

Interpretable by Design: Learning Predictors by Composing Interpretable Queries

Aditya Chattopadhyay, Stewart Slocum, Benjamin D. Haeffele,
René Vidal, *Fellow, IEEE* and Donald Geman, *Life Senior Member, IEEE*

Abstract—There is a growing concern about typically opaque decision-making with high-performance machine learning algorithms. Providing an explanation of the reasoning process in domain-specific terms can be crucial for adoption in risk-sensitive domains such as healthcare. We argue that machine learning algorithms should be interpretable by design and that the language in which these interpretations are expressed should be domain- and task-dependent. Consequently, we base our model’s prediction on a family of user-defined and task-specific binary functions of the data, each having a clear interpretation to the end-user. We then minimize the expected number of queries needed for accurate prediction on any given input. As the solution is generally intractable, following prior work, we choose the queries sequentially based on information gain. However, in contrast to previous work, we need not assume the queries are conditionally independent. Instead, we leverage a stochastic generative model (VAE) and an MCMC algorithm (Unadjusted Langevin) to select the most informative query about the input based on previous query-answers. This enables the online determination of a query chain of whatever depth is required to resolve prediction ambiguities. Finally, experiments on vision and NLP tasks demonstrate the efficacy of our approach and its superiority over post-hoc explanations.

Index Terms—Explainable AI, Interpretable ML, Computer Vision, Generative Models, Information Theory



1 INTRODUCTION

IN recent years, interpreting large machine learning models has emerged as a major priority, particularly for transparency in making decisions or predictions that impact human lives [1], [2], [3]. In such domains, understanding *how* a prediction is made may be as important as achieving high predictive accuracy. For example, medical regulatory agencies have recently emphasized the need for computational algorithms used in diagnosing, predicting a prognosis, or suggesting treatment for a disease, to explain why a particular decision was made [4], [5].

On the other hand, it is widely believed that there exists a fundamental trade-off in machine learning between interpretability and predictive performance [6], [7], [8], [9], [10]. Simple models like decision trees and linear classifiers are often regarded as *interpretable*¹ but at the cost of potentially reduced accuracy compared with larger *black box* models such as deep neural networks. As a result, considerable effort has been given to developing methods that provide *post-hoc* explanations of black box model predictions, i.e., given a prediction from a (fixed) model provide additional annotation or elaboration to explain how the prediction was made. As a concrete example, for image classification problems, one common family of post-hoc explanation methods produces attribution maps which seek to estimate the regions of the image that are *most important* for prediction. This is typically approached by attempting to capture the effect or sensitivity of perturbations to the input (or intermediate

features) on the model output [11], [12], [13], [14], [15], [16], [17], [18]. However, post-hoc analysis has been critiqued for a variety of issues [2], [19], [20], [21], [22], [23] (see also §2) and often fails to provide explanations in terms of concepts that are intuitive or interpretable for humans [24].

This naturally leads to the question of what an *ideal* explanation of a model prediction would entail; however, this is potentially highly *task-dependent* both in terms of the task itself as well as what the user seeks to obtain from an explanation. For instance, a model for image classification is often considered interpretable if its decision can be explained in terms of patterns occurring in salient parts of the image [25] (e.g., the image is a car because there are wheels, a windshield, and doors), whereas in a medical task explanations in terms of causality and mechanism could be desired (e.g., the patient’s chest pain and shortness of breath is likely not a pulmonary embolism because the blood D-dimer level is low, suggesting thrombosis is unlikely). Note that some words or patterns may be *domain-dependent* and therefore not interpretable to non-experts, and hence what is interpretable ultimately depends on the end user, namely the person who is trying to understand or deconstruct the decision made by the algorithm [26].

In addition to this *task-dependent* nature of model interpretation, there are several other desirable intuitive aspects of interpretable decisions that one can observe. The first is that meaningful interpretations are often *compositional* and can be constructed and explained from a set of *elementary units* [27]. For instance, words, parts of an image, or domain-specific concepts [28], [29], [30] could all be a suitable basis to form an explanation of a model’s prediction depending on the task. Moreover, the basic principle that simple and *concise* explanations are preferred (i.e., Occam’s razor) suggests that interpretability is enhanced when an explana-

• The authors are with the Mathematical Institute for Data Science and the Center for Imaging Science of The Johns Hopkins University, MD, 21218. E-mail: {achatto1, sslocum3, bhaeffele, rovidal, geman}@jhu.edu

1. Although later in the paper we will discuss situations in which even these simple models need not be interpretable.

tion can be composed from the smallest number of these elementary units as possible. Finally, we would like this explanation to be *sufficient* for describing model predictions, meaning that there should be no external variables affecting the prediction that are not accounted for by the explanation.

Inspired by these desirable properties, we propose a framework for learning predictors that are *interpretable by design*. The proposed framework is based on composing a subset of user-defined *concepts*, i.e., functions of the input data which we refer to as *queries*, to arrive at the final prediction. Possible choices for the set of queries Q based on the style of interpretation that is desired include:

- 1) **Salient image parts:** For vision problems, if one is interested in explanations in terms of salient image regions then this can be easily accomplished in our framework by defining the query set to be a collection of small patches (or even single pixels) within an image. This can be thought of as a generalization of the pixel-wise explanations generated by attribution maps.
- 2) **Concept-based explanations:** In domains such as medical diagnosis or species identification, the user might prefer explanations in terms of concepts identified by the community to be relevant for the task. For instance, a “Crow” is determined by the shape of the beak, color of the feathers, etc. In our framework, by simply choosing a query for each such concept, the user can easily obtain concept-based explanations (see Fig. 1(b)).
- 3) **Visual scene interpretation:** In visual scene understanding, one seeks a rich semantic description of a scene by accumulating the answers to queries about the existence of objects and relationships, perhaps generating a scene graph [31]. One can design a query set Q by instantiating these queries with trained classifiers. The answers to chosen queries in this context would serve as a semantic interpretation of the scene.
- 4) **Deep neuron-based explanations:** The above three examples are query sets based on domain knowledge. Recent techniques [30], [32], [33] have shown the ability of different neurons in a trained deep network to act as concept detectors. These are learnt from data by solving auxiliary tasks without any explicit supervisory signal. One could then design a Q in which each query corresponds to the activation level of a specific concept neuron. Such a query set will be useful for tasks in which it is difficult to specify interpretable functions/queries beforehand.

Given a user-specified set of queries Q , our framework makes its prediction by selecting a short sequence of queries such that the sequence of query-answer pairs provides a complete explanation for the prediction. More specifically, the selection of queries is done by first learning a generative model for the joint distribution of queries and output labels and then using this model to select the “most informative” queries for a given input. The final prediction is made using the Maximum A Posteriori (MAP) estimate of the output given these query-answer pairs. Fig. 1(a) gives an illustration of our proposed framework, where the task is to predict the bird species in an image and the queries are based on color, texture and shape attributes of birds. We argue that the sequence of query-answer pairs provides a meaningful

