

Master Thesis

**Uncertainty Calibration with Online Conformal
Prediction in Neural Architecture Search:
An Evaluation under the BANANAS Framework**

Cheng Chen
(matriculation number 1662473)

July 31, 2025

Submitted to
Data and Web Science Group
Prof. Dr. Margret Keuper
University of Mannheim

Abstract

Some contents

Contents

Abstract	ii
1. Introduction	1
1.1. Related Work	1
1.2. Contributions and Limitations	1
1.3. Outline	1
2. Background	2
2.1. Neural Architecture Search	2
2.1.1. Overview	2
2.1.2. Bayesian Optimization and BANANAS	5
2.2. Uncertainty Quantification Methods	8
2.3. Conformal Prediction	10
2.3.1. Theoretical Background	10
2.3.2. Full Conformal Prediction	13
2.3.3. Extensions of Conformal Prediction	14
3. Methodology	17
3.1. The BANANAS-CP Framework	17
3.2. Uncertainty Calibration Algorithms	19
3.2.1. Split Conformal Prediction	19
3.2.2. Conformal Prediction with Cross-validation	22
3.2.3. Conformal Prediction with Bootstrapping	23
3.3. Distribution Estimation	25
3.4. Acquisition Function and Search Strategy	29
3.4.1. Acquisition Functions	29
3.4.2. Acquisition Optimization Strategy	29
4. Experiment Design	31
4.1. Dataset	31
4.2. Setups and Implementation	33
5. Results	35
5.1. Baseline	35
6. Conclusion	36
6.1. Discussion	36
6.2. Limitations and Future Work	36

Contents

Bibliography	37
Acronyms	46
A. Program Code and Data Resources	47
B. Additional Experimental Results	48
Ehrenwörtliche Erklärung	49

1. Introduction

BANANAS with Conformal Prediction (BANANAS-CP)

Bayesian Optimization with Neural Architectures for Neural Architecture Search (BANANAS)

Neural Architecture Search (NAS)

Conformal Prediction (CP) Split Conformal Prediction (SCP) Conformal Prediction with Cross-validation (CrossVal-CP) Conformal Prediction with Bootstrapping (Bootstrap-CP)

1.1. Related Work

1.2. Contributions and Limitations

1.3. Outline

Having gained an overview of the research question and the background, the remainder of this thesis is organized as follows. First, Chapter 2 reviews the related works on neural architecture search, uncertainty quantification, and in particular, conformal prediction. In Chapter 3, after proposing a novel framework to incorporate uncertainty calibration into the architecture search process in Section 3.1, we describe its methodological steps in more detail. In Section 3.2, we identify different types of conformal prediction algorithms that are applicable for NAS, and consider the use of the underlying surrogate models. In Section 3.3 and Section 3.4, we further examine how the calibrated predictions can be incorporated into a Bayesian optimization process. In Chapter 4, we present an overview of the general experiment setups and the strategy for progressively tuning configurations, along with a description of the benchmark dataset used for research. In Chapter 5, we present the experimental results and compare the performance of the algorithms with state-of-the-art methods. Finally, Chapter 6 and Chapter 7 conclude this work and discuss potential future directions.

2. Background

This chapter offers the technical background related to the research question of this work. We start by providing a comprehensive overview of NAS and introduce the three dimensions that characterize a NAS algorithm, followed by an anatomy of the high-performing search algorithm BANANAS. Then, we review the existing uncertainty quantification techniques, with a focus on CP algorithms, particularly those related to the novel framework we propose in Chapter 3.

2.1. Neural Architecture Search

2.1.1. Overview

In the recent decades, deep learning has achieved remarkable success in a variety of areas, including computer vision, natural language understanding, and machine translation. This success is partly attributed to the meticulously hand-crafted neural network architectures. With the rising demand for efficient architecture engineering in complex domains, NAS has emerged as a technique for automating the design of neural architectures for specific tasks.

NAS has been a rapidly progressing research domain in the past years. Since the seminal work that achieves competitive performance on CIFAR-10 [60], numerous NAS algorithms built on different techniques have been proposed. In general, NAS algorithms can be characterized by three key dimensions: search space, search strategy, and performance evaluation strategy [15, 54, 55]. Figure 2.1 illustrates a typical architecture search process.

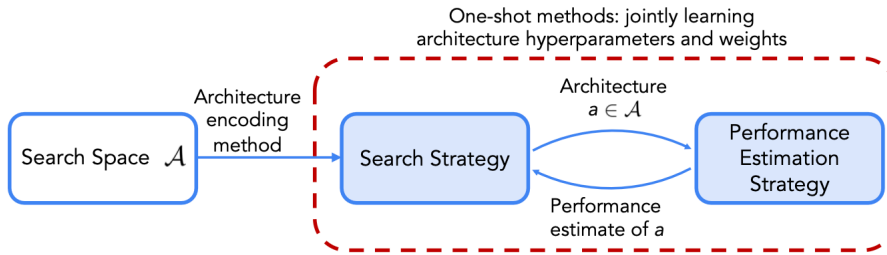


Figure 2.1: Overview of an architecture search process. The search strategy iteratively selects architectures from a predefined search space \mathcal{A} . The performance estimation strategy evaluates the model performance on the target dataset and returns the performance to the search strategy.

2. Background

Next, we provide definitions of the terms and review the research progress of each domain.

Search Space A search space defines a set of architectures that the search algorithm is allowed to select. The search space is often the first step when setting up NAS and perhaps is also the most essential step, because the design of the search space represents an important trade-off between human bias and efficiency of search: a smaller search space incorporating more prior human knowledge and involving more manual decisions will enable NAS algorithms to find high-performing architectures more easily, in contrast a larger space with more primitive building blocks provides higher odds of discovering truly novel architectures [54]. Common search spaces range in size from a few thousand to over 10^{20} .

There are four major categories of search spaces in the NAS literature [54]. We start with two types of search spaces that have relatively simple architecture topologies. The macro search spaces [5, 20, 60] encode the entire neural architecture at a high level. Typically, an entire architecture is often represented by a Directed Acyclic Graph (DAG), with nodes defining the operation types and edges representing data flows. Each node is allowed to have distinct structures, such as convolution, pooling. As a result, macro search spaces are highly flexible and possess high representation power. Another type is the chain-structured search spaces. As suggested by the name, chain-structured search spaces consist of neural networks that can be written as a sequence of operation layers. These search spaces often take state-of-the-art manual designs as the backbone. For example, there are several chain-structured search spaces based on the convolutional networks [10] or the transformer architectures [57].

The third group is the cell-based search spaces, which perhaps are the most popular type of search spaces in NAS research. The cell-based search spaces are inspired by the fact that state-of-the-art human-designed architectures often consist of repeated blocks. For instance, the high-performing Transformer [49] contains 6 identical stacked encoder and decoder layers. Thus, instead of searching for the entire network architecture from scratch, [61] propose to only search over relatively small cells, and stack the cells according to a predefined skeleton to form the overall architecture. Building on this idea, [61] proposes the first modern cell-based search space, NASNet, which comprises of two types of cells: the normal cell that preserves the dimensionality and the reduction cell that reduces the spatial dimension, as illustrated in Figure 2.2. Since its emergence, many other cell-based search spaces have been developed. In general, these cell search spaces share a high-level similarity, but differ in the design of the fixed macro structure, the layout and constraints in the cells, and the choices of operations within the cells [14, 30, 58]. The cell-based approach significantly reduces the size and the complexity of the search space. However, it has been criticized for limiting the expressiveness of NAS, potentially hindering the discovery of highly novel architectures [54].

The last main category is the hierarchical search spaces. Different from the aforementioned types of search spaces that mostly have a flat representation, hierarchical search spaces involve designing patterns at different levels, where each higher-level pattern is often represented as a DAG of lower-level patterns [12, 29].

2. Background

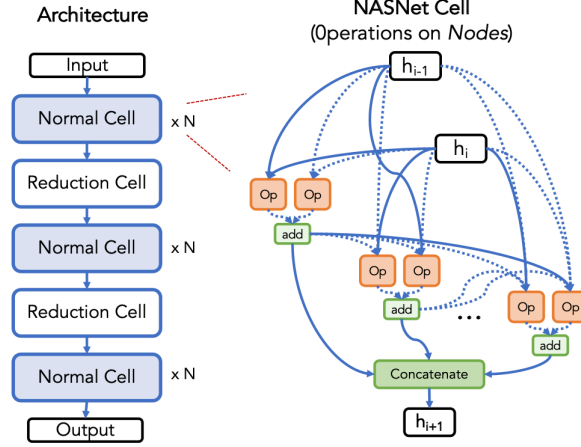


Figure 2.2: Overview of a cell-based search space NasNet. The outer skeleton across cells (left) is fixed, and the operations, represented by nodes, within the cells are searchable (right) [54].

In addition to the architecture topology, another important design accompanying a search space is the architecture encodings, because many NAS algorithms require representations of the architectures to, for example, mutate an architecture or train a predictive model to extrapolate its performance. For search spaces that can be represented by a DAG, adjacency matrix is a commonly used encoding method. In addition, other encoding techniques, including graph-based encoding [38], path-based encoding [53] and conditionally-structured encoding methods tailored for hierarchical search spaces have been proposed. [52] has shown that the effect of the encoding methods varies across different NAS subroutines.

Search Strategy According to [54], there are generally two main categories of search strategies: black-box optimization based techniques and one-shot techniques.

The black-box optimization based techniques largely overlap with another sub-area of AutoML: the hyperparameter tuning. Common techniques for hyperparameter tuning have been proven to be efficient for NAS as well, including reinforcement learning [60, 61], evolutionary algorithms [32, 41], gradient descent [30], and etc. In particular, we take a close look at the search strategies based on Bayesian optimization, since they are closely related to the research question of this work. Specifically, initial Bayesian optimization based approaches typically use the Gaussian Process (GP) as the surrogate model [20]. However, these algorithms often demonstrate under-performance compared to their competitors due to several limitations: 1) search spaces are usually high-dimensional, non-continuous, and graph-like; 2) GPs requires custom distance metrics among architectures, which involves a time-consuming matrix inversion step. Besides, GPs are difficult to scale since the computation complexity grows cubically with the number of observations. To address these challenges, a new framework that using a neural predictor as the surrogate model for Bayesian optimization has been proposed and demonstrated strong

2. Background

performance [31, 47, 53]. We review this framework in details in Section 2.1.2.

The one-shot techniques are introduced to avoid training each architecture from scratch. The key idea is to train a *supernet* that comprises all possible architectures in the search space as subnetworks. Once a supernet is trained, each architecture from the search space can be evaluated by inheriting the weights from the corresponding subnet within the supernet [7, 30].

Performance Evaluation The performance evaluation refers to the process of estimating the performance of architectures. The estimated performance is communicated back to the search algorithm to guide the next search. The simplest performance estimation strategy is to fully train an architecture on the training data and then evaluate its performance on the validation data. However, training each architecture demands substantial computation resources and typically takes several hours or days on a GPU. Consequently, many methods for speeding up the performance evaluation process for architectures have been proposed. One popular line of work is to predict the performance of neural networks before they are fully trained using the zero-cost proxies [33].

In this work, we primarily run experiments on the benchmark dataset NAS-Bench-201 [14], which offers queryable validation and test accuracies for all architectures in the search space and thereby eliminates the need to train neural networks when simulating NAS experiments. Hence, we provide only a brief overview of this aspect and refer the readers to [54] for a comprehensive introduction to the performance evaluation techniques.

2.1.2. Bayesian Optimization and BANANAS

As briefly mentioned in Section 2.1.1, NAS search strategies based on Bayesian optimization with a neural network as the surrogate model have demonstrated strong performance. In particular, [53] identifies five critical components of this framework. After performing a thorough analysis on each component’s effect towards the search performance, [53] develops the BANANAS algorithm based on both theoretical and empirical findings. This method is proven to be efficient and has achieved state-of-the-art performance on popular NAS benchmarks.

In this section, we present a detailed review of the work [53]. We start with the theoretical background and give an introduction into the Bayesian optimization. Next, we walk through the five identified components and provide a summary of the experimental findings.

