
Meta-Statistical Learning: Supervised Learning of Statistical Inference

Maxime Peyrard¹ Kyunghyun Cho^{2,3}

Abstract

This work demonstrates that the tools and principles driving the success of large language models (LLMs) can be repurposed to tackle distribution-level tasks, where the goal is to predict properties of the data-generating distribution rather than labels for individual datapoints. These tasks encompass statistical inference problems such as parameter estimation, hypothesis testing, or mutual information estimation. Framing these tasks within traditional machine learning pipelines is challenging, as supervision is typically tied to individual datapoint. We propose *meta-statistical learning*, a framework inspired by multi-instance learning that reformulates statistical inference tasks as supervised learning problems. In this approach, entire datasets are treated as single inputs to neural networks, which predict distribution-level parameters. Transformer-based architectures, without positional encoding, provide a natural fit due to their permutation-invariance properties. By training on large-scale synthetic datasets, meta-statistical models can leverage the scalability and optimization infrastructure of Transformer-based LLMs. We demonstrate the framework’s versatility with applications in hypothesis testing and mutual information estimation, showing strong performance, particularly for small datasets where traditional neural methods struggle.

accepted mechanism through which scientists relate noisy observations to theoretical models; it underpins the design of experiments, the validation of theories, and the interpretation of empirical results (Barlow, 1993; Altman, 1990; James, 2006; Dienes, 2008; Salganik, 2019).

However, the practice of statistical inference is notoriously difficult. Real-world data is noisy often deviating from idealized assumptions (Gurland and Tripathi, 1971; Hoekstra et al., 2012; Knief and Forstmeier, 2021; Czyż et al., 2023). In particular, inference in low-sample regimes presents a persistent challenge, yet it is crucial across many applied sciences. In such settings, estimators must balance universality with bias and variance, whereas more robust estimators require strong assumptions on the underlying distribution (Casella and Berger, 2024). Some statistical quantities simply lack universally unbiased estimators – like the standard deviation – necessitating context-dependent correction strategies (Gurland and Tripathi, 1971; Bengio and Grandvalet, 2003). In general, designing statistical estimators requires making choices regarding the bias-variance and robustness-universality trade-offs, requiring manual effort to craft estimators to specific goals (Silvey, 2013).

Machine learning, itself a form of statistical inference, can provide a flexible approach to these challenges. Instead of manually designing statistical estimators, we propose to *learn* them from data, leveraging amortized learning strategies to train models that generalize across diverse data distributions and adapt their estimation strategy contextually to new inputs. We call this approach *meta-statistical learning*, wherein entire datasets are treated as input objects and statistical inference tasks are directly framed as supervised learning problems (see Figure 1).

Meta-statistical learning shifts the unit of analysis from individual data points to entire datasets. Unlike traditional supervised learning, where the goal is to predict a label y for an individual sample x drawn from a joint distribution $P_{X,Y}$, meta-statistical models learn to map datasets to their target statistical properties from large amounts of synthetic datasets. This formulation aligns naturally with modern deep learning tools, which can easily handle various input modalities such as images (Voulodimos et al., 2018), graphs (Ma and Tang, 2021), time-series (Gamba, 2017; Lim and Zohren, 2021; Torres et al., 2021),

1. Introduction

Statistical inference is the backbone of many scientific inquiries, providing a rigorous framework for quantifying evidence, testing hypotheses, and estimating uncertainty (Walker and Lev, 1953; Casella and Berger, 2024). Across scientific disciplines, statistical inference is an ac-

*Equal contribution ¹Université Grenoble Alpes, CNRS, Grenoble INP, LIG ²New York University ³Genentech. Correspondence to: Maxime Peyrard <maxime.peyrard@univ-grenoble-alpes.fr>.

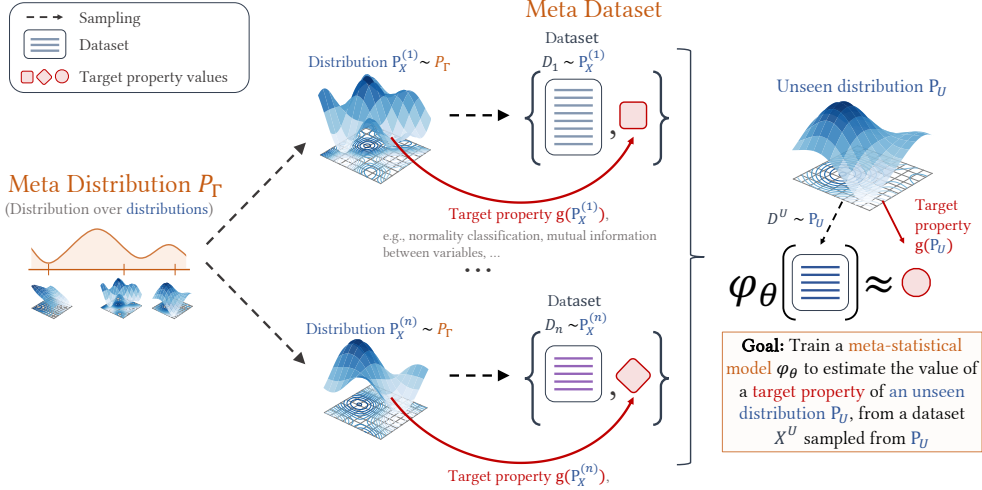


Figure 1: **Illustration of the Meta-Statistical learning setup:** a meta-distribution P_T dictates the sampling of meta-datapoints, couples of datasets X and the label y a property of data-generating distribution P_X . The meta-statistical model learns to predict y from entire datasets effectively converting statistical inference in a supervised learning problem.

language (Vaswani et al., 2017; Devlin et al., 2019; Radford et al., 2019), and tabular data (Gulati and Roysdon, 2023; Hollmann et al., 2025). The approach remains firmly within supervised learning where datasets are just another input modality processed by neural networks. However, this perspective enables us to flexibly tackle statistical inference tasks by supervised learning.

The success of large language models (LLMs) (Vaswani et al., 2017; Radford et al., 2019; Achiam et al., 2023) serves as both an inspiration and a blueprint for the meta-statistical research direction. First, attention-based architectures inherently satisfy the permutation invariance property necessary for processing datasets as unordered collections (Lee et al., 2019; Zhang et al., 2022), making them default candidate architectures. Second, the substantial infrastructure developed for LLMs—including algorithmic frameworks, hardware optimizations, and software ecosystems—can be readily repurposed for meta-statistical modeling. Finally, large-scale training datasets can be synthetically generated for most statistical inference tasks, enabling robust generalization across data distributions and strong performance in low-sample settings where classical estimators often struggle.

Contributions. In this work: (i) We introduce the *meta-statistical framework* and establish its connections to related paradigms such as meta-learning and amortized inference (Section 2). (ii) We evaluate various meta-statistical architectures, demonstrating that dataset encoders based on Set Transformer (Lee et al., 2019) variants achieve strong performance and generalization while avoiding the quadratic computational cost of standard attention mech-

anisms (Section 4). (iii) We showcase the versatility of the framework through three statistical inference problems: estimating standard deviation, conducting normality tests, and estimating mutual information. Our results highlight the robustness and generalization capabilities of meta-statistical models across diverse data distributions and, especially in low-sample-size settings (Section 5). (iv) We release for practitioners meta-statistical normality test and MI estimator, strong in low-sample regimes.

We believe the meta-statistical framework offers a promising path to leverage the principles that made LLMs successful, re-purposing them to tackle challenging statistical inference problems.

2. Meta-Statistical Learning

2.1. Background: Supervised Learning

Supervised learning aims to find a function $f: \mathcal{X} \rightarrow \mathcal{Y}$ that maps input data to output labels based on a finite dataset of observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, where $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ and $i \in \{1, \dots, n\}$. Here, \mathcal{X} denotes the input space (e.g., \mathbb{R}^d for d -dimensional data), and \mathcal{Y} denotes the output space, which is continuous for regression or discrete for classification. The data points are assumed to be i.i.d. samples from an unknown joint distribution $P_{X,Y}$ over $\mathcal{X} \times \mathcal{Y}$.

The function f is modeled by a parameterized family $\{f_\theta: \theta \in \Theta\}$, where θ represents the parameters (e.g., weights in a neural network). The quality of f_θ is evaluated using a loss function $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, which measures the discrepancy between predicted and true outputs.

Generalization. The goal is to minimize the expected risk: $R(\theta) = \mathbb{E}_{(x,y) \sim P_{X,Y}} [L(f_\theta(x), y)]$, but since $P_{X,Y}$ is unknown, the empirical risk is minimized as a proxy. Generalization is achievable because machine learning algorithms perform *induction*, based on assumptions about the underlying structure of the data and the expectation of how new data relate to observed ones. Typically, we expect the model to generalize *in distribution*, where new instances are sampled from $P_{X,Y}$. However, we often also care about generalization *out of distribution*, where new instances are sampled from a different, but related distribution.

2.2. Meta-Statistical Learning

Instead of learning a mapping from individual data points to their labels, **meta-statistical learning** maps entire datasets to their labels. Meta-statistical learning remains within standard supervised learning with the dataset being just another modality representable by a neural network.

Setup and notation. Meta-statistical learning aims to find a function $\varphi : \Gamma \rightarrow \mathcal{Y}$ that maps input datasets to labels based on a finite meta-dataset $\mathcal{S} = \{(\mathcal{D}_i, y_i)\}_{i=1}^M$, where $\mathcal{D}_i \in \Gamma$ is itself a dataset $\mathcal{D}_i = \{(x_{i,j})\}_{j=1}^{n_i}$. As in standard supervised learning, \mathcal{Y} denotes the output space, which is continuous for regression or discrete for classification. The meta-datapoints are assumed to be sampled i.i.d. from an unknown joint meta-distribution $P_{\Gamma,Y}$, a distribution over datasets (their data-generating distribution) and their target labels. The function φ is modeled by a parameterized family $\{\varphi_\theta : \theta \in \Theta\}$ that can process entire datasets as input (e.g., a recurrent neural network, convolutional neural network, or Transformer). The quality of φ_θ is still evaluated using a loss function $L_\Gamma : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. The learning objective remains to minimize the expected risk, but taken over the meta-distribution:

$$R(\theta) = \mathbb{E}_{(\mathcal{D}, y) \sim P_{\Gamma,Y}} [L_\Gamma(\varphi_\theta(\mathcal{D}), y)].$$

2.3. Structure of the Meta-Generalization Problem

To provide additional structure to the generative process that produces a meta-datapoint $(\mathcal{D}_i, y_i) \sim P_{\Gamma,Y}$, we decompose it into two steps: (i) sample a distribution P_X , and (ii) sample a dataset $\mathcal{D}_i \sim P_X$. This process is illustrated in Figure 1. The label can either be a property of the dataset itself, $y_i = A(\mathcal{D}_i)$, or a property of the distribution, $y_i = g(P_X)$. When the label is a property A of the dataset, we refer to it as a **descriptive** label, such as the column-wise average. When the label is a property g of the distribution, we refer to it as an **inferential** label, such as determining whether the dataset was sampled from a normal distribution or estimating the mutual information between two variables.