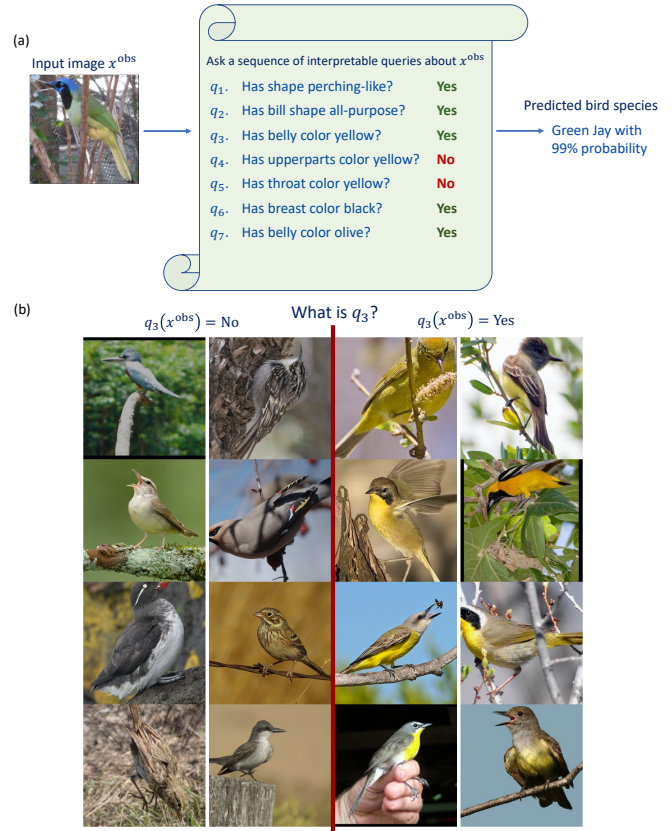


Fig. 1. (a) An illustration of our proposed learning framework. The prediction of a bird species is explained through a short sequence of interpretable queries, (q_1, q_2, \dots, q_7) , derived from a user-defined query set of domain-specific attribute for birds. (b) Interpretable queries. Each query in this case corresponds to a well-defined bird attribute. For instance, q_3 asks “Does the bird have belly color yellow?”. We visualize some example images which evaluate to “Yes” and observe that all of them correspond to birds with a yellow belly. Similarly, all images which evaluate to “No” corresponds to birds which do not have a yellow belly.

explanation to the user that captures the subjective nature of interpretability depending on the task at hand, and that is, by construction, compositional, concise and sufficient.

At first glance, one might think that classical decision trees [34], [35] based on Q could also produce interpretable decisions by design. However, the classical approach to determining decision tree branching rules based on the empirical distribution of the data is prone to over-fitting due to data fragmentation. Whereas random forests [36], [37] are often much more competitive than classical decision trees in accuracy [38], [39], [40], they sacrifice interpretability, the very property we want to hardwire into our decision algorithm. Similarly, the accuracy of a single tree can be improved by using deep networks to learn queries directly from data, as in Neural Decision Trees (NDTs) [41]. However, the opaqueness of the interpretation of these learnt queries makes the explanation of the final output, in terms of logical operations on the queries at the internal nodes, unintelligible. Figure 2 illustrates this with an example.

In this paper we make the following contributions;

- We propose a novel framework for prediction that is *interpretable by design*. We allow the end-user to specify a set Q of queries about input X and formulate learning as the problem of selecting a minimal set of queries

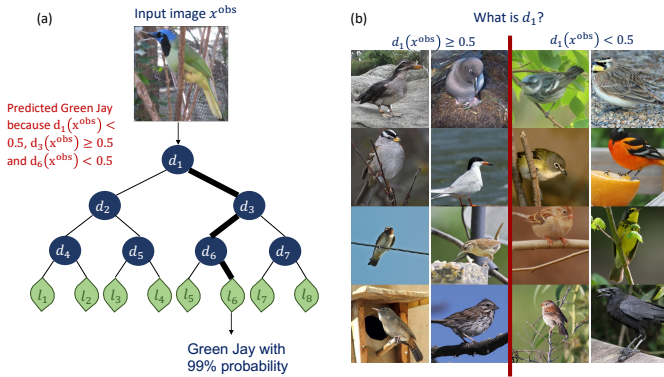


Fig. 2. **The interpretability of an explanation depends on how interpretable the queries are.** (a) An illustration of a Deep Neural decision tree [41] trained on the CUB-2011 dataset of bird images. The bold path denotes the trajectory the input image x^{obs} takes through the tree. Each d_i corresponds to an internal node of the tree and is a black-box function/query learnt from data. Each l_i denotes a leaf and computes the final classification for x^{obs} . The prediction can be explained as a conjunction of internal node functions, but is it really interpretable? (b) Example images that get routed to the left sub-tree ($d_1 \geq 0.5$) and right sub-tree ($d_1 < 0.5$). Notice that the interpretation of d_1 is not clear from these examples. Compare this to Fig. 1 where the semantics of each query is unambiguous to the end-user.

from Q whose answers are sufficient for predicting output Y . We formulate this query selection problem as an optimization problem over strategies that minimize the number of queries needed on average to predict Y from X . A prediction for Y is then made based on the selected query-answer pairs, which provide an explanation for the prediction that is by construction interpretable. The set of selected query-answer pairs can be viewed as a *code* for the input. However, a major difference between our framework and coding theory is that, due to the constraint of interpretability, Q is a vanishingly small collection of the functions of X , whereas coding theory typically considers Q to be all possible binary functions of X .

- Since computing the exact solution to our optimization problem is computationally challenging, we propose to greedily select a minimal set of queries by using the *Information Pursuit* (IP) algorithm [42]. IP sequentially selects queries in order of maximum *information gain* until enough evidence is gathered from the query-answer pairs to predict Y . This sequence of query-answer pairs serves as the *explanation* for predicting Y from X . To ameliorate the computational challenge of computing information gain for high-dimensional input and query spaces, prior work [42] had assumed that query answers were conditionally independent given Y , an assumption that is largely inadequate for most prediction tasks we encounter in practice. In this paper, we propose a latent variable graphical model for the joint distribution of queries and outputs, $p(Q(X), Y)$, and learn the required distributions using Variational Autoencoders (VAEs). We then use the Unadjusted Langevin Algorithm (ULA) to generate samples required to carry out IP. This gives us a tractable algorithm for any task and query set. To the best of our knowledge, ours is the first implementation of IP that

uses deep generative models and does *not* assume that query answers are conditionally independent given Y .

- Finally, we demonstrate the utility of our framework on various vision and NLP tasks. In binary image classification using MNIST, Fashion-MNIST & KMNIST, and bird species identification using CUB-200, we observe that IP finds succinct explanations which are highly predictive of the class label. We also show, across various datasets, that the explanations generated by our method are shorter and more predictive of the class label than state-of-the-art post-hoc explanation methods like Integrated Gradients and DeepSHAP.

2 RELATED WORK

Methods for interpretable deep learning can be separated into those that seek to explain existing models (post-hoc methods) and those that build models that are interpretable by design. Because they do not negatively impact performance and are convenient to use, post-hoc explanations have been the more popular approach, and include a great diversity of methods.

Saliency maps estimate the contribution of each feature through first-order derivatives [11], [12], [16], [17], [43]. Linear perturbation-based methods like LIME [44] train a linear model to locally approximate a deep network around a particular input, and use the coefficients of this model to estimate the contribution of each feature to the prediction. Another popular set of methods use game-theoretic Shapley values as attribution scores, estimating feature contributions by generating predictions on randomly sampled subsets of the input [45]. We provide quantitative comparisons between IP and these methods in Section 5.1.2. Recently, there has been interest in concept-based analogues of these methods that leverage similar approaches to measure the sensitivity of a prediction to high-level, human-friendly concepts as opposed to raw features [46], [47], [48].

Despite certain advantages, what all the above post-hoc methods have in common is that they come with little guarantee that the explanations they produce actually reflect how the model works [2]. Indeed, several recent studies [18], [19], [20], [21], [22] call into question the veracity of these explanations towards the trained model. Adebayo *et al.* [19] show that several popular attribution methods act similar to edge detectors and are insensitive to the parameters of the model they attempt to explain! Yang *et al.* [20] find that these methods often produce false-positive explanations, assigning importance to features that are irrelevant to the prediction of the model. It is also possible to adversarially manipulate post-hoc explanations to hide any spurious biases the trained model might have picked up from data [23].

Interpretability by design. These issues have motivated recent work on deep learning models which are *interpretable by design*, i.e., constrained to produce explanations that are faithful to the underlying model, albeit with varying conceptions of “faithfulness”. Several of these models are constructed so they behave similarly to or can be well-approximated by a classically interpretable model, such as a linear classifier [49], [50] or a decision tree [51]. This allows for an approximately faithful explanation in raw feature space. In a similar vein, Pillai & Pirsiavash [52] fix

a post-hoc explanation method (e.g. Grad-CAM [16]), and regularize a model to generate consistent explanations with the chosen post-hoc method. However, our method does not just behave *like* a fully interpretable model or generate *approximately* faithful explanations, but rather it produces explanations that are guaranteed to be faithful and fully explain a given prediction.