Bayesian optimization [34] is a sequential decision-making process that seeks to find the global maximum (minimum is the negation of the maximum) of an unknown black-box objective function $f : X \rightarrow R$ over an input space $X \subseteq R^D$. In a Bayesian optimization process, the unknown objective function f is treated as a random function and the prior belief over f is encoded by a surrogate model, usually a Gaussian Process or a Parzen-Tree Estimator [8]. At each iteration, the surrogate model updates the prior with the observations and forms a posterior probabilistic distribution of f . Then, the acquisition function, another key component that trades off exploration and exploitation

2. Background

Algorithm 1 Bayesian Optimization

- Input:** surrogate model \mathcal{M} , acquisition function ϕ , objective function $f(\cdot)$, number of iterations T .
- 1: Initialize the set of observations: $D \leftarrow \emptyset$
 - 2: **for** t in $1, \dots, T$ **do**
 1. Fit surrogate model \mathcal{M} to current observations set \mathcal{D}_{t-1} .
 2. Evaluate acquisition function and select the next point for query:
 $x_t = \operatorname{argmax}_{x \in \mathcal{X}} \phi(x, \mathcal{M})$
 3. Query the objective function: $y_t = f(x_t)$
 4. Update the observations set: $\mathcal{D}_t \leftarrow \mathcal{D}_{t-1} \cup \{(x_t, y_t)\}$
 - 3: **end for**
 - 4: **Output:** $x^* = \operatorname{argmax}_{t=1, \dots, T} f(x_t)$
-

in the process, evaluates a set of candidates based on the posterior distribution and picks the data point with the largest acquisition score for next query. Algorithm 1 outlines this procedure.

The acquisition function adopted in the original paper [34] is the Expected Improvements (EI). Other popular alternatives include: Thompson Sampling (TS), Independent Thompson Sampling (ITS), Upper Confident Bound (UCB), and Probability of Improvements (PI). Different acquisition functions typically favor exploration and exploitation differently. Nevertheless, [1] shows that EI is competitive in reaching the optimum value with comparably few iterations.

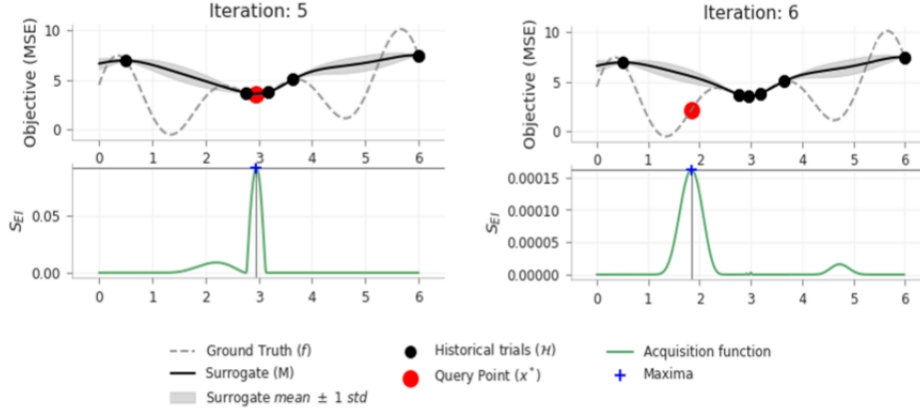


Figure 2.3: Example of Bayesian optimization with Gaussian Process as the surrogate and EI as the acquisition function to explore the minimum of the objective function [1].

Next, we return to the NAS framework "Bayesian optimization + neural predictor". Now it becomes obvious that this framework is essentially an optimization task searching the maximum with a neural network as the surrogate and neural architectures in a search

2. Background

space being the inputs. Specifically, [53] identifies five critical components within the framework, which are listed as follows:

Architecture Encoding This item refers to obtaining a vector representation of the architectures. Previous work often encode architectures using an adjacency matrix based approach, where nodes are assigned with an arbitrary ordering and then binary features for all edges in the DAG are set to form the final representation. Notably, the resulting representation of a specific architecture is not deterministic since the encoding relies on an arbitrary indexing of the nodes.

In contrast, [53] proposes a novel path-based encoding mechanism with optional path truncation. This method simply checks if a path from the input node to the output node, expressed in terms of operations (e.g., input \rightarrow 1×1 conv \rightarrow 3×3 pool \rightarrow output), is present in an architecture. The final encoding is a binary vector indicating which of the possible paths within a cell are present in the architecture. In this way, an architecture is always mapped to the same (though not necessarily unique) path encoding. Experiments show that the path-based encoding substantially increases the performance of neural predictors.

Neural Predictor This item is about choosing an appropriate neural network for surrogate. A set of neural architectures and their corresponding validation accuracies are randomly sampled from the search space for training and comparing different neural predictors. Among all tested neural predictors, which include VAEs, GCNs, and FNNs with either the adjacency matrix or path-based encoding, FNNs with path encoding demonstrates the strongest performance.

Uncertainty Estimation Uncertainty estimates are required to form the probabilistic distribution. In particular, Bayesian neural networks (BNNs) and an ensemble of Feed-forward Neural Networks (FNNs) are investigated. For BNNs, the posterior distribution is inferred over the network weights. For FNNs, the predictive distribution is simply modeled as a Gaussian, parameterized by the mean and standard deviation of the ensemble base learners' predictions. The results show that an ensemble of even only 3 to 5 neural networks in general yields more reliable uncertainty estimates than BNNs.

Acquisition Function In the experiments, five commonly used acquisition functions are examined: TS, ITS, UCB, PI, and EI. Each function is adapted to the Gaussian assumption, thereby requiring only the mean and standard deviation estimates to compute the acquisition scores. The experiments show that overall ITS yields the best performance among all the options, although the marginal outperformance is subtle. The results indicate that the acquisition function does not have as significant impact on the search performance as the other examined components in the framework.

Acquisition Optimization In each iteration of the Bayesian optimization, the goal is to select a candidate from the search space that maximizes the acquisition score. Evaluating the acquisition function for every architecture available in the search space is computationally infeasible, therefore [53] proposes to create a set of 100 to 1000 candidates and

2. Background

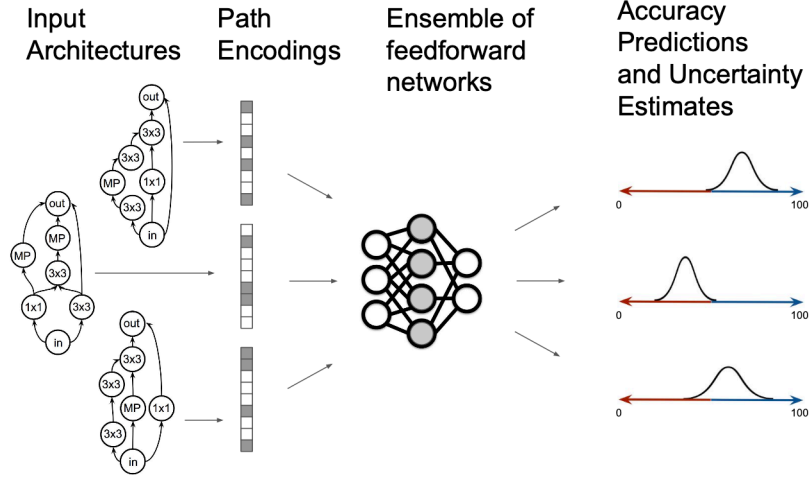


Figure 2.4: Diagram of the BANANAS framework [53].

then choose the architecture with the maximal acquisition score in this set. Specifically, [53] explores various approaches for creating this candidate set. The simplest and most natural way is to draw architectures at random. Consider that architectures close in edit distance to those used for training the surrogate model are likely to have more accurate estimates, an alternative is a mutation-based sampling approach, where the candidate set is created via local search by randomly modifying an operation or an edge of the best-performing architectures that have been evaluated so far. In addition, [53] also examines a hybrid approach that combines random search with mutation-based search. Their experiments show that the mutation-based approach outperforms its competitors and suggest it is better to search locally rather than globally.

Finally, the best components found in the aforementioned analyses are transformed into the BANANAS algorithm, which composed of an ensemble of FNNs using the path encoding, and ITS with a mutation strategy for acquisition (Figure 2.4). Notably, uncertainty is estimated based on the Gaussian assumption.

2.2. Uncertainty Quantification Methods

While BANANAS employs a relatively simple method for uncertainty estimation, there are many commonly used Uncertainty Quantification (UQ) techniques that are generally more theoretically grounded and sophisticated. In this section, we review some of these UQ methods and explain the rationale for choosing CP for the development of our approach. Despite not directly tied to this work, this thesis aspires to offer a comprehensive overview of common UQ techniques applied in deep learning.

UQ is the science of quantitative characterization and estimation of uncertainties in both computational and real world applications. Formally, the source of uncertainty

2. Background

can be categorized into two types:

Aleatory Uncertainty: also known as data uncertainty, refers to uncertainty that arises due to inherent noises or randomness in a system and can not be reduced.

Epistemic Uncertainty: also known as model uncertainty, refers to uncertainty that arises due to lack of knowledge, and can be reduced by better modeling or collecting more data.

A family of methods for quantifying uncertainty directly model the full predictive distribution of the output, i.e., $p(y | x)$. One example is the Gaussian Process, which is typically parameterized by a kernel function and the predictive inference can be obtained in an analytical manner based on the observations. Another important example is the Bayesian Neural Network [21, 36], in which each weight in the network is seen as a random variable rather than a single number. In practice, however, the posterior distributions over weights are intractable due to the integral operation and the high dimensionality. Therefore, the posterior can only be approximated using e.g., variational inference or Monte-Carlo sampling methods [36, 47], both requiring multiple forward passes through the network.

Like the inference stage in BNNs, some methods also rely on aggregating statistics from sub-networks, such as Monte-Carlo Dropout [16] and Deep Ensemble [27]. The idea of Monte-Carlo Dropout is to apply *dropout*, a regularization technique often used for preventing overfitting during the training time by randomly deactivating neurons, also at the inference time. Specifically, the model run multiple times on the same input with stochastically deactivated neurons and thereby get different predictions, which approximate the distribution of the output. Deep Ensemble involves using a network that outputs two values in the final layer, corresponding to the predicted mean $\mu(x)$ and the variance $\sigma^2(x)$, respectively. This neural network is typically trained by minimizing a custom loss function in which $\mu(x)$ and $\sigma^2(x)$ have opposing effects. To avoid overfitting and lower estimation variance, multiple instances of the neural network are initialized with different weights and trained independently. Then, the final estimates of $\mu(x)$ and $\sigma^2(x)$ are obtained by aggregating the predictions from all these sub-networks.

Another popular approach is Quantile Regression (QR) [23]. Instead of modeling the full distribution, QR only models a sequence of discrete quantiles of the output. We present QR with more details in Section 3.2.

However, these UQ techniques face several challenges in practice. Some methods require multiple forward propagations either during the training or inference time, such as BNNs, Deep Ensemble. This is in general computationally intensive and may cause the models difficult to scale and consequently limit their applications for real-time prediction or in an online setting. In addition, the inference quantile of BNNs depends in if the prior of weights is correctly specified, which demands expertise and can potentially introduce human bias. On the other hand, some models are pre-trained and are only accessible via API, making intervening the training process practically infeasible [35]. Moreover, pre-trained models developed using certain datasets may struggle to generalize across different domains or contexts.

2. Background

Post-hoc uncertainty calibration techniques can serve as an effective approach in the restricted scenarios described above. Instead of assessing and measuring uncertainties on the model level, calibration targets on adjusting a model’s predicted probabilities to make the reported uncertainties better aligned with the actual likelihoods.

Platt Scaling [40] probably is the first calibration technique applied in modern machine learning. This method is introduced in the context of Support Vector Machines and is intended for classification problems. For regression problems, one widely used calibration technique is Isotonic Regression [37], a technique of fitting a step-wise line to a sequence of observations such that the fitted line is non-decreasing (or non-increasing) everywhere. The monotonicity of the resulting line allows to preserve the order of predicted quantiles, making this method appealing for calibration. For example, [25] show that accurate uncertainty estimates can be obtained by training a recalibrator on a holdout dataset using isotonic regression.

