldsymbol{\theta}} \mathbf{f}_i(\boldsymbol{\theta}_*) \in \mathbb{R}^{O \times D}$, $\mathbf{H}_{\text{GN}} = \sum_{i=1}^N \boldsymbol{\Phi}_i^\top \boldsymbol{\Phi}_i + \delta \mathbf{I}$, $\hat{\mathbf{e}}_{N+1}(\mathbf{y}) = \mathbf{y} - \mathbf{f}_{N+1}(\boldsymbol{\theta}_*)$ and $\hat{\mathbf{V}}_{N+1} = \boldsymbol{\Phi}_{N+1} \mathbf{H}_{\text{GN}}^{-1} \boldsymbol{\Phi}_{N+1}^\top$. There is a change in notation in the last expression of the multi-output leverage score to avoid confusion with the Hessian. Using this we can derive analogous expressions to Eqs. (13) and (14):

$$\mathbf{y}_i - \mathbf{f}_i(\boldsymbol{\theta}_*^+(\mathbf{y})) \approx \mathbf{y}_i - \left[\mathbf{f}_i(\boldsymbol{\theta}_*) + \nabla_{\boldsymbol{\theta}} \mathbf{f}_i(\boldsymbol{\theta}_*) (\hat{\boldsymbol{\theta}}_*^+(\mathbf{y}) - \boldsymbol{\theta}_*) \right] \quad (36)$$

$$= \mathbf{y}_i - \mathbf{f}_i(\boldsymbol{\theta}_*) - \boldsymbol{\Phi}_i \mathbf{H}_{\text{GN}}^{-1} \boldsymbol{\Phi}_{N+1}^\top (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \hat{\mathbf{e}}_{N+1}(\mathbf{y}) \quad (37)$$

$$= \underbrace{\mathbf{y}_i - \mathbf{f}_i(\boldsymbol{\theta}_*) + \hat{\mathbf{V}}_{i,N+1} (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \mathbf{f}_{N+1}(\boldsymbol{\theta}_*)}_{\mathbf{a}_i} - \underbrace{\hat{\mathbf{V}}_{i,N+1} (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \mathbf{y}}_{\mathbf{B}_i} \quad (38)$$

and

$$\mathbf{y} - \mathbf{f}_{N+1}(\boldsymbol{\theta}_*^+(\mathbf{y})) \approx \underbrace{-(\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \mathbf{f}_{N+1}(\boldsymbol{\theta}_*)}_{\mathbf{a}_{N+1}} + \underbrace{(\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \mathbf{y}}_{\mathbf{B}_{N+1}} \quad (39)$$

where $\hat{\mathbf{V}}_{i,N+1} = \boldsymbol{\Phi}_i \mathbf{H}_{\text{GN}}^{-1} \boldsymbol{\Phi}_{N+1}^\top$. The normalized nonconformity scores can also be extended to the multi-output setting. In the case of deleted-CRR, we have:

$$\mathbf{a}_i \leftarrow (\mathbf{I} - \bar{\mathbf{V}}_i)^{-1} \mathbf{a}_i, \quad \mathbf{B}_i \leftarrow (\mathbf{I} - \bar{\mathbf{V}}_i)^{-1} \mathbf{B}_i \quad \forall i = 1, \dots, N+1 \quad (40)$$

where,

$$\bar{\mathbf{V}}_i = \hat{\mathbf{V}}_i - \hat{\mathbf{V}}_{i,N+1} (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \hat{\mathbf{V}}_{i,N+1}^\top \quad \forall i = 1, \dots, N \quad (41)$$

$$\bar{\mathbf{V}}_{N+1} = \hat{\mathbf{V}}_{N+1} (\mathbf{I} + \hat{\mathbf{V}}_{N+1})^{-1} \quad (42)$$

and we introduced $\hat{\mathbf{V}}_i = \boldsymbol{\Phi}_i \mathbf{H}_{\text{GN}}^{-1} \boldsymbol{\Phi}_i^\top$.

D DERIVATION OF APPROXIMATE ADD-ONE-IN POSTERIOR WITH LAPLACE-GGN

Khan et al. (2019) show that the Laplace-GGN posterior can be equivalently stated as exact inference in the following linear regression model (see their Theorem 1):

$$q_*(\boldsymbol{\theta}) \propto \prod_{i=1}^N e^{-\frac{1}{2}(\tilde{y}_i - \boldsymbol{\phi}_i^\top \boldsymbol{\theta})^2} p(\boldsymbol{\theta}) \quad (43)$$

where $p(\boldsymbol{\theta}) \propto \exp(\frac{1}{2} \delta \|\boldsymbol{\theta}\|^2)$, $\boldsymbol{\phi}_i := \nabla_{\boldsymbol{\theta}} \mathbf{f}_i(\boldsymbol{\theta}_*)^\top$ and $\tilde{y}_i := \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_* + e_i$ with $e_i = y_i - f_i(\boldsymbol{\theta}_*)$. This can be viewed as approximating the original non-conjugate terms by conjugate factors (see Sec. 5.4 in (Khan & Rue, 2023)): $e^{-\ell(y_i, f_i(\boldsymbol{\theta}))} \approx e^{-\frac{1}{2}(\tilde{y}_i - \boldsymbol{\phi}_i^\top \boldsymbol{\theta})^2}$, which take a similar interpretation

to site functions in expectation propagation. This is used along with the standard formula for online Bayesian updating to derive the approximate AOI posterior $\hat{q}_*^+(\boldsymbol{\theta}) \approx p(\boldsymbol{\theta}|\mathcal{D}_{N+1})$,

$$p(\boldsymbol{\theta}|\mathcal{D}_{N+1}) \propto \prod_{i=1}^{N+1} e^{-\ell(y_i, f_i(\boldsymbol{\theta}))} p(\boldsymbol{\theta}) \quad (44)$$

▷ split off the $(N+1)^{\text{th}}$ likelihood from the others

$$= e^{-\ell(y, f_{N+1}(\boldsymbol{\theta}))} \prod_{i=1}^N e^{-\ell(y_i, f_i(\boldsymbol{\theta}))} p(\boldsymbol{\theta}) \quad (45)$$

▷ apply the Laplace-GGN posterior approximation

$$\approx e^{-\ell(y, f_{N+1}(\boldsymbol{\theta}))} q_*(\boldsymbol{\theta}) \quad (46)$$

▷ approximate the $(N+1)^{\text{th}}$ likelihood by its site function

$$\approx e^{-\frac{1}{2}(\tilde{y} - \phi_{N+1}^\top \boldsymbol{\theta})^2} q_*(\boldsymbol{\theta}) \quad (47)$$

▷ this corresponds to an unnormalized Gaussian distribution

$$\propto \mathbb{N}(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_*^+(y), \hat{\boldsymbol{\Sigma}}_*^+) \quad (48)$$

E MULTI-OUTPUT REGRESSION PREDICTION INTERVALS

For a Bayesian posterior predictive distribution (under Gaussian assumptions) in the multi-output case, $p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \mathbb{N}(\mathbf{y}_*|\hat{\mathbf{y}}_*, \boldsymbol{\Sigma}_{y_*})$, the confidence *region* is given by an ellipsoid,

$$C_\alpha(\mathbf{x}_*) = \{\mathbf{y} \in \mathbb{R}^O : (\mathbf{y} - \hat{\mathbf{y}}_*)^\top \boldsymbol{\Sigma}_{y_*}^{-1} (\mathbf{y} - \hat{\mathbf{y}}_*) \leq \chi_{O,\alpha}^2\} \quad (49)$$

where $\chi_{O,\alpha}^2$ is the quantile function for the chi-squared distribution with O degrees of freedom. Using this definition, it is straightforward to evaluate empirical coverage. Efficiency is then given by the volume of the corresponding ellipsoid,

$$\text{Vol}[C_\alpha(\mathbf{x}_*)] = (\chi_{O,\alpha}^2)^{\frac{O}{2}} \det(\boldsymbol{\Sigma}_{y_*})^{\frac{1}{2}} \text{Vol}[B_O] \quad (50)$$

where B_O is the unit ball with O dimensions.

F TIME COMPLEXITY: FCP VS. ACP-GN

Table 3: We have number of test points (M), number of grid points (K), number of train points (N), parameter count (D), total epochs (E), and cost of forward pass/gradient computation/Jacobian computation ($[FP]$). We highlight the additional complexity of FCP over ACP-GN (purple), ACP over ACP-GN (green), ACP-GN over ACP (blue) and ACP/ACP-GN over FCP (red).

	Train	Predict
FCP	—	$MKN^2E[FP]$
ACP	$NE[FP] + ND^2 + D^3$	$MK([FP] + ND^2)$
ACP-GN	$NE[FP] + ND^2 + D^3$	$M([FP] + ND^2 + N \log N)$

We write the time complexity for our ACP-GN and compare it against full conformal prediction (FCP) and approximate full conformal prediction (ACP) (Martinez et al., 2023) in Table 3. This is shown for the deleted nonconformity score similar to Sec. A.2 in Martinez et al. (2023) and for scalar targets only. For the standard score, ACP and ACP-GN are unchanged but the factor of N is dropped in FCP. “Train” refers to the upfront time complexity that can be re-used when constructing prediction intervals (“predict”) for new batches of test points.

We state the best time complexity of the ridge regression confidence machine routine given the coefficients $\{(a_i, b_i)\}_{i=1}^{N+1}$ as $\mathcal{O}(N \log N)$ (see Sec. 2.3 in (Vovk et al., 2005)). This is the same as

the asymmetric implementation of Burnaev & Vovk (2014) shown in Alg. 2 (sorting the N change-points). The cost of network prediction (single forward pass) is architecture-dependent but since this is constant across the methods we do not expand on this here. We also regard the cost of gradient and Jacobian computation as comparable to a forward pass (see the discussion in App. D in Novak et al. (2022)).

FCP has a complexity multiplicative in the number of test points, grid points, train points (twice) and epochs. Whereas for ACP-GN, it is just multiplicative in the number of test points, given an upfront cost which is cubic in the number of parameters. We can further reduce the time complexity for ACP-GN when scalable approximations to the Gauss-Newton Hessian are used (as investigated in App. H.1). These extensions could also be adapted for ACP which was in fact left as future work in Martinez et al. (2023). This is shown in Table 4 for the cases of Kronecker-factored approximate curvature (KFAC) and last-layer approximation (LL). Most crucially, both these approximations relax the cubic dependence on the parameter count P which typically exceeds millions of parameters for modern architectures. Instead the cubic dependence shifts to the input and output dimensionality of the network layers in the case of KFAC or just the last layer which are typically far smaller than P .