Several generalization questions arise from this setup:

(i) **Within-distribution generalization:** The function φ_θ should generalize across different datasets resampled from the same distribution P_X . In the inferential case, where the label depends only on P_X , φ_θ should produce the same prediction for all datasets of fixed size sampled from P_X . The predictions of φ_θ should not systematically overestimate or underestimate the label. This is measured by the **variance** and the **bias** of φ_θ as a statistical estimator of $y = g(P_X)$.

(ii) **Length generalization:** The function φ_θ should generalize to datasets of varying lengths. Statistical inference is harder for smaller datasets, so we expect performance to improve with larger datasets. This is measured by the **consistency** of φ_θ as a statistical estimator of $y = g(P_X)$.

(iii) **In-meta-distribution generalization:** Similar to standard supervised learning, φ_θ should generalize to new meta-datapoints sampled from the same meta-distribution $P_{\Gamma,Y}$. For example, if φ_θ is trained to predict the standard deviation of datasets sampled from exponential distributions, it should generalize to exponential distributions with unseen rate parameters.

(iv) **Out-of-meta-distribution generalization:** Analogous to out-of-distribution generalization, φ_θ could be expected to generalize to distributions and datasets sampled from a different meta-distribution than $P_{\Gamma,Y}$. For instance, if φ_θ is trained on datasets from Normal, Uniform, and Exponential distributions, it can be tested on datasets sampled from Log-normal, Cauchy, or Weibull distributions.

2.4. Related Work

The idea of processing multiple data points simultaneously originates from multi-instance learning, where models receive sets of instances and assign labels at the group level (Maron and Lozano-Pérez, 1997; Dietterich et al., 1997; Ilse et al., 2018). Once datasets can be meaningfully represented by neural networks, amortized learning techniques allowing models to generalize quickly to new datasets naturally emerge (Ganguly et al., 2023; Lopez-Paz et al., 2015; Kim et al., 2024).

A notable example is the *neural statistician* framework (Edwards and Storkey, 2017), which employs variational autoencoders (VAEs) to learn dataset representations in an unsupervised manner. Similarly, Hewitt et al. (2018) applied VAEs to infer generative models from few data points. The concept of learning dataset-level representations has also been explored through meta-features (Jomaa et al., 2021; Kotlar et al., 2021; Hartmann et al., 2023), where models extract high-level statistics tailored for specific tasks. For instance, Kotlar et al. (2021) learned meta-features for anomaly detection, while Wu et al. (2022) trained models to predict dataset-level statistics such as the number of distinct values. Recently, Hollmann et al.

(2025) employed transformers trained on synthetic datasets for missing value imputation, which we recognize as an instance of meta-statistical learning in low-sample-size settings.

Approaches of a meta-statistical nature have also been successfully applied in causal discovery (Lopez-Paz et al., 2015; Löwe et al., 2022; Lorch et al., 2022; Wu et al., 2024). These methods generate synthetic data with known causal structures and train neural networks to infer causal properties from a set of observations (Ke et al.). For example, Kim et al. (2024) proposed an attention-based model trained on simulated datasets to identify causal parents of target variables. Meta-statistical learning is a type of amortized learning focused on estimating statistical parameters; it builds upon and generalizes these previous works.

Relationship to Meta-Learning. Meta-learning, or *learning to learn*, is a paradigm focused on generalizing across tasks drawn from different distributions (Schmidhuber et al., 1996; Hospedales et al., 2021; Huisman et al., 2021). Meta-learning seeks to acquire transferable meta-knowledge, enabling rapid adaptation to new tasks (Schmidhuber, 1987; Thrun, 1998; Schmidhuber, 1993; Vanschoren, 2019). A broad range of approaches exist (Vinyals et al., 2016; Santoro et al., 2016; Finn et al., 2017; Snell et al., 2017), some emphasizing dataset-level processing to extract useful representations (Mishra et al., 2017; Ravi and Larochelle, 2017; Munkhdalai and Yu, 2017; Shyam et al., 2017). This is particularly relevant in few-shot learning (Finn et al., 2017; Snell et al., 2017; Wang et al., 2023; Wu et al., 2020; Rivolli et al., 2022). Notably, neural processes represent a class of meta-learners that use a meta-distribution over functions, adapting their prior to new datasets using observed input-output pairs (Garnelo et al., 2018b;a; Kim et al., 2019). Meta-statistical learning shares conceptual similarities with meta-learning, as both focus on generalization across distributions. However, while the target of meta-learning remains instance-level predictions, meta-statistical learning emphasizes distributional properties. These paradigms are complementary: insights from dataset-level analysis can directly improve generalization in meta-learning (Jomaa et al., 2021; Kotlar et al., 2021).

3. Experimental Setup

Our experiments demonstrate the versatility of meta-statistical learning by achieving strong performance across diverse tasks with minimal task-specific effort. Here, we describe the template used to run experiments with various descriptive and inferential tasks.

Meta-Dataset Generation. We construct meta-datasets by sampling datasets \mathcal{D} and labels y from a predefined meta-

distribution $P_{\Gamma, y}$. The generation process involves two stages: first, a distribution family (e.g., Normal, Uniform) is randomly selected, and its parameters are sampled from predefined priors to yield a data-generating distribution P_X . A dataset \mathcal{D} of size n is then sampled from P_X , with n also drawn from a prior. Thus, P_{Γ} defines the set of base distributions, parameter priors, and dataset size priors.

In- vs. Out-of-Meta-Distribution. With in-meta-distribution (IMD) settings, both training and testing datasets are sampled from P_{Γ} . For out-of-meta-distribution (OoMD) testing, we modify P_{Γ} by changing the set of base distributions (e.g., replacing Normal with Cauchy). This tests the robustness of meta-statistical estimators to unseen distributions.

3.1. Meta-Statistical Models

Models should predict dataset-level properties y from datasets \mathcal{D} of varying sizes n . The architecture we consider consists of a dataset encoder ϕ and a prediction head ρ , defined as $\varphi(\mathcal{D}; \rho, \phi) = \rho \circ \phi(\mathcal{D})$, where ϕ transforms \mathcal{D} into a fixed-dimensional representation, and ρ is a Multi-Layer Perceptron that predicts the target. The model is trained with MSE loss for regression tasks and cross-entropy loss for classification tasks.

LSTM Encoder. The Long Short-Term Memory (LSTM) network (Hochreiter and Schmidhuber, 1997) processes datasets sequentially. For a dataset $\mathcal{D} \in \mathbb{R}^{n \times k}$, the dataset representation is the average of all hidden states of the last layer: $\phi(\mathcal{D}) = \frac{1}{n} \sum_{t=1}^n h_t$, where h_t is the hidden state at timestep t . LSTMs lack the inductive bias of permutation invariance, making them a baseline model.

Vanilla Transformer Encoder. The Vanilla Transformer (VT) (Vaswani et al., 2017) uses multi-head self-attention without positional encodings to ensure permutation invariance. The dataset representation is the output of a special token, analogous to the CLS token in BERT (Devlin et al., 2019): $\phi(\mathcal{D}) = z_{\text{CLS}}$.

Set Transformer. The Set Transformer (Lee et al., 2019) is designed for set-structured data and ensures permutation invariance. Furthermore, it reduces the quadratic cost of attention by performing attention on a fixed set of m inducing points, where m is a hyperparameter. The inducing points are learned as a projection of the full sequence at each layer. The enhanced Set Transformer 2 (ST2) (Zhang et al., 2022) incorporates **SetNorm**, a normalization technique that improves over LayerNorm (Ba, 2016) by preserving permutation invariance while improving the convergence properties of the Set Transformer.

4. Experiments on Descriptive Tasks

In descriptive tasks, the label y of a dataset \mathcal{D} is the output of an algorithm A applied to \mathcal{D} , i.e., $y = A(\mathcal{D})$. Simple tasks like median or correlation serve as unit testing of meta-statistical models. However, for more computationally intensive algorithms, such as optimal transport, meta-statistical models could serve as fast approximations. For datasets $\mathcal{D} \in \mathbb{R}^{n \times m}$, we consider four descriptive tasks: the **per-column median** label $y \in \mathbb{R}^m$ consists of the medians of each column. The **Pearson correlation** coefficient $y \in \mathbb{R}$ is computed between the two columns. The **win rate** (Bradley-Terry) is the fraction of rows where the value in the first column exceeds that in the second: $y = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\mathcal{D}_{i,1} > \mathcal{D}_{i,2})$, where $\mathbb{I}(\cdot)$ is the indicator function. Finally, the 1D **optimal transport** (OT) label $y \in \mathbb{R}$ is the optimal transport cost between the empirical distributions of the two columns.

Meta-Dataset Generation. To construct the meta-dataset, we sample datasets D from predefined probability distributions as described in Section 3. Once a dataset is sampled we simply compute the target label y by applying the target algorithm. We experiment with various numbers of columns m . By having $k > 1$, we produce k computation in parallel with one forward pass (independently of the batch dimension). We observe no significant difference when varying k and fix $k = 2$ in the experiments. The meta-dataset contains 30K training meta datapoints per task, with dataset sizes sampled from $n \in [5, 300]$. Details about meta-datasets and which distribution families are in- or out-of-meta-distribution are provided in Appendix A.

Meta-Statistical Models. After optimizing hyperparameters and architecture choices (e.g., pooling mechanisms and head-to-dimensionality ratio) on a small validation set of 1K meta datapoints, we compare four meta-statistical model variants: LSTM, Vanilla Transformer (VT), and two ST2 variants with 16 or 32 inducing points. ST2(16) is the fastest model for both training and inference. In Appendix A.4, we show that VT scales quadratically, while LSTM and ST2 scale linearly, with better slopes for ST2. Additionally, ST2(16) achieves a 12x faster training time per batch normalized by parameters compared to VT, meaning an ST2(16) model with 12 times more parameters can be trained in the same time as VT. However, for consistency in reporting, we compare models with approximately the same number of parameters ($\sim 10K$ in this section).