Recently, Conformal Prediction (CP) [45, 51] is emerged as a new, distribution-free framework for uncertainty quantification. This technique has been widely adopted for both classification and regression problems despite of its relatively short history. Fundamentally, CP serves as a model-agnostic wrapper and can be applied with any arbitrary prediction algorithms. It works by constructing a prediction interval (for regression) or a set of possible values (for classification) that will cover the true value with a predefined probability. The fact that CP is applied in a post-hoc fashion with minimal assumptions enables it to capture both aleatoric and epistemic uncertainties [35]. Moreover, like most other post-hoc calibration methods, CP offers the advantage of being lightweighted with little overhead to implement. We present a thorough introduction into the CP framework in the following section.

2.3. Conformal Prediction

2.3.1. Theoretical Background

We start by introducing the concept of *exchangeable data*, which is an key prerequisite for understanding CP.

Definition (Exchangeability): A sequence of random variables (Z_1, Z_2, \dots, Z_n) is said to be *exchangeable* if its joint probability distribution is invariant under any permutation of the indices. That is, for every permutation π of the set $\{1, 2, \dots, n\}$,

$$P(Z_1, \dots, Z_n) = P(Z_{\pi(1)}, \dots, Z_{\pi(n)})$$

Note that exchangeability allows for dependencies among data points, as long as the joint distribution is invariant under any permutation. Accordingly, it is weaker than the i.i.d. assumption. Then, the goal of CP can be formally defined as:

Definition (Conformal Prediction): Let $B = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be a bag of observed examples from an exchangeable joint distribution $\mathbb{P}_{X,Y}$, with $x_i \subset \mathbb{R}^d$ representing some features and Y the target variable. Given a new unseen input x_{n+1} ,

2. Background

conformal prediction constructs a prediction region $\mathcal{C}_n^\tau(x_{n+1})$ using a *conformity score function* such that

$$\mathbb{P}(y_{n+1} \in \mathcal{C}_n^\tau(x_{n+1})) \geq 1 - \tau \quad (2.1)$$

where $\tau \in (0, 1)$ is the predefined nominal mis-coverage rate. This probability is also known as the **finite-sample validity** property of CP and is taken over the joint distribution of all $n + 1$ samples. Moreover, the prediction regions for different τ should be nested. That means, if $\tau_1 \geq \tau_2$, then $1 - \tau_1$ is a lower confidence level than $1 - \tau_2$ and we have $\mathcal{C}_n^{\tau_1}(x_{n+1}) \subseteq \mathcal{C}_n^{\tau_2}(x_{n+1})$ [45].

Conformity score function is a real-valued function measures how dissimilar the unseen example x_{n+1} is compared to the existing examples in the bag B . Now generalizing this setting to a supervised learning framework, a predictive algorithm f that maps features to the target variable naturally induces a conformity measure. Specifically, we denote \hat{f}_n the predictor trained on the n existing observations in the bag and $d(z_{n+1}, B)$, $z_{n+1} = (x_{n+1}, y_{n+1})$ the conformity function, such that a higher value indicates a greater deviation of the new sample from the existing ones. When using \hat{f}_n to generate a point prediction for x_{n+1} , since the prediction rule is learnt from examples in the bag, intuitively, the distance between the predicted value $\hat{f}_n(x_{n+1})$ and the true value y_{n+1} informs how different z_{n+1} is compared to the rest examples in the bag. With a conformity score function of this kind, the prediction set is constructed by collecting all values that lie within a certain distance from the true value, i.e., all values such that $d(\cdot, B)$ smaller than a threshold.

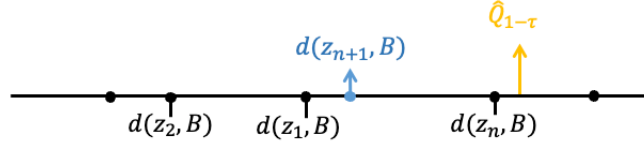


Figure 2.5: Illustration of Symmetrical Data in Conformal Prediction

An important intuition behind selecting the threshold that eventually defines the prediction region is *symmetry*, that is, all data points are treated symmetrically within the CP framework [3, 4]. Suppose $d(z_1, B), d(z_2, B), \dots, d(z_n, B)$ is a sequence of observations in the form of conformity scores, the exchangeability assumption implies that the rank of a new data point $d(z_{n+1}, B)$ is uniformly distributed over the observed values. Figure 2.5 illustrates the concept of symmetry, in which the new data point is equally likely to take any value on the number line. As a result, the $(1 - \tau)$ -th quantile of all the conformity scores is a natural answer to achieve the coverage probability stated in Formula 2.1. Recall that the coverage probability is achieved on the $n + 1$ samples, the

2. Background

quantile value should be adjusted with a finite-sample correction, resulting in a threshold that is exactly equal to

$$\frac{\lceil (1 - \tau)(n + 1) \rceil}{n}\text{-th quantile of conformity scores} \quad (2.2)$$

Finite-sample Validity Note that the coverage guarantee provided by the CP framework is marginal, meaning it holds on average over all data points. This is a significantly weaker requirement than the conditional coverage, which requires the guarantee to hold for each individual input (Figure 2.6). Nevertheless, CP remains a powerful tool for uncertainty quantification, because the finite-sample coverage is valid regardless of the size of the observations bag or the choice made for the underlying prediction algorithm or the conformity score function. Under the distance interpretation presented earlier, a smaller score indicates that the new sample conforms more closely with the existing observations. In such case, the function is also referred as a negatively-oriented score and, to be precise, it is actually a non-conformity measure. A positively-oriented one also works in practice, provided with correct modifications.

Although coverage is guaranteed for any conformity scoring function, this scoring method is actually an important engineering decision and has significant impact on the effectiveness of the constructed prediction set. Imagine a classification problem, if we produce a set that contains all possible classes, then surely the true label is covered by the set. However, such prediction set is neither informational nor actionable, offering little value in application. Naturally, this introduces a second evaluation criterion for CP: while ensuring coverage guarantees, ideally the yielded prediction set should be as compact as possible.

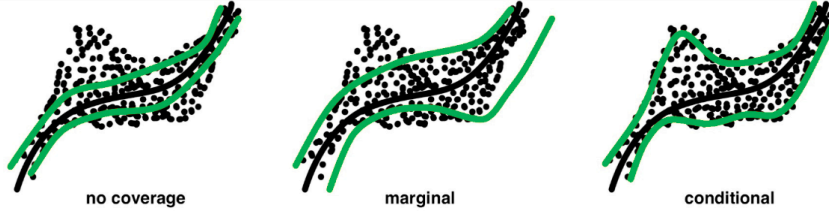


Figure 2.6: Example of Different Coverage Types

A Link to Hypothesis Testing It is noticeable that CP exhibits a couple of structural similarities to hypothesis testing, which is a fundamental tool in statistical inference and perhaps sounds more familiar to most ears. For example, both methods involves producing an interval and relying on a certain threshold to make decisions. In fact, these terms are often used interchangeably in many research work. This motivates us to seek an interpretation of the CP framework through the hypothesis testing perspective, such that we gain a deeper and clearer understanding of CP.

2. Background

Table 2.1: Comparisons between conformal prediction and hypothesis testing

	Conformal Prediction	Hypothesis Testing
Input	nominal mis-coverage level τ	significance level α
Decision rule	compare the conformity score of a data point against the quantile value	compare p -value, which is an empirical measure, to the significance level
Output	prediction interval offering marginal coverage probability $1 - \tau$	confidence interval that contains the true parameter value with probability $(1 - \alpha)$

Both methods take a predefined value ranging from 0 to 1 as input. In the context of CP, this value is often interpreted as mis-coverage rate or error rate, referring to the probability that a prediction region fails to cover the true value. The corresponding concept in hypothesis testing is the significance level, referring to the probability of rejecting the null hypothesis when it is actually true. Table 2.1 offers a summary of the terms used in CP alongside with their parallel in hypothesis testing.

Intuitively, we can think CP as running a sequence of hypothesis tests for a given new point x_{n+1} , with each testing a candidate value for y_{n+1} . The hypotheses are as:

H_0 : The new example (x_{n+1}, y) conforms to the existing observations.

H_1 : The new example (x_{n+1}, y) does not conform to the existing observations.

In this context, the fraction of conformity scores that is larger than that for (x_{n+1}, y) , or equivalently, the complement of the quantile of the conformity score for (x_{n+1}, y) , can be seen as a p -value. If the p -value exceeds the predefined significance level, we fail to reject the null hypothesis, and consequently, this candidate y is included in the prediction set. In fact, this interpretation from the hypothesis testing perspective reflects the core idea of the approach known as full conformal prediction, which we formally introduce in the following section.

2.3.2. Full Conformal Prediction

In fact, conformal prediction is originally designed for a transductive setting, in which the prediction for a test point is generated using the entire dataset available at the inference time [17]. Hence, this transductive CP framework is also known as Full Conformal Prediction (FCP). Typically, new observations are revealed sequentially, which means, after the prediction is generated for a test point, its true value is observed before the next prediction is made.

The core idea for constructing a prediction set for a given test point x_{n+1} is to iterate over all possible values in the target variable space. As performed in hypothesis testing, for each candidate y , a conformity score and its corresponding p -value will be calculated. Specifically, the conformity score for candidate \tilde{y} is computed in a way where a prediction rule $\hat{f}_{\tilde{y}}$ is learnt using all pairs in the bag $B = \{(x_1, y_1), \dots, (x_n, y_n), (x_{n+1}, \tilde{y})\}$, then

2. Background

the conformity score is computed as a distance measure between $\hat{f}_{\tilde{y}}(x_{n+1})$ and \tilde{y} . This process is repeated for all available candidates and the prediction set is constructed by collecting every \tilde{y} with which the null hypothesis in the previous sub-section can not be rejected. Typically, if the target variable space is continuous, the space is first discretized to form a finite grid, then each element in the set can be checked, e.g., using grid search [4].

Moreover, after the true value y_{n+1} is revealed, the underlying bag of observations gets updated by including this data point, and the above procedure needs to start from scratch when generating prediction for each new test point. Suppose the target variable space has a size of K , to generate prediction regions for m instances, the underlying prediction algorithm should fit for $K \cdot m$ times. Clearly, FCP is in general extremely computationally expensive and renders the application of FCP highly unsuitable for training-intensive works, such as with neural networks.

2.3.3. Extensions of Conformal Prediction

Since the transductive version of CP that was first proposed in [17], several variants of CP have been developed with different computational complexities, formal guarantees, and practical applications.

To address the aforementioned inefficient computation problem of FCP, Split Conformal Prediction (SCP), also known as Inductive Conformal Prediction (ICP), was first introduced in [39] by replacing the transductive inference with inductive inference. It aims to learn a general prediction rule about the data using the observed records. Then, this rule can be applied directly to obtain predictions when new data arrives in sequence, without re-using the training data and retraining the model repeatedly. The main concept involves splitting the data into two non-overlapping subsets, designated for training and calibration, respectively. A predictive model is fit exclusively on the training set, then (non-)conformity scores are computed using the calibration set to obtain the quantile value that determines the width of the prediction interval. Due to its simplicity and computational efficiency, SCP is one of the most commonly used techniques in the CP family. We introduce the methodological steps of SCP in details with pseudo-code in Section 3.2.1.

Apart from high computational cost, CP still faces several challenges in practice:

Distribution/Covariate Shift: The finite-sample validity offered by CP depends on exchangeable data. However, the key assumption of exchangeability is often violated in real-world applications, in which the underlying data generation process might vary over time. For example, in finance market behavior can shift drastically in response to major world events.

Adaptivity: The conformity score adopted in the original SCP work is based on the absolute residual $|y - \hat{y}|$, which leads to a prediction interval with fixed width and does not adapt to the intrinsic complexity of the specific test example.

Various CP algorithms have been developed to address limitations and broaden application domains. These extended algorithms can be mainly categorized into three

2. Background

types in accordance with the structural components in a CP framework.