Table 4: We have number of network layers (L), layer input/output dimensionality ($I_{l,\text{in}}/I_{l,\text{out}}$), cost of Gauss-Newton Hessian evaluation/inversion (H_N) along with its KFAC (H_N^{KFAC}) and last-layer (H_N^{LL}) approximation, and cost of operations related to Gauss-Newton influence in the KFAC case $[INF]^{\text{KFAC}}$. Other terms are defined in Table 3.

	Train	Predict
ACP-GN	$NE[FP] + H_N^{\text{a}}$	$M([FP] + ND^2 + N \log N)$
ACP-GN(kfac)	$NE[FP] + H_N^{\text{KFAC}^{\text{b}}}$	$M([FP] + [INF]^{\text{KFAC}^{\text{d}}} + N \log N)$
ACP-GN(LL)	$NE[FP] + H_N^{\text{LL}^{\text{c}}}$	$M([FP] + NI_{L,\text{in}}^2 + N \log N)$

^a $H_N = N[FP] + ND^2 + D^3$
^b $H_N^{\text{KFAC}} = N[FP] + N \sum_{l=1}^L (I_{l,\text{in}}^2 + I_{l,\text{out}}^2) + \sum_{l=1}^L (I_{l,\text{in}}^3 + I_{l,\text{out}}^3)$
^c $H_N^{\text{LL}} = N[FP] + NI_{L,\text{in}}^2 + I_{L,\text{in}}^3$
^d $[INF]^{\text{KFAC}} = N \left(D + \sum_{l=1}^L (I_{l,\text{out}} I_{l,\text{in}}^2 + I_{l,\text{in}} I_{l,\text{out}}^2) \right)$

G EXPERIMENTAL DETAILS

G.1 UCI REGRESSION

For every dataset, both the input features and targets are standardized to have zero mean and unit variance. We use a batch size of 256 in SGD training with an initial learning rate of 10^{-2} that is decayed to 10^{-5} using a cosine schedule. All methods require tuning the L2 regularizer/prior precision. Additionally, for Linearized Laplace we have the observation noise. For the small and medium datasets, these are tuned using online marginal likelihood optimization (Immer et al., 2021a) that alternates between standard neural network training and gradient-based updates to the hyperparameters using the differentiable marginal likelihood estimate. We use a layerwise structure in the prior precision. The marginal likelihood estimate is also used for early stopping. The prior precision and observation noise are initialized to 1. Overall, using Adam optimizer we run for 5000 epochs with a hyperparameter learning rate of 10^{-2} decayed to 10^{-3} using a cosine schedule, 100 burn-in epochs, and take 50 hyperparameter steps on single marginal likelihood evaluation every 50 epochs.

For the large datasets, the (scalar) prior precision is tuned via grid-search for each order of magnitude from 10^{-2} to 10^4 . 10% of the training set is used for validation and once the best prior precision is found, the network is retrained on full training set. The observation noise is fit to the training data via maximum likelihood after training.

We follow the above procedure for all methods except “ACP-GN (split + refine)” where we first train with L2 regularizer fixed to $\delta/N = 10^{-4}$. For the small datasets, the hyperparameters are tuned via post-hoc marginal likelihood training with 5000 steps and the same learning rate schedule described

Table 5: We repeat the bounding box localization experiment in Table 2 but with 50% calibration split.

	Avg. Volume ($\times 10^{-2}$)			Avg. Coverage		
	85%	90%	95%	85%	90%	95%
LA	0.710 \pm 0.009	0.945 \pm 0.013	1.405 \pm 0.019	92.11 \pm 0.18 (X)	94.49 \pm 0.18 (X)	96.65 \pm 0.16 (X)
SCP	1.177 \pm 0.023	1.723 \pm 0.037	3.137 \pm 0.081	87.36 \pm 0.27 (X)	91.87 \pm 0.27 (X)	95.91 \pm 0.19 (✓)
CRF	1.248 \pm 0.030	1.900 \pm 0.044	3.781 \pm 0.092	87.48 \pm 0.34 (X)	91.96 \pm 0.25 (X)	96.02 \pm 0.20 (✓)
CQR	1.627 \pm 0.027	2.739 \pm 0.052	7.971 \pm 0.159	87.46 \pm 0.31 (X)	91.42 \pm 0.28 (✓)	95.97 \pm 0.14 (✓)
ACP-GN	0.384 \pm 0.004	0.605 \pm 0.007	1.225 \pm 0.017	90.97 \pm 0.20 (X)	94.09 \pm 0.20 (X)	97.41 \pm 0.16 (X)
SCP-GN	1.149 \pm 0.025	1.615 \pm 0.034	2.643 \pm 0.054	86.84 \pm 0.29 (✓)	91.35 \pm 0.31 (✓)	95.56 \pm 0.21 (✓)
ACP-GN (split + refine)	0.365 \pm 0.006	0.556 \pm 0.010	1.104 \pm 0.028	87.78 \pm 0.37 (X)	91.75 \pm 0.31 (✓)	95.62 \pm 0.17 (✓)

earlier. In Fig. 2 we demonstrate the trade-off between coverage and efficiency in the split variant of ACP-GN.

We set the sensitivity hyperparameter (β) of CRF (see Eq. 16 in (Papadopoulos & Haralambous, 2011)) to 1 as used in Romano et al. (2019). The additional network is identical to the original one and re-uses the same hyperparameters and training configuration. It is worth mentioning that the original proposal of CRF in the context of neural networks (Papadopoulos et al., 2002) used a ridge regression model to predict residuals.

The quantile regression network in CQR uses the same architecture as that used in the other methods except there are two output units. As done in CRF, hyperparameters from SCP are re-used along with the training configuration. Rather than retraining for each desired significance level, the output layer is appended with additional units and trained jointly. We do not perform tuning of the quantiles so it is expected the intervals can be made more efficient (Romano et al., 2019).

Results on additional UCI datasets are reported in Table 8. Furthermore, in Table 9 and 10 we repeat all experiments with 25% calibration split.

G.2 BOUNDING BOX LOCALIZATION

Images are resized to 224×224 followed by scaling the pixel values to $[0, 1]$ and then normalizing to statistics computed from the ImageNet dataset as required by the VGG-19 backbone. Bounding box coordinates denote the top-left and bottom-right corners and are scaled to $[0, 1]$. The reported volumes use the standardized targets.

For SGD training we use a batch size of 128 for 200 epochs and an initial learning rate of 10^{-2} . A learning rate schedule is used that takes incremental steps towards the base learning rate for the first 5 epochs (warm-up) and then decays to 0 using a cosine schedule. The SGD optimizer uses nesterov momentum 0.9 and weight decay 5×10^{-4} . At inference with squared-error loss, we perform *post-hoc* finetuning of hyperparameters (regularization coefficient, observation noise) using the marginal likelihood. This is optimized using Adam for 250 epochs with a learning rate of 10^{-1} . For LA, we use the expressions outlined in App. E to evaluate the volume and coverage.

Due to the last-layer approximation, the off-diagonal entries of the $O \times O$ multi-output leverage scores are zero. We can understand this by realizing the Jacobian of an output with respect to a parameter tied to a different output is zero. Due to this inherent output independence, we can simply use the expressions given in the single-output case in parallel over the output dimensions rather than use the expressions derived for the multi-output setting in App. C.

When training the additional network in CRF and the quantile regression network in CQR, the backbone parameters are initialized to those of the original network and not updated. They are trained in a similar fashion except for 100 epochs, without warm-up in the learning rate schedule and without data augmentation. The sensitivity hyperparameter (β) of CRF is set to 0.01 after trying a few different values on a single seed.

In Table 5, we repeat the experiment but with 50% calibration split. In this case the model achieves 99.4% classification accuracy and 37.0% localization error.

H ADDITIONAL EXPERIMENTS

H.1 ABLATION WITH SCALABLE APPROXIMATIONS TO THE GAUSS-NEWTON MATRIX

For most deep architectures, storing and inverting the full Gauss-Newton matrix is infeasible. We investigate two choices for scalable approximations to the Gauss-Newton matrix in the context of our experiments involving the large UCI datasets. These are inspired by two popular choices for scalable Laplace approximations.

The first is the use of Kronecker-factored approximate curvature (KFAC) (Martens & Grosse, 2015) as an approximation to the GN. By approximating each layer’s Gauss Newton independently as a Kronecker product, this leads to a block-diagonal factorization for the overall GN matrix enabling efficient storage and computation of inverses. We use the specific form proposed in Immer et al. (2021b) that performs an eigendecomposition of the Kronecker factors and avoids the “dampening” approximation of Ritter et al. (2018). With regards to refinement, we observed that the one-step solve with KFAC approximation performs much worse than without refinement. Hence, we instead take a gradient-based approach as outlined in Antorán et al. (2022) optimizing the linearized network’s objective using a similar configuration to the original NN training. The linearized network’s loss gradients are evaluated without explicitly instantiating Jacobians via vector-Jacobian and Jacobian-vector products.

The results are shown in Table 11. For LA, we observe that KFAC leads to slightly tighter intervals but the miscoverage increases. In the case of SCP-GN, we find KFAC to be very competitive to the full Gauss-Newton with little change to the interval width or coverage. For ACP-GN we show results for the standard (*i.e.* AOI) nonconformity score as opposed to the studentized nonconformity score reported in Table 1 and 8. This is because we found that the deleted and studentized nonconformity scores combined with the KFAC approximation performed poorly. However for the reported standard nonconformity score, we observe KFAC leads to much tighter intervals with just a slight increase in miscoverage. For the “ACP-GN (split + refine)”, on all datasets except bike, we observe tighter intervals with KFAC however the miscoverage is often slightly greater.

The second choice is a last-layer approximation or neural linear model approach (Ober & Rasmussen, 2019) that can be considered a special case of subnetwork inference for Linearized Laplace (Daxberger et al., 2021b). This makes the AOI estimation exact with respect to a linear model whose basis features are given by the activations of the penultimate layer. This can be recovered as a special case of Gauss-Newton influence. This leads to storing and inverting a much smaller matrix whose size corresponds to the number of parameters in the final layer. In the context of Linearized Laplace, the last-layer approximation can be combined with KFAC for increased scalability but we do not investigate this configuration here.