In-meta-distribution performance. Table 1 shows the MSE of the four meta-statistical models on a test set sampled from the same meta-distribution as the training data. All models approximate the descriptive tasks well, but the LSTM-based model, lacking permutation invariance, performs worse than attention-based models. Notably, ST2,

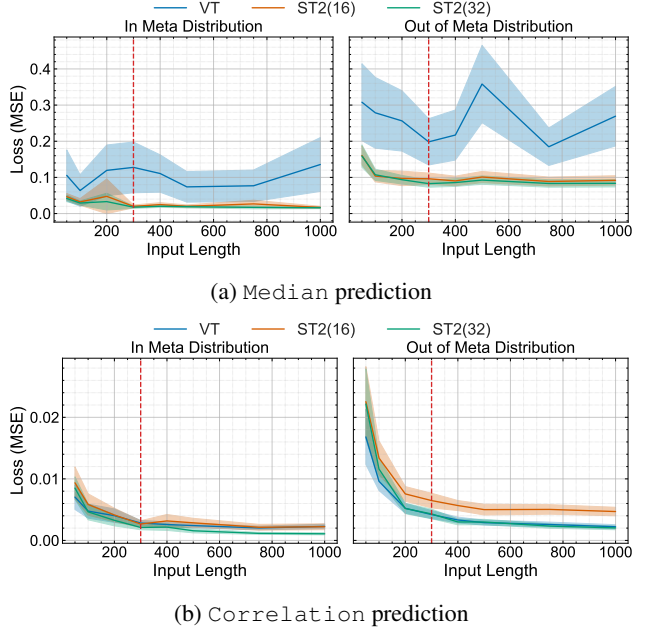


Figure 2: **Generalization across dataset lengths and meta-distributions.** The left panel shows MSE as a function of dataset length for in-meta-distribution datasets, while the right panel displays the same for out-of-meta-distribution datasets. The vertical red line marks the largest dataset seen during training ($n = 300$). LSTM is excluded due to its errors being an order of magnitude higher. Additional tasks can be found in Appendix A.3.

despite being much faster than VT, narrowly outperforms it. Given its strength and efficiency, ST2(16) is our main model in the rest of the paper, with VT considered as an alternative baseline.

Generalization Performance. We evaluate meta-statistical models’ generalization capabilities on two aspects: (i) **Out-of-Meta-Distribution (OoMD)**: Datasets from unseen distributions. (ii) **Length Generalization**: Datasets with lengths outside the training range. Figure 2 shows strong length generalization, where models maintain their performance for larger datasets than seen during training, both IMD and OoMD. They are also robust to OoMD datasets despite a small performance degradation. Manual inspection reveals that the degradation mainly comes from cases where the magnitude of the input values exceeds the range seen during training. This is discussed further in Section 6. Additional results and generalization plots are provided in Appendix A.

	Median	Corr	WinRate (BT)	OT (ID)
LSTM	$2.9e^{-1} \pm 0.8$	$5.9e^{-2} \pm 1.5$	$4.4e^{-2} \pm 0.9$	$8.5e^{-2} \pm 2.9$
VT	$6.0e^{-2} \pm 1.9$	$9.2e^{-3} \pm 4.6$	$7.1e^{-3} \pm 1.5$	$5.5e^{-2} \pm 1.4$
ST2(16)	$4.2e^{-2} \pm 1.7$	$7.5e^{-3} \pm 2.8$	$2.9e^{-3} \pm 1.2$	$4.5e^{-2} \pm 1.9$
ST2(32)	$4.4e^{-2} \pm 0.9$	$9.1e^{-3} \pm 5.1$	$1.6e^{-2} \pm 0.5$	$3.0e^{-2} \pm 1.5$

Table 1: Performance comparison meta-statistical models across tasks, measured by Mean Squared Error with respect to correct output on the test set. **Bold** indicates no significant difference with the best model.

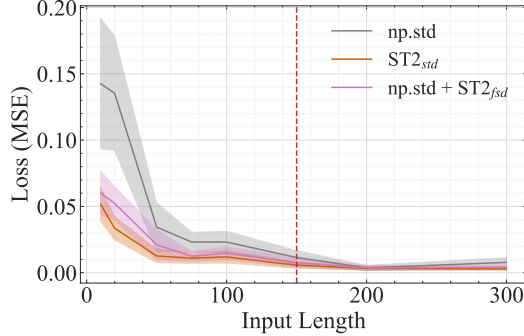


Figure 3: MSE of σ estimators as a function of dataset sizes, for dataset sampled **out-of-meta-distribution**.

5. Experiments on Inferential Tasks

In inferential tasks, the label y represents a property g of the underlying distribution P_X from which a dataset \mathcal{D} is sampled: $y = g(P_X)$. We illustrate the meta-statistical framework with three such tasks: standard deviation estimation, normality testing, and mutual information estimation. Details on meta-dataset creation and models are in Appendix B. For all tasks in this section, the dataset sizes during training are sampled from $n \in [5, 150]$, depicted by vertical red lines in the plots.

5.1. Standard Deviation Estimation

The standard deviation ($\sigma = \sqrt{\mathbb{E}[(X - \mathbb{E}[X])^2]}$) quantifies the spread of a distribution P_X . Unlike the mean or variance, estimating σ is non-trivial due to the square root’s non-linearity (Gurland and Tripathi, 1971; Gupta, 1952). In fact, no universal unbiased estimator exists across all distributions (Gurland and Tripathi, 1971; Fenstad et al., 1980). We use this task to show meta-statistical learning in action.

Meta-Dataset. To create the meta-dataset, we follow the procedure outlined in Section 3, keeping different distribution families for in- and out-of-meta-distribution. We use 100K meta datapoints for training.

Meta-Statistical Model. We train two ST2-based models: $ST2_{std}$, which predicts the standard deviation σ , and $ST2_{fsd}$,

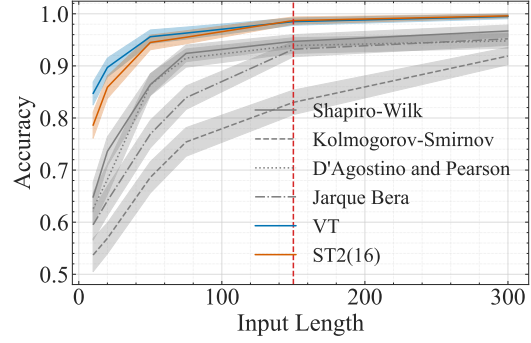


Figure 4: Accuracy of normality classifiers as a function of dataset sizes. The non-normal distributions are sampled **out-of-meta distribution** for meta-statistical models.

which estimates the finite sample correction to apply to the sample standard deviation, defined as $y = \sigma - np.std(X)$. This allows constructing a corrected estimator by adjusting $np.std$ with $ST2_{fsd}$ ’s predictions. Both models share the same architecture: 16 inducing points, five hidden layers (128 dimensions), and 12 attention heads per layer, totaling around 950K parameters.

Results. We compare the ST2-based estimator to the sample standard deviation ($np.std$ with Bessel’s correction) across dataset lengths for out-of-meta-distribution scenarios in Figure 3. $ST2_{std}$ achieves strong MSE performance, converging to a low error in high-sample sizes. Also, the learned correction from $ST2_{fsd}$ further reduces the bias of $np.std$, effectively capturing finite sample errors. Notably, $ST2_{fsd}$ also lowers variance of $np.std$, suggesting the correction is data-dependent rather than a fixed offset. Full tables of bias, variance, and MSE across distributions and dataset lengths are provided in Appendix B, confirming these observations.

5.2. Normality Testing

The task is now to determine whether a dataset $\mathcal{D} \sim P_X$ originates from a normal distribution, formulated as a binary classification task: $y = 1$ if P_X is normal, $y = 0$ otherwise. Normality testing is crucial in hypothesis testing, model selection, and preprocessing (Shapiro and Wilk, 1965; Razali et al., 2011), particularly before applying t-tests, linear regression, or ANOVA with small samples (Altman, 1990; Das and Imon, 2016; Kwak Sang Gyu, 2019). However, standard tests struggle in low-sample settings (Razali et al., 2011). We propose to train meta-statistical models for normality classification, aiming for robust generalization in such regimes. Details on meta-dataset creation and model properties are in Appendix C.

Meta-dataset creation. We construct a balanced meta-dataset of normally and not normally distributed datasets

following the process described in Section 3. For non-normality, we choose an alternative distribution from a pre-defined set detailed in Appendix C. We use 40K meta datapoints for training.

Estimators. We transform traditional normality tests into binary classifiers by thresholding their p -values, optimizing the threshold on the training meta-dataset for maximum classification accuracy. We consider four widely used tests: the *Shapiro-Wilk test* (Shapiro and Wilk, 1965), known to be effective for small samples (Razali et al., 2011); the *D’Agostino-Pearson test* (D’agostino and Pearson, 1973), which combines skewness and kurtosis; the *Kolmogorov-Smirnov test* (Massey Jr, 1951), a non-parametric test based on cumulative distribution differences; and the *Jarque-Bera test* (Jarque and Bera, 1987), which assesses skewness and kurtosis deviations from theoretical expectations.

We then train two meta-statistical models: one based on VT and another on ST2 with 16 inducing points. Both use four layers, a hidden dimensionality of 32, and 12 attention heads. The classification head is a single-layer MLP with 32 neurons, totaling approximately 50K parameters per model.

	Accuracy \uparrow	AuROC \uparrow	Brier \downarrow	BT \uparrow
KS	0.88 \pm 0.01	0.93 \pm 0.01	0.09 \pm 0.01	0.12 \pm 0.02
SW	0.89 \pm 0.01	0.95 \pm 0.01	0.18 \pm 0.01	0.16 \pm 0.03
JB	0.88 \pm 0.01	0.93 \pm 0.01	0.16 \pm 0.01	0.13 \pm 0.02
AP	0.90 \pm 0.01	0.95 \pm 0.01	0.18 \pm 0.01	0.17 \pm 0.03
ST2	0.92 \pm 0.01	0.97 \pm 0.01	0.06 \pm 0.01	0.25 \pm 0.04
VT	0.91 \pm 0.01	0.97 \pm 0.01	0.07 \pm 0.01	0.17 \pm 0.03

Table 2: Normality test classifiers with datasets drawn from Gaussian or Uniform distributions of sizes $n \in [10, 300]$. AuROC refers to the area under the ROC curve, Brier loss is the calibration error, and BT measures the relative strengths of classifiers in a paired evaluation.

Results. Figure 4 summarizes the accuracy of the proposed meta-statistical models, in settings where negative labels correspond to datasets sampled from distribution families unseen during training. Consistent with prior comparisons of normality tests, Shapiro-Wilk and D’Agostino-Pearson perform best among the baselines (Razali et al., 2011). Across all dataset sizes, meta-statistical models consistently and largely outperform baselines, with particularly strong gains in small-sample settings ($n < 100$), making them highly relevant for biomedical applications (Kwak Sang Gyu, 2019). Meta-statistical models achieve near-perfect accuracy (> 0.98) as n increases demonstrating their consistency. Overall, this task seems relatively easy for meta-statistical models, which generalize smoothly out-of-meta distribution. However, note that the training of

meta-statistical models could be harder if the input datasets are standardized during training (see Section 6).