One line of work focus on the coverage rate. Instead of using a static predefined mis-coverage rate τ for computing the quantile of conformity scores, [18] develops a method called *ACI* where the applied coverage level is dynamically adapted based on empirical mis-coverage frequency of previous prediction sets. This approach is applied in an online setting. At each time t , the applied coverage rate τ_t is derived by decreasing (resp. increasing) the current value if the prediction sets were historically under-covering (resp. over-covering). The amplitude of rate adjustment is jointly controlled by the pre-specified learning rate and the empirical mis-coverage level. As an example, if the preset target coverage rate is 90%, historical coverage rate suggests the applied coverage should set to 93%, then this adjusted value is used for computing quantile. In this setting, the empirical mis-coverage frequency of previous examples serves as a signal of distribution shift or violation of exchangeability. [18] show that *ACI* is capable of forming prediction sets that are robust to changes in the underlying data distribution. [59] extends *ACI* to time-series data and introduce *AgACI*, which is a parameter-free variant of *ACI* that uses online expert aggregation to adaptively combine multiple *ACI* experts. In contrast to *ACI*, where the learning rate should be carefully chosen in advance, *AgACI* leverages a number of k experts working with different learning rates, and automatically learns the optimal learning rate by aggregating across experts such that each expert’s contribution is proportional to their corresponding performances over previous iterations. [59] show *AgACI* demonstrates strong performance and produces tighter, well-calibrated prediction intervals consistently.

Another line of work strives to find suitable and efficient conformity score functions tailored for specific tasks. For example, [28] introduces a simple score for regression that accounts for heteroskedasticity, thus offering local adaptivity. [43] proposes a novel conformity score crafted for multi-label classification tasks. While providing marginal coverage, this score also demonstrates full adaptiveness to data complexities and enhances the approximated conditional coverage. In parallel, [42] develops a conformity score that involves both upper and lower bounds corresponding to a given coverage rate, enabling to combine the strengths of quantile regression and conformal prediction. Furthermore,

The underlying prediction algorithm also plays an important role in CP, since the trained predictor serves as an approximation of the real data generation process and determines the base from which intervals are constructed. SCP fits the predictor only once, however potentially at a cost of training accuracy and but statistical efficiency, since both training set and calibration set only see a subset of samples. One way to address this challenge is combining conformal prediction with techniques like cross-validation[50], Jackknife+ [6] and bootstrapping [22], such that all data points are used for both training and calibration, with only a limited number of model fits.

In addition, there are works beyond the above three categories. Some works boost prediction reliability through enhancing the datasets. [48] proposes a CP algorithm targeted on a special classification case where the ground-truth is ambiguous and consequently cause a biased distribution in manually-annotated labels. This method ap-

2. Background

proximates the true distribution of labels by resampling data points using Monte-Carlo sampling. [46] proposes to integrate a test-time data augmentation into CP to reduce prediction set size and improve stability. Moreover, Conformal Risk Control [2] also emerges as an important extension of CP that not only provides distribution-free coverage guarantees but also explicitly controls risk. It is achieved by replacing the mis-coverage probability in normal CP with the expected value of a custom loss function. The prediction set is then bounding to the predefined loss tolerance instead of the mis-coverage rate. Several works have adopted this risk-control approach such as [35].

3. Methodology

Despite of its provable strength, BANANAS assumes a Gaussian distribution for measuring uncertainty. However, this assumption does not necessarily hold in real world. To mitigate the potential limitations caused by inaccurate uncertainty estimates, this work proposes a new framework that integrates a conformal prediction-based uncertainty calibration process into BANANAS in an online setting.

An algorithm outlining the overall procedure of BANANAS-CP is presented in Section 3.1, followed by detailed descriptions of each methodological step. Section 3.2 presents different conformal predictions algorithms to be explored. Next in Section 3.3, methods for the estimation and evaluation of the conditional distribution of each candidate architecture are discussed. Finally, in Section 3.4 we introduce how the calibrated distribution can be combined with different acquisition functions and acquisition search strategies in the Bayesian optimization process.

3.1. The BANANAS-CP Framework

Refer to Section 2.1.2 for a detailed introduction of the original BANANAS algorithm. In this section, we emphasis the key ideas of the uncertainty calibration mechanism, which corresponds to Step 1 to 6 of the inner iteration in Algorithm 2.

Bayesian optimization is a form of sequential decision-making task. In the applications of neural architecture search, the typical goal is to find the architecture that has the best evaluation performance on a fixed dataset under a given search budget. At each iteration t , a surrogate model is trained on all architectures evaluated at step $\{0, 1, 2, \dots, t-1\}$ and their associated validation accuracies, to predict the scores of unseen architectures for the next search.

In the standard BANANAS setting, the surrogate model is an ensemble of m feed-forward neural networks, typically $m = 5$. At iteration t , a set of candidate architectures is sampled, and a conditional Gaussian distribution is estimated for each candidate based on the ensemble predictions, as expressed below:

$$\hat{f}(a) \sim \mathcal{N} \left(\frac{1}{m} \sum_{i=1}^m f_i(a), \sqrt{\frac{1}{m} \sum_{i=1}^m \left(f_i(a) - \frac{1}{m} \sum_{j=1}^m f_j(a) \right)^2} \right) \quad (3.1)$$

where a denotes an architecture sampled from the search space, and $f_i(a)$ is the predicted accuracy from the i -th base learner of the ensemble for architecture a .

3. Methodology

Algorithm 2 The BANANAS-CP Framework

Input - NAS parameters: search space \mathcal{A} , evaluation dataset \mathcal{D} , exploration budget T , the number of initially sampled architectures t_0 , acquisition function ϕ , surrogate model \mathcal{M} that approximates the true objective function, function $f(\cdot)$ returning validation error of an architecture after training.

Input - Calibration parameters: a function $C(\cdot)$ to create calibration set, a non-conformity score function $s(\cdot)$, and an array of desired quantile levels q .

- 1: Draw t_0 architectures $\{a_1, a_2, \dots, a_{t_0}\}$ uniformly at random from \mathcal{A} and train each individual architecture on \mathcal{D} .
 - 2: $\mathcal{A}_{t_0} \leftarrow \{a_1, a_2, \dots, a_{t_0}\}$,
 - 3: **for** t in $t_0 + 1, \dots, T$ **do**
 1. Apply $C(\cdot)$ and split all evaluated architectures into two disjoint datasets; use them as a training set $\mathcal{A}_{t,train}$, and a calibration set $\mathcal{A}_{t,cal}$.
 2. Train the surrogate model \mathcal{M}_t on $\{a, f(a)\}, a \in \mathcal{A}_{t,train}$ using the path encoding to represent each architecture.
 3. Compute the conformity scores s on $\mathcal{A}_{t,cal}$.
 4. Generate a set of candidate architectures from \mathcal{A} .
 5. **for** each a_i in candidates **do**
 - a) Estimate the value for each quantile level q_i in q and calibrate using conformity scores computed in the previous step, with q_i implying a mis-coverage rate $2q_i$ or $2(1 - q_i)$ for conformal prediction.
 - b) Fit a distribution F_i based on the estimated quantile values.
 - c) Compute the acquisition score $\phi(a_i)$.
 6. **end for**
 7. Denote a_t as the candidate architecture with maximum $\phi(a)$; evaluate $f(a_t)$.
 8. $\mathcal{A}_t \leftarrow \mathcal{A}_{t-1} \cup \{a_t\}$
 - 6: **end for**
 - 7: **Output:** $a^* = \operatorname{argmax}_{t=1, \dots, T} f(a_t)$
-

3. Methodology

In the BANANAS-CP framework, a key distinction is that all architectures evaluated at step $\{0, 1, 2, \dots, t-1\}$ are divided disjointly into a training set and a calibration set. Then, the surrogate model is trained exclusively using samples in the training set, while the calibration set is used to compute conformity scores for quantile calibration. In practice, at each iteration t , the surrogate model estimates a conditional distribution \hat{F} for an unseen architecture over its validation accuracy on the target dataset, either based on a specific distribution assumption or a probabilistically-interpretable modeling approach, e.g. quantile regression. Following the definition in [13, 25], calibration means that for any quantile level $p \in [0, 1]$, the empirical fraction of data-points below the p -th percentile of the predicted distribution \hat{F} should converge to p as the sample size goes to infinity. For example, if $p = 80\%$, then the 80th percentile of \hat{F} is set to the threshold value such that 80% of previously evaluated architectures fall below, thereby aligning with the empirical coverage. In an online setting, the objective of the calibration process can be defined as:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{I}\{y_t \leq Q_t(p)\} \rightarrow p \quad \text{for all } p \in [0, 1] \quad (3.2)$$

as $t \rightarrow \infty$, where \mathbb{I} is the indicator function and $Q_t(p)$ represents the distribution \hat{F} in the format of quantile function [13, 25].

Next, as in the standard Bayesian optimization process, the acquisition function picks the architecture for the next evaluation based on the conditional distribution of all sampled candidates.

3.2. Uncertainty Calibration Algorithms

As reviewed in Section 2.3, numerous conformal prediction algorithms have been proposed in recent research. This work identifies several approaches applicable in NAS for building a calibration set and computing conformity scores. This section provides an overview of these splitting strategies, as well as the conformity scoring functions that are commonly used for regression problems.

3.2.1. Split Conformal Prediction

To begin, a natural choice for a baseline calibration strategy is the SCP. In this section, we start by introducing the standard SCP procedure, then proceed with the adaptations required to incorporate it into the BANANAS-CP framework.

Implementation steps of SCP are summarized in Algorithm 3. Imagine a regression task where the non-conformity level is measured by the absolute residual, i.e. $|y_i - \hat{y}(x_i)|$. In this case, the algorithm produces a prediction interval for the test point with a width of $[\hat{y}_{test} - \hat{q}, \hat{y}_{test} + \hat{q}]$, where \hat{q} is the conformity threshold as defined in line 6 in.

In this work, we explore SCP in combination with different prediction algorithms. First, we follow the settings in BANANAS and use an ensemble of five FNNs as the

3. Methodology

Algorithm 3 Split Conformal Prediction

Input: A set of observations $\{(x_i, y_i)\}_{i=1}^n$, a prediction algorithm $h(\cdot)$, a non-conformity measure $s(\cdot)$, nominal mis-coverage rate τ , fraction of data assigned to the training set p_{train} , test data x_{n+1} .

Output: a prediction set $\mathcal{C}_\tau(x_{n+1})$ that covers y_{n+1} with probability $1 - \tau$.

- 1: Allocate at random a proportion of p_{train} of the observations to the training set \mathcal{D}_{train} and use the rest for calibration \mathcal{D}_{cal} .
 - 2: Train the point predictor $h(\cdot)$ on \mathcal{D}_{train} .
 - 3: Initialise a scoring set $S = \emptyset$
 - 4: **for** (x_i, y_i) in \mathcal{D}_{cal} **do**
 $S \leftarrow S \cup \{s(h(x_i), y_i)\}$
 - 5: **end for**
 - 6: Return $\mathcal{C}_\tau(x_{n+1}) \leftarrow \{y \mid s(h(x_{n+1}), y) \leq q\}$, where q is the $\lceil (1 - \tau)(n_s + 1) \rceil$ -th smallest value of S , with $n_s = |S|$.
-

underlying surrogate model. In this case, note that the bounds of the prediction set as identified in Algorithm 3 should not be simply interpreted as the quantile values of a distribution, since the prediction algorithm does not directly model the τ -quantile of the variable Y , i.e., $Q_Y(\tau) = F_Y^{-1}(\tau) = \inf \{y : F_Y(y) \geq \tau\}$, with $\tau \in [0, 1]$ denoting a quantile level and F_Y its cumulative distribution function. Thus, the ensemble predictor must be used in conjunction with a valid distribution assumption to obtain valid quantile values. Motivated by the goal of achieving a completely distribution-agnostic solution, we next replace the ensemble model with a quantile regressor that directly models the quantiles of a distribution. In the remainder of this section, we discuss the configurations designated for each prediction algorithm.