The results are shown in Table 12. The last-layer approximation has the effect of making the existing undercoverage in LA much worse. In the case of SCP-GN, the approximation leads to the intervals being no more efficient than SCP, that is without normalization, whilst maintaining correct coverage. For ACP-GN, the intervals undercover by a large margin on the bike, community and protein datasets. The split and refine variant is still able to successfully correct this attaining the desired coverage whilst being competitive to the full Gauss-Newton on tightness. We also observed the different varieties of nonconformity score all performed similarly – results reported are those of the studentized variety as in Table 1 and 8.

H.2 COMPARISON AGAINST APPROXIMATE FULL CONFORMAL PREDICTION

We compare ACP-GN and ACP-GN (split+refine) to Approximate full Conformal Prediction (ACP) (Martinez et al., 2023) on several UCI datasets taken from Table 1 and 8. ACP uses the AOI nonconformity score (referred to as the ordinary scheme in (Martinez et al., 2023)). Whilst the exact Hessian is computed in (Martinez et al., 2023) with a damping term to ensure positive eigenvalues, we approximate the Hessian by a (damped) Gauss-Newton matrix in order to keep approximations consistent in the comparison. The damping term is tuned in the same way as ACP-GN, as described in App. G.1. The Gauss-Newton matrix has been used in previous works when evaluating influence function (Bae et al., 2022; Nickl et al., 2023). Furthermore, we use “direct approach” (see Eq. 4 in (Martinez et al., 2023)). Martinez et al. (2023) only considered classification tasks so ACP needs to

be adapted for regression. We define a grid of candidate targets using a simple discretization strategy that constructs a fine, uniform grid of 200 points delimited by the training targets.

Table 6 shows that in some datasets and target coverage levels, ACP-GN is more efficient but in others ACP is more efficient. However, we used 200 grid points which exceeds the upper end of what was investigated in prior work (Chen et al., 2018). We also emphasise that ACP-GN has lower space complexity since it only computes changepoints for each combination of test and train point, whereas ACP computes residuals for each combination of test point, train point and candidate target value. For the datasets considered we did not observe a considerable difference in the running time between the two methods but we do expect for larger datasets ACP to be slower due to the need for batched computation.

Table 6: We compare against Approximate full Conformal Prediction (ACP) that uses a target discretization strategy.

		Avg. Width				Avg. Coverage	
		90%	95%	99%	90%	95%	99%
energy N=768 I=8	ACP	1.706±0.045	2.219±0.055	3.694±0.070	87.41±0.28 (✓)	93.44±0.31 (✓)	98.58±0.16 (✓)
	ACP-GN	1.467 ±0.010	1.882 ±0.013	3.097 ±0.015	88.32±0.40 (✓)	93.87±0.29 (✓)	99.01±0.09 (✓)
	ACP-GN (split + refine)	1.745±0.016	2.174±0.021	3.300±0.045	90.54±0.25 (✓)	94.96±0.22 (✓)	99.18±0.10 (✓)
concrete N=1,030 I=8	ACP	17.854±0.104	21.966±0.119	31.068±0.139	90.36±0.15 (✓)	94.97±0.13 (✓)	98.43±0.07 (✓)
	ACP-GN	15.690 ±0.112	19.856 ±0.141	29.179 ±0.190	89.95±0.18 (✓)	94.96±0.13 (✓)	98.91±0.06 (✓)
	ACP-GN (split + refine)	18.907±0.103	24.012±0.129	35.949±0.387	90.26±0.20 (✓)	95.23±0.27 (✓)	99.17±0.06 (✓)
wine N=1,599 I=11	ACP	2.158±0.003	2.676±0.005	3.630 ±0.008	90.46±0.13 (✓)	94.88±0.08 (✓)	98.22±0.07 (✓)
	ACP-GN	2.091 ±0.005	2.651 ±0.006	3.797±0.013	90.91±0.11 (✓)	95.50±0.10 (✓)	99.14±0.05 (✓)
	ACP-GN (split + refine)	2.474±0.007	3.054±0.008	4.324±0.020	89.56±0.27 (✓)	94.81±0.13 (✓)	98.99±0.09 (✓)
kin8nm N=8,192 I=8	ACP	0.207±0.001	0.255 ±0.001	0.355 ±0.001	88.83±0.28 (✗)	94.16±0.19 (✓)	98.77±0.10 (✓)
	ACP-GN	0.213 ±0.001	0.262±0.001	0.365±0.002	89.68±0.28 (✓)	94.58±0.20 (✓)	98.85±0.08 (✓)
	ACP-GN (split + refine)	0.232±0.001	0.284±0.001	0.406±0.003	90.45±0.17 (✓)	95.02±0.22 (✓)	99.10±0.07 (✓)
power N=9,568 I=4	ACP	12.131 ±0.022	14.840 ±0.017	20.804 ±0.050	89.28±0.25 (✓)	94.68±0.18 (✓)	98.82±0.10 (✓)
	ACP-GN	12.526±0.024	15.248±0.020	21.592±0.062	90.16±0.26 (✓)	95.13±0.19 (✓)	98.89±0.08 (✓)
	ACP-GN (split + refine)	12.744±0.045	15.442±0.039	22.043±0.151	90.17±0.29 (✓)	95.26±0.17 (✓)	98.98±0.09 (✓)
community N=1,994 I=100	ACP	0.401 ±0.001	0.502 ±0.003	0.700±0.005	89.71±0.48 (✓)	94.00±0.36 (✓)	97.71±0.30 (✗)
	ACP-GN	0.459±0.003	0.594±0.004	0.936 ±0.007	90.68±0.55 (✓)	95.21±0.38 (✓)	99.18±0.14 (✓)
	ACP-GN (split + refine)	0.521±0.006	0.654±0.008	1.011±0.017	90.95±0.64 (✓)	95.26±0.43 (✓)	99.16±0.17 (✓)
bike N=10,886 I=18	ACP	93.647 ±1.300	121.850 ±1.603	186.090±2.393	89.03±0.31 (✓)	94.36±0.22 (✓)	98.43±0.14 (✗)
	ACP-GN	97.966±1.193	130.732±1.565	216.677±2.917	89.12±0.21 (✗)	94.32±0.18 (✗)	98.72±0.09 (✗)
	ACP-GN (split + refine)	130.933±5.221	174.190±7.065	288.603 ±12.066	90.09±0.25 (✓)	94.93±0.19 (✓)	99.06±0.07 (✓)

H.3 COMPARISON AGAINST FULL CONFORMAL PREDICTION

We compare full conformal prediction (FCP) against ACP-GN and ACP (Martinez et al., 2023) as well as LA and SCP. This is evaluated on a synthetic dataset with outliers (500 train points and 100 test points) taken from Papadopoulos (2024) (see Sec. 5.2). This generates data from a Gaussian Process prior with RBF kernel. There is a coin flip of probability 0.1 for which the observation noise standard deviation is increased by a factor of 10. The experiment is repeated 20 times with different seeds also controlling the data generation. A MLP with a single hidden layer, 100 units and Tanh activation function is used throughout. This is trained using Adam for 500 epochs (full-batch training) with an initial learning rate of 10^{-2} that decays following a cosine schedule to 10^{-5} . We use a uniform grid of 50 points for the target discretization in FCP which is on the lower end as suggested by previous works (Chen et al., 2018) and enough to give valid coverage. ACP follows the configuration outlined in App. H.2 with 200 grid points. The AOI nonconformity score is used throughout for FCP, ACP and ACP-GN. SCP uses a 50% calibration split. Hyperparameters are tuned using the online marginal likelihood procedure as described in App. G.1 which we exclude from the running time.

As Table 7 shows, all conformal methods give the correct coverage but FCP is indeed the most efficient. However, the running time is a factor of 10^4 slower than all other methods including ACP and ACP-GN. This experiment was run on a NVIDIA A100 GPU.

Table 7: We compare against Full Conformal Prediction (FCP) on a synthetically-generated dataset.

	90%	Avg. Width 95%	99%	90%	Avg. Coverage 95%	99%	Time ($\times 10^2$)
LA	1.061 ± 0.025	1.264 ± 0.030	1.661 ± 0.039	94.45 ± 0.56 (✗)	95.15 ± 0.52 (✓)	96.60 ± 0.45 (✗)	0.034 ± 0.001
SCP	0.623 ± 0.024	1.466 ± 0.090	3.632 ± 0.131	91.10 ± 0.82 (✓)	96.10 ± 0.45 (✓)	99.50 ± 0.15 (✓)	0.023 ± 0.001
FCP	0.451 ± 0.011	1.287 ± 0.053	3.090 ± 0.092	87.00 ± 1.06 (✓)	95.25 ± 0.52 (✓)	98.90 ± 0.32 (✓)	169.527 ± 3.560
ACP	0.532 ± 0.013	1.370 ± 0.054	3.177 ± 0.094	90.35 ± 0.79 (✓)	95.50 ± 0.49 (✓)	99.00 ± 0.32 (✓)	0.035 ± 0.001
ACP-GN	0.581 ± 0.021	1.450 ± 0.053	3.529 ± 0.105	91.00 ± 0.81 (✓)	95.60 ± 0.48 (✓)	99.25 ± 0.27 (✓)	0.034 ± 0.000