While classification lacks a direct bias-variance formulation, we analyze false positive and false negative rates as well as precision and recall in Appendix C.3, showing more balanced error profiles for meta-statistical estimators. In Table 2, we present key metrics for evaluating classifier performance: the Area Under the Receiver Operating Characteristic Curve (AuROC), the Brier Score, and the Bradley-Terry (BT) scores from a paired evaluation. The AuROC measures a classifier’s ability to discriminate between positive and negative classes across different decision thresholds. A higher AuROC indicates better separability. Unlike accuracy, AuROC provides a threshold-independent measure of performance. The Brier loss (Brier, 1950) quantifies the calibration of a model’s predicted probabilities. Lower values indicate better calibration. The Bradley-Terry (BT) score (Bradley and Terry, 1952; Huang et al., 2004) ranks models based on pairwise comparisons, assessing how often one classifier outperforms another across test instances (Peyrard et al., 2021; Colombo et al., 2023). The accuracy scores are lower than those of Figure 4 because the uniform is among the hardest out-of-meta-distribution to recognize as non-Gaussian. Still, across metrics, the meta-statistical estimators perform strongly. In particular, we find it interesting that they are particularly well-calibrated.

5.3. Mutual Information Estimation

Mutual information (MI) quantifies the dependency between two random variables X and Y and is defined as:

$$\text{MI}(X;Y) = \int \int P_{X,Y}(x,y) \log \frac{P_{X,Y}(x,y)}{P_X(x)P_Y(y)} dx dy.$$

Here, P_X and P_Y denote the marginal distributions of X and Y , respectively.

MI possesses key properties such as invariance to homeomorphisms and adherence to the *Data Processing Inequality*, making it fundamental in machine learning and related fields (Li et al., 2021; Belghazi et al., 2018; van den Oord et al., 2018; Tishby et al., 2000). However, MI estimation remains challenging, particularly for small sample sizes and non-Gaussian distributions (Song and Ermon, 2020; McAllester and Stratos, 2020; Czyż et al., 2023).

We adopt a meta-statistical approach, training models to predict $y = \text{MI}(X;Y)$ between two dataset columns. Focusing on low-sample, non-Gaussian settings, but we restrict experiments to the one-dimensional case for simplicity. Details on meta-dataset creation, models, and extra results are provided in Appendix D.

Meta-dataset Creation. We construct a meta-dataset inspired by the benchmark methodology in (Czyż et al., 2023), where distributions with ground-truth MI are generated in two steps: (i) by sampling a distribution with known MI, (ii) optionally applying MI-preserving transformations. This process creates complex distributions and datasets with known MI. For generating meta-dataset in this way, we again follow the process described in Section 3 using different base-distribution and MI-preserving transformations between in-meta-distribution and out-of-meta-distribution. We use 50K meta datapoints for training.

Estimators. We compare our approach with the best-performing 1D estimators from (Czyż et al., 2023), including Kraskov-Stögbauer-Grassberger (KSG) (Kraskov et al., 2004), Canonical Correlation Analysis (CCA) (Murphy, 2023), and three neural estimators: MINE (Belghazi et al., 2018), InfoNCE (van den Oord et al., 2018), and NWJE (Nguyen et al., 2007; Nowozin et al., 2016; Poole et al., 2019). We train two meta-statistical models: one based on Vanilla Transformer (VT) and the other on Set Transformer 2 (ST2). Both models consist of five layers, with a hidden dimensionality of 256 and 12 attention heads. The regression head is a single hidden-layer MLP with 128 neurons, resulting in models with approximately 1M parameters.

Estimation Performance. The mean squared error (MSE) results for both in- and out-of-meta-distribution testing are shown in Table 3. Meta-statistical models outperform baseline estimators across all sample sizes, with significant advantages in low-sample scenarios. Baseline models, particularly neural ones, struggle with small sample sizes, while only KSG and CCA begin to match meta-statistical models for sample sizes greater than 100 in the out-of-meta-distribution regime.

Bias and Variance of MI Estimators. We examine the bias and variance of MI estimators by resampling datasets from fixed distributions and measuring the variance and bias of the estimates. In Figure 5, we visualize the bias and variance for a challenging distribution identified by previous works (Czyż et al., 2023) (additive noise). Even at a sample size of $n = 100$, meta-statistical models show clear improvements in both bias (estimates centered around 0) and variance. A more detailed analysis of bias and variance is available in Appendix D (Table 7). Compared to baseline estimators, meta-statistical models demonstrate significantly lower bias, close to zero, and lower or comparable variance. These results are promising, suggesting that further scaling could create even more robust meta-statistical MI estimators. Currently, the ST2 model can be trained in less than an hour on a single GPU, with inference orders of magnitude faster than existing neural baselines.

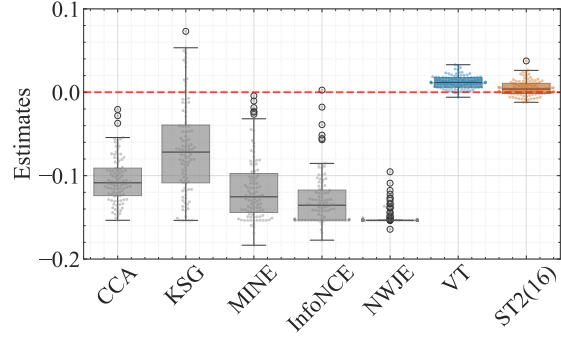


Figure 5: We estimate statistics for MI estimators over 150 resampled datasets of size $n = 100$ from a fixed distribution (additive noise (Czyż et al., 2023)). Each dot represents the difference between an estimate and the true mutual information (MI).

$n \in$	IMD		OoMD	
	[10, 100]	[100, 200]	[10, 100]	[100, 200]
CCA	$7.4e^{-2} \pm 9.3$	$1.4e^{-2} \pm 1.1$	$1.3e^{-1} \pm 1.2$	$4.9e^{-2} \pm 3.3$
KSG	$2.9e^{-2} \pm 1.5$	$7.8e^{-3} \pm 3.5$	$1.2e^{-2} \pm 0.4$	$7.2e^{-3} \pm 2.3$
MINE	$2.5e^0 \pm 2.4$	$2.8e^{-2} \pm 1.2$	$5.4e^0 \pm 7.7$	$1.6e^{-1} \pm 2.1$
NWJE	-	$7.3e^{-2} \pm 4.7$	$6.6e^0 \pm 7.5$	$6.3e^{-2} \pm 5.4$
InfoNCE	$1.5e^1 \pm 2.2$	$1.9e^{-2} \pm 0.7$	$2.3e^1 \pm 3.4$	$3.4e^{-1} \pm 4.5$
VT	$4.6e^{-3} \pm 2.4$	$2.5e^{-3} \pm 1.3$	$1.5e^{-2} \pm 0.8$	$7.7e^{-3} \pm 3.2$
ST2(16)	$6.2e^{-3} \pm 3.0$	$2.4e^{-3} \pm 1.1$	$1.3e^{-2} \pm 0.7$	$8.5e^{-3} \pm 3.1$

Table 3: MSE loss of mutual information estimators both in- and out-of-meta-distribution. **Bold** indicates no significant difference with the best estimator.

6. Discussion

With the meta-statistical framework, statistical *inference* becomes synonymous with *inference* in machine learning, and what is hard for statistical inference reveals itself as hard to learn for our models. Our experiments reveal interesting difficulties in statistical inference. Predicting normality was the *easiest* task for meta-statistical models, requiring only 50K parameters for strong generalization. Estimating mutual information, as expected, demanded significantly larger models (1M parameters). Surprisingly, predicting the standard deviation was particularly difficult: while small models (<10K parameters) could easily approximate sample standard deviation (descriptive), nearly 1M parameters were needed to predict the standard deviation (inferential) better than np.std . Training a model to predict only the corrective term also required nearly 1M parameters and yielded an estimator equivalent to directly estimating the true standard deviation, suggesting that finite-sample errors is the main driver of difficulty in this problem. This raises intriguing questions about what makes meta-statistical models work and fail: do these models im-

plicitly perform Bayesian inference with input-dependent priors?

Limitations and Future Work. While meta-statistical learning brings the advantages of machine learning to statistical inference, it also imports its challenges. A key question is the choice of meta-distribution during training—what constitutes a good meta-distribution to sample from? Additionally, the evaluation of estimators becomes more difficult; a model trained on a narrow meta-distribution might generalize poorly outside its training regime.

Interpretability is another challenge. The precise algorithm computed at inference to perform the statistical inference becomes unknown and difficult to interpret (Teney et al., 2022). Also, like LLMs, meta-statistical models could exhibit unexpected failure cases and lack strict guarantees of validity. However, they also offer a promising testbed for mechanistic interpretability research: they process structured numerical inputs without tokenization, operate in a single forward pass, and construct mathematical representations rather than linguistically ambiguous ones.

Failure cases also merit further study. Models struggled when input scales exceeded training ranges and we found one case of poor generalization to one unseen distribution family (log-normal) in the standard deviation estimation task (documented in Appendix B.4). In the normality test setting, we believe that standardizing the datasets would make training harder but encourage better generalization OoMD by preventing the meta-statistical estimators from picking up on spurious associations between the meta-distribution and the labels. Overall, like LLMs, these models would benefit from larger and more diverse training data. Future directions include learned row embeddings to accommodate varying input row dimensions and magnitudes, as well as scaling laws to guide the training of larger models with optimized data mixes. One limitation of this work is the focus on one-dimensional datasets to explore inference tasks in a controlled setting, but real-world inference involves high-dimensional data, where traditional estimators often struggle. Scaling to higher dimensions is a key next step, and we anticipate meta-statistical models to generalize well. In general, this work focuses on demonstrating promises of the meta-statistical perspective encouraging further efforts in crafting and evaluating learned statistical estimators with methods inspired by natural language processing.