Ensemble Predictor Following the settings in the original BANANAS, an ensemble by default consists of five neural networks, where each neural network is a fully-connected multi-layer perceptron with 20 layers of width 20. The neural networks are trained by minimizing the mean absolute error (MAE), using the Adam optimizer with a learning rate of 0.01. In parallel to BANANAS, we assume that the validation accuracy of each unseen candidate architecture a follows a Gaussian distribution, which is parameterized by the predictive mean ($\hat{\mu}$) and standard deviation ($\hat{\sigma}$) provided by the ensemble model, as demonstrated in equation 3.1. For a specific significance level α (suppose $\alpha < 0.5$), the central quantile interval can be written as:

$$\left[\hat{\mu} - \Phi_{1-\alpha/2}^{-1} \cdot \hat{\sigma}, \hat{\mu} + \Phi_{1-\alpha/2}^{-1} \cdot \hat{\sigma} \right] \quad (3.3)$$

where $\Phi_{1-\alpha/2}^{-1}$ denotes the $(1 - \frac{\alpha}{2})$ -th quantile of the standard normal distribution.

Now, take a closer look at the formula 3.3 and recall the example based on the absolute residuals, which is presented earlier in this section. We observe that the confidence

3. Methodology

interval under the Gaussian assumption takes a close form to the prediction interval produced by CP when the conformity scoring function is exactly chosen as:

$$s(\cdot) = \frac{|y_i - \hat{y}(x_i)|}{\hat{\sigma}(x_i)} \quad (3.4)$$

Hence, the bounds of the CP-derived prediction interval can be *approximately* interpreted as empirically calibrated quantile estimates under the Gaussian assumption, provided that the conformity scoring function is chosen appropriately. Note that the absolute residual can be seen as a special case of equation 3.4 as well, where the empirical standard deviation estimate is disregarded and fixed at one. In fact, this scaled absolute residual (equation 3.4) is a popular choice for measuring conformity in practice. Ideally, we would like the CP-derived prediction interval also demonstrates local adaptivity, i.e., the prediction interval should have a larger width if the prediction task is difficult and smaller otherwise. The scaled absolute residual accounts for heteroskedasticity and is able to adjust the width of the prediction band by multiplying the standard deviation estimate. In contrast, the band produced with a pure residual score has constant-width everywhere regardless of the input, which limits its effectiveness in application. Therefore, in this work, we use the scaled absolute residual as the conformity scoring function for ensemble predictors, unless otherwise specified.

Quantile Regressor We now explain how a quantile regressor can be leveraged to build a probabilistic surrogate for Bayesian optimization. We follow the methods established previously in [42, 44].

We start with a brief introduction into the quantile regression [23]. Suppose $(x, y) \sim F$ denote data drawn from a joint distribution that is characterized by its cumulative distribution function F , the aim of the conditional quantile regression is to estimate a given quantile of the conditional distribution of Y given $X = x$. The conditional quantile function for α -quantile is:

$$Q(\alpha) = \inf \{y \in \mathbb{R} : \mathbb{P}(Y \leq y \mid X) \geq \alpha\} \quad (3.5)$$

and can be estimated by minimizing the Pinball loss on the training data [23]:

$$\ell_\alpha(y, \hat{y}) = \begin{cases} \alpha(y - \hat{y}), & \text{if } y \geq \hat{y} \\ (1 - \alpha)(\hat{y} - y), & \text{otherwise} \end{cases} \quad (3.6)$$

where \hat{y} is the predicted quantile value. As illustrated in Figure 3.1, the Pinball loss is asymmetric and the intuition behind is that under-estimate and over-estimate receive different penalties across quantiles. For instance, if $\alpha = 0.9$, then we would expect that empirically 90% of observations should fall below the prediction. In this case, the loss function places a higher penalty for underestimate.

3. Methodology

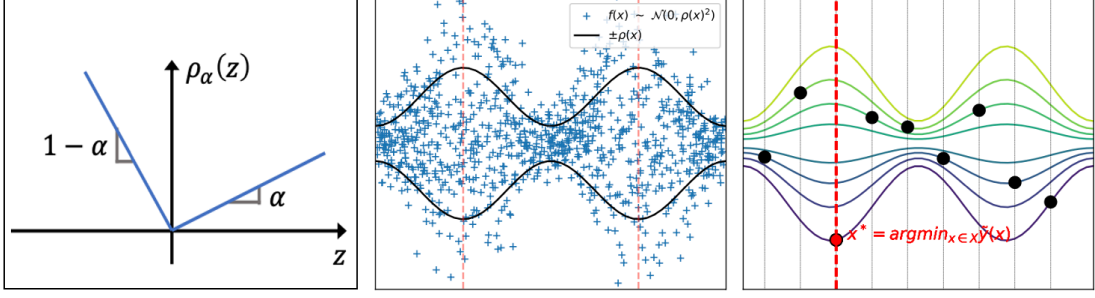


Figure 3.1: Visualization of the Pinball loss function, where $z = y - \hat{y}$ [42] (left); Samples from a synthetic heteroskedastic function (middle) and the sampling procedure based on $|q| = 8$ predicted quantiles [44] (right).

Quantile regression in the BANANAS-CP framework is implemented by training a dedicated neural network for each quantile level q_i in the array q as defined in Algorithm 2 using the corresponding Pinball loss $\ell_{q_i}(y, \hat{y})$.

While quantile regression can model the shape of any continuous distribution given enough data, the predictions are not guaranteed to be well calibrated in practice. In fact, it is not uncommon that quantile regression generates non-monotonic predictions, a phenomenon referred as quantile crossing. To address this issue, we apply a post-hoc calibration upon the predicted quantiles using the Conformal Quantile Regression (CQR) from [42]. This method consists of a novel conformity score tailored for quantile estimation and the key idea of calibration is to apply quantile-aware offsets, which are computed on the calibration set, on the original predicted quantiles.

A close work is [44] that employs CQR to obtain quantiles with robust coverage during hyperparameter tuning via Bayesian optimization. Specifically, the calibrated quantiles are used to select the candidate for the next search, where a set of candidates is first sampled uniformly at random, and then for each of those candidates a random quantile is simply picked and is treated as the acquisition score (Figure 3.1). We follow their notation and interpretation in defining the conformity score for a quantile surrogate:

$$E_i = \max \{ \hat{q}_{\alpha_j}(x_i) - y_i, y_i - \hat{q}_{1-\alpha_j}(x_i) \} \quad (3.7)$$

where $\hat{q}_\alpha(x_i)$ denotes the predicted α -quantile at x_i . Note that the sign of the score is positive when the target y_i is outside of the interval and negative when the target falls inside the predicted interval. This allows the conformity score to account for both overcoverage and undercoverage cases. In addition, the score amplitude always measures the distance to the closer quantile between $\hat{q}_{\alpha_j}(x_i)$ and $\hat{q}_{1-\alpha_j}(x_i)$ [42, 44].

3.2.2. Conformal Prediction with Cross-validation

Solving a NAS problem is usually computationally expensive, as each neural architecture evaluation incurs the cost of fully training and validating the underlying model on

3. Methodology

the target dataset. Motivated by the fact that NAS based on Bayesian optimization is typically allocated with a budget of 100 to 200 epochs, an additional heuristic for constructing the calibration set via cross-validation (hereafter: CrossVal-CP) is employed to avoid reducing the sample size for obtaining a holdout set as performed in SCP.

The CrossVal-CP method is a natural extension of SCP and is first formally presented in [50]. At each step, the evaluated architectures are divided at random into K folds. A dedicated surrogate model is trained on $K - 1$ folds, while the remaining one is used as the calibration set to calculate the conformity scores. This process is repeated for K times over each individual fold. Finally, conformity scores from all calibration folds are combined to form the overall calibration set, on which the quantile is computed to determine the calibration offset. For an unseen data point, the prediction is obtained by aggregating the predictions of the K trained models, e.g., by taking the average of the model outputs. Algorithm 4 summarizes this procedure.

Algorithm 4 Conformal Prediction with Cross-validation

Input: A set of observations $\{(x_i, y_i)\}_{i=1}^n$, number of folds K , a prediction algorithm $h(\cdot)$, a non-conformity measure $s(\cdot)$, nominal mis-coverage rate τ , test data x_{n+1} .

Output: a prediction set $\mathcal{C}_\tau(x_{n+1})$ that covers y_{n+1} with probability $1 - \tau$.

- 1: Initialise a conformity scoring set $S = \emptyset$
 - 2: Split the observations $\{(x_i, y_i)\}_{i=1}^n$ into K folds at random. I_k denotes the index set containing indices of samples in the k -th fold.
 - 3: **for** k in 1, 2, ..., K **do**
 1. Train $\hat{h}_{-k}(\cdot)$ on $\{(x_i, y_i) \mid i \notin I_k\}$
 2. Compute conformity score on the k -th fold $S_k = \{s((\hat{h}_{-k}(x_i), y_i) \mid i \in I_k\}$
 3. $S \leftarrow S \cup S_k$
 - 4: **end for**
 - 5: Predict x_{n+1} : $h(x_{n+1}) \leftarrow \text{aggregate}(\{\hat{h}_{-1}(x_{n+1}), \dots, \hat{h}_{-K}(x_{n+1})\})$
 - 6: Return $\mathcal{C}_\tau(x_{n+1}) \leftarrow \{y \mid s((h(x_{n+1}), y) \leq q\}$, where q is the $\lceil (1 - \tau)(n_s + 1) \rceil$ -th smallest value of S , with $n_s = |S|$.
-

Note that the only distinction between SCP and CrossVal-CP is how the calibration set is constructed. Since it does not place any additional restriction on the choices of the underlying surrogate, CrossVal-CP can be applied in conjunction with either ensemble predictor or quantile regressor in the same way as SCP. See Section 3.2.1 for detailed configurations.

3.2.3. Conformal Prediction with Bootstrapping

Inspired by the fact that BANANAS is built on an ensemble surrogate, we further explore incorporating Jackknife+-after-bootstrap [22], a wrapper for predictive inference

3. Methodology

Algorithm 5 Conformal Prediction with Bootstrapping

Input: A set of observations $\{(x_i, y_i)\}_{i=1}^n$, number of bootstraps B , a prediction algorithm $h(\cdot)$, a non-conformity measure $s(\cdot)$, nominal mis-coverage rate τ , test data x_{n+1} .

Output: a prediction set $\mathcal{C}_\tau(x_{n+1})$ that covers y_{n+1} with probability $1 - \tau$.

- 1: Sample all available data with replacement and create B subsets. I_b denotes the indices of data points included in the b -th sample.
- 2: Train $\hat{h}_b(\cdot)$ on $\{(x_i, y_i) \mid i \in I_b\}$ for b in 1, 2, ..., B
- 3: Initialise a conformity scoring set $S = \emptyset$
- 4: **for** i in 1, 2, ..., n **do**
 1. Initialize an empty for leave-one-out estimates $LOO_i = \emptyset$
 2. For b in 1, 2, ..., B , if $i \notin I_b$: $LOO_i \leftarrow LOO_i \cup \hat{h}_b(x_i)$
3. $S \leftarrow S \cup s(\text{aggregate}(LOO_i), y_i)$
- 8: **end for**
- 9: Predict x_{n+1} : $h(x_{n+1}) \leftarrow \text{aggregate}(\{\hat{h}_1(x_{n+1}), \dots, \hat{h}_B(x_{n+1})\})$
- 10: Return $\mathcal{C}_\tau(x_{n+1}) \leftarrow \{y \mid s(h(x_{n+1}), y) \leq q\}$, where q is the $\lceil (1 - \tau)(n_s + 1) \rceil$ -th smallest value of S , with $n_s = |S|$.

designed specifically for use with ensemble learners, into the calibration step (hereafter: Bootstrap-CP).

In contrast to fitting m neural networks on the same training data with different random weights initializations, as applied in the BANANAS framework, a different technique to build an ensemble model is via bootstrapping. Specifically, the ensemble method starts by creating multiple training datasets by resampling the available data points with replacement. In the next step, multiple models are trained on each of the bootstrapped subsets, and their predictions are aggregated to produce the single final prediction [9]. This technique offers more accurate and stable estimates than a single model and has shown superior performance in application.