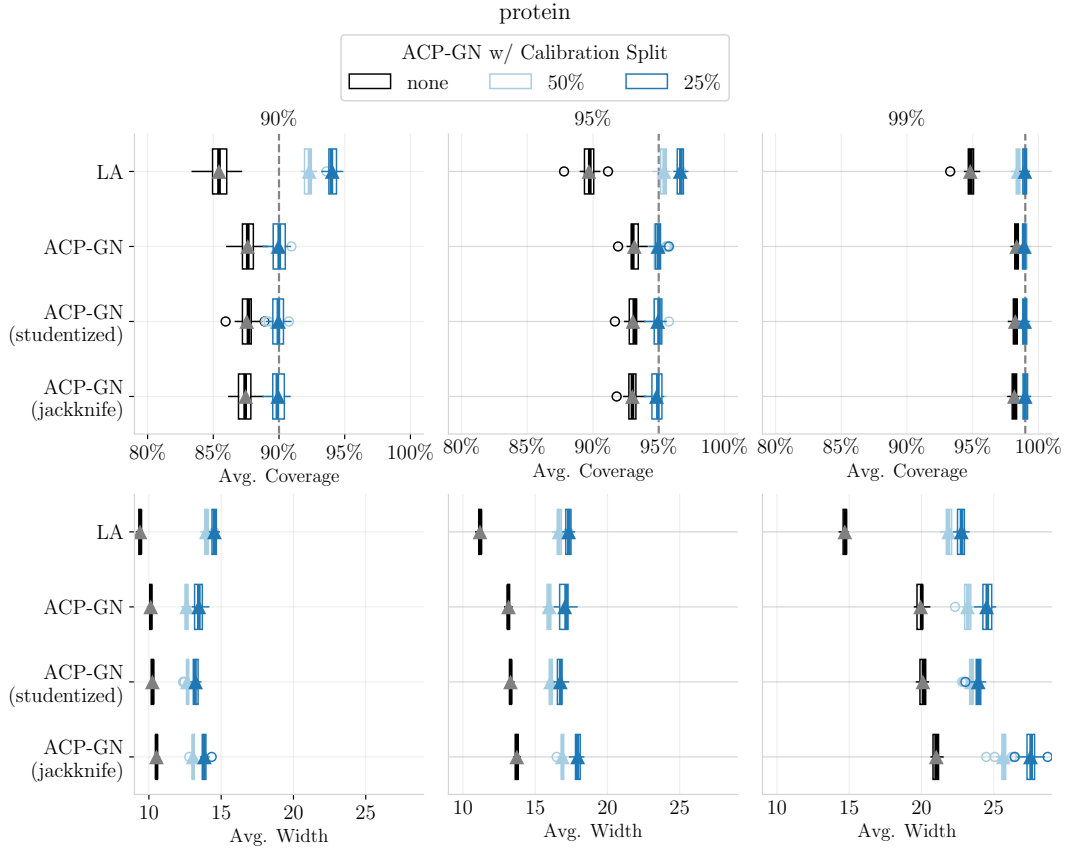


Figure 2: Our extension to ACP-GN that employs a separate calibration set to evaluate nonconformity scores combined with *refinement* of the linearized model shows considerable improvement to the coverage. Unsurprisingly it leads to larger average interval widths due to a smaller training set. This is shown for protein dataset at $\{90\%, 95\%, 99\%\}$ confidence levels.

Table 8: Additional results on UCI regression datasets.

		90%	Avg. Width 95%	99%	90%	Avg. Coverage 95%	99%
concrete $N=1,030$ $I=8$	LA	15.523 \pm 0.087	18.497 \pm 0.103	24.310 \pm 0.136	89.30 \pm 0.17 (✓)	93.53 \pm 0.10 (✓)	97.37 \pm 0.11 (✗)
	SCP	19.216 \pm 0.124	24.526 \pm 0.194	42.110 \pm 0.802	89.92 \pm 0.29 (✓)	94.97 \pm 0.14 (✓)	99.11 \pm 0.05 (✓)
	CRF	18.623 \pm 0.117	23.737 \pm 0.167	40.390 \pm 0.506	89.75 \pm 0.36 (✓)	94.97 \pm 0.13 (✓)	99.05 \pm 0.06 (✓)
	CQR	22.486 \pm 0.081	27.145 \pm 0.091	39.827 \pm 0.212	90.51 \pm 0.27 (✓)	95.24 \pm 0.16 (✓)	99.01 \pm 0.08 (✓)
	ACP-GN	15.727 \pm 0.114	19.810 \pm 0.147	29.192 \pm 0.179	89.60 \pm 0.20 (✓)	94.77 \pm 0.14 (✓)	98.85 \pm 0.08 (✓)
	SCP-GN	18.851 \pm 0.118	23.780 \pm 0.121	36.686 \pm 0.313	89.68 \pm 0.28 (✓)	94.87 \pm 0.14 (✓)	99.15 \pm 0.07 (✓)
	ACP-GN (split + refine)	18.907 \pm 0.103	24.012 \pm 0.129	35.949 \pm 0.387	90.26 \pm 0.20 (✓)	95.23 \pm 0.27 (✓)	99.17 \pm 0.06 (✓)
wine $N=1,599$ $I=11$	LA	2.099 \pm 0.001	2.501 \pm 0.002	3.287 \pm 0.002	91.12 \pm 0.13 (✓)	94.58 \pm 0.10 (✓)	98.34 \pm 0.05 (✓)
	SCP	2.183 \pm 0.009	2.768 \pm 0.012	3.941 \pm 0.020	90.41 \pm 0.19 (✓)	95.07 \pm 0.12 (✓)	99.12 \pm 0.04 (✓)
	CRF	2.304 \pm 0.009	2.921 \pm 0.015	4.381 \pm 0.031	90.14 \pm 0.18 (✓)	95.16 \pm 0.15 (✓)	99.02 \pm 0.05 (✓)
	CQR	1.977 \pm 0.005	2.490 \pm 0.016	3.867 \pm 0.027	90.24 \pm 0.12 (✓)	95.05 \pm 0.06 (✓)	99.07 \pm 0.06 (✓)
	ACP-GN	2.103 \pm 0.003	2.665 \pm 0.006	3.827 \pm 0.007	91.16 \pm 0.12 (✓)	95.80 \pm 0.05 (✓)	99.17 \pm 0.04 (✓)
	SCP-GN	2.134 \pm 0.007	2.676 \pm 0.015	3.753 \pm 0.026	90.06 \pm 0.24 (✓)	95.09 \pm 0.16 (✓)	99.01 \pm 0.05 (✓)
	ACP-GN (split + refine)	2.474 \pm 0.007	3.054 \pm 0.008	4.324 \pm 0.020	89.56 \pm 0.27 (✓)	94.81 \pm 0.13 (✓)	98.99 \pm 0.09 (✓)
kin8nm $N=8,192$ $I=8$	LA	0.213 \pm 0.001	0.254 \pm 0.001	0.334 \pm 0.001	89.66 \pm 0.29 (✓)	93.88 \pm 0.23 (✗)	98.09 \pm 0.10 (✗)
	SCP	0.231 \pm 0.001	0.285 \pm 0.001	0.409 \pm 0.003	90.54 \pm 0.29 (✓)	95.31 \pm 0.20 (✓)	99.18 \pm 0.08 (✓)
	CRF	0.229 \pm 0.001	0.282 \pm 0.001	0.401 \pm 0.002	90.59 \pm 0.31 (✓)	95.41 \pm 0.23 (✓)	99.18 \pm 0.06 (✓)
	CQR	0.254 \pm 0.001	0.304 \pm 0.001	0.426 \pm 0.003	90.24 \pm 0.27 (✓)	95.21 \pm 0.25 (✓)	99.02 \pm 0.09 (✓)
	ACP-GN	0.213 \pm 0.001	0.262 \pm 0.001	0.365 \pm 0.002	89.68 \pm 0.28 (✓)	94.58 \pm 0.20 (✓)	98.85 \pm 0.08 (✓)
	SCP-GN	0.230 \pm 0.001	0.282 \pm 0.001	0.400 \pm 0.002	90.57 \pm 0.29 (✓)	95.20 \pm 0.22 (✓)	99.18 \pm 0.07 (✓)
	ACP-GN (split + refine)	0.232 \pm 0.001	0.284 \pm 0.001	0.406 \pm 0.003	90.45 \pm 0.17 (✓)	95.02 \pm 0.22 (✓)	99.10 \pm 0.07 (✓)
power $N=9,568$ $I=4$	LA	13.345 \pm 0.028	15.901 \pm 0.034	20.898 \pm 0.045	92.08 \pm 0.23 (✗)	95.86 \pm 0.17 (✗)	98.86 \pm 0.09 (✓)
	SCP	12.732 \pm 0.040	15.359 \pm 0.041	21.688 \pm 0.143	90.27 \pm 0.24 (✓)	94.94 \pm 0.17 (✓)	98.98 \pm 0.10 (✓)
	CRF	12.573 \pm 0.039	15.130 \pm 0.039	21.548 \pm 0.145	90.26 \pm 0.27 (✓)	94.93 \pm 0.18 (✓)	98.98 \pm 0.09 (✓)
	CQR	12.860 \pm 0.026	15.180 \pm 0.042	21.946 \pm 0.154	90.37 \pm 0.24 (✓)	94.93 \pm 0.14 (✓)	98.89 \pm 0.06 (✓)
	ACP-GN	12.526 \pm 0.024	15.248 \pm 0.020	21.592 \pm 0.062	90.16 \pm 0.26 (✓)	95.13 \pm 0.19 (✓)	98.89 \pm 0.08 (✓)
	SCP-GN	12.736 \pm 0.041	15.362 \pm 0.044	21.680 \pm 0.146	90.25 \pm 0.25 (✓)	94.95 \pm 0.18 (✓)	98.98 \pm 0.10 (✓)
	ACP-GN (split + refine)	12.744 \pm 0.045	15.442 \pm 0.039	22.043 \pm 0.151	90.17 \pm 0.29 (✓)	95.26 \pm 0.17 (✓)	98.98 \pm 0.09 (✓)
community $N=1,994$ $I=100$	LA	0.548 \pm 0.074	0.653 \pm 0.088	0.858 \pm 0.116	90.90 \pm 0.59 (✓)	93.83 \pm 0.50 (✓)	97.05 \pm 0.33 (✗)
	SCP	0.534 \pm 0.010	0.731 \pm 0.020	1.159 \pm 0.028	90.30 \pm 0.42 (✓)	95.53 \pm 0.24 (✓)	99.12 \pm 0.16 (✓)
	CRF	0.526 \pm 0.010	0.721 \pm 0.020	1.154 \pm 0.029	90.20 \pm 0.43 (✓)	95.33 \pm 0.25 (✓)	99.12 \pm 0.13 (✓)
	CQR	0.554 \pm 0.020	0.701 \pm 0.034	1.077 \pm 0.038	90.90 \pm 0.57 (✓)	95.70 \pm 0.30 (✓)	99.35 \pm 0.15 (✓)
	ACP-GN	0.570 \pm 0.108	0.755 \pm 0.158	1.224 \pm 0.285	90.90 \pm 0.62 (✓)	95.30 \pm 0.47 (✓)	99.25 \pm 0.14 (✓)
	SCP-GN	0.474 \pm 0.013	0.656 \pm 0.024	1.104 \pm 0.034	90.58 \pm 0.37 (✓)	94.95 \pm 0.27 (✓)	99.00 \pm 0.15 (✓)
	ACP-GN (split + refine)	0.519 \pm 0.006	0.652 \pm 0.008	1.010 \pm 0.016	90.97 \pm 0.61 (✓)	95.28 \pm 0.41 (✓)	99.12 \pm 0.16 (✓)
facebook_l $N=40,948$ $I=53$	LA	67.580 \pm 2.637	80.527 \pm 3.142	105.831 \pm 4.130	97.63 \pm 0.09 (✗)	98.15 \pm 0.08 (✗)	98.76 \pm 0.06 (✗)
	SCP	20.771 \pm 2.452	44.192 \pm 2.799	178.890 \pm 4.596	90.00 \pm 0.10 (✓)	95.03 \pm 0.08 (✓)	99.05 \pm 0.04 (✓)
	CRF	18.689 \pm 1.944	35.765 \pm 1.870	121.181 \pm 4.148	89.98 \pm 0.12 (✓)	95.10 \pm 0.09 (✓)	99.08 \pm 0.04 (✓)
	CQR	19.625 \pm 0.878	25.970 \pm 1.322	59.922 \pm 12.774	90.16 \pm 0.24 (✓)	94.91 \pm 0.11 (✓)	98.97 \pm 0.04 (✓)
	ACP-GN	17.986 \pm 0.480	41.063 \pm 0.770	199.331 \pm 10.821	90.36 \pm 0.16 (✓)	95.56 \pm 0.17 (✗)	99.37 \pm 0.08 (✗)
	SCP-GN	20.460 \pm 2.463	39.712 \pm 3.148	125.078 \pm 8.075	90.00 \pm 0.10 (✓)	94.94 \pm 0.09 (✓)	99.04 \pm 0.03 (✓)
	ACP-GN (split + refine)	29.006 \pm 4.472	54.099 \pm 4.576	172.252 \pm 6.758	90.06 \pm 0.10 (✓)	95.20 \pm 0.08 (✓)	99.14 \pm 0.04 (✓)