Another approach to building interpretable models by design is to generate explanations in terms of high-level, interpretable concepts rather than in raw feature space, often by applying a linear classifier to a final latent space of concepts [25], [49], [53]. However these concepts are learned from data, and may not align with the key concepts identified by the user. For example, Prototypical Part Networks [25] take standard convolutional architectures and insert a “prototype layer” before the final linear layer, learning a fixed number of visual concepts that are used to represent the input. This allows the network to explain a prediction in terms of these “prototype” concepts. Since these prototypes are learned embeddings, there is no guarantee that their interpretation will coincide with the user’s requirements. Furthermore, these explanations may require a very large number of concepts, while in contrast, we seek minimal-length explanations to preserve interpretability.

Attention-based models are another popular family of models that are sometimes considered interpretable by design [54], [55]. However, attention is only a small part of the overall computation and can be easily manipulated to hide model biases [56]. Moreover, the attention coefficients are not necessarily a sufficient statistic for the model prediction.

Perhaps most similar to our work are Concept Bottleneck Networks [24], which first predict an intermediate set of human-specified concepts c and then use c to predict the final class label. Nevertheless, the learnt mapping from concepts to labels is still a black-box. To remedy this, the authors suggest using a linear layer for this mapping but this can be limiting since linearity is often an unrealistic assumption [27]. In contrast, our framework makes no linearity assumptions about the final classifier and the classification is explainable as a sequence of interpretable query-answer pairs obtained about the input (see Fig. 1(a)).

Neural networks and decision trees. Unlike the above methods, which can be thought of as deep interpretable linear classifiers, our method can be described as a deep decision tree that branches on responses to an interpretable query set. Spanning decades, there has been a variety of work building decision trees from trained neural networks [29], [57], [58], [59] and using neural networks within nodes of decision trees [41], [60], [61], [62]. Our work differs from these in three important aspects. First, rather than allowing arbitrary splits, we branch on responses to an interpretable query set. Second, instead of using empirical estimates of information gains based on training data (which inevitably encounter data-fragmentation [63] and hence overfitting), or using heuristics like agglomerative clustering on deep representations [29], we calculate information gain from a generative model, leading to strong generalization. Third, for a given input, say x^{obs} , we use a generative model to compute the queries along the branch traversed by x^{obs} in an online manner. The entire tree is never constructed. This

allows for much very deep terminal nodes when necessary to resolve ambiguities in prediction. As an example, for the task of topic classification using the HuffPost dataset (§5.0.3), our framework asks about 199 queries (on average) before identifying the topic. Such large depths are impossible in standard decision trees due to memory limitations.

Information bottleneck and minimal sufficient statistics. The problem of finding minimal-length, task-sufficient codes is not new. For example, the *information bottleneck* method [64] seeks a minimum-length encoding for X that is (approximately) sufficient to solve task Y . Our concept of description length differs in that we constrain the code to consist of interpretable query functions rather than *all functions* of the input, as in the information bottleneck and classical information theory. Indeed, arbitrary subsets of the input space (e.g. images) are overwhelmingly *not* interpretable to humans.

Sequential active testing and hard attention. The *information pursuit* (IP) algorithm we use was introduced in [42] under the name “active testing,” which sequentially observes parts of an input (rather than the whole input at once), using mutual information to determine “where to look next,” which is calculated online using on a scene model. Sequentially guiding the selection of partial observations has also been independently explored in Bayesian experimental design [65]. Subsequent works in these two areas include many ingredients of our approach (e.g. generative models [31], [66] and MCMC algorithms [67]). Of particular interest is the work of Branson *et al.* [68] which used the CUB dataset to identify bird species by sequentially asking pose and attribute queries to a human user. They employ IP to generate the query sequence based on answers provided by the user, much like our experiments in §5.0.2. However, for the sake of tractability, all the above works assume that query answers are independent conditioned on Y . We do not. Rather, to the best of our knowledge, ours is the first implementation of the IP algorithm that uses deep generative models and only assumes that queries are independent given Y and some latent variable Z . This greatly improves performance, as we show in §5.

The strategy of inference through sequential observations of the input has been recently re-branded in the deep learning community as *Hard Attention* [69], [70], [71]. However, high variance in gradient estimates and scalability issues have prevented widespread adoption. In the future, we wish to explore how our work could inform more principled and better-performing reward functions for Hard Attention models.

Visual question answering. Although it may appear that our work is also related to the Visual Question Answering (VQA) literature [72], [73], [74], [75], [76], [77], we note that our work addresses a very different problem. VQA focuses on training deep networks for *answering* a large set of questions about a visual scene. In contrast, our framework is concerned with *selecting* a small number of queries to ask about a given image to solve a task, say classification. As we move on to more complex tasks, an interesting avenue for future work would involve using VQA systems to supply answers to the queries used in our framework. However,

this would require significantly more complex generative models than the ones considered here.

3 LEARNING INTERPRETABLE PREDICTORS BY COMPOSING QUERIES VIA INFORMATION PURSUIT

Let X and Y be the input data and the corresponding output/hypothesis, both random variables assuming values in \mathcal{X} and \mathcal{Y} respectively. In supervised learning, we seek to infer Y from X using a finite set of samples drawn from the joint distribution $p_{XY}(x, y)$.² As motivated in Section 1, useful explanations for prediction should be *task-dependent*, *compositional*, *concise* and *sufficient*. We capture such properties through a suitably rich set Q of binary functions $q(x)$, or *queries*, whose answers $\{q(x)\}_{q \in Q}$ collectively determine the task Y . More precisely, a query set Q is *sufficient* for Y if

$$p(y | x) = p(y | \{x' \in \mathcal{X} : q(x') = q(x) \forall q \in Q\}). \quad (1)$$

In other words, Q is sufficient for Y if whenever two inputs x and x' have identical answers for all queries in Q , their corresponding posteriors are equal, i.e., $p(y | x) = p(y | x')$.

Given a fixed query set Q , how do we compose queries into meaningful representations that are predictive of Y ? We answer this by first formally defining an explanation strategy π and then formulating the task of composing queries as an optimization problem.

Explanation strategies based on composing queries. An *explanation strategy*, or just *strategy*, is a function, $\pi : K^* \rightarrow Q$, where K^* is the set of all finite-length sequences generated using elements from the set $K = \{(q, q(x)) \mid q \in Q, x \in \mathcal{X}\}$ of query-answer pairs. We require that Q contains a special query, q_{STOP} , which signals the strategy to stop asking queries and output $expl_Q^\pi(x)$, the set of query-answer pairs asked before q_{STOP} . More formally, a strategy π is recursively defined as follows; given input sample x^{obs}

- 1) $q_1 = \pi(\emptyset)$. The first query is independent of x^{obs} .
- 2) $q_{k+1} = \pi(\{q_i, q_i(x^{obs})\}_{1:k})$. All subsequent queries depend on the query-answer pairs observed so far for x^{obs} .
- 3) If $q_{L+1} = q_{STOP}$ terminate, and return

$$expl_Q^\pi(x^{obs}) := \{q_i, q_i(x^{obs})\}_{1:L}. \quad (2)$$

Notice that each q_i depends on x^{obs} , but we drop this dependency in the notation for brevity. We call the number of pre-STOP queries for a particular x^{obs} as the explanations' description length and denote it by $t^\pi(x^{obs}) := |expl_Q^\pi(x^{obs})|$. Computing a strategy on x^{obs} is thus akin to traversing down the branch of a decision tree dictated by x^{obs} . Each internal node encountered along this branch computes the query proposed by the strategy based on the path (query-answer pairs) observed so far.