Jackknife+ is a type of CP algorithm that is closely related to the leave-one-out (LOO) method [6]. Given a set of observations $\{(x_i, y_i)\}_{i=1}^n$, the idea is to fit an LOO estimator \hat{h}_{-i} using all available data except for the i -th sample, and this process iterates over all individual samples. Then, the predictive interval around the i -th point is obtained by offsetting the prediction from $\hat{h}_{-i}(x_i)$ with the quantile of all LOO conformity scores. Equivalently, Jackknife+ can be viewed as a special case of CrossVal-CP when the number of folds is exactly set to $K = n$.

Jackknife+-after-bootstrap [22] integrates both approaches and provides a cost-efficient wrapper by leveraging only the available bootstrapped sets and their corresponding fitted models, thereby avoiding re-fitting ensembles on each individual bootstrapped sample. [56] extends this method to online setting and proves its efficiency for time-series data. In contrast to the CP algorithms described in earlier sections,

3. Methodology

Bootstrap-CP requires no data splitting because sampling with replacement automatically creates holdout sets. Training the bootstrap ensemble on random subsets from the full data also reduces the chance of overfitting.

Implementation of Bootstrap-CP is shown in Algorithm 5. Notably, if a particular data point appears in all bootstrapped samples, it is then excluded from the computation of conformity scores since it has no associated LOO estimator. In Bootstrap-CP, the absolute residual is mainly used for measuring conformity, due to potentially insufficient LOO outputs for standard deviation estimation, i.e., LOO_i in Algorithm 5 might have fewer than two points. Specifically, Bootstrap-CP is only applied with the ensemble model. This concludes our experiment setups and the BANANAS-CP framework is finally evaluated under five various predictor+CP configurations.

3.3. Distribution Estimation

As described in Section 2.1.2, Bayesian optimization generally relies on a continuous posterior distribution at $X = x$ to obtain the acquisition score. Here, we denote by $F_t(x)$ the Cumulative Distribution Function (CDF) of the posterior distribution of the data point x at step t . In the context of NAS, where the target variable is assumed to be continuous and real-valued, the distribution can be represented by the inverse of its CDF, or quantile function without loss of generality, i.e. $Q_t(x) = F_t^{-1}(x)$.

In the BANANAS-CP framework, as outlined in Algorithm 2, we are able to generate calibrated quantile estimates for a finite set of discrete quantile levels. Intuitively, assuming the quantiles estimates are accurate, increasing the granularity of quantile levels should lead to a more accurate approximation of the underlying continuous distribution. However, it is computationally prohibited to estimate an infinite number of quantiles in practice, especially with significantly limited training data. Therefore, we propose an approach for constructing a continuous distribution from discrete quantile estimates with mild assumption. Specifically, the distribution estimator is defined as:

Definition Let $\{q_i\}_{i=1}^n$ be a quantile of percentile levels such that $0 < q_1 < q_2, \dots, < q_n < 1$, and $\{v_i\}_{i=1}^n$ are the corresponding quantiles, i.e., $F^{-1}(q_i) = v_i$, the empirical CDF of the distributions \hat{F} is constructed by applying linear interpolation between adjacent quantiles. Consequently, the Probability Density Function (PDF) of a specific interval (v_a, v_b) , $a, b \in \{1, \dots, n\}$ and $a < b$ is:

$$PDF(x) = \frac{q_b - q_a}{v_b - v_a}, \quad \text{if } x \in (v_a, v_b)$$

Diagnosis Analysis To assess the validity of this approach, we first perform a diagnosis analysis using synthetic datasets generated by a left-skewed Gaussian distribution, which we believe resembles the true underlying distribution of the validation performances of architectures in a search space. The experiments on the synthetic data are intended to

3. Methodology

reflect the comparisons in a real NAS application, therefore two kinds of distribution estimators are examined: a Gaussian estimator and a linear-interpolation based quantile estimator. The analysis is repeated with different parameterizations, e.g, the number of quantiles, or the size of the samples, etc.

Table 3.1 reports the performance of three distribution estimators on synthetic datasets with various sample sizes. In particular, the linear-interpolation based quantile estimator is evaluated under 10 and 20 quantile levels, and the quantile estimates for interpolation are obtained by taking the empirical quantiles of the sample data. The estimation performance is assessed using the mean, standard deviation and Kullback–Leibler (KL) divergence [26]. For each sample size, we run 50 trials and aggregate the metrics over the trials to reduce the effects of randomness. Results in Table 3.1 indicate that the linear-interpolation based quantile estimator offers a better approximation to an asymmetric distribution than a Normal distribution, and increasing the number of quantile levels leads to improved approximation performance, provided with sufficient data. However, a caveat is that quantile-based estimation tends to produce biased standard deviation estimates, which may lead to undesired effect.

Having considered the behaviors of the various acquisition functions employed in the real BANANAS-CP application (see Section 3.4), we additionally present several visualizations based on a synthetic dataset with 150 samples (Figure 3.2, Figure 3.3) to confirm the shapes of the estimated distributions are not significantly deviated from that of the true underlying distribution, thereby assuring the effectiveness of the acquisition functions.

Evaluation Metrics Within the BANANAS-CP framework, the calibration quality at a specific epoch is measured by the Root Mean Squared Calibration Error (RMSCE) [25]. Suppose \hat{F}_t^{-1} is the CDF of the distribution estimated at the t -th step and y_t is the true value revealed after the estimation, consider a sequence of $\{(\hat{F}_t^{-1}, y_t)\}_{t=1}^T$ that represents a neural architecture search process after T epochs, the calibration error at the T -th epoch is defined as:

$$\text{RMSCE}(\hat{F}_1^{-1}, y_1, \dots, \hat{F}_T^{-1}, y_T) = \sum_{j=1}^m w_j (p_j - \hat{p}_j)^2 \quad (3.8)$$

$$\text{with } \hat{p}_j = \frac{\left| \left\{ y_t \mid \hat{F}_t^{-1}(y_t) \leq p_j, t = 1, 2, \dots, T \right\} \right|}{T}$$

where m represents the number of discrete quantile levels and the scalars w_j are quantile weights. For simplicity, we adopt $w_j = 1, \forall j \in [0, 1]$ in this work. Alternatively, calibration errors of the quantiles can be weighted by the number of observations falling into the corresponding intervals. Note that calibration errors calculated with different numbers of quantiles are not directly comparable. In general, increasing the number of quantiles tends to lead to larger calibration errors.

3. Methodology

Table 3.1: Statistical metrics of distributions estimated by three methods on synthetic datasets with sample size ranging from 50 to 500. For each method and sample size, the mean, standard deviation, and KL divergence are reported to assess the estimation performance.

Sample Size	Estimator	Mean	Standard Deviation	KL Divergence
50	Gaussian	-0.7985	0.6005	0.5727
	Quantile ($ q = 10$)	-0.7481	0.5503	0.1646
	Quantile ($ q = 20$)	-0.7440	0.5536	0.2928
100	Gaussian	-0.7989	0.6145	0.6181
	Quantile ($ q = 10$)	-0.7744	0.5941	0.0953
	Quantile ($ q = 20$)	-0.7654	0.5798	0.1297
150	Gaussian	-0.7986	0.6100	0.5977
	Quantile ($ q = 10$)	-0.7860	0.6133	0.0811
	Quantile ($ q = 20$)	-0.7791	0.5950	0.0948
200	Gaussian	-0.7969	0.6162	0.6227
	Quantile ($ q = 10$)	-0.7948	0.6347	0.0625
	Quantile ($ q = 20$)	-0.7833	0.6050	0.0676
250	Gaussian	-0.7921	0.6127	0.6166
	Quantile ($ q = 10$)	-0.7936	0.6384	0.0526
	Quantile ($ q = 20$)	-0.7815	0.6084	0.0576
300	Gaussian	-0.7940	0.6140	0.6168
	Quantile ($ q = 10$)	-0.8002	0.6521	0.0540
	Quantile ($ q = 20$)	-0.7876	0.6153	0.0482
350	Gaussian	-0.7919	0.6131	0.6176
	Quantile ($ q = 10$)	-0.8002	0.6570	0.0507
	Quantile ($ q = 20$)	-0.7866	0.6188	0.0432
400	Gaussian	-0.7919	0.6131	0.6174
	Quantile ($ q = 10$)	-0.8041	0.6650	0.0491
	Quantile ($ q = 20$)	-0.7871	0.6215	0.0388
450	Gaussian	-0.7908	0.6114	0.6121
	Quantile ($ q = 10$)	-0.8054	0.6668	0.0468
	Quantile ($ q = 20$)	-0.7888	0.6240	0.0357
500	Gaussian	-0.7922	0.6081	0.5956
	Quantile ($ q = 10$)	-0.8070	0.6657	0.0459
	Quantile ($ q = 20$)	-0.7912	0.6219	0.0328
-	Ground Truth	-0.7939	0.6080	0.0000

3. Methodology

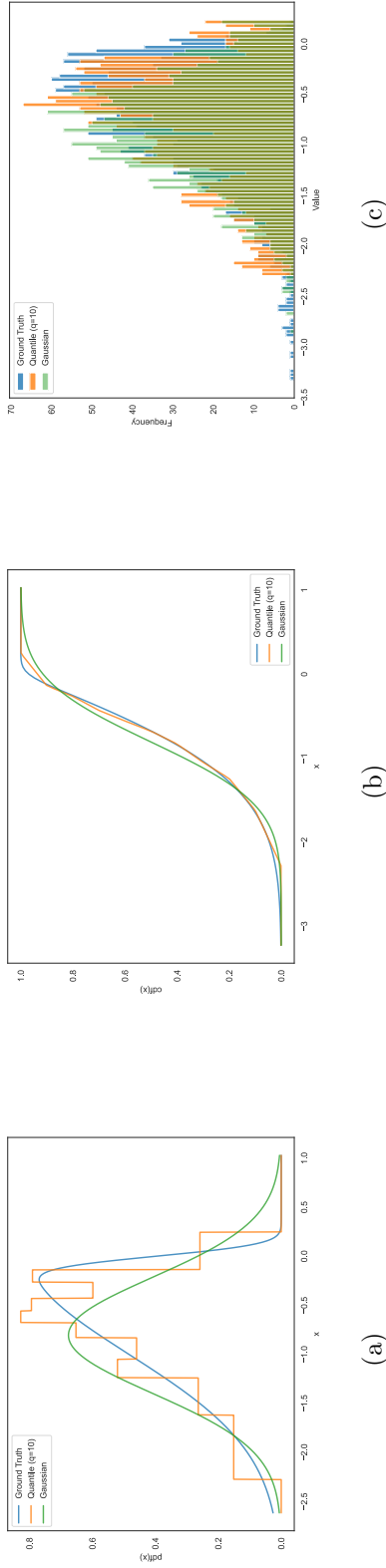


Figure 3.2: Visualization of the distributions estimated from 10 quantile levels using a synthetic dataset with 150 samples. The figure includes the estimated PDF (a), CDF (b), and the histogram of 2000 samples drawn from the estimated distribution (c).