Impact Statements

This paper presents work whose goal is to advance the field of Machine Learning. Meta-statistical learning aims to enhance inference in low-sample settings, benefiting applied

fields of Science like medicine and economics by improving estimator reliability. Learned estimators may inherit biases from the data they are trained on, potentially leading to misleading conclusions if not carefully validated. Further, as with any data-driven methodology, interpretability remains a challenge; understanding why a model makes a particular statistical inference is crucial for scientific rigor.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Douglas G Altman. 1990. *Practical statistics for medical research*. Chapman and Hall/CRC.
- Jimmy Lei Ba. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Roger J Barlow. 1993. *Statistics: a guide to the use of statistical methods in the physical sciences*, volume 29. John Wiley & Sons.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 531–540. PMLR.
- Yoshua Bengio and Yves Grandvalet. 2003. No unbiased estimator of the variance of k-fold cross-validation. *Advances in Neural Information Processing Systems*, 16.
- Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.
- Glenn W Brier. 1950. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- George Casella and Roger Berger. 2024. *Statistical inference*. CRC press.
- Pierre Colombo, Maxime Peyrard, Nathan Noiry, Robert West, and Pablo Piantanida. 2023. The glass ceiling of automatic evaluation in natural language generation. In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 178–183, Nusa Dua, Bali. Association for Computational Linguistics.
- Paweł Czyż, Frederic Grabowski, Julia Vogt, Niko Beerenwinkel, and Alexander Marx. 2023. Beyond normal: On

- the evaluation of mutual information estimators. In *Advances in Neural Information Processing Systems*, volume 36, pages 16957–16990. Curran Associates, Inc.
- Ralph D’agostino and Egon S Pearson. 1973. Tests for departure from normality. empirical results for the distributions of b_2 and b_3 . *Biometrika*, 60(3):613–622.
- Keya Rani Das and AHMR Imon. 2016. A brief review of tests for normality. *American Journal of Theoretical and Applied Statistics*, 5(1):5–12.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Zoltan Dienes. 2008. *Understanding psychology as a science: An introduction to scientific and statistical inference*. Bloomsbury Publishing.
- Thomas G Dietterich, Richard H Lathrop, and Tomás Lozano-Pérez. 1997. Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, 89(1-2):31–71.
- Harrison Edwards and Amos Storkey. 2017. Towards a neural statistician. In *5th International Conference on Learning Representations*, pages 1–13.
- Grete U Fenstad, Morten Kjaernes, and Lars WallØe. 1980. Robust estimation of standard deviation. *Journal of statistical computation and simulation*, 10(2):113–132.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR.
- John Cristian Borges Gamboa. 2017. Deep learning for time-series analysis. *arXiv preprint arXiv:1701.01887*.
- Ankush Ganguly, Sanjana Jain, and Ukrit Watchareeruechai. 2023. Amortized variational inference: A systematic review. *Journal of Artificial Intelligence Research*, 78:167–215.
- Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. 2018a. Conditional neural processes. In *International conference on machine learning*, pages 1704–1713. PMLR.
- Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. 2018b. Neural processes. *arXiv preprint arXiv:1807.01622*.
- Manbir Gulati and Paul Roysdon. 2023. Tabmt: Generating tabular data with masked transformers. In *Advances in Neural Information Processing Systems*, volume 36, pages 46245–46254. Curran Associates, Inc.
- AK Gupta. 1952. Estimation of the mean and standard deviation of a normal population from a censored sample. *Biometrika*, 39(3/4):260–273.
- John Gurland and Ram C Tripathi. 1971. A simple approximation for unbiased estimation of the standard deviation. *The American Statistician*, 25(4):30–32.
- Valentin Hartmann, Léo Meynert, Maxime Peyrard, Dimitrios Dimitriadis, Shruti Tople, and Robert West. 2023. Distribution inference risks: Identifying and mitigating sources of leakage. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 136–149.
- Luke B Hewitt, Maxwell I Nye, Andreea Gane, Tommi Jaakkola, and Joshua B Tenenbaum. 2018. The variational homoencoder: Learning to learn high capacity generative models from few examples. *arXiv preprint arXiv:1807.08919*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Rink Hoekstra, Henk AL Kiers, and Addie Johnson. 2012. Are assumptions of well-known statistical techniques checked, and why (not)? *Frontiers in psychology*, 3:137.
- Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeyer, and Frank Hutter. 2025. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326.
- Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. 2021. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169.
- Tzu-kuo Huang, Chih-jen Lin, and Ruby Weng. 2004. A generalized bradley-terry model: From group competition to individual skill. In *Advances in Neural Information Processing Systems*, volume 17. MIT Press.
- Mike Huisman, Jan N Van Rijn, and Aske Plaat. 2021. A survey of deep meta-learning. *Artificial Intelligence Review*, 54(6):4483–4541.

- Maximilian Ilse, Jakub Tomczak, and Max Welling. 2018. Attention-based deep multiple instance learning. In *International conference on machine learning*, pages 2127–2136. PMLR.
- Frederick James. 2006. *Statistical methods in experimental physics*. World Scientific Publishing Company.
- Carlos M Jarque and Anil K Bera. 1987. A test for normality of observations and regression residuals. *International Statistical Review/Revue Internationale de Statistique*, pages 163–172.
- Hadi S Jomaa, Lars Schmidt-Thieme, and Josif Grabocka. 2021. Dataset2vec: Learning dataset meta-features. *Data Mining and Knowledge Discovery*, 35(3):964–985.
- Nan Rosemary Ke, Silvia Chiappa, Jane X Wang, Jorg Bornschein, Anirudh Goyal, Melanie Rey, Theophane Weber, Matthew Botvinick, Michael Curtis Mozer, and Danilo Jimenez Rezende. Learning to induce causal structure. In *International Conference on Learning Representations*.
- Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. 2019. Attentive neural processes. *arXiv preprint arXiv:1901.05761*.
- Jang-Hyun Kim, Claudia Skok Gibbs, Sangdoo Yun, Hyun Oh Song, and Kyunghyun Cho. 2024. Targeted cause discovery with data-driven learning. *arXiv preprint arXiv:2408.16218*.
- Ulrich Knief and Wolfgang Forstmeier. 2021. Violating the normality assumption may be the lesser of two evils. *Behavior Research Methods*, 53(6):2576–2590.
- Miloš Kotlar, Marija Punt, Zaharije Radivojević, Miloš Cvetanović, and Veljko Milutinović. 2021. Novel meta-features for automated machine learning model selection in anomaly detection. *IEEE Access*, 9:89675–89687.
- Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. 2004. Estimating mutual information. *Phys. Rev. E*, 69:066138.
- Park Sung-Hoon Kwak Sang Gyu. 2019. Normality test in clinical research. *jrd*, 26(1):5–11.
- Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. 2019. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pages 3744–3753. PMLR.
- Bo Li, Yifei Shen, Yezhen Wang, Wenzhen Zhu, Colorado J. Reed, Jun Zhang, Dongsheng Li, Kurt Keutzer, and Han Zhao. 2021. Invariant information bottleneck for domain generalization. *CoRR*, abs/2106.06333.
- Bryan Lim and Stefan Zohren. 2021. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194):20200209.
- David Lopez-Paz, Krikamol Muandet, Bernhard Schölkopf, and Iliya Tolstikhin. 2015. Towards a learning theory of cause-effect inference. In *International Conference on Machine Learning*, pages 1452–1461. PMLR.
- Lars Lorch, Scott Sussex, Jonas Rothfuss, Andreas Krause, and Bernhard Schölkopf. 2022. Amortized inference for causal structure learning. *Advances in Neural Information Processing Systems*, 35:13104–13118.
- Sindy Löwe, David Madras, Richard Zemel, and Max Welling. 2022. Amortized causal discovery: Learning to infer causal graphs from time-series data. In *Conference on Causal Learning and Reasoning*, pages 509–525. PMLR.
- Yao Ma and Jiliang Tang. 2021. *Deep learning on graphs*. Cambridge University Press.
- Oded Maron and Tomás Lozano-Pérez. 1997. A framework for multiple-instance learning. *Advances in neural information processing systems*, 10.
- Frank J Massey Jr. 1951. The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78.
- David McAllester and Karl Stratos. 2020. Formal limitations on the measurement of mutual information. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 875–884. PMLR.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. 2017. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*.
- Tsendsuren Munkhdalai and Hong Yu. 2017. Meta networks. In *International conference on machine learning*, pages 2554–2563. PMLR.
- Kevin P. Murphy. 2023. *Probabilistic Machine Learning: Advanced Topics*. MIT Press.
- XuanLong Nguyen, Martin J Wainwright, and Michael Jordan. 2007. Estimating divergence functionals and the likelihood ratio by penalized convex risk minimization. In *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc.

- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. 2016. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Maxime Peyrard, Wei Zhao, Steffen Eger, and Robert West. 2021. Better than average: Paired evaluation of NLP systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2301–2315, Online. Association for Computational Linguistics.
- Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. 2019. On variational bounds of mutual information. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5171–5180. PMLR.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Sachin Ravi and Hugo Larochelle. 2017. Optimization as a model for few-shot learning. In *International conference on learning representations*.
- Nornadiiah Mohd Razali, Yap Bee Wah, et al. 2011. Power comparisons of shapiro-wilk, kolmogorov-smirnov, lilliefors and anderson-darling tests. *Journal of statistical modeling and analytics*, 2(1):21–33.
- Adriano Rivolli, Luís PF Garcia, Carlos Soares, Joaquin Vanschoren, and André CPLF de Carvalho. 2022. Meta-features for meta-learning. *Knowledge-Based Systems*, 240:108101.
- Matthew J Salganik. 2019. *Bit by bit: Social research in the digital age*. Princeton University Press.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850. PMLR.
- Juergen Schmidhuber, Jieyu Zhao, and Marco Wiering. 1996. Simple principles of metalearning.
- Jürgen Schmidhuber. 1987. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. Ph.D. thesis, Technische Universität München.
- Jürgen Schmidhuber. 1993. A neural network that embeds its own meta-levels. In *IEEE International Conference on Neural Networks*, pages 407–412. IEEE.
- Samuel Sanford Shapiro and Martin B Wilk. 1965. An analysis of variance test for normality (complete samples). *Biometrika*, 52(3-4):591–611.
- Pranav Shyam, Shubham Gupta, and Ambedkar Dukkipati. 2017. Attentive recurrent comparators. In *International conference on machine learning*, pages 3173–3181. PMLR.
- Samuel Silvey. 2013. *Optimal design: an introduction to the theory for parameter estimation*, volume 1. Springer Science & Business Media.
- Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30.
- Jiaming Song and Stefano Ermon. 2020. Understanding the limitations of variational mutual information estimators.
- Damien Teney, Maxime Peyrard, and Ehsan Abbasnejad. 2022. Predicting is not understanding: Recognizing and addressing underspecification in machine learning. In *European Conference on Computer Vision*, pages 458–476. Springer.
- Sebastian Thrun. 1998. Lifelong learning algorithms. In *Learning to learn*, pages 181–209. Springer.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. 2000. The information bottleneck method.
- José F Torres, Dalil Hadjout, Abderrazak Sebaa, Francisco Martínez-Álvarez, and Alicia Troncoso. 2021. Deep learning for time series forecasting: a survey. *Big Data*, 9(1):3–21.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748.
- Joaquin Vanschoren. 2019. Meta-learning. *Automated machine learning: methods, systems, challenges*, pages 35–61.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. *Advances in neural information processing systems*, 29.
- Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. 2018. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*, 2018(1):7068349.

- Helen M Walker and Joseph Lev. 1953. Statistical inference.
- Xixi Wang, Xiao Wang, Bo Jiang, and Bin Luo. 2023. Few-shot learning meets transformer: Unified query-support transformers for few-shot classification. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(12):7789–7802.
- Menghua Wu, Yujia Bao, Regina Barzilay, and Tommi Jaakkola. 2024. Sample, estimate, aggregate: A recipe for causal discovery foundation models. *arXiv preprint arXiv:2402.01929*.
- Mike Wu, Kristy Choi, Noah Goodman, and Stefano Ermon. 2020. Meta-amortized variational inference and learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6404–6412.
- Renzhi Wu, Bolin Ding, Xu Chu, Zhewei Wei, Xiening Dai, Tao Guan, and Jingren Zhou. 2022. Learning to be a statistician: learned estimator for number of distinct values. *arXiv preprint arXiv:2202.02800*.
- Lily Zhang, Veronica Tozzo, John Higgins, and Rajesh Ranganath. 2022. Set norm and equivariant skip connections: Putting the deep in deep sets. In *International Conference on Machine Learning*, pages 26559–26574. PMLR.