Notice also that we restrict our attention to *sequential* strategies so that the resulting explanations satisfy the property of being *prefix-free*.³ This means that explanations generated for predictions made on an input signal x_1 cannot be a sub-part for explanations generated for predictions on a different input signal x_2 ; otherwise, the explanation

2. We denote random variables by capital letters and their realizations with small letters.

3. The term prefix-free comes from the literature on instantaneous codes in information theory.

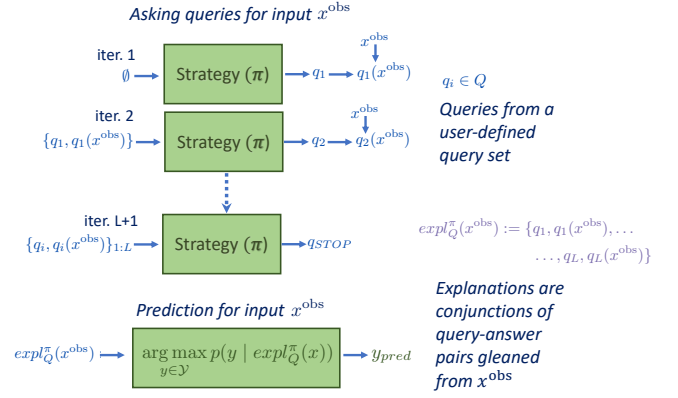


Fig. 3. Schematic view of the overall framework for quantifying explanations for predicting y from x^{obs} . For details see Sec. 3.

procedure is ambiguous because a terminal node carrying one label could be an internal node of a continuation leading to a different label. Sequential strategies generate prefix-free explanations by design. For non-sequential strategies, which are just functions mapping an input X to a set of queries in Q , it is not clear how to effectively encode the constraint of generating prefix-free explanations.

Concise and approximately sufficient strategies. In machine learning, we are often interested in solving a task *approximately* rather than *exactly*. Let Q be sufficient for Y , choose a distance-like metric d on probability distributions and let $\epsilon > 0$. We propose the following optimization problem to efficiently compose queries for prediction,

$$\begin{aligned} \min_{\pi} \mathbb{E}_X [|expl_Q^\pi(X)|] &=: H_Q^\epsilon(X; Y) \\ \text{s.t. } \mathbb{E}_X [d(p(Y | X), p(Y | expl_Q^\pi(X)))] &\leq \epsilon \quad (\epsilon\text{-Sufficiency}), \end{aligned} \quad (3)$$

where the minimum is taken over all strategies π . The solution π^* to (3) provides a criterion for an optimal strategy for the task of inferring Y approximately from X . The *minimal expected description length* objective, $H_Q^\epsilon(X; Y)$, ensures the conciseness of the explanations, while the constraint ensures approximate sufficiency of the explanation. The "metric" d on distributions could be KL-divergence, total variation, Wasserstein distance, etc. The hyper-parameter ϵ controls how approximate the explanations are. The posterior $p(y | expl_Q^\pi(x^{obs}))$ should be interpreted as the conditional probability of y given the event

$$[x^{obs}]_{\pi, Q} := \{x \in \mathcal{X} \mid expl_Q^\pi(x) = expl_Q^\pi(x^{obs})\}. \quad (4)$$

$expl_Q^\pi(X)$ can also be interpreted as a random variable which maps input X to its equivalence class $[X]_{\pi, Q}$.

The final prediction/inference for the input x^{obs} is then taken to be the usual MAP estimator, namely

$$y_{pred} = \arg \max_{y \in \mathcal{Y}} p(y | expl_Q^\pi(x^{obs})). \quad (5)$$

The sequence of query-answers streams obtained by π on x^{obs} serves as the explanation for y_{pred} . One could also monitor the posterior over the labels Y evolving as successive queries get asked to gain more insight into the strategy's decision-making process. Fig. 3 illustrates the overall framework in detail.

Information Pursuit: a greedy approximation. Unfortunately, solving (3) is known to be NP-Complete and hence generally intractable [78]. As an approximate solution to (3) we propose to use a greedy algorithm called Information Pursuit (IP). IP was introduced by Geman & Jedynek in 1996 [42] as a model-based, online construction of a single but deep branch. The IP strategy, that is, $\pi = \text{IP}$, is recursively defined as follows,

$$q_1 = \text{IP}(\emptyset) = \arg \max_{q \in Q} I(q(X); Y) \quad (6)$$

$$q_{k+1} = \text{IP}(\{q_i, q_i(x^{\text{obs}})\}_{1:k}) = \arg \max_{q \in Q} I(q(X); Y \mid S_k^{\text{IP}}(x^{\text{obs}}))$$

where I denotes mutual information and $S_k^{\text{IP}}(x^{\text{obs}})$ corresponds to the event $\{x \in \mathcal{X} \mid \{q_i, q_i(x^{\text{obs}})\}_{1:k} = \{q_i, q_i(x)\}_{1:k}\}$. Ties in choosing q_{k+1} are broken arbitrarily if the maximum is not unique.

The algorithm stops when there are no more informative queries left in Q , that is, it satisfies the following condition:

$$q_{L+1} = q_{STOP} \quad \text{if} \quad \max_{q \in Q} I(q(X); Y \mid S_m^{\text{IP}}(x^{\text{obs}})) \leq \epsilon \\ \forall m \in \{L, L+1, \dots, L+T\}, \quad (7)$$

where hyper-parameter $T > 0$ is chosen via cross-validation. This termination criteria corresponds to taking the distance-like metric d in (3) as the KL-divergence between the two distributions. Further details about the relation between this termination criteria and the ϵ -Sufficiency constraint in (3) are provided in Appendix A.3. For tasks in which Y is a function of X , a common scenario in many supervised learning problems, we use a simpler alternative,

$$q_{L+1} = q_{STOP} \quad \text{if} \quad \arg \max_{y \in \mathcal{Y}} p(y \mid S_m^{\text{IP}}(x^{\text{obs}})) \geq 1 - \epsilon \\ \forall m \in \{L, L+1, \dots, L+T\}. \quad (8)$$

The key distinction between the information gain criteria used in standard decision tree induction and IP is that the former uses the empirical distributions to compute (6) while the latter is based on generative models (as we will see in Section 4). The use of generative models guards against data fragmentation [63] and thus allows for asking longer sequences of queries without grossly over-fitting.

How does IP compare to the optimal strategy π^* ? We begin by characterizing the constraint in (3) in terms of mutual information, the quantity that drives IP.

Proposition 1. *Let $S_k^\pi(X)$ be a random variable where any realization $S_k^\pi(x^{\text{obs}})$, $x^{\text{obs}} \in \mathcal{X}$, denotes the event*

$$S_k^\pi(x^{\text{obs}}) := \{x' \in \mathcal{X} \mid \{q_i, q_i(x^{\text{obs}})\}_{1:k} = \{q_i, q_i(x')\}_{1:k}\},$$

where q_i is the i^{th} query selected by π for input x^{obs} . Here we use the convention that $S_0^\pi(X) = \Omega$ (the entire sample space) and $S_l^\pi(X) = S_{t^\pi(X)}^\pi(X) \quad \forall l > t^\pi(X)$. If Q is finite⁴ and d is taken to be the KL-divergence, then objective (3) can be rewritten as

$$H_Q^\epsilon(X; Y) := \min_{\pi} \mathbb{E}_X [\text{expl}_Q^\pi(X)] \\ \text{s.t.} \quad \sum_{k=1}^{\tau^\pi} I(Y; S_k^\pi(X) \mid S_{k-1}^\pi(X)) \geq I(X; Y) - \epsilon, \quad (9)$$

4. The assumption of Q being a finite set is benign. Many interested applications can be addressed with a finite Q as we show in our experiments.

where $\tau^\pi = \max\{t^\pi(x) : x \in \mathcal{X}\}$ and $t^\pi(X)$ is defined as the number of queries selected by π for input X until q_{STOP} .

See Appendix A.1 for a detailed proof. The objective in (9) can be alternatively stated as,

$$\max_{\pi} \sum_{k=1}^{\tau^\pi} I(Y; S_k^\pi(X) \mid S_{k-1}^\pi(X)) \\ \text{s.t.} \quad \mathbb{E}_X [\text{expl}_Q^\pi(X)] \leq \gamma, \quad (10)$$

where $\gamma > 0$ is a user-defined hyper-parameter. From (10) it is clear that the optimal strategy π^* would ask a sequence of queries about X that would maximize the cumulative sum of the mutual information each additional query provides about Y , conditioned on the history of query-answers observed so far, subject to a constraint on the average number of queries that can be asked. As stated before, solving for π^* is infeasible but a greedy approximation that makes locally optimal choices is much more amenable.