Table 9: We repeat the UCI regression experiment in Table 1 but with 25% calibration split. In the case of yacht, the desired coverage of 99% was too small a significance level to accurately evaluate the quantile of the calibration scores.

		Avg. Width		Avg. Coverage			
		90%	95%	99%	90%	95%	99%
yacht N=308 I=6	LA	1.690±0.017	2.014±0.020	2.647±0.027	88.73±0.61 (✓)	90.78±0.59 (✗)	93.89±0.60 (✗)
	SCP	2.306±0.117	3.866±0.133	—	90.11±0.70 (✓)	95.88±0.43 (✓)	—
	CRF	2.281±0.112	3.818±0.130	—	90.11±0.67 (✓)	95.82±0.46 (✓)	—
	CQR	3.236±0.151	4.895±0.211	—	90.33±0.65 (✓)	95.85±0.39 (✓)	—
	ACP-GN	1.594 ±0.016	2.385 ±0.029	6.915 ±0.067	87.36±0.58 (✓)	92.56±0.68 (✓)	99.03±0.11 (✓)
	SCP-GN	2.070±0.084	3.111±0.110	—	90.53±0.54 (✓)	95.95±0.51 (✓)	—
	ACP-GN (split + refine)	3.931±0.083	5.926±0.144	—	91.66±0.46 (✓)	96.66±0.41 (✓)	—
boston N=506 I=13	LA	9.398±0.046	11.199 ±0.055	14.718±0.072	91.24±0.31 (✓)	94.34±0.22 (✓)	97.53±0.11 (✗)
	SCP	10.086±0.184	14.746±0.345	38.574±1.914	90.12±0.51 (✓)	96.03±0.34 (✓)	99.19±0.12 (✓)
	CRF	9.874±0.194	14.204±0.270	37.421±2.132	90.03±0.45 (✓)	95.52±0.24 (✓)	99.28±0.11 (✓)
	CQR	10.918±0.121	14.618±0.208	31.130±0.823	90.09±0.31 (✓)	95.65±0.21 (✓)	99.17±0.10 (✓)
	ACP-GN	9.182 ±0.046	12.111±0.038	20.512 ±0.057	90.64±0.26 (✓)	95.49±0.16 (✓)	99.11±0.08 (✓)
	SCP-GN	9.787±0.152	13.372±0.247	27.274±1.221	90.21±0.44 (✓)	95.49±0.29 (✓)	99.13±0.14 (✓)
	ACP-GN (split + refine)	17.372±0.252	22.852±0.371	42.081±1.170	90.68±0.44 (✓)	95.79±0.28 (✓)	98.95±0.15 (✓)
energy N=768 I=8	LA	1.502±0.006	1.790±0.007	2.353±0.009	88.96±0.35 (✓)	92.92±0.33 (✗)	96.95±0.23 (✗)
	SCP	1.824±0.037	2.365±0.046	4.313±0.129	90.05±0.45 (✓)	95.04±0.25 (✓)	99.39±0.09 (✓)
	CRF	1.807±0.037	2.341±0.046	4.316±0.133	89.97±0.47 (✓)	95.02±0.23 (✓)	99.36±0.10 (✓)
	CQR	4.779±0.032	5.181±0.035	7.692±0.171	90.55±0.26 (✓)	95.51±0.22 (✓)	99.35±0.06 (✓)
	ACP-GN	1.462 ±0.006	1.884 ±0.008	3.076 ±0.015	88.28±0.33 (✓)	93.69±0.33 (✓)	98.88±0.11 (✓)
	SCP-GN	1.812±0.034	2.348±0.047	4.226±0.120	90.16±0.43 (✓)	95.19±0.24 (✓)	99.31±0.11 (✓)
	ACP-GN (split + refine)	2.498±0.056	3.251±0.072	5.601±0.168	90.19±0.53 (✓)	95.28±0.37 (✓)	99.43±0.07 (✓)
bike N=10,886 I=18	LA	100.451±2.394	119.694±2.853	157.305±3.749	89.82±0.39 (✓)	93.29±0.33 (✗)	96.83±0.16 (✗)
	SCP	113.594±1.370	159.114±2.297	287.486±5.425	90.48±0.22 (✓)	95.11±0.20 (✓)	98.87±0.08 (✓)
	CRF	112.104±1.653	156.161±2.486	283.911±6.247	90.38±0.22 (✓)	95.17±0.19 (✓)	98.93±0.09 (✓)
	CQR	107.141 ±3.598	130.756 ±3.259	212.513 ±5.131	90.34±0.26 (✓)	95.18±0.18 (✓)	98.97±0.09 (✓)
	ACP-GN	98.813 ±2.485	130.893±3.231	213.131±5.630	89.36±0.43 (✓)	94.41±0.27 (✗)	98.67±0.09 (✗)
	SCP-GN	105.020±1.209	139.391±1.711	228.656±4.493	90.17±0.21 (✓)	95.10±0.16 (✓)	98.99±0.07 (✓)
	ACP-GN (split + refine)	137.114±0.759	178.639±1.307	285.997±3.834	90.20±0.18 (✓)	95.05±0.16 (✓)	99.16±0.10 (✓)
protein N=45,730 I=9	LA	9.385±0.022	11.183±0.027	14.697±0.035	85.43±0.18 (✗)	89.69±0.15 (✗)	94.81±0.10 (✗)
	SCP	12.188±0.036	16.069±0.050	24.828±0.126	89.85±0.11 (✓)	94.91±0.07 (✓)	98.92±0.05 (✓)
	CRF	11.532 ±0.047	15.398±0.071	25.006±0.147	89.89±0.09 (✓)	94.93±0.08 (✓)	98.92±0.03 (✓)
	CQR	13.174±0.162	14.422 ±0.148	17.924 ±0.078	90.10±0.09 (✓)	95.04±0.05 (✓)	98.92±0.04 (✓)
	ACP-GN	10.243±0.019	13.294±0.027	20.101±0.053	87.54±0.15 (✗)	93.04±0.11 (✗)	98.24±0.05 (✗)
	SCP-GN	11.743±0.035	15.312±0.035	23.243±0.097	89.87±0.09 (✓)	94.91±0.07 (✓)	98.95±0.04 (✓)
	ACP-GN (split + refine)	13.229±0.042	16.741±0.055	23.927±0.068	89.96±0.11 (✓)	94.93±0.10 (✓)	98.94±0.05 (✓)
facebook_2 N=81,311 I=53	LA	66.088±2.760	78.749±3.289	103.493±4.322	97.47±0.12 (✗)	98.01±0.09 (✗)	98.65±0.06 (✗)
	SCP	15.739 ±0.240	34.114±0.646	145.543±1.431	89.85±0.09 (✓)	94.96±0.06 (✓)	99.01±0.03 (✓)
	CRF	15.257 ±0.663	29.889±1.398	100.709±6.136	89.82±0.08 (✓)	94.97±0.05 (✓)	99.01±0.03 (✓)
	CQR	17.929 ±0.703	22.314 ±1.070	32.257 ±1.604	89.85±0.06 (✓)	95.00±0.04 (✓)	99.02±0.03 (✓)
	ACP-GN	18.396±0.546	40.088±1.091	166.792±5.632	90.47±0.11 (✗)	95.45±0.08 (✗)	99.35±0.04 (✗)
	SCP-GN	15.633 ±0.220	33.192±0.708	122.079±4.142	89.81±0.07 (✓)	94.96±0.06 (✓)	99.00±0.02 (✓)
	ACP-GN (split + refine)	20.487±0.502	41.051±0.844	152.090±6.966	90.18±0.08 (✓)	95.08±0.07 (✓)	99.11±0.03 (✓)

Table 10: We repeat the UCI regression experiment in Table 8 but with 25% calibration split.