A. Details about the Descriptive tasks experiments

A.1. Details of Meta-Dataset Creation

To ensure reproducibility of the experiments, we describe the synthetic data generation process in detail. The datasets were generated using a custom-built class, `DescMetaDatasetGenerator`, which allows for the creation of datasets with various distributions and customizable descriptive target variables. The key components and configurations are outlined below.

In-Meta-Distribution. The set of distributions used to generate datasets during training is parameterized as follows:

- **normal::** It has two parameters: the mean and the variance. Mean values are sampled from $[-3, 3]$, and variances are sampled from $[0.1, 1.5]$.
- **uniform::** It has two parameters: the lower bound and the upper bound. The lower bounds are sampled from $[-3.5, -0.5]$ and the upper bounds from $[0.5, 3.5]$.
- **beta::** It has two parameters: a and b . Parameters a and b are sampled from $[1, 3]$ and $[2, 5]$, respectively.
- **exponential::** It has one parameter: scale sampled from $[1, 2]$.

Out-of-Meta-Distribution. The set of distribution used to test models for unseen distribution families is parametrized as follows:

- **gamma::** It has two parameters: shape and scale. Shape parameters are sampled from $[1, 5]$, and scale parameters from $[1, 2]$.
- **log-normal::** It has two parameters: mean and variance. Means are sampled from $[0, 1]$, and standard deviations from $[0.5, 0.75]$.

Dataset Characteristics. Once a distribution P_X has been sampled, we use it to sample one dataset. In general, we could sample several dataset per distributions but we prefer to sample only one to maximize the diversity of distributions seen during training. Each dataset is defined by the following parameters:

- **Number of variables (n_{var}):** The number of features (columns) in the dataset. For our experiments, we set $n_{\text{var}} = 2$.

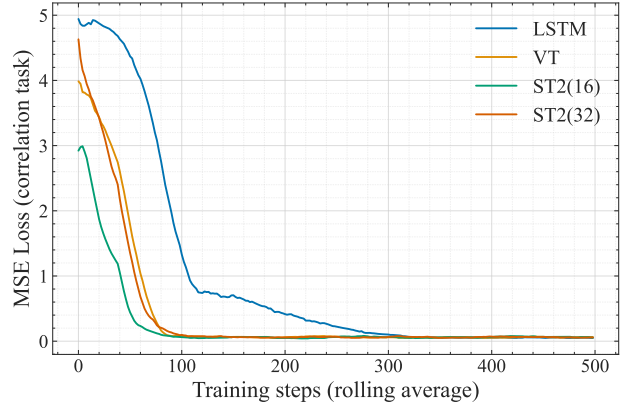


Figure 6: **Training curves:** Comparison of training convergence of meta-statistical models on the correlation task.

- **Number of rows ($n_{\text{row_range}}$):** The number of samples (rows) in the dataset, sampled uniformly from the range $[5, 300]$. During testing, we explore longer lengths to test the generalization of meta-statistical models.
- To generate the target values y , each dataset is passed through the target descriptive functions: per-column-mean, per-column-median, correlation, win rate, optimal transport (1D).

This results in a meta-datapoint. We then sample many meta-datapoints to build a meta-dataset with the following split sizes: 30K training, 300 validation, 3K for testing in-meta-distribution and 3K for testing out-of-meta-distribution.

A.2. Examples of Training Curves

For meta-statistical models of approximately the same size ($\approx 10K$ parameters), we compare their convergence during training on the task of predicting the correlation between variable A and variable B, the two columns of the dataset. We consider the same meta-statistical models and the same meta-dataset generation parameters as considered in the results of the main paper. We report the results in Figure 6.

A.3. More Generalization Plots

In Figure 7, we report the same generalization plots as the main paper for other descriptive tasks. Similar conclusion holds: models generalize very well with lengths and tend to suffer from an offset of performance out-of-meta-distribution.

	Training Time	Inference Time
VT	$2.2e^{-5} \pm 1.1$	$2.7e^{-3} \pm 1.2$
LSTM	$6.5e^{-6} \pm 1.0$	$8.5e^{-4} \pm 3.9$
ST2(32)	$5.9e^{-6} \pm 6.9$	$1.5e^{-3} \pm 2.4$
ST2(16)	$1.7e^{-6} \pm 0.2$	$2.2e^{-4} \pm 0.9$

Table 4: Comparison of training and inference times for various models, normalized per batch per number of parameters.

A.4. Details about Efficiency

In Figure 8, we present the inference time of meta-statistical models as a function of the input dataset size n . As expected, the VT scales quadratically, whereas LSTM and ST2 variants scale linearly with slopes in favor of ST2. We also compare the efficiency per parameter. For this we compute both the training and inference time of each model per batch averaged over 1K batches, and normalized by the number of parameters in the model. The results are reported in Table 4. Given the strong performance of ST2 and the clear computational advantage we see it as strong meta-statistical architecture.

B. Details about Standard Deviation Experiments

B.1. Details of Meta-Dataset Creation

We construct a meta-dataset by generating datasets labeled with the ground truth standard deviation, using a set of distributions for which the standard deviation is well-defined. To create each meta-datapoint, we first sample the base distribution uniformly at random from a set of pre-defined distribution families (see below). Then, we sample the parameters of the distribution, resulting in a distribution P_X . A dataset size n is then drawn uniformly at random from the range $[10, 150]$, and the dataset D is sampled with n rows. We generate 50K meta-datapoints for training and 3K for validation.

In-Meta-Distribution. These are the distributions **seen** during training. The base distributions are the following, with the priors on their parameters:

- **normal**: the mean is sampled from $\mathcal{U}(-1, 1)$, and the variance is sampled from $\mathcal{U}(0.5, 2.0)$.
- **uniform**: the lower bound is sampled from $\mathcal{U}(0, 0.5)$, and the upper bound is sampled from $\mathcal{U}(0.5, 1.5)$.
- **exponential**: the scale parameter is sampled from $\mathcal{U}(1, 2)$.

- **gamma**: the shape parameter is sampled from $\mathcal{U}(1, 5)$, and the scale parameter is sampled from $\mathcal{U}(1, 2)$.

Out-of-Meta-Distribution. These are the distributions **seen** during training. The base distributions are the following, with the priors on their parameters:

- **beta**: the α parameter is sampled from $\mathcal{U}(1, 5)$, and the β parameter is sampled from $\mathcal{U}(1, 5)$.
- **lognormal**: the mean of the underlying normal distribution is sampled from $\mathcal{U}(0, 1)$, and the standard deviation from $\mathcal{U}(0.1, 1)$.
- **weibull**: the shape parameter is sampled from $\mathcal{U}(1, 5)$, and the scale parameter is sampled from $\mathcal{U}(1, 2)$.

B.2. Details about Meta-Statistical Models

For these experiments, we train two meta-statistical models based on **Set Transformer 2 (ST2)**. Variants with different numbers of inducing points were tested, such as ST2 (16), which uses 16 inducing points (`num_inds = 16`).

- **ST2_{std}**: an ST2 (16) encoder with a regression MLP trained to predict the standard deviation σ_{P_X} of the distribution.
- **ST2_{fstd}**: an ST2 (16) encoder with a regression MLP trained to predict the finite sample error made by the sample standard deviation, i.e., it predicts $y = \sigma_{P_X} - \text{np.std}(X)$.

This design probes whether meta-statistical models can reliably estimate finite sampling errors from data, offering insights into their expected performance. We can craft a new standard deviation estimator by combining the sample standard deviation (`np.std`) with the corrections predicted by ST2_{fstd}. Both ST2-based models use the same architecture, consisting of 16 inducing points, five hidden layers with 128 dimensions, and 12 attention heads per layer, resulting in approximately 960K parameters. The prediction head is implemented as an MLP with a single hidden layer of 64 units.

Training Configuration. The models were trained on a regression task using the binary MSE loss. We employed a batch size of 64 and optimized the model parameters using the Adam optimizer with a learning rate of 1×10^{-5} . The training process spanned 10 epochs.

B.3. Bias and Variance

In Figure 9, we report an experiment targeted at measuring bias and variance by resampling 150 datasets from a fixed `exponential` distribution. Let the true standard deviation be denoted as σ and the estimates for the i -th dataset be $\hat{\sigma}_i$, for $i = 1, \dots, 150$.

The bias of the estimator is computed as:

$$\text{Bias} = \frac{1}{n} \sum_{i=1}^n \hat{\sigma}_i - \sigma, \quad (1)$$

where $n = 50$ is the number of resampled datasets.

The variance of the estimator is computed as:

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n \left(\hat{\sigma}_i - \frac{1}{n} \sum_{j=1}^n \hat{\sigma}_j \right)^2. \quad (2)$$

It shows that the meta-statistical models can improve over the sample standard deviation `np.std` both in terms of bias and variance. Especially, the learned correction is capable of reducing the bias of `np.std` indicating that it has not just learned a constant offset. We also report bias, variance, and MSE across dataset sizes and over many meta-datapoints sampled from various distributions (both in- and out-of-meta-distribution) in Table 5.

B.4. Interesting Failure Case

While meta-statistical models generally demonstrate strong robustness when presented with datasets sampled from distribution unseen during training (OoMD), we found a interesting failure case. For the log-normal family that was unseen during training, meta-statistical estimators of standard deviations failed to provide improvements over `np.std` and even performed worse. This is illustrated in Figure 10 across dataset sizes. Log-normal is a skewed distribution which is particularly challenging estimators of standard deviation which can explain why meta-statistical model poorly generalize in this case. These models can benefit from a larger and more diverse training meta-dataset to exhibit even more robust generalization. Note, however, that the models do not fail for other unseen distribution families as it can be seen in Figure 3 of the main paper.

C. Details about Normality Tests Experiments

C.1. Details of Meta-Dataset Creation

We construct a meta-dataset by generating datasets labeled with the ground truth binary indicator of normality, using a diverse set of alternative distributions. In previous studies, the uniform distribution was used the contrast distribution Razali et al. (2011). To create each meta-datapoint,

we first determine whether the dataset is sampled from a normal distribution by flipping a fair coin. If the outcome is normal, we sample the mean and variance of the Gaussian distribution; otherwise, we sample uniformly at random one distribution from the set of non-Gaussian distributions and then draw its parameters.. A dataset size n is then drawn uniformly at random from the range $[5, 150]$, and the dataset D is sampled with n rows. We generate 40K meta-datapoints for training, 1K for validation, and 1K meta-datapoints for testing.