Suppose that one has been given the answers to k queries about a given input, the locally optimal choice would then be to ask the most informative query about Y conditioned on the history of these k query-answers observed. This greedy choice at each stage gives rise to the IP strategy. Obtaining approximation guarantees for IP is still an open problem; however in the special case where Q is taken to be the set of all possible binary functions of X , it is possible to show that IP asks at most 1 query more than π^* on average. More formally, we have the following result, whose proof can be found in Appendix A.2.

Proposition 2. *Let Y be discrete. Let $\tilde{H}_Q(X; Y)$ be the expected description length obtained by the IP strategy. If $H(Y|X) = 0$ and Q is the set of all possible binary functions of X such that $H(q(X) \mid Y) = 0 \quad \forall q \in Q$, then*

$$H(Y) \leq \tilde{H}_Q(X; Y) \leq H(Y) + 1 \quad (11)$$

Having posed the problem of finding explanations as an optimization problem and proposed a greedy approximation to solving it, in the next section we propose a tractable implementation of IP based on deep generative models.

4 INFORMATION PURSUIT USING VARIATIONAL AUTOENCODERS AND UNADJUSTED LANGEVIN

IP requires probabilistic models relating query-answers and data to compute the required mutual information terms in (6). Specifically, computing q_{k+1} in (6) (for any iteration number k) requires computing the mutual information between $q(X)$ and Y given the history $S_k^{\text{IP}}(x^{\text{obs}})$ till time k . As histories become longer, we quickly run out of samples in our dataset which belong to the event $S_k^{\text{IP}}(x^{\text{obs}})$. As a result, non-parametric sample-based methods to estimate mutual information (such as [79]) would be impractical. In this section, we propose a model-based approach to address this challenge for a general supervised learning task and query set Q . In §5 we adapt this model to the specific cases where Q is taken to be image patches or task-based concepts.

Information Pursuit Generative Model. To make learning tractable, we introduce latent variables Z to account for all

the dependencies between different query-answers, and we posit the following factorization of $Q(X), Y, Z$

$$\begin{aligned} p_{Q(X)ZY}(Q(x), z, y) & \\ &= \prod_{q \in Q} p_{q(X)|ZY}(q(x) | z, y) p_Y(y) p_Z(z), \end{aligned} \quad (12)$$

where $Q(X) = \{q(X) : q \in Q\}$, and z and $q(x)$ denote realizations of Z and $q(X)$ respectively. In other words, we assume that the query-answers are conditionally independent given the label y and a latent vector z . The independence assumption in (12) shows up ubiquitously in many machine learning applications, such as the following.

- 1) **$q(X)$ as object presence indicators evaluated at non-overlapping windows:** Let Q be a set of non-overlapping windows in the image X with $q(X)$ being a random variable indicating the presence of an object at the q^{th} location. The correlation between the qs is entirely due to latent image generating factors Z , such as lighting, camera position, scene layout, and texture along with the scene description signal Y .
- 2) **$q(X)$ as snippets of speech utterances:** A common assumption in speech recognition tasks is that the audio frame features ($q(X)$) are conditionally independent given latent phonemes Z (which is often modeled as a Hidden Markov Model).

The latent space Z is often a lower-dimensional space compared to the original high-dimensional X . We learn Z from data in an unsupervised manner using variational inference. Specifically, we parameterize the distributions $\{p_{\omega}(q(x) | z, y) \forall q \in Q\}$ with a **Decoder Network** with shared weights ω . These weights are learned using stochastic Variational Bayes [80] by introducing an approximate posterior distribution $p'_{\phi}(z | y, Q(x))$ parameterized by another neural network with weights ϕ called the **Encoder Network** and priors $p_Y(y)$ and $p_Z(z)$. More specifically, the parameters ϕ and ω are learned by maximizing the Evidence Lower Bound (ELBO) objective. Appendix A.7 gives more details on this optimization procedure. The learned Decoder Network $p_{\omega^*}(q(x) | z, y)$ is then used as a plug-in estimate for the true distribution $p_{q(X)|ZY}(q(x) | z, y)$, which is in turn used to estimate (12).

Implementing IP using the generative model. Once the Decoder Network has been learned using variational inference, the first query $q_1 = \text{IP}(\emptyset)$ is the one that maximizes the mutual information with Y as per (6). The mutual information term for any query q is completely determined by $p(q(x), y)$, which is obtained by numerically marginalizing the nuisances Z from (12) using Monte Carlo integration. In particular, we carry out the following computation $\forall q \in Q$,

$$\begin{aligned} p_{q(X)Y}(q(x), y) &= \int_z p_{Q(X)ZY}(Q(x), z, y) dz \\ &= \int_z p_{q(X)|ZY}(q(x) | z, y) p_Y(y) p_Z(z) dz \\ &\approx \frac{1}{N} \sum_{i=1}^N p_{\omega^*}(q(x) | y, z^{(i)}) p_Y(y) \\ &=: \tilde{p}(q(x), y). \end{aligned} \quad (13)$$

In the last approximation, $p_{\omega^*}(q(x) | y, z^{(i)})$ is the distribution obtained using the trained decoder network. N is the

number of i.i.d. samples drawn and $z^i \sim p_Z(z)$. We then estimate mutual information numerically via the following formula,

$$I(Y; q(X)) = \sum_{q(x), y} \tilde{p}(q(x), y) \log \frac{\tilde{p}(q(x), y)}{\tilde{p}(q(x)) \tilde{p}(y)}. \quad (14)$$

The computation of subsequent queries q_{k+1} requires the mutual information conditioned on observed history $S_k^{\text{IP}}(x^{\text{obs}})$, which can be calculated from the distribution

$$\begin{aligned} p(q(x), y | S_k^{\text{IP}}(x^{\text{obs}})) & \\ &= \int p(q(x), z, y | S_k^{\text{IP}}(x^{\text{obs}})) dz \\ &= \int p(q(x) | z, y, S_k^{\text{IP}}(x^{\text{obs}})) p(z | y, S_k^{\text{IP}}(x^{\text{obs}})) p(y | S_k^{\text{IP}}(x^{\text{obs}})) dz \\ &= \int p(q(x) | z, y) p(z | y, S_k^{\text{IP}}(x^{\text{obs}})) p(y | S_k^{\text{IP}}(x^{\text{obs}})) dz. \end{aligned} \quad (15)$$

The first equality is an application of the law of total probability. The last equality appeals to the assumption that $\{q(X), q \in Q\}$ are conditionally independent given Y, Z (12).

To estimate the right-hand side of (15) via Monte Carlo integration, one needs to sample $z^i \sim p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$ and compute

$$\begin{aligned} p(q(x), y | S_k^{\text{IP}}(x^{\text{obs}})) &\approx \tilde{p}(q(x), y | S_k^{\text{IP}}(x^{\text{obs}})) \\ &:= \frac{1}{N} \sum_{i=1}^N p_{\omega^*}(q(x) | z^{(i)}, y) p(y | S_k^{\text{IP}}(x^{\text{obs}})), \end{aligned} \quad (16)$$

where the term $p(y | S_k^{\text{IP}}(x^{\text{obs}}))$ is estimated recursively via the Bayes' theorem. This computation is as follows,

$$\begin{aligned} p(y | S_k^{\text{IP}}(x)) &\propto p(y, S_k^{\text{IP}}(x)) \\ &= p(q_k(x), y, S_{k-1}^{\text{IP}}(x)) \\ &\propto p(q_k(x) | y, S_{k-1}^{\text{IP}}(x)) p(y | S_{k-1}^{\text{IP}}(x)) \end{aligned} \quad (17)$$

$S_0^{\text{IP}}(x) = \emptyset$ (since no evidence via queries has been gathered from x yet) and so $p(y | S_0^{\text{IP}}(x)) = p_Y(y)$. The posterior $p(y | S_k^{\text{IP}}(x))$ is obtained by normalizing the last equation in (17) such that $\sum_y p(y | S_k^{\text{IP}}(x)) = 1$. This recursive updating of the posterior is similar to the posterior updates used in Bayesian sequential filtering [81]. The term $p(q_k(x) | y, S_{k-1}^{\text{IP}}(x))$ is estimated using (16).