		90%	Avg. Width 95%	99%	90%	Avg. Coverage 95%	99%
concrete $N=1,030$ $I=8$	LA	15.523 \pm 0.087	18.497 \pm 0.103	24.310 \pm 0.136	89.30 \pm 0.17 (✓)	93.53 \pm 0.10 (✓)	97.37 \pm 0.11 (✗)
	SCP	17.563 \pm 0.209	22.787 \pm 0.277	39.972 \pm 1.300	89.35 \pm 0.43 (✓)	94.93 \pm 0.27 (✓)	99.02 \pm 0.12 (✓)
	CRF	16.946 \pm 0.176	21.898 \pm 0.271	38.480 \pm 1.287	89.44 \pm 0.33 (✓)	94.88 \pm 0.30 (✓)	99.03 \pm 0.12 (✓)
	CQR	20.063 \pm 0.092	24.620 \pm 0.144	37.035 \pm 0.329	89.88 \pm 0.27 (✓)	94.90 \pm 0.19 (✓)	99.11 \pm 0.07 (✓)
	ACP-GN	15.727 \pm 0.114	19.810 \pm 0.147	29.192 \pm 0.179	89.60 \pm 0.20 (✓)	94.77 \pm 0.14 (✓)	98.85 \pm 0.08 (✓)
	SCP-GN	17.373 \pm 0.177	22.170 \pm 0.233	34.104 \pm 0.566	89.47 \pm 0.45 (✓)	94.94 \pm 0.25 (✓)	99.11 \pm 0.12 (✓)
	ACP-GN (split + refine)	25.331 \pm 0.182	32.279 \pm 0.244	47.494 \pm 0.795	90.30 \pm 0.40 (✓)	95.45 \pm 0.21 (✓)	99.13 \pm 0.11 (✓)
wine $N=1,599$ $I=11$	LA	2.099 \pm 0.001	2.501 \pm 0.002	3.287 \pm 0.002	91.12 \pm 0.13 (✓)	94.58 \pm 0.10 (✓)	98.34 \pm 0.05 (✓)
	SCP	2.112 \pm 0.010	2.686 \pm 0.018	4.032 \pm 0.028	90.11 \pm 0.15 (✓)	95.00 \pm 0.11 (✓)	99.27 \pm 0.05 (✓)
	CRF	2.186 \pm 0.013	2.753 \pm 0.013	4.397 \pm 0.075	89.92 \pm 0.15 (✓)	94.92 \pm 0.13 (✓)	99.20 \pm 0.06 (✓)
	CQR	1.966 \pm 0.004	2.421 \pm 0.016	3.920 \pm 0.031	89.97 \pm 0.22 (✓)	94.95 \pm 0.09 (✓)	99.12 \pm 0.08 (✓)
	ACP-GN	2.103 \pm 0.003	2.665 \pm 0.006	3.827 \pm 0.007	91.16 \pm 0.12 (✓)	95.80 \pm 0.05 (✓)	99.17 \pm 0.04 (✓)
	SCP-GN	2.068 \pm 0.009	2.621 \pm 0.018	3.760 \pm 0.035	90.07 \pm 0.18 (✓)	95.00 \pm 0.11 (✓)	99.09 \pm 0.06 (✓)
	ACP-GN (split + refine)	2.908 \pm 0.016	3.527 \pm 0.023	5.089 \pm 0.037	90.14 \pm 0.20 (✓)	94.93 \pm 0.21 (✓)	99.09 \pm 0.09 (✓)
kin8nm $N=8,192$ $I=8$	LA	0.213 \pm 0.001	0.254 \pm 0.001	0.334 \pm 0.001	89.66 \pm 0.29 (✓)	93.88 \pm 0.23 (✗)	98.09 \pm 0.10 (✗)
	SCP	0.223 \pm 0.001	0.275 \pm 0.001	0.394 \pm 0.003	90.10 \pm 0.27 (✓)	95.35 \pm 0.20 (✓)	99.15 \pm 0.09 (✓)
	CRF	0.221 \pm 0.001	0.271 \pm 0.001	0.381 \pm 0.003	90.15 \pm 0.27 (✓)	95.36 \pm 0.19 (✓)	99.12 \pm 0.09 (✓)
	CQR	0.241 \pm 0.001	0.286 \pm 0.001	0.392 \pm 0.003	90.41 \pm 0.24 (✓)	95.18 \pm 0.20 (✓)	98.98 \pm 0.09 (✓)
	ACP-GN	0.213 \pm 0.001	0.262 \pm 0.001	0.365 \pm 0.002	89.68 \pm 0.28 (✓)	94.58 \pm 0.20 (✓)	98.85 \pm 0.08 (✓)
	SCP-GN	0.222 \pm 0.001	0.273 \pm 0.001	0.386 \pm 0.003	90.05 \pm 0.27 (✓)	95.36 \pm 0.20 (✓)	99.12 \pm 0.08 (✓)
	ACP-GN (split + refine)	0.243 \pm 0.001	0.298 \pm 0.001	0.423 \pm 0.004	90.60 \pm 0.27 (✓)	95.27 \pm 0.16 (✓)	99.17 \pm 0.09 (✓)
power $N=9,568$ $I=4$	LA	13.345 \pm 0.028	15.901 \pm 0.034	20.898 \pm 0.045	92.08 \pm 0.23 (✗)	95.86 \pm 0.17 (✗)	98.86 \pm 0.09 (✓)
	SCP	12.620 \pm 0.051	15.212 \pm 0.069	21.720 \pm 0.245	90.29 \pm 0.22 (✓)	95.01 \pm 0.19 (✓)	98.99 \pm 0.10 (✓)
	CRF	12.457 \pm 0.049	15.019 \pm 0.069	21.553 \pm 0.189	90.32 \pm 0.24 (✓)	95.06 \pm 0.17 (✓)	98.96 \pm 0.10 (✓)
	CQR	12.818 \pm 0.041	15.081 \pm 0.058	21.903 \pm 0.170	90.52 \pm 0.22 (✓)	94.90 \pm 0.18 (✓)	99.02 \pm 0.08 (✓)
	ACP-GN	12.526 \pm 0.024	15.248 \pm 0.020	21.592 \pm 0.062	90.16 \pm 0.26 (✓)	95.13 \pm 0.19 (✓)	98.89 \pm 0.08 (✓)
	SCP-GN	12.627 \pm 0.050	15.215 \pm 0.070	21.724 \pm 0.242	90.32 \pm 0.22 (✓)	95.01 \pm 0.18 (✓)	99.00 \pm 0.09 (✓)
	ACP-GN (split + refine)	12.834 \pm 0.054	15.561 \pm 0.067	22.548 \pm 0.255	90.18 \pm 0.26 (✓)	95.18 \pm 0.19 (✓)	99.05 \pm 0.10 (✓)
community $N=1,994$ $I=100$	LA	0.548 \pm 0.074	0.653 \pm 0.088	0.858 \pm 0.116	90.90 \pm 0.59 (✓)	93.83 \pm 0.50 (✓)	97.05 \pm 0.33 (✗)
	SCP	0.471 \pm 0.009	0.642 \pm 0.013	1.033 \pm 0.027	90.12 \pm 0.55 (✓)	95.20 \pm 0.46 (✓)	98.97 \pm 0.20 (✓)
	CRF	0.464 \pm 0.010	0.617 \pm 0.015	0.992 \pm 0.029	90.35 \pm 0.53 (✓)	95.03 \pm 0.45 (✓)	98.90 \pm 0.18 (✓)
	CQR	0.659 \pm 0.034	0.926 \pm 0.058	1.349 \pm 0.067	91.28 \pm 0.49 (✓)	96.00 \pm 0.37 (✓)	99.65 \pm 0.18 (✓)
	ACP-GN	0.460 \pm 0.002	0.593 \pm 0.005	0.932 \pm 0.006	90.45 \pm 0.46 (✓)	95.05 \pm 0.42 (✓)	99.21 \pm 0.14 (✓)
	SCP-GN	0.448 \pm 0.008	0.621 \pm 0.013	1.069 \pm 0.048	90.10 \pm 0.57 (✓)	95.35 \pm 0.42 (✓)	99.03 \pm 0.19 (✓)
	ACP-GN (split + refine)	0.561 \pm 0.009	0.700 \pm 0.010	1.092 \pm 0.029	89.97 \pm 0.69 (✓)	94.92 \pm 0.45 (✓)	98.85 \pm 0.20 (✓)
facebook_1 $N=40,948$ $I=53$	LA	67.580 \pm 2.637	80.527 \pm 3.142	105.831 \pm 4.130	97.63 \pm 0.09 (✗)	98.15 \pm 0.08 (✗)	98.76 \pm 0.06 (✗)
	SCP	18.756 \pm 0.625	42.110 \pm 1.381	176.332 \pm 5.531	89.98 \pm 0.12 (✓)	95.08 \pm 0.09 (✓)	99.03 \pm 0.04 (✓)
	CRF	16.845 \pm 0.513	34.078 \pm 0.936	114.769 \pm 4.111	90.04 \pm 0.13 (✓)	95.12 \pm 0.09 (✓)	98.98 \pm 0.05 (✓)
	CQR	21.120 \pm 1.049	27.024 \pm 1.543	43.095 \pm 2.232	90.05 \pm 0.14 (✓)	95.04 \pm 0.08 (✓)	98.95 \pm 0.03 (✓)
	ACP-GN	17.986 \pm 0.480	41.063 \pm 0.770	199.331 \pm 10.821	90.36 \pm 0.16 (✓)	95.56 \pm 0.17 (✗)	99.37 \pm 0.08 (✗)
	SCP-GN	18.696 \pm 0.612	39.460 \pm 1.552	130.967 \pm 6.779	90.04 \pm 0.13 (✓)	95.04 \pm 0.10 (✓)	99.03 \pm 0.04 (✓)
	ACP-GN (split + refine)	24.639 \pm 1.195	51.864 \pm 2.687	181.585 \pm 11.155	90.27 \pm 0.13 (✓)	95.29 \pm 0.09 (✓)	99.14 \pm 0.04 (✓)