28

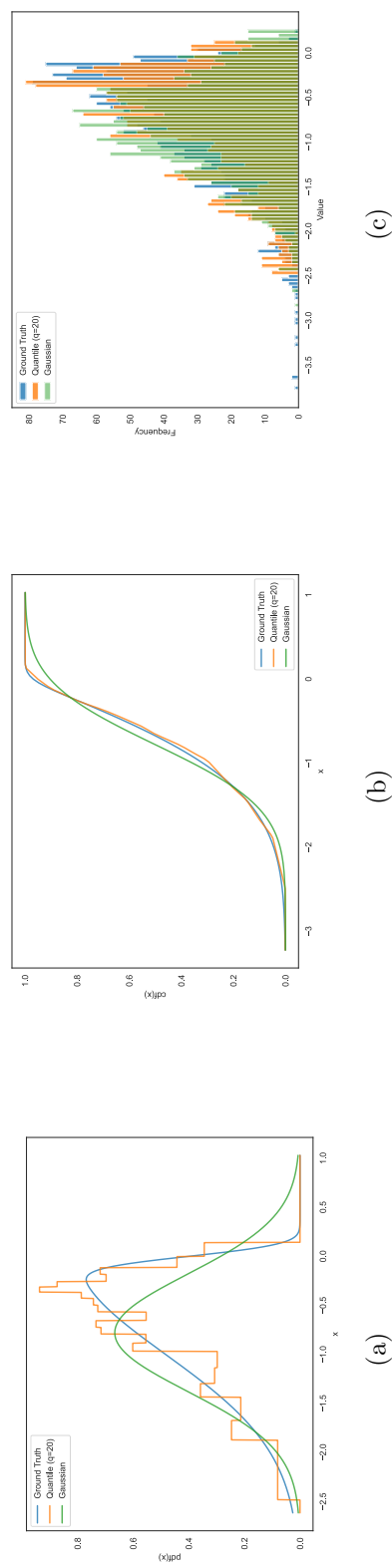


Figure 3.3: Visualization of the distributions estimated from 20 quantile levels using a synthetic dataset with 150 samples. The figure includes the estimated PDF (a), CDF (b), and the histogram of 2000 samples drawn from the estimated distribution (c).

3.4. Acquisition Function and Search Strategy

Finally, we describe the acquisition functions and the acquisition optimization strategies used within the BANANAS-CP framework.

3.4.1. Acquisition Functions

We consider four commonly used acquisition functions. Consistent with the notation in the previous sections, $\hat{f}(a)$ is the predicted performance of architecture a and \hat{F}_a represents the CDF of the estimated distribution of $f(a)$. Then, depending on the acquisition function used, the specific acquisition score can be calculated by:

ITS: A sample is drawn from the distribution at random and its value is seen as the acquisition score for the candidate architecture a .

UCB: The acquisition score is given by $\mu + \gamma \cdot \sigma$ for a Gaussian distribution, where γ is the exploration factor. For non-Gaussian distributions, we follow [13] and generalize the function to a quantile function, i.e., $\hat{F}_a^{-1}(\gamma)$. As in the original formulation, higher values of γ promotes exploration. In our experiments, we set $\gamma = 0.75$ due to these concerns: a) the distribution of model performance is believed to be left-skewed; b) estimates of extreme quantiles are generally based on sparse observations and therefore might be less reliable. A value of 0.75 is likely to offer a reasonable trade-off between exploration and accurate estimation.

PI: The probability of improvements corresponds to $1 - \hat{F}_a^{-1}(f_{max})$, with f_{max} being the highest model performance ever observed.

EI: The expected value of improvements can be written as $\mathbb{E}[\max(0, f(a) - f_{max})]$, with f_{max} being the highest model performance ever observed.

3.4.2. Acquisition Optimization Strategy

In parallel to the settings in [53] (see Section 2.1.2), we compare three different approaches for constructing a set with 100¹ candidates in order to compute the acquisition scores. The motivations and the specific approaches are described below:

Mutation We investigate the mutation-based search strategy because this approach demonstrates the best performance in [53]. In line with their approach, the candidates are selected by randomly changing one operation or one edge of the k best models that have been found so far, where k is a search hyperparameter with the default value of 2.

¹ We have conducted preliminary experiments using 1,000 candidate samples as well. The results indicate that increasing the sample size does not lead to significant improvements in performance. Considering the size of the search space (NAS-Bench-201), and the required computational time, the candidate set size is fixed at 100 for all subsequent primary experiments.

3. Methodology

Random Sampling This approach is explored under the assumption that globally sampled architectures may improve the quality of the calibration process. As indicated by the name, the candidate set is created by randomly sampling architectures from the entire search space.

Dynamic This approach aims to resemble the "random + mutation" search strategy in [53]. The search process begins with random sampling until utilizing the first half of the search budget, then switches to a mutation-based strategy that progressively reduces the number of best ever-found models considered for mutation. To be precise, the number of models to be mutated decreases by 2 every 20 epoch. For instance, in a NAS task with 160 epochs, candidates are picked via random sampling for the first 80 epochs. Starting from epochs 80/100/120/140, the candidates are selected by mutating the best 8/6/4/2 models, respectively.

4. Experiment Design

To compare the performance of BANANAS-CP with the original BANANAS method and assess the role of uncertainty calibration, we choose the widely used tabular benchmark dataset NAS-Bench-201 [14] for experiments. In this chapter, we first provide a description of the dataset (Section 4.1). Then, we introduce the general setups that are shared across all experiments, along with the strategy for step-wise configuration tuning (Section 4.2).

4.1. Dataset

NAS-Bench-201 is a cell-based architecture search space. Each cell is expressed as a densely connected DAG with in total 4 nodes and 6 edges. The nodes within a cell represents the sum of all feature maps transformed through the associated operations of the edges pointing to this node, and the edges represent the architectures operation that are chosen from the predefined operation set. Specifically, the operation set comprises 5 representative operations: (1) zeroize, (2) skip connection, (3) 1-by-1 convolution, (4) 3-by-3 convolution, and (5) 3-by-3 average pooling layer. This search space contains all possible architectures generated by 4 nodes and 5 associated operation options, which results in 15,625 cell candidates in total. The macro structure of an architecture is defined as a chain of blocks, which is initiated with one 3-by-3 convolution with 16 output channels and a batch normalization layer, and consists of three stacks of cells that are connected by a residual layer. Figure 4.1 illustrates the structure of an architecture in this search space.

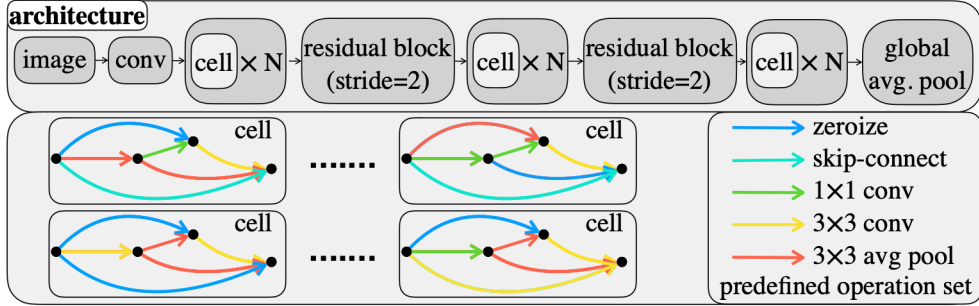


Figure 4.1: Illustration of the skeleton (top) and the design of individual cells (bottom) of architectures in NAS-Bench-201 [14].

4. Experiment Design

Architectures in the search space are evaluated on three datasets that are widely used for image classification tasks: CIFAR-10, CIFAR-100 [24] and ImageNet-16-120 [11]. Each dataset is split into the training, validation, and test sets using a standard evaluation pipeline. NAS-Bench-201 provides the training, validation, and test losses as well as accuracies for all architectures in the search space. The following gives a brief introduction to these datasets:

CIFAR10 : The dataset consists of 60K 32×32 color images in 10 classes. In NAS-Bench-201, 25K images with 10 classes are assigned into the training and the validation sets, respectively. The test set contains 10K images, with 1K images per class.

CIFAR100 : This dataset has the same images as CIFAR-10 but in 100 classes. The training set has 50K images, and each of the validation and the test sets has 5K images.

ImageNet-16-120 : This dataset contains 151.7K training images, 3K validation images, and 3K test images with 120 classes. Each image has 16×16 pixels.

Since its introduction, NAS-Bench-201 has contributed to the NAS community in several aspects. First, it provides full training and test results (e.g., accuracy, training time, etc.) for every possible architecture in the space on three datasets, allowing reproducible NAS experiments without training models from scratch. Because all architectures are pre-evaluated, NAS methods can be benchmarked extremely efficiently and subsequent NAS research can just focus on the search algorithms without any model evaluation. Moreover, using a unified dataset splitting strategy, NAS-Bench-201 reduces the variability caused by different implementation details or training setups, which is a limitation of previous benchmarks, and thus enables consistent comparisons across different NAS algorithms.

NAS-Bench-201 also serves as a foundation for extending benchmark datasets. In particular, [19] evaluates all 6,466 non-isomorphic architectures in the space for robustness against adversarial attacks and common image corruptions, and introduces a dataset that includes both clean and robust accuracy values. This dataset covers adversarial attacks and corruptions of different severity levels, enabling a systematical study on how architectural variations impact robustness.

In this work, we leverage the API offered by NAS-Bench-201 and directly query the pre-computed train and test metrics of architectures ¹. Notably, NAS-Bench-201 also provides several analytical metrics, such as model rankings and accuracy correlations across the three datasets, which further guide our post-hoc analysis of the experimental results.

¹Performance metrics for training and testing on the three datasets are downloaded from NASLib: <https://github.com/automl/NASLib/tree/Develop>

4.2. Setups and Implementation

We now present the common experimental setups used throughout this study. As introduced in Chapter 3, we primarily run experiments for five calibrated NAS algorithms varying in calibration techniques and/or surrogate predictors: SCP with ensemble predictor, SCP with quantile regressor, CrossVal-CP with ensemble predictor, CrossVal-CP with quantile regressor, and Bootstrap-CP with ensemble predictor. In line with [53], each algorithm is given a search budget of 150 epochs to ensure consistent benchmarking with BANANAS. As in previous works, the validation accuracy is used as the supervision signal to guide the search. Each algorithm is repeated for 50 trials with different random seeds and the final results are obtained by aggregating the performance across all trials.

Although training neural networks is avoided thanks to the benchmark dataset, each search algorithm still involves a large number of hyperparameters, making it unrealistic to tune them all. Also, not all hyperparameters have the same impact on search performance; some might be more important than the others. Therefore, we select a subset of hyperparameters that we believe, either base on experience or preliminary testing results, are less important or already well-set, and fix their values throughout the experiments. For instance, we believe path-encoding is stronger than other architecture encoding techniques, thus architectures are always encoded using paths present in the cell in all experiments. As for the acquisition optimization strategy, the number of candidates that the acquisition function evaluates in each iteration is fixed at 100 and the maximal mutation allowed for each model is 1 in case the mutation strategy is adopted. In addition, we follow [53] and output 10 architectures to mimic the parallelized evaluation procedure.

The hyperparameters are tuned progressively. In the first stage, we focus on configurations common to all algorithms. We start by conducting a thorough analysis on the baseline method, i.e., SCP with ensemble predictor, to find the optimal general setting, like the number of quantiles, the size of the initial dataset (the number of model evaluations before fitting the surrogate), the acquisition functions and the sampling strategies, etc. This optimal setting will be applied to other more advanced approaches in the next stage. Then, we turn to hyperparameter specific to each search algorithm and conduct separate experiments that are discussed in the respective sections.

We borrow the implementation of BANANAS from NASLib², which is a modern Python-based framework for NAS developed by the AutoML Freiburg group. NASLib is well modularized, enabling a relatively easy integration of the new calibration block. In NASLib, each NAS algorithm typically comprises a *trainer* for initiating search and evaluation iterations and an *optimizer* for encapsulating specific search logics, including an inherent *predictor* if applicable. Specifically, *trainer* serves as a generic engine and is shared across all NAS algorithms. All predictors should conform to a particular interface so that they can be invoked inside the *trainer*.

²<https://github.com/automl/NASLib>

4. *Experiment Design*

Building on this structure, we add new modules for the quantile regressor, the distribution estimators (including compatible acquisition functions), and the calibration algorithms. We leverage the existing implementation for *trainer* with mild modifications on the export functionalities, aiming for a better access to intermediate outputs, such as the estimated distribution at each iteration. In addition, we also provide tools for analyzing and interpreting the experimental result.

5. Results

This chapter describes the performance of the methods presented in Chapter 3 applied under the experimental setups described in Chapter 4.

5.1. Baseline

We consider SCP with the ensemble predictor as the baseline calibrated strategy.

6. Conclusion

This chapter presents the central findings of this work as well as their critical discussions (Section 6.1). Finally, it highlights limitations and corresponding opportunities for further research (Section 6.2).