Having estimated $p(q(x), y | S_k^{\text{IP}}(x^{\text{obs}}))$, we then numerically compute the mutual information between query-answer $q(X)$ and Y given history for every $q \in Q$ via the formula

$$\begin{aligned} I(Y; q(X) | S_k^{\text{IP}}(x^{\text{obs}})) &= \\ &\sum_{q(x), y} \tilde{p}(q(x), y | S_k^{\text{IP}}(x^{\text{obs}})) \log \frac{\tilde{p}(q(x), y | S_k^{\text{IP}}(x^{\text{obs}}))}{\tilde{p}(q(x) | S_k^{\text{IP}}(x^{\text{obs}})) \tilde{p}(y | S_k^{\text{IP}}(x^{\text{obs}}))}. \end{aligned} \quad (18)$$

Estimating $p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$ with the Unadjusted Langevin Algorithm. Next we describe how to sample from this posterior $p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$ using the Unadjusted Langevin Algorithm (ULA). ULA is an iterative algorithm used to approximately sample from any distribution with a density known only up to a normalizing factor. It has been successfully applied to many high-dimensional Bayesian

inference problems [82], [83], [84]. Given an initialization $z^{(0)}$, ULA proceeds by

$$z^{(i+1)} = z^{(i)} + \eta \nabla U(z^{(i)}) + \sqrt{2\eta} \zeta^{(i+1)}. \quad (19)$$

Here $(\zeta^{(i)})_{i \geq 1} \sim \mathcal{N}(0, I)$ and η is the step-size. Asymptotically, the chain $(z^{(i)})_{i \geq 1}$ converges to a stationary distribution that is “approximately” equal to a measure with density $\propto e^{U(z)}$ [85].

For IP, we need samples from $p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$. This is achieved by initializing $z^{(0)}$ using the last iterate of the ULA chain used to simulate $p(z | y, S_{k-1}^{\text{IP}}(x^{\text{obs}}))$.⁵ We then run ULA for N iterations by recursively applying (19) with $U(z) := \log p(z, S_k^{\text{IP}}(x^{\text{obs}}) | y) = \log p(S_k^{\text{IP}}(x^{\text{obs}}) | z, y) p(z) p(y)$.

The number of steps N is chosen to be sufficiently large to ensure the ULA chain converges “approximately” to the desired $z \sim p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$. We use the trained decoder network $\prod_{i=1}^k p_\omega(q_i(x) | z, y)$, with q_i being the i^{th} query asked by IP for input x , as a proxy for $p(S_k^{\text{IP}}(x^{\text{obs}}) | z, y)$. We then obtain stochastic approximations of (15) by time averaging the iterates,

$$p(q(x), y | S_k^{\text{IP}}(x^{\text{obs}})) \approx \frac{1}{N} \sum_{i=1}^N p_\omega(q(x) | z^{(i)}, y) p(y | S_k^{\text{IP}}(x^{\text{obs}})), \quad (20)$$

where $(z^{(i)})_{1:N}$ are the iterates obtained using the ULA chain whose stationary distribution is “approximately” $p(z | y, S_k^{\text{IP}}(x^{\text{obs}}))$.

Algorithmic complexity for IP. For any given input x , the per-iteration cost of the IP algorithm is $\mathcal{O}(N + |Q|m)^6$, where $|Q|$ is the total number of queries, N is the number of ULA iterations, and m is cardinality of the product sample space $q(X) \times Y$. For simplicity we assume that the output hypothesis Y and query-answers $q(X)$ are finite-valued and also that the number of values query answers can take is the same. However, our framework can handle more general cases. See Appendix A.6 for more details.

5 EXPERIMENTS

In this section, we empirically evaluate the effectiveness of our method. We begin by analyzing the explanations provided by IP for classifying individual input data, in terms of words, symbols, or patterns (the queries). We find in each case that IP discovers concise explanations which are amenable to human interpretation. We then perform quantitative comparisons which show that (i) IP explanations are more faithful to the underlying model than existing attribution methods; and (ii) the predictive accuracy of our method using a given query set is competitive with black-box models trained on features provided by the same set.

5.0.1 Binary Image Classification with Patch Queries

Task and query set. We start with the simple task of binary image classification. We consider three popular datasets – MNIST [86], Fashion-MNIST [87] and KMNIST [88]. We choose a threshold for binarizing these datasets since they

are originally grayscale. We choose the query set Q as the set of all $w \times w$ overlapping patch locations in the image. The answer $q(X)$ for any $q \in Q$ is the w^2 pixel intensities observed at the patch indexed by location q . This choice of Q reflects the user’s desire to know which parts of the input image are most informative for a particular prediction, a common practice for explainability in vision tasks [25]. We conduct experiments for multiple values of w and conclude that $w = 3$ provides a good trade-off between the required number of queries and the interpretability of each query. Note that when $w > 1$ the factorization in (12) that we use to model $p(Q(x), y, z)$ and compute mutual information no longer holds as the overlapping queries $q(X)$ are now *causally related* (and therefore dependent even when conditioned on Z , making them unable to be modeled by a VAE). So instead of training a VAE to directly model the query set $p(Q(x) | y, z)$, we train a VAE to model the pixel distribution $p(x | y, z)$, and then compute the probability distribution over the patch query $p(q(x) | z, y)$ as the product of the probabilities of all pixels in that patch.⁷

IP in action. Fig. 4(a) illustrates the decision-making process of IP using 3×3 patch queries on an image x^{obs} of a 6 from the MNIST test set. The first query is near the center of the image; recall from (6) that this choice is independent of the particular input image and represents the patch whose pixel intensities have maximum mutual information with Y (the class label). The updated posterior, $p(Y | S_1^{\text{IP}}(x^{\text{obs}}))$, concentrates most of its mass on the digit “1”, perhaps because most of the other digits do not commonly have a vertical piece of stroke at the center patch. However, the next query (about three pixels below the center patch) reveals a horizontal stroke and the posterior mass over the labels immediately shifts to $\{2, 3, 6, 8\}$. The next two queries are well-suited to discerning between these four possibilities and we see that after asking 4 questions, IP is more than 90% confident that the image is a 6. Such rich explanations in terms of querying informative patches based on what is observed so far and seeing how the belief $p(Y | S_k^{\text{IP}}(x^{\text{obs}}))$ of the model evolves over time is missing from post-hoc attribution methods which output static importance scores for every pixel towards the black-box model’s final prediction.

Explanation length vs. task complexity. Fig. 6 shows that IP requires an average of 5.2, 12.9 and 14.5 queries of size 3×3 to predict the label with 99% confidence ($\epsilon = 0.01$ in (8)) on MNIST, KMNIST and FashionMNIST, respectively. This reflects the intuition that more complex tasks require longer explanations. For reference, state-of-the-art deep networks on these datasets obtain test accuracies in order $\text{MNIST} \geq \text{KMNIST} \geq \text{FashionMNIST}$ (see last row in Table 1).

Effect of patch size on interpretability. We also run IP on MNIST with patch sizes of 1×1 (single pixels), 2×2 , 3×3 , and 4×4 . We observed that IP terminates at 99% confidence after 21.1, 9.6, 5.2, and 4.6 queries on average, respectively. While this suggests that larger patches lead to shorter explanations, we note that explanations with larger patches use more pixels (e.g. on MNIST, IP uses 21.1 pixels

7. Since the patches overlap in our query set, when computing the conditional probability of a patch query given history we only consider the probability of the pixels in the patch that have not yet been observed in our history.

5. $z^{(0)} \sim \mathcal{N}(0, I)$ for the first iteration of IP.