The normal distribution parameters are sampled according to: the mean is sampled from $\mathcal{U}(-3, 3)$, and the standard deviation from $\mathcal{U}(0.5, 3)$.

In-Meta-Distribution. These are the distributions **seen** during training as non-normal distributions. Not however that because the distribution sampled parameters and the datasets are then sampled from the distributions, in-meta-distribution evaluation still produces different meta-datapoints.

- **gamma:** the shape parameter is sampled from $\mathcal{U}(1, 5)$, and the scale parameter from $\mathcal{U}(0.5, 2)$.
- **triangular:** the lower bound is sampled from $\mathcal{U}(-3, 0)$, the mode from $\mathcal{U}(\text{lower bound}, 3)$, and the upper bound from $\mathcal{U}(\text{mode}, 5)$.
- **cauchy:** the location parameter is sampled from $\mathcal{U}(-1, 1)$, and the scale parameter from $\mathcal{U}(0.5, 2)$.
- **laplace:** the location parameter is sampled from $\mathcal{U}(-1, 1)$, and the scale parameter from $\mathcal{U}(0.5, 2)$.
- **weibull:** the shape parameter is sampled from $\mathcal{U}(0.5, 5)$.
- **vonmises:** the mean direction μ is sampled from $\mathcal{U}(-\pi, \pi)$, and the concentration parameter κ from $\mathcal{U}(0.5, 5)$.
- **arcsine:** the lower bound is sampled from $\mathcal{U}(-3, 0)$, and the upper bound from $\mathcal{U}(0, 3)$. Data is generated by transforming uniform samples with a sine function scaled to the specified bounds.
- **bimodal:** two Gaussian components are used, where the means of the components are sampled from $\mathcal{U}(-3, -1)$ and $\mathcal{U}(1, 3)$, respectively. The standard deviations are sampled from $\mathcal{U}(0.5, 1)$, and the mixture ratio is sampled from $\mathcal{U}(0.3, 0.7)$.

Out-of-Meta-Distribution. These are the distributions **unseen** during training as non-normal distributions.

	Bias			Variance			MSE		
	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$
np.std	$-3.4e^{-2} \pm 2.3$	$-2.5e^{-2} \pm 1.6$	$-3.3e^{-3} \pm 14.3$	$1.1e^{-1} \pm 0.3$	$3.5e^{-2} \pm 0.9$	$1.7e^{-2} \pm 0.4$	$1.2e^{-1} \pm 0.2$	$4.3e^{-2} \pm 0.6$	$2.2e^{-2} \pm 0.3$
ST2 _{std}	$-8.1e^{-3} \pm 27.7$	$1.9e^{-2} \pm 2.1$	$1.8e^{-2} \pm 1.4$	$1.6e^{-2} \pm 0.3$	$1.1e^{-2} \pm 0.2$	$6.4e^{-3} \pm 1.2$	$3.7e^{-2} \pm 0.4$	$2.2e^{-2} \pm 0.3$	$1.1e^{-2} \pm 0.1$
np.std + ST2 _{fsd}	$-5.4e^{-3} \pm 28.6$	$-1.5e^{-2} \pm 2.2$	$-1.1e^{-2} \pm 1.4$	$3.6e^{-2} \pm 1.0$	$1.5e^{-2} \pm 0.4$	$8.4e^{-3} \pm 1.8$	$5.7e^{-2} \pm 0.8$	$2.5e^{-2} \pm 0.3$	$1.4e^{-2} \pm 0.2$
np.std	$-5.6e^{-2} \pm 2.9$	$-3.7e^{-2} \pm 2.1$	$-1.4e^{-2} \pm 1.6$	$1.7e^{-1} \pm 1.1$	$7.9e^{-2} \pm 4.9$	$1.9e^{-2} \pm 1.0$	$2.0e^{-1} \pm 0.8$	$9.1e^{-2} \pm 4.3$	$2.6e^{-2} \pm 0.7$
ST2 _{std}	$7.1e^{-2} \pm 9.9$	$-2.1e^{-2} \pm 6.5$	$1.2e^{-2} \pm 5.2$	$1.3e^{-2} \pm 0.3$	$8.0e^{-3} \pm 1.9$	$5.0e^{-3} \pm 1.6$	$2.7e^{-1} \pm 0.3$	$1.1e^{-1} \pm 0.2$	$7.9e^{-2} \pm 1.5$
np.std + ST2 _{fsd}	$4.9e^{-2} \pm 6.0$	$-4.8e^{-2} \pm 4.3$	$-3.2e^{-2} \pm 3.4$	$9.1e^{-2} \pm 8.3$	$4.8e^{-2} \pm 4.0$	$1.0e^{-2} \pm 0.6$	$1.9e^{-1} \pm 0.5$	$9.8e^{-2} \pm 4.2$	$4.0e^{-2} \pm 1.1$

Table 5: **Bias and Variance of Standard Deviation Estimators** The top part shows the in-meta-distribution results, while the bottom part reports the out-of-meta-distribution results.

- **uniform**: the lower bound is sampled from $\mathcal{U}(-3, 0)$, and the upper bound from $\mathcal{U}(0, 3)$.
- **exponential**: the scale parameter is sampled from $\mathcal{U}(0.5, 2)$.
- **beta**: the shape parameters a and b are sampled independently from $\mathcal{U}(0.5, 5)$.
- **log-normal**: the mean is sampled from $\mathcal{U}(-1, 1)$, and the standard deviation from $\mathcal{U}(0.5, 1.5)$.

C.2. Details of Meta-Statistical Models

The experiments compared the performance of multiple model variants, including:

- **Vanilla Transformer (VT)**: This model utilized the `vanilla_transformer` architecture.
- **Set Transformer 2 (ST2)**: a set Transformer with 16 inducing points (`num_inds = 16`).

Both models have four layers, hidden dimensionality of 32, and 12 attention heads. To build a full meta-statistical model from the meta-statistical encoders, we add a prediction head made of a MLP with one layer of 32 neurons before predicting the probability of being normally distributed. The meta-model based on VT has a total of 51K parameters, and the meta-model based on ST2 has a total of 54K. ST2 has more parameters to learn the projected attention but ends being much faster to train and use at inference because of the constant cost of attention compared to quadratic in the length of the input for VT.

Training Configuration. The models were trained on a binary classification task using the binary cross-entropy loss. We employed a batch size of 24 and optimized the model parameters using the Adam optimizer with a learning rate of 0.0005. The training process spanned 7 epochs. These hyper-parameters were selected to balance computational efficiency and convergence stability. In the main paper Figure 4, the meta-statistical models are evaluated OoMD. For completeness, we report their results on unseen test set but in-meta-distribution in Table 11

C.3. Details about Precision and Recall

While classification lacks a direct bias-variance formulation, we analyze false positive and false negative rates as well as precision and recall in Table 6. Interestingly, most baselines, when optimizing their p-value threshold for maximum accuracy end up maximizing recall at the expense of precision with high false positive rates. On the contrary, the meta-statistical models perform well in both precision and recall with balanced error profiles, i.e., similar amount of false negatives and false positives.

D. Details about Mutual Information Experiments

D.1. Details of Meta-Dataset Creation

We construct a meta-dataset inspired by the benchmark methodology in (Czyż et al., 2023), where distributions with ground-truth MI are generated in two steps: (i) by sampling a distribution with known MI, (ii) potentially applying MI-preserving transformation. This process creates complex distributions and datasets with known MI. For generating meta-dataset in this way, we again follow the process described in Section 3 using different base-distribution and MI-preserving transformation between in-meta-distribution and out-of-meta-distribution.

In-Meta-Distribution. The in-meta-distributions focus on bivariate relationships and transformations:

- **binormal-[base, wigglyfy, halfcube, asinh, normal_cdfise]**: Standard bivariate normal distribution with correlation sampled from $\mathcal{U}(-1, 1)$, following by any of the MI-preserving transformation (or none).
- **bimodal_gaussians-base**: Bivariate Gaussian mixture model with correlation sampled from $\mathcal{U}(-1, 1)$.
- **bistudent-[base, asinh]**: Bivariate Student’s t-distribution with degrees of freedom sampled

	TP	FP	TN	FN	Prec.	Recall	F1-Score
Shapiro-Wilk	439	138	383	40	0.7608	0.9165	0.8314
Kolmogorov-Smirnov	479	364	157	0	0.5682	1.0000	0.7247
D’Agostino and Pearson	423	147	374	56	0.7421	0.8831	0.8065
Jarque Bera	449	296	225	30	0.6027	0.9374	0.7337
VT	471	38	441	50	0.9253	0.9040	0.9146
ST2(16)	474	50	429	47	0.9046	0.9098	0.9072

Table 6: Reporting True Positive (TP), False Positive (FP), True Negative (TN), False Negative (FN), the associated Precision, Recall, and F1-Score of normality classifiers over 1,000 datasets of size $n = 50$.

from $\mathcal{U}(1,10)$, potentially followed by asinh transform.

Out-of-Meta-Distribution. The out-of-meta-distributions extend the variety of data by incorporating additional transformations and configurations:

- **additive_noise**=[base, wigglyfy, halfcube, asinh, normal_cdfise]: A model where X is sample from $\mathcal{U}(-1,1)$ and $Y = X + \epsilon$, a Gaussian noise sampled from $\mathcal{U}(0.01,2)$ followed by any of the MI-preserving transformation.
- **bistudent**=[wigglyfy, halfcube, normal_cdfise]: Bivariate Student’s t-distribution (degrees of freedom sampled from $\mathcal{U}(1,10)$) followed by any of the MI-preserving transformation except asinh.

For the computation of the true MI from these distributions and their transformations, we refer to (Czyż et al., 2023); we sample using the tool they provided: <https://github.com/cbg-ethz/bmi>.

D.2. Details of the Training of Meta-Statistical Models

The meta-statistical models were trained using the code framework outlined in the previous section. This section provides detailed information about the training process.

Default Configuration. The default configuration specifies the following key components:

- **Dimensionality:** In these experiments, we focus on 1D dimensional variables.
- **Dataset Parameters:** Dataset sizes are sampled uniformly from the range [10, 150].
- **Meta-Dataset Properties:** The meta-dataset contains 50,000 training, 500 validation, and 1,000 testing meta-datapoints.

- **Training Parameters:** A regression task was specified with a batch size of 64, learning rate of 0.0001, and 20 epochs of training.

- **Model Architecture:** The base model used a 5-layer encoder (`n_enc_layers`), with hidden dimensions of 256 and 128 for the `phi` and `theta` components, respectively. Multi-head attention mechanisms used 12 heads.

Model Variants. The experiments compared the performance of multiple model variants, including:

- **Vanilla Transformer (VT):** This model utilized the `vanilla_transformer` architecture.
- **Set Transformer 2 (ST2):** a set Transformer with 16 inducing points (`num_inds = 16`).