6.1. Discussion

6.2. Limitations and Future Work

Bibliography

- [1] Apoorv Agnihotri and Nipun Batra. Exploring bayesian optimization. *Distill*, 2020. <https://distill.pub/2020/bayesian-optimization>.
- [2] Anastasios N. Angelopoulos, Stephen Bates, Adam Fisch, Lihua Lei, and Tal Schuster. Conformal risk control. *arXiv preprint arXiv:2208.02814*, 2022.
- [3] Anastasios Nikolas Angelopoulos. Conformal prediction and distribution-free uncertainty quantification. https://www.youtube.com/watch?v=nq1000Lu_iE&t=547s, 2023.
- [4] Anastasios Nikolas Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021. Version 6 (Dec 8, 2022).
- [5] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- [6] Rina Foygel Barber, Emmanuel J. Candes, Aaditya Ramdas, and Ryan J. Tibshirani. Predictive inference with the jackknife+, 2020.
- [7] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 550–559. PMLR, July 2018.
- [8] James S. Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, volume 24, pages 2546–2554. Curran Associates, Inc., 2011.
- [9] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [10] Han Cai, Ligeng Zhu, and Song Han. ProxylessNAS: Direct neural architecture search on target task and hardware. In *International Conference on Learning Representations*, 2019.
- [11] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the CIFAR datasets. *arXiv preprint arXiv:1707.08819*, 2017. 36,000 3232 ImageNet images over 1,000 classes (also variants at 1616 and 6464).

Bibliography

- [12] Aristeidis Christoforidis, George Kyriakides, and Konstantinos Margaritis. A novel evolutionary algorithm for hierarchical neural architecture search. arXiv preprint arXiv:2107.08484, 2021. Published July 18, 2021.
- [13] Shachi Deshpande, Charles Marx, and Volodymyr Kuleshov. Online calibrated and conformal prediction improves bayesian optimization. In *Proceedings of the 27th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 238, pages 7262–7273. PMLR, 2024.
- [14] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *International Conference on Learning Representations (ICLR)*, 2020.
- [15] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(55):1–21, 2019. Published online August 16, 2018.
- [16] Yarın Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, USA, June 2016. PMLR.
- [17] Alexander Gammerman, Volodya Vovk, and Vladimir Vapnik. Learning by transduction. In *Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 148–155, Madison, WI, USA, 1998. Morgan Kaufmann.
- [18] Isaac Gibbs and Emmanuel J. Candès. Adaptive conformal inference under distribution shift. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, pages 1660–1672, 2021.
- [19] Steffen Jung, Jovita Lukasik, and Margret Keuper. Neural architecture design and robustness: A dataset. In *International Conference on Learning Representations (Poster Track)*, 2023.
- [20] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabás Póczos, and Eric Xing. Neural architecture search with bayesian optimisation and optimal transport. In *Advances in Neural Information Processing Systems*, pages 2020–2029, 2018.
- [21] Alex Kendall and Yarın Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, pages 5574–5584. Curran Associates, Inc., 2017.
- [22] Byol Kim, Chen Xu, and Rina Foygel Barber. Predictive inference is free with the jackknife+-after-bootstrap. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

Bibliography

- [23] Roger Koenker and Gilbert Bassett. Regression quantiles. *Econometrica: Journal of the Econometric Society*, pages 33–50, 1978.
- [24] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical Report 2009-TR-XXX, University of Toronto.
- [25] Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2796–2804. PMLR, Jul 2018.
- [26] Solomon Kullback and Richard A Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, 1951.
- [27] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pages 6402–6413, 2017.
- [28] 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.
- [29] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. In *International Conference on Learning Representations (Poster)*, April 2018. Originally published on arXiv Nov 1, 2017 as arXiv:1711.00436.
- [30] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. arXiv preprint arXiv:1806.09055, 2018. Published 2019 at ICLR.
- [31] Lizheng Ma, Jiaxu Cui, and Bo Yang. Deep neural architecture search with deep graph bayesian optimization. In *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence*, pages 500–507, IEEE, 2019. arXiv preprint arXiv:1905.06159.
- [32] Krzysztof Maziarczyk, Mingxing Tan, Andrey Khorlin, Marin Georgiev, and Andrea Gesmundo. Evolutionary-neural hybrid agents for architecture search. *CoRR*, abs/1811.09828, 2018.
- [33] Joe Mellor, Jack Turner, Amos Storkey, and Elliot J. Crowley. Neural architecture search without training. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 7588–7598. PMLR, July 18–24 2021.
- [34] Jonas Mockus, Vytautas Tiesis, and Antanas Zilinskas. The application of Bayesian methods for seeking the extremum. *Towards Global Optimization*, 2(117–129):2, 1978.

Bibliography

- [35] Luca Mossina, Joseba Dalmau, and Léo Andéol. Conformal semantic image segmentation: Post-hoc quantification of predictive uncertainty. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 3574–3584, June 2024.
- [36] Radford M. Neal. *Bayesian Learning for Neural Networks*, volume 118 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1996.
- [37] Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learning. In *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, pages 625–632. ACM, 2005.
- [38] Xuefei Ning, Yin Zheng, Tianchen Zhao, Yu Wang, and Hua-Zhong Yang. A generic graph-based neural architecture encoding scheme for predictor-based nas. In *Computer Vision – ECCV 2020*, volume 12358, pages 189–204. Springer, 2020.
- [39] Harris Papadopoulos, Kostas Proedrou, Vladimir Vovk, and Alexander Gammerman. Inductive confidence machines for regression. In Tapio Elomaa, Heikki Mannila, and Hannu Toivonen, editors, *Machine Learning: ECML 2002*, volume 2430 of *Lecture Notes in Computer Science*, pages 345–356. Springer, 2002.
- [40] John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In Alex J. Smola, Peter L. Bartlett, Bernhard Schölkopf, and Dale Schuurmans, editors, *Advances in large margin classifiers*, pages 61–74. MIT Press, 1999.
- [41] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V. Le. Regularized evolution for image classifier architecture search. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, pages 4780–4789, Honolulu, Hawaii, USA, 2019. AAAI Press.
- [42] Yaniv Romano, Evan Patterson, and Emmanuel J Candès. Conformalized quantile regression. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [43] Yaniv Romano, Matteo Sesia, and Emmanuel J. Candès. Classification with valid and adaptive coverage. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- [44] David Salinas, Jacek Golebiowski, Aaron Klein, Matthias Seeger, and Cédric Archambeau. Optimizing hyperparameters with conformal quantile regression. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023.
- [45] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(Mar):371–421, 2008.
- [46] Divya M. Shanmugam, Helen Lu, Swami Sankaranarayanan, and John Gutttag. Test-time augmentation improves efficiency in conformal prediction. In *Proceedings of*

Bibliography

- the *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20622–20631, 2025.
- [47] Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization with robust bayesian neural networks. In *Advances in Neural Information Processing Systems*, volume 29 of *NeurIPS*, pages 4134–4142, 2016.
- [48] David Stutz, Abhijit Guha Roy, Tatiana Matejovicova, Patricia Strachan, Ali Taylan Cemgil, and Arnaud Doucet. Conformal prediction under ambiguous ground truth. *Transactions on Machine Learning Research*, 4:1–25, 2023.
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [50] Vladimir Vovk. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74(1-2):9–28, 2015.
- [51] Vladimir Vovk, Alex Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, New York, 2005.
- [52] Colin White, Willie Neiswanger, Sam Nolen, and Yash Savani. A study on encodings for neural architecture search. In *Advances in Neural Information Processing Systems*. NeurIPS, 2020. Originally released as arXiv:2007.04965 in July 2020.
- [53] Colin White, Willie Neiswanger, and Yash Savani. Bananas: Bayesian optimization with neural architectures for neural architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [54] Colin White, Mahmoud Safari, Rhea Sukthanker, Binxin Ru, Thomas Elsken, Arber Zela, Debadepta Dey, and Frank Hutter. Neural architecture search: Insights from 1000 papers. *CoRR*, abs/2301.08727, 2023.
- [55] Martin Wistuba, Ambrish Rawat, and Tejaswini Pedapati. A survey on neural architecture search. *CoRR*, abs/1905.01392, 2019.
- [56] Chen Xu and Yao Xie. Conformal prediction interval for dynamic time-series. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 11559–11569. PMLR, 18–24 Jul 2021.
- [57] Jin Xu, Xu Tan, Renqian Luo, Kaitao Song, Jian Li, Tao Qin, and Tie-Yan Liu. Nas-bert: Task-agnostic and adaptive-size bert compression with neural architecture search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1933–1943, 2021.

Bibliography

- [58] Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter. NAS-bench-101: Towards reproducible neural architecture search. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7105–7114. PMLR, June 2019.
- [59] Margaux Zaffran, Aymeric Dieuleveut, Olivier Féron, Yannig Goude, and Julie Josse. Adaptive conformal predictions for time series. In *Proceedings of the 39th International Conference on Machine Learning (ICML 2022)*, 2022.
- [60] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- [61] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8697–8710, 2018.

List of Algorithms

1.	Bayesian Optimization	6
2.	The BANANAS-CP Framework	18
3.	Split Conformal Prediction	20
4.	Conformal Prediction with Cross-validation	23
5.	Conformal Prediction with Bootstrapping	24

List of Figures

Figure 2.1: Overview of Neural Architecture Search	2
Figure 2.2: Overview of Cell-based Search Space	4
Figure 2.3: Illustration of Bayesian Optimization	6
Figure 2.4: Diagram of the BANANAS framework	8
Figure 2.5: Illustration of Symmetrical Data	11
Figure 2.6: Example of Different Coverage Types	12
Figure 3.1: Pinball loss and CQR within Bayesian optimization	22
Figure 3.2: Visualization of Estimated Distribution Based on 10 Quantiles. . .	28
Figure 3.3: Visualization of Estimated Distribution Based on 20 Quantiles. . .	28
Figure 4.1: Illustration of the overall network architecture structure in NAS- Bench-201	31

List of Tables

Table 2.1: Comparison between Conformal Prediction and Hypothesis Testing .	13
Table 3.1: Statistical Comparison of Distribution Estimators	27

Acronyms

BANANAS	Bayesian Optimization with Neural Architectures for Neural Architecture Search
BANANAS-CP	BANANAS with Conformal Prediction
Bootstrap-CP	Conformal Prediction with Bootstrapping
CDF	Cumulative Distribution Function
CP	Conformal Prediction
CQR	Conformal Quantile Regression
CrossVal-CP	Conformal Prediction with Cross-validation
DAG	Directed Acyclic Graph
EI	Expected Improvements
FCP	Full Conformal Prediction
FNNs	Feedforward Neural Networks
ITS	Independent Thompson Sampling
LOO	leave-one-out
NAS	Neural Architecture Search
PDF	Probability Density Function
PI	Probability of Improvements
RMSCE	Root Mean Squared Calibration Error
SCP	Split Conformal Prediction
TS	Thompson Sampling
UCB	Upper Confident Bound
UQ	Uncertainty Quantification

A. Program Code and Data Resources

The source code and a documentation are available at the GitHub repository: <https://github.com/chengc823/Thesis>. The datasets used for experiments and algorithm evaluations are sourced from the [NASLib repository](#).

In case of access or permission issues to the private repository, please reach out at: chechen@mail.uni-mannheim.de.

B. Additional Experimental Results

Ehrenwörtliche Erklärung

Ich versichere, dass ich die beiliegende Bachelor-, Master-, Seminar-, oder Projektarbeit ohne Hilfe Dritter und ohne Benutzung anderer als der angegebenen Quellen und in der untenstehenden Tabelle angegebenen Hilfsmittel angefertigt und die den benutzten Quellen wörtlich oder inhaltlich entnommenen Stellen als solche kenntlich gemacht habe. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen. Ich bin mir bewusst, dass eine falsche Erklärung rechtliche Folgen haben wird.

Declaration of Used AI Tools

Tool	Purpose	Where?	Useful?
ChatGPT	Rephrasing	Throughout	+
ChatGPT	Debugging LaTeX syntax errors	Equation	+
ChatGPT	Rendering LaTeX tables from Python frame	Tables	+

Unterschrift

Mannheim, den 30.07.2025