6. In this computation we have assumed, for simplicity, a unit cost for any operation that was computed in a batch concurrently on a GPU.

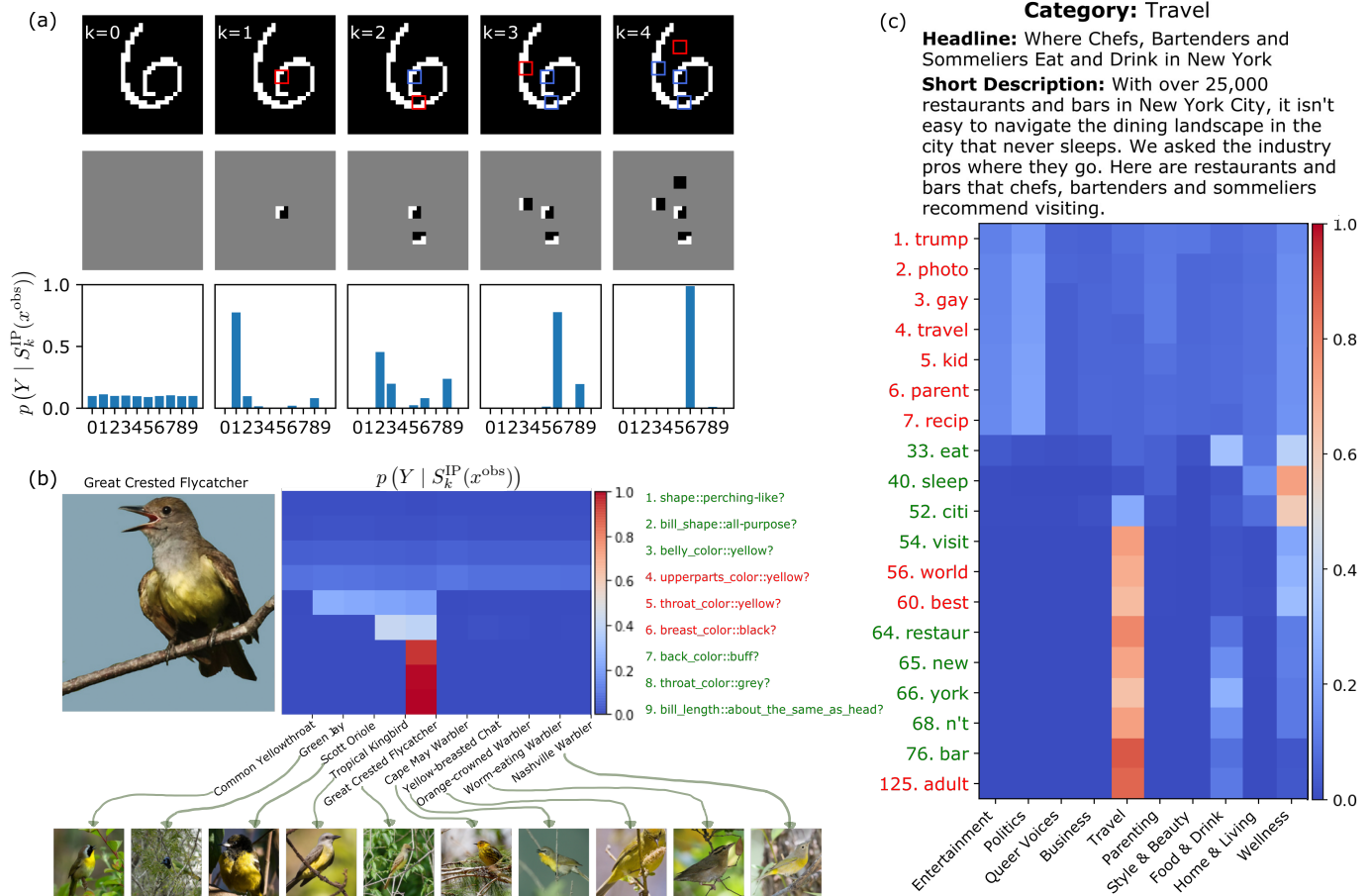


Fig. 4. **(a) IP on MNIST.** The top row displays the test image with red boxes denoting the current queried patch and blue boxes denoting previous patches. The second row shows the revealed portion of the image that IP gets to use at each query. The final row shows the model's estimated posteriors at each query, beginning at a nearly uniform prior before converging on the true digit "6" after 4 queries. **(b) IP on CUB Bird Species Classification.** On the left we show the input image and on the right we have a heatmap of the estimated class probabilities at each iteration. We only show the top 10 most probable classes out of the 200. To the right, we display the queries asked at each iteration, with red indicating a "no" response and green a "yes" response. **(c) IP on HuffPost News.** We show the input news item and a heatmap depicting the evolution of topic probabilities as IP asks queries and gathers answers. Words colored in red are absent from the sentence while words in green are present. For our visualization, we compute the KL divergence between each successive posterior and plot only the top 20 queries that led to the greatest change in posterior class probabilities.

on average for 1×1 patches and 54.7 pixels on average for 4×4 patches). That being said, very small patch queries are hard to interpret (see Fig. 5) and very large patch queries are also hard to interpret since each patch contains many image features. Overall, we found that 3×3 patches represented the right trade-off between interpretability in terms of edge patterns and minimality of the explanations. Specifically, single pixels are not very interpretable to humans but the explanations generated are more efficient in terms of *number of pixels needed to predict the label*. On the other extreme, using the entire image as a query is not interesting from an interpretability point of view since it does not help us understand which parts of the image are salient for prediction. We refer the reader to Appendix B.3.1 for additional patch size examples and quantitative analysis.

5.0.2 Concept-Based Queries

Task and query set. What if the end-user is interested in a more semantic explanation of the decision process in terms of high-level concepts? This can be easily incorporated into our framework by choosing an appropriate query set Q .

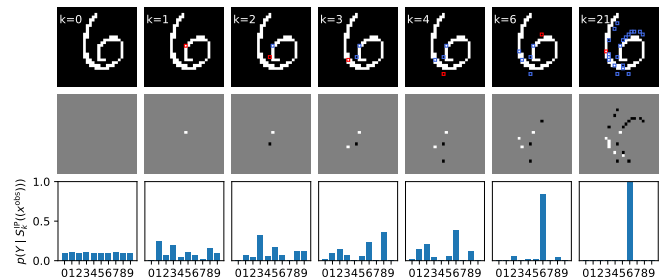


Fig. 5. **IP with 1×1 patches on MNIST.** Through the first 6 iterations, IP asks queries in the same center vertical region as in Fig. 4(a) (which uses 3×3 queries), outlining the distinctive loop in the bottom of the "6". However, reaching 99% confidence requires a total of 21 1×1 queries as opposed to just 4 3×3 ones. For conciseness, we show only the 6 queries that led to the greatest KL divergence between successive posterior class probabilities.

As an example we consider the challenging task of bird species classification on the Caltech-UCSD Birds-200-2011 (CUB) dataset [89]. The dataset contains images of 200

different species of birds. Each image is annotated with 312 binary attributes representing high-level concepts, such as the colour and shape of the beak, wings, and head. Unfortunately, these attribute annotations are very noisy. We follow [24] in deciding attribute labels by majority voting. For example, if more than 50% of images in a class have black wings, then we set all images in that class to have black wings. We construct Q by choosing a query for asking the presence/absence of each of these 312 binary attributes. Unfortunately, attribute annotations are not available at test time. To remedy this, we train a CNN (see [24] for details) to answer each query using the training set annotations, which is then used to answer queries at test time. Subsequently, we learn a VAE to model the joint distribution of query-answers supplied by this CNN (instead of the ground truth annotations) and Y , so our generative model can account for any estimation errors incurred by the CNN. Finally, we carry out IP as explained in §4.