To build a full meta-statistical model from the meta-statistical encoders, we add a prediction head made of a MLP with one layer of 128 neurons before predicting the 1 number target output. The meta-model based on VT has a total of 1,008,385 parameters, and the meta-model based on ST2 has a total of 1,280,065. ST2 has more parameters to learn the projected attention but ends being much faster to train and use at inference because of the constant cost of attention compared to quadratic in the length of the input for VT.

D.3. Details about the Bias and Variance of MI Estimators

To estimate the bias and variance of a mutual information (MI) estimator, we follow a systematic procedure. First, we sample a base distribution and an MI-preserving transformation from the out-of-meta-distribution set. From this fixed setup, we resample $n = 50$ datasets, leading to 50 MI estimates from the given estimator. Let the true mutual information be denoted as I_{true} and the MI estimates for the i -th dataset be \hat{I}_i , for $i = 1, \dots, 50$.

The bias of the estimator is computed as:

$$\text{Bias} = \frac{1}{n} \sum_{i=1}^n \hat{I}_i - I_{\text{true}}, \quad (3)$$

where $n = 50$ is the number of resampled datasets.

The variance of the estimator is computed as:

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n \left(\hat{I}_i - \frac{1}{n} \sum_{j=1}^n \hat{I}_j \right)^2. \quad (4)$$

To capture a broader picture of estimator behavior, we repeat this process for 100 random choices of base distributions and transformations. For each random choice, we report the bias and variance of the estimator as calculated above. Finally, we summarize these results across different sample sizes in Table 7.

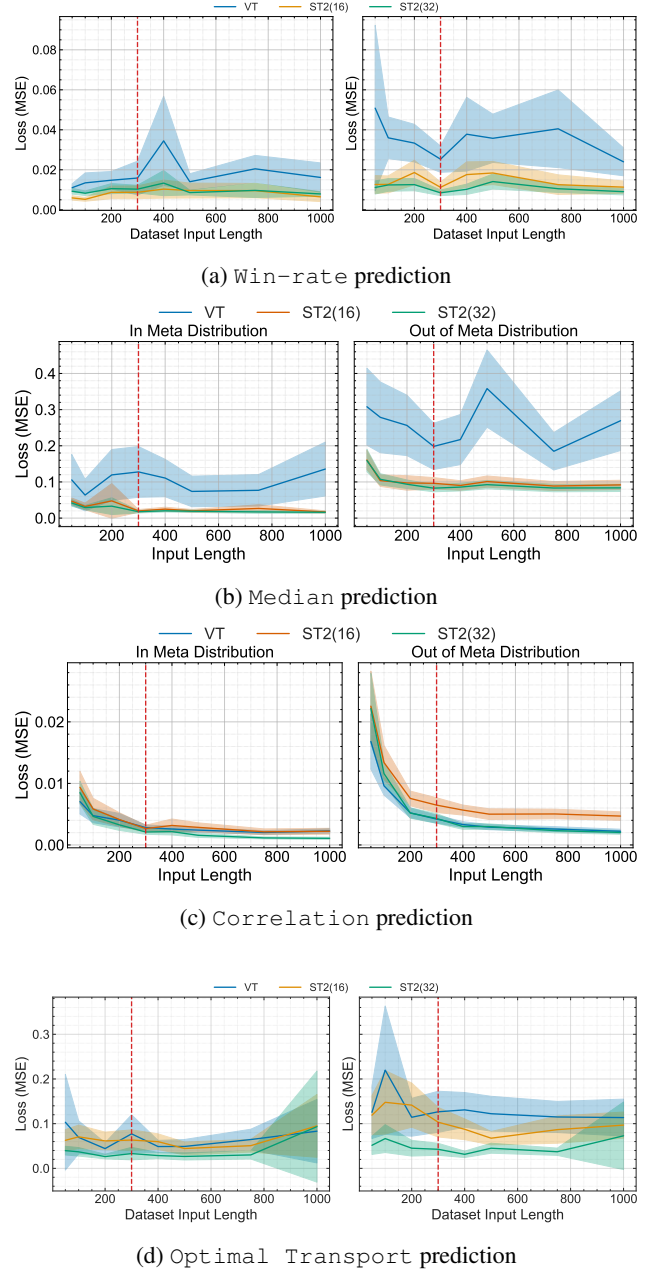


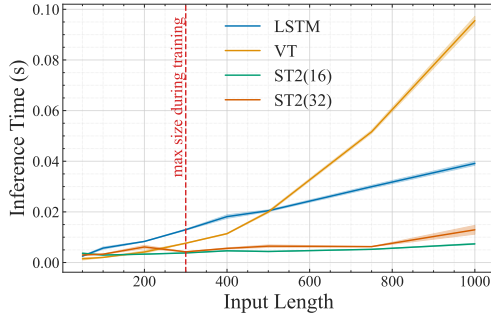
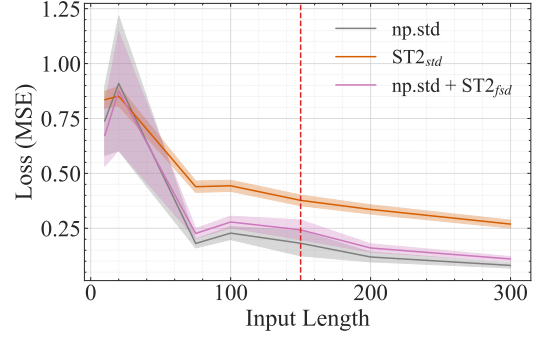
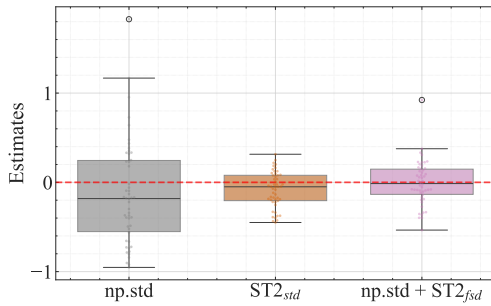
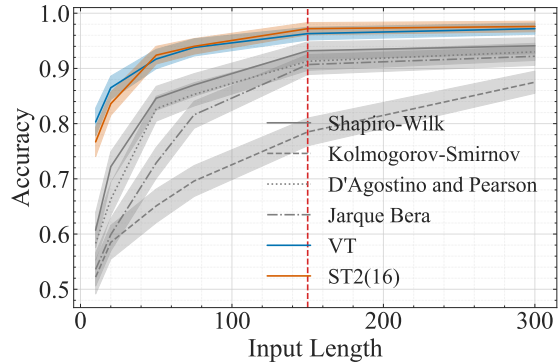
Figure 7: Generalization Across Dataset Lengths and Meta-Distributions. For each subplot, the left panel illustrates the performance of meta-statistical models on test datasets that vary in input length, including lengths not observed during training, while remaining within the training meta-distribution. For each subplot, the right panel presents the same comparison but for test datasets sampled from entirely new meta-distributions, with distributions unseen during training. Note that LSTM is excluded because its errors are an order of magnitude higher.

$n \in$	Bias			Variance		
	[10, 50]	[50, 100]	[100, 150]	[10, 50]	[50, 100]	[100, 150]
CCA	$-5.0e^{-2} \pm 3.7$	$-1.5e^{-2} \pm 1.2$	$-3.7e^{-2} \pm 1.8$	$8.6e^{-3} \pm 3.3$	$3.5e^{-3} \pm 1.2$	$1.3e^{-3} \pm 0.5$
KSG	$-1.0e^{-1} \pm 0.6$	$-2.6e^{-2} \pm 2.9$	$-2.9e^{-2} \pm 1.2$	$5.6e^{-3} \pm 1.5$	$5.6e^{-3} \pm 1.7$	$1.8e^{-3} \pm 0.7$
VT	$4.4e^{-3} \pm 10.8$	$4.0e^{-3} \pm 9.7$	$-7.8e^{-4} \pm 57.2$	$9.8e^{-3} \pm 7.1$	$3.1e^{-3} \pm 1.2$	$1.1e^{-3} \pm 0.5$
ST2(16)	$1.3e^{-2} \pm 1.0$	$9.8e^{-3} \pm 9.8$	$7.5e^{-4} \pm 60.7$	$7.7e^{-3} \pm 4.5$	$3.6e^{-3} \pm 1.4$	$9.8e^{-4} \pm 3.9$

Table 7: Bias and Variance of Mutual Information Estimators

	Bias			Variance			MSE		
	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$	$n \in [10, 50]$	$n \in [50, 100]$	$n \in [100, 150]$
np.std	$-1.0e^{-2} \pm 1.1$	$1.3e^{-2} \pm 0.2$	$-2.2e^{-2} \pm 0.7$	$1.9e^{-2} \pm 0.4$	$8.1e^{-3} \pm 1.6$	$3.8e^{-3} \pm 0.6$	$2.0e^{-2} \pm 0.2$	$8.2e^{-3} \pm 0.6$	$3.9e^{-3} \pm 0.3$
ST2 _{std}	$-1.2e^{-1} \pm 0.8$	$2.2e^{-3} \pm 9.7$	$-2.4e^{-2} \pm 0.9$	$8.6e^{-3} \pm 1.3$	$5.1e^{-3} \pm 0.6$	$3.4e^{-3} \pm 0.5$	$1.5e^{-2} \pm 0.1$	$7.5e^{-3} \pm 0.5$	$3.9e^{-3} \pm 0.3$
ST2 _{np.std} + ST2 _{fsd}	$2.0e^{-2} \pm 0.8$	$4.5e^{-2} \pm 1.6$	$9.2e^{-3} \pm 8.5$	$8.9e^{-3} \pm 1.6$	$5.7e^{-3} \pm 1.0$	$3.7e^{-3} \pm 0.6$	$1.4e^{-2} \pm 0.1$	$6.9e^{-3} \pm 0.5$	$3.9e^{-3} \pm 0.3$
np.std + ST2 _{fsd}	$1.4e^{-3} \pm 12.0$	$3.4e^{-2} \pm 0.7$	$4.7e^{-3} \pm 4.6$	$1.0e^{-2} \pm 0.1$	$5.8e^{-3} \pm 0.7$	$3.6e^{-3} \pm 0.5$	$1.4e^{-2} \pm 0.1$	$6.9e^{-3} \pm 0.5$	$3.6e^{-3} \pm 0.3$

Table 8: Bias and Variance of Standard Deviation Estimators


 Figure 8: **Inference time** comparison of meta-statistical models per batch as a function of input dataset length. Models have similar parameter counts $\approx 10K$).

 Figure 10: MSE of σ estimators as a function of dataset sizes, for the log-normal distribution.

 Figure 9: Estimate statistics for standard deviation estimators over 150 resampled datasets of size $n \in [10, 50]$ for the exponential distributions. Each dot represents the difference between an estimate and the true standard deviation. An unbiased estimator should be centered around zero.

 Figure 11: Accuracy of meta-statistical models compared to standard normality tests converted into classifiers using optimized p -value thresholds. The non-normal distributions are sampled in-meta-distribution for meta-statistical models.