Table 11: KFAC approximation to the Gauss-Newton evaluated on large UCI datasets

		Avg. Width		Avg. Coverage			
		90%	95%	90%	95%		
bike $N=10,886$ $I=18$	LA	100.451 \pm 2.394	119.694 \pm 2.853	157.305 \pm 3.749	89.82 \pm 0.39 (✓)	93.29 \pm 0.33 (✗)	96.83 \pm 0.16 (✗)
	LA(<i>kfac</i>)	91.501 \pm 3.058	109.030 \pm 3.644	143.290 \pm 4.789	87.24 \pm 1.01 (✗)	91.07 \pm 0.84 (✗)	95.20 \pm 0.55 (✗)
	SCP	131.138 \pm 0.812	180.477 \pm 1.244	324.756 \pm 4.635	90.33 \pm 0.21 (✓)	95.17 \pm 0.15 (✓)	99.00 \pm 0.07 (✓)
	SCP-GN	122.245 \pm 1.073	160.505 \pm 1.761	254.409 \pm 3.767	90.34 \pm 0.24 (✓)	95.26 \pm 0.15 (✓)	99.02 \pm 0.08 (✓)
	SCP-GN(<i>kfac</i>)	121.401 \pm 0.708	160.214 \pm 1.328	259.909 \pm 3.247	90.46 \pm 0.21 (✓)	95.22 \pm 0.18 (✓)	98.99 \pm 0.07 (✓)
	ACP-GN	102.401 \pm 0.637	133.781 \pm 7.394	212.144 \pm 11.677	89.27 \pm 0.35 (✗)	94.45 \pm 0.24 (✗)	98.69 \pm 0.10 (✗)
	ACP-GN(<i>kfac</i>)	94.964 \pm 2.095	124.501 \pm 3.011	198.018 \pm 5.012	88.67 \pm 0.68 (✗)	94.06 \pm 0.52 (✗)	98.49 \pm 0.18 (✗)
	ACP-GN (split + refine)	125.090 \pm 4.180	162.697 \pm 5.654	254.893 \pm 9.363	90.20 \pm 0.22 (✓)	94.95 \pm 0.13 (✓)	99.01 \pm 0.08 (✓)
	ACP-GN(<i>kfac</i>) (split + refine)	162.693 \pm 2.312	211.953 \pm 2.906	334.975 \pm 5.205	91.31 \pm 0.19 (✗)	95.81 \pm 0.11 (✗)	99.24 \pm 0.05 (✓)
	LA	0.548 \pm 0.074	0.653 \pm 0.088	0.858 \pm 0.116	90.90 \pm 0.59 (✓)	93.83 \pm 0.50 (✓)	97.05 \pm 0.33 (✗)
community $N=1,994$ $I=100$	LA(<i>kfac</i>)	0.472 \pm 0.003	0.562 \pm 0.004	0.739 \pm 0.005	90.45 \pm 0.42 (✓)	93.55 \pm 0.41 (✗)	96.95 \pm 0.30 (✗)
	SCP	0.534 \pm 0.010	0.735 \pm 0.020	1.164 \pm 0.029	90.17 \pm 0.41 (✓)	95.33 \pm 0.22 (✓)	99.17 \pm 0.17 (✓)
	SCP-GN	0.473 \pm 0.013	0.660 \pm 0.024	1.116 \pm 0.034	90.55 \pm 0.39 (✓)	95.10 \pm 0.25 (✓)	99.12 \pm 0.14 (✓)
	SCP-GN(<i>kfac</i>)	0.476 \pm 0.012	0.657 \pm 0.023	1.095 \pm 0.030	90.62 \pm 0.41 (✓)	95.30 \pm 0.23 (✓)	98.90 \pm 0.18 (✓)
	ACP-GN	0.611 \pm 0.151	0.784 \pm 0.191	1.222 \pm 0.289	90.92 \pm 0.62 (✓)	95.28 \pm 0.45 (✓)	99.25 \pm 0.13 (✓)
	ACP-GN(<i>kfac</i>)	0.476 \pm 0.018	0.612 \pm 0.023	0.964 \pm 0.035	90.78 \pm 0.48 (✓)	95.20 \pm 0.37 (✓)	99.25 \pm 0.12 (✓)
	ACP-GN (split + refine)	0.523 \pm 0.005	0.661 \pm 0.007	1.040 \pm 0.016	91.38 \pm 0.55 (✓)	95.45 \pm 0.44 (✓)	99.20 \pm 0.14 (✓)
	ACP-GN(<i>kfac</i>) (split + refine)	0.530 \pm 0.005	0.674 \pm 0.008	1.067 \pm 0.019	91.20 \pm 0.64 (✓)	95.62 \pm 0.44 (✓)	99.35 \pm 0.12 (✓)
	LA	9.385 \pm 0.022	11.183 \pm 0.027	14.697 \pm 0.035	85.43 \pm 0.18 (✗)	89.69 \pm 0.15 (✗)	94.81 \pm 0.10 (✗)
	LA(<i>kfac</i>)	9.132 \pm 0.027	10.881 \pm 0.032	14.301 \pm 0.042	84.33 \pm 0.18 (✗)	88.63 \pm 0.16 (✗)	94.04 \pm 0.12 (✗)
protein $N=45,730$ $I=9$	SCP	13.041 \pm 0.088	17.161 \pm 0.098	26.181 \pm 0.119	89.78 \pm 0.08 (✓)	94.83 \pm 0.06 (✓)	98.94 \pm 0.04 (✓)
	SCP-GN	12.426 \pm 0.085	16.102 \pm 0.096	24.032 \pm 0.138	89.78 \pm 0.10 (✓)	94.86 \pm 0.08 (✓)	98.94 \pm 0.03 (✓)
	SCP-GN(<i>kfac</i>)	12.679 \pm 0.084	16.570 \pm 0.096	24.981 \pm 0.133	89.76 \pm 0.10 (✓)	94.85 \pm 0.08 (✓)	98.94 \pm 0.03 (✓)
	ACP-GN	10.127 \pm 0.020	13.156 \pm 0.027	19.943 \pm 0.061	87.62 \pm 0.15 (✗)	93.14 \pm 0.10 (✗)	98.35 \pm 0.04 (✗)
	ACP-GN(<i>kfac</i>)	9.695 \pm 0.034	12.595 \pm 0.041	19.093 \pm 0.064	86.02 \pm 0.18 (✗)	91.92 \pm 0.15 (✗)	97.77 \pm 0.06 (✗)
	ACP-GN (split + refine)	12.614 \pm 0.034	15.961 \pm 0.044	23.188 \pm 0.069	89.88 \pm 0.10 (✓)	94.90 \pm 0.09 (✓)	98.97 \pm 0.05 (✓)
	ACP-GN(<i>kfac</i>) (split + refine)	12.129 \pm 0.105	15.363 \pm 0.128	22.260 \pm 0.152	88.39 \pm 0.14 (✗)	93.82 \pm 0.13 (✗)	98.49 \pm 0.07 (✗)
	LA	67.580 \pm 2.637	80.527 \pm 3.142	105.831 \pm 4.130	97.63 \pm 0.09 (✗)	98.15 \pm 0.08 (✗)	98.76 \pm 0.06 (✗)
	LA(<i>kfac</i>)	67.534 \pm 3.806	80.471 \pm 4.535	105.757 \pm 5.960	97.39 \pm 0.15 (✗)	97.88 \pm 0.12 (✗)	98.57 \pm 0.09 (✗)
	facebook_1 $N=40,948$ $I=53$	SCP	20.771 \pm 2.452	44.192 \pm 2.799	178.890 \pm 4.596	90.00 \pm 0.10 (✓)	95.03 \pm 0.08 (✓)
SCP-GN	20.460 \pm 2.463	39.712 \pm 3.148	125.078 \pm 8.075	90.01 \pm 0.10 (✓)	94.94 \pm 0.09 (✓)	99.04 \pm 0.03 (✓)	
SCP-GN(<i>kfac</i>)	20.303 \pm 2.476	40.531 \pm 2.973	128.583 \pm 5.746	90.01 \pm 0.11 (✓)	95.04 \pm 0.09 (✓)	99.03 \pm 0.03 (✓)	
facebook_2 $N=81,311$ $I=53$	ACP-GN	30.466 \pm 7.858	62.895 \pm 13.942	223.609 \pm 45.677	90.56 \pm 0.19 (✗)	95.61 \pm 0.16 (✗)	99.26 \pm 0.05 (✗)
	ACP-GN(<i>kfac</i>)	19.275 \pm 0.542	40.991 \pm 0.964	149.319 \pm 5.539	90.24 \pm 0.12 (✓)	95.38 \pm 0.09 (✗)	99.24 \pm 0.06 (✗)
	ACP-GN (split + refine)	33.780 \pm 5.321	57.292 \pm 5.730	141.074 \pm 6.581	90.13 \pm 0.11 (✓)	95.23 \pm 0.06 (✓)	99.08 \pm 0.04 (✓)
	ACP-GN(<i>kfac</i>) (split + refine)	23.732 \pm 1.180	45.239 \pm 2.220	121.631 \pm 3.699	89.81 \pm 0.10 (✓)	94.86 \pm 0.08 (✓)	98.90 \pm 0.04 (✓)
	LA	66.088 \pm 2.760	78.749 \pm 3.289	103.493 \pm 4.322	97.47 \pm 0.12 (✗)	98.01 \pm 0.09 (✗)	98.65 \pm 0.06 (✗)
	LA(<i>kfac</i>)	63.591 \pm 3.085	75.774 \pm 3.675	99.584 \pm 4.830	97.26 \pm 0.15 (✗)	97.84 \pm 0.12 (✗)	98.54 \pm 0.08 (✗)
	SCP	16.387 \pm 0.208	35.387 \pm 0.462	152.706 \pm 1.591	89.97 \pm 0.07 (✓)	95.00 \pm 0.06 (✓)	99.06 \pm 0.03 (✓)
	SCP-GN	16.287 \pm 0.202	33.655 \pm 0.563	118.489 \pm 4.305	89.99 \pm 0.07 (✓)	94.98 \pm 0.06 (✓)	99.00 \pm 0.03 (✓)
	SCP-GN(<i>kfac</i>)	16.213 \pm 0.195	33.686 \pm 0.530	116.672 \pm 2.618	89.96 \pm 0.08 (✓)	94.98 \pm 0.06 (✓)	99.02 \pm 0.02 (✓)
	ACP-GN	26.536 \pm 3.513	56.470 \pm 7.754	197.251 \pm 27.935	90.54 \pm 0.11 (✗)	95.46 \pm 0.08 (✗)	99.23 \pm 0.04 (✗)
ACP-GN(<i>kfac</i>)	20.497 \pm 1.859	43.246 \pm 3.955	147.832 \pm 13.223	90.34 \pm 0.08 (✗)	95.22 \pm 0.05 (✗)	99.18 \pm 0.03 (✗)	
ACP-GN (split + refine)	23.095 \pm 1.120	42.653 \pm 1.161	123.337 \pm 2.954	90.17 \pm 0.08 (✓)	95.15 \pm 0.05 (✓)	99.08 \pm 0.03 (✓)	
ACP-GN(<i>kfac</i>) (split + refine)	20.487 \pm 0.257	39.636 \pm 0.582	117.994 \pm 2.407	90.11 \pm 0.09 (✓)	95.15 \pm 0.06 (✓)	99.11 \pm 0.03 (✓)	

Table 12: Last-layer (LL) approximation to the Gauss-Newton evaluated on large UCI datasets

		Avg. Width			Avg. Coverage		
		90%	95%	99%	90%	95%	99%
bike $N=10,886$ $I=18$	LA	100.451 \pm 2.394	119.694 \pm 2.853	157.305 \pm 3.749	89.82 \pm 0.39 (✓)	93.29 \pm 0.33 (✗)	96.83 \pm 0.16 (✗)
	LA(LL)	79.923 \pm 3.087	95.234 \pm 3.679	125.159 \pm 4.835	82.80 \pm 1.32 (✓)	86.98 \pm 1.21 (✗)	92.19 \pm 0.95 (✗)
	SCP	131.138 \pm 0.812	180.477 \pm 1.244	324.756 \pm 4.635	90.33 \pm 0.21 (✓)	95.17 \pm 0.15 (✓)	99.00 \pm 0.07 (✓)
	SCP-GN	122.245 \pm 1.073	160.505 \pm 1.761	254.409 \pm 3.767	90.34 \pm 0.24 (✓)	95.26 \pm 0.15 (✓)	99.02 \pm 0.08 (✓)
	SCP-GN(LL)	130.965 \pm 0.796	180.227 \pm 1.213	324.427 \pm 4.404	90.26 \pm 0.22 (✓)	95.17 \pm 0.15 (✓)	99.00 \pm 0.07 (✓)
	ACP-GN	98.813 \pm 2.485	130.893 \pm 3.231	213.131 \pm 5.630	89.36 \pm 0.43 (✓)	94.41 \pm 0.27 (✗)	98.67 \pm 0.09 (✗)
	ACP-GN(LL)	76.808 \pm 2.933	100.920 \pm 4.078	160.448 \pm 6.648	81.69 \pm 1.31 (✗)	88.20 \pm 1.23 (✗)	95.17 \pm 0.78 (✗)
	ACP-GN (split + refine)	128.336 \pm 4.336	170.782 \pm 5.859	281.632 \pm 10.176	89.98 \pm 0.22 (✓)	94.94 \pm 0.16 (✓)	99.01 \pm 0.06 (✓)
community $N=1,994$ $I=100$	ACP-GN(LL) (split + refine)	128.853 \pm 1.101	178.079 \pm 1.464	325.171 \pm 3.916	90.12 \pm 0.23 (✓)	95.12 \pm 0.20 (✓)	99.02 \pm 0.07 (✓)
	LA	0.548 \pm 0.074	0.653 \pm 0.088	0.858 \pm 0.116	90.90 \pm 0.59 (✓)	93.83 \pm 0.50 (✓)	97.05 \pm 0.33 (✗)
	LA(LL)	0.455 \pm 0.013	0.542 \pm 0.015	0.712 \pm 0.020	89.28 \pm 1.00 (✓)	92.33 \pm 0.94 (✗)	96.05 \pm 0.79 (✗)
	SCP	0.534 \pm 0.010	0.735 \pm 0.020	1.164 \pm 0.029	90.17 \pm 0.41 (✓)	95.33 \pm 0.22 (✓)	99.17 \pm 0.17 (✓)
	SCP-GN	0.473 \pm 0.013	0.660 \pm 0.024	1.116 \pm 0.034	90.55 \pm 0.39 (✓)	95.10 \pm 0.25 (✓)	99.12 \pm 0.14 (✓)
	SCP-GN(LL)	0.533 \pm 0.010	0.730 \pm 0.020	1.157 \pm 0.028	90.33 \pm 0.42 (✓)	95.50 \pm 0.24 (✓)	99.15 \pm 0.16 (✓)
	ACP-GN	0.570 \pm 0.108	0.755 \pm 0.158	1.224 \pm 0.285	90.90 \pm 0.62 (✓)	95.30 \pm 0.47 (✓)	99.25 \pm 0.14 (✓)
	ACP-GN(LL)	0.438 \pm 0.012	0.564 \pm 0.016	0.887 \pm 0.026	88.83 \pm 1.06 (✓)	93.78 \pm 0.99 (✗)	98.42 \pm 0.65 (✗)
protein $N=45,730$ $I=9$	ACP-GN (split + refine)	0.519 \pm 0.006	0.652 \pm 0.008	1.010 \pm 0.016	90.97 \pm 0.61 (✓)	95.28 \pm 0.41 (✓)	99.12 \pm 0.16 (✓)
	ACP-GN(LL) (split + refine)	0.495 \pm 0.005	0.648 \pm 0.006	1.040 \pm 0.014	90.75 \pm 0.40 (✓)	95.70 \pm 0.29 (✓)	99.35 \pm 0.12 (✓)
	LA	9.385 \pm 0.022	11.183 \pm 0.027	14.697 \pm 0.035	85.43 \pm 0.18 (✗)	89.69 \pm 0.15 (✗)	94.81 \pm 0.10 (✗)
	LA(LL)	8.667 \pm 0.026	10.328 \pm 0.031	13.573 \pm 0.041	82.51 \pm 0.17 (✗)	87.08 \pm 0.17 (✗)	92.81 \pm 0.12 (✗)
	SCP	13.041 \pm 0.088	17.161 \pm 0.098	26.181 \pm 0.119	89.78 \pm 0.08 (✓)	94.83 \pm 0.06 (✓)	98.94 \pm 0.04 (✓)
	SCP-GN	12.426 \pm 0.085	16.102 \pm 0.096	24.032 \pm 0.138	89.78 \pm 0.10 (✓)	94.86 \pm 0.08 (✓)	98.94 \pm 0.03 (✓)
	SCP-GN(LL)	13.035 \pm 0.088	17.150 \pm 0.098	26.167 \pm 0.119	89.77 \pm 0.08 (✓)	94.82 \pm 0.06 (✓)	98.94 \pm 0.04 (✓)
	ACP-GN	10.243 \pm 0.019	13.294 \pm 0.027	20.101 \pm 0.053	87.54 \pm 0.15 (✗)	93.04 \pm 0.11 (✗)	98.24 \pm 0.05 (✗)
facebook_1 $N=40,948$ $I=53$	ACP-GN(LL)	8.731 \pm 0.032	11.341 \pm 0.038	17.192 \pm 0.060	82.62 \pm 0.18 (✗)	89.16 \pm 0.15 (✗)	96.29 \pm 0.08 (✗)
	ACP-GN (split + refine)	12.660 \pm 0.028	16.073 \pm 0.031	23.445 \pm 0.057	89.83 \pm 0.09 (✓)	94.90 \pm 0.09 (✓)	98.97 \pm 0.05 (✓)
	ACP-GN(LL) (split + refine)	12.616 \pm 0.029	16.078 \pm 0.036	23.672 \pm 0.075	89.81 \pm 0.11 (✓)	94.77 \pm 0.09 (✓)	98.98 \pm 0.03 (✓)
	LA	67.580 \pm 2.637	80.527 \pm 3.142	105.831 \pm 4.130	97.63 \pm 0.09 (✗)	98.15 \pm 0.08 (✗)	98.76 \pm 0.06 (✗)
	LA(LL)	66.025 \pm 4.620	78.673 \pm 5.506	103.394 \pm 7.235	96.90 \pm 0.19 (✗)	97.45 \pm 0.17 (✗)	98.10 \pm 0.14 (✗)
	SCP	20.771 \pm 2.452	44.192 \pm 2.799	178.890 \pm 4.596	90.00 \pm 0.10 (✓)	95.03 \pm 0.08 (✓)	99.05 \pm 0.04 (✓)
	SCP-GN	20.460 \pm 2.463	39.712 \pm 3.148	125.078 \pm 8.075	90.01 \pm 0.10 (✓)	94.94 \pm 0.09 (✓)	99.04 \pm 0.03 (✓)
	SCP-GN(LL)	20.735 \pm 2.454	44.237 \pm 2.780	177.595 \pm 4.404	90.02 \pm 0.11 (✓)	95.03 \pm 0.08 (✓)	99.04 \pm 0.03 (✓)
facebook_2 $N=81,311$ $I=53$	ACP-GN	17.986 \pm 0.480	41.063 \pm 0.770	199.331 \pm 10.821	90.36 \pm 0.16 (✓)	95.56 \pm 0.17 (✗)	99.37 \pm 0.08 (✗)
	ACP-GN(LL)	17.305 \pm 0.514	36.951 \pm 1.220	134.876 \pm 6.983	89.54 \pm 0.12 (✗)	94.53 \pm 0.09 (✗)	98.63 \pm 0.09 (✗)
	ACP-GN (split + refine)	29.006 \pm 4.472	54.099 \pm 4.576	172.252 \pm 6.758	90.06 \pm 0.10 (✓)	95.20 \pm 0.08 (✓)	99.14 \pm 0.04 (✓)
	ACP-GN(LL) (split + refine)	19.198 \pm 0.585	42.501 \pm 0.841	175.527 \pm 2.964	90.02 \pm 0.11 (✓)	95.12 \pm 0.08 (✓)	99.06 \pm 0.03 (✓)
	LA	66.088 \pm 2.760	78.749 \pm 3.289	103.493 \pm 4.322	97.47 \pm 0.12 (✗)	98.01 \pm 0.09 (✗)	98.65 \pm 0.06 (✗)
	LA(LL)	64.783 \pm 4.014	77.194 \pm 4.783	101.450 \pm 6.286	97.02 \pm 0.15 (✗)	97.56 \pm 0.13 (✗)	98.23 \pm 0.09 (✗)
	SCP	16.387 \pm 0.208	35.387 \pm 0.462	152.706 \pm 1.591	89.97 \pm 0.07 (✓)	95.00 \pm 0.06 (✓)	99.06 \pm 0.03 (✓)
	SCP-GN	16.287 \pm 0.202	33.655 \pm 0.563	118.489 \pm 4.305	89.99 \pm 0.07 (✓)	94.98 \pm 0.06 (✓)	99.00 \pm 0.03 (✓)
	SCP-GN(LL)	16.323 \pm 0.199	35.192 \pm 0.446	151.199 \pm 1.862	90.00 \pm 0.07 (✓)	94.99 \pm 0.06 (✓)	99.05 \pm 0.03 (✓)
	ACP-GN	18.396 \pm 0.546	40.088 \pm 1.091	166.792 \pm 5.632	90.47 \pm 0.11 (✗)	95.45 \pm 0.08 (✗)	99.35 \pm 0.04 (✗)
	ACP-GN(LL)	17.671 \pm 0.815	37.040 \pm 1.617	127.615 \pm 5.783	89.94 \pm 0.09 (✓)	94.74 \pm 0.06 (✗)	98.73 \pm 0.05 (✗)
	ACP-GN (split + refine)	21.469 \pm 0.906	42.184 \pm 0.788	152.460 \pm 2.499	90.14 \pm 0.08 (✓)	95.10 \pm 0.06 (✓)	99.13 \pm 0.02 (✗)
	ACP-GN(LL) (split + refine)	16.895 \pm 0.238	36.727 \pm 0.477	149.723 \pm 2.197	89.98 \pm 0.07 (✓)	95.03 \pm 0.06 (✓)	99.03 \pm 0.02 (✗)