IP in action. Consider the image of a *Great Crested Flycatcher* in Fig. 4(b). IP proceeds by asking most informative queries about various bird attributes progressively making the posterior over the species labels more and more peaked. After 5 queries, IP has gathered that the image is of a bird that has a perching-like shape, all-purpose beak and yellow belly, but does not have a yellow throat nor yellow upperparts. This results in a posterior concentrated on just 4 species that exhibit these characteristics. IP then proceeds to discount *Green Jay* and *Scott Oriole* which have black breasts with query 6. Likewise, *Tropical Kingbirds* have grayish back and is segregated from *Great Crested Flycatchers* which have buff-coloured backs with query 7. Finally after 9 queries, IP is 99% confident about the current class. Such concept-based explanations are more accessible to non-experts, especially on fine-grained classification datasets, which typically require domain expertise. On average IP takes 14.7 queries to classify a given bird image with $\epsilon = 0.007$ as the stopping criteria (See (7)).

5.0.3 Word-based Queries

Task and query set. Our framework can also be successfully applied to other domains like NLP. As an example we consider the task of topic identification from newspaper extended headlines (headline + short description field) using the the Huffington Post News Category Dataset [90]. We adopt a simple query set that consists of binary queries probing the existence of words in the extended headline. The words are chosen from a pre-defined vocabulary obtained by stemming all words in the HuffPost dataset and choosing the top-1,000 according to their tf-idf scores [91]. We process the dataset to merge redundant categories (such as *Style & Beauty* and *Beauty & Style*), remove semantically ambiguous, HuffPost-specific categories (e.g. *Impact* or *Fifty*) and remove categories with few samples, arriving at 10 final categories (see Appendix B.1).

IP in action. Fig. 4(c) shows an example run of IP on the HuffPost dataset. Note that positive responses to queries are very sparse, since each extended headline only contains 8.6 words on average out of the 1,000 in the vocabulary. As a result, IP asks 125 queries before termination. As discussed in §2, such long decision paths would be impossible in decision

trees due to data fragmentation and memory limitations. For clarity of presentation we only show the 20 queries with the greatest impact on the estimated posterior (as measured by KL-divergence from previous posterior). Upon reaching the first positive query “eat”, the probability mass concentrates on the categories *Food & Drink* and *Wellness* with little mass on *Travel*. However, as the queries about the existence of “citi”, “visit”, “york”, and “bar” in the extended headline come back positive, the model becomes more and more confident that “Travel” is the correct class. IP requires about 199.3 queries on average to predict the topic of the extended headline with $\epsilon = 10^{-3}$ as the stopping criteria (See (7)). Additional details on the HuffPost query set are in Appendix B.1.

Further examples of IP performing inference on all tasks can be found in Appendix B.3.

5.1 Quantitative Evaluation

5.1.1 Classification Accuracy

We compare the classification accuracy of our model’s prediction based on the query-answers gathered by IP until termination with several other baseline models. For each of the models considered, we first give a brief description and then comment on their performance with respect to IP. All the results are summarized in Table 1.

DECISION TREE refers to standard classification trees learnt using the popular CART algorithm [34]. In the Introduction, we mentioned that classical decision trees learnt using Q to supply the node splitting functions will be interpretable by construction but are not competitive with state-of-the-art methods. This is illustrated in our results in Table 1. Across all datasets, IP obtains superior performance since it is based on an underlying generative model (VAE) and only computes the branch of the tree traversed by the input data in an online manner, thus it is not shackled by data fragmentation and memory limitations.

MAP USING Q refers to the Maximum A Posteriori estimate obtained using the posterior distribution over the labels given the answers to all the queries in Q (for a given input). Recall, IP asks queries until the stopping criteria is reached (Equation (7) & Equation (8)). Naturally, there is a trade-off between the length of the explanations and the predictive performance observed. If we ask all the queries then the resulting explanations of length $|Q|$ might be too long to be desirable. The results for IP reported in Table 1 use different dataset-specific stopping criteria according to the elbow in their respective accuracy vs. explanation length curves (see Fig. 6). On the binary image datasets, (MNIST, KMNIST, and FashionMNIST) IP obtains an accuracy within 3% of the best achievable upon seeing all the query-answers with only about 2% of the total queries in Q . Similarly for the CUB and Huffpost datasets, IP achieves about the same accuracy as MAP USING Q but asks less than 5% and 20% of total possible queries respectively.

BLACK-BOX USING Q refers to the best performing deep network model we get by training on features supplied by evaluating all $q \in Q$ on input data from the various training datasets. For the binary image datasets, this is just a 4-layer CNN with ReLU activations. For CUB we use the results reported by the sequential model in [24]. For

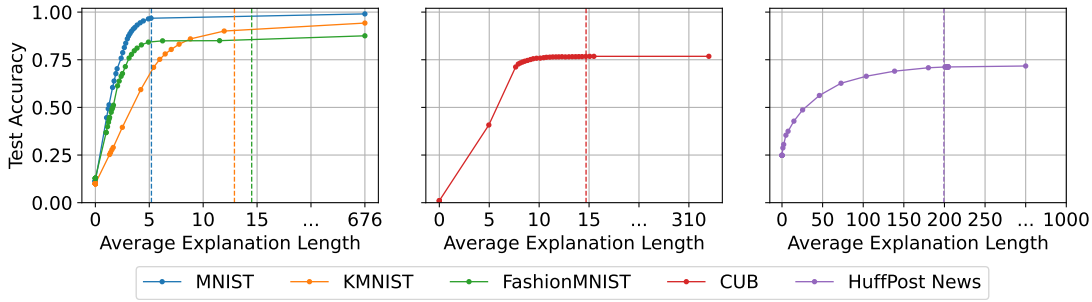


Fig. 6. **Trade-off between predictive performance and explanation length** Different points along the curves correspond to different values of ϵ as the stopping criteria (7) is varied. The colored dotted vertical line in each plot indicates the avg. explanation length v/s test accuracy at the ϵ value used as the stopping criteria for reporting results for the IP strategy in this work. For each plot, the x-axis ranges from 0 to the size of the query set, $|Q|$, chosen for that task.

HuffPost, we found a single hidden layer with ReLU non-linearity give the best performance. Further architectural and training details are in Appendix B.2. In Table 1 we show that across all datasets, the predictive performance obtained by MAP USING Q is on par with the best performance we obtained using black-box expressive non-interpretable networks BLACK-BOX USING Q . Thus, our generative models, which form the backbone for IP, are competitive with state-of-the-art prediction methods.

BLACK-BOX refers to the best performing black-box model on these datasets in terms of classification accuracy as reported in literature; to the best of our knowledge. In Table 1, we see a performance gap in each dataset when compared with MAP USING Q which uses an interpretable query set. This is expected since explainability can be viewed as an additional constraint on learning. For example, on FashionMNIST we see an almost 8.5% relative fall in accuracy due to binarization. This is because it is harder to decipher between some classes like shirts and pullovers at the binary level. On the other hand, binary patches are easily interpretable as edges, foregrounds and backgrounds. Similarly, there is a relative drop of accuracy of about 17% for the HuffPost dataset since our queries regarding the existence of different words ignore their arrangement in sentences. Thus we lose crucial contextual information used by state-of-the-art transformer models [92]. Ideally, we would like query sets to be easily interpretable, lead to short explanations and be sufficient to solve the task. Finding such query sets is nontrivial and will be explored in future work.

TABLE 1
Classification accuracy of our model (Information Pursuit) relative to baselines on different test sets. See 5.1.1 for details on each model.

Model	MNIST	KMNIST	Fashion	CUB	HuffPost
INFORMATION PURSUIT	96.78%	91.02%	85.60%	76.73%	71.21%
DECISION TREE [34]	90.23%	78.00%	80.80%	68.80%	63.00%
MAP USING Q	99.05%	94.25%	87.56%	76.80%	71.72%
BLACK-BOX USING Q	99.15%	95.10%	88.43%	76.30%	71.48%
BLACK-BOX	99.83% [93]	98.83% [88]	96.70% [94]	82.70% [24]	86.45% ⁸

8. We fine-tuned a Bert Large Uncased Transformer model [92] with the last layer replaced with a linear one. See Appendix B.2.3 for details.

5.1.2 Comparison to current attribution methods

At first glance, it might seem that using attribution methods/saliency maps can provide the same insights as to which parts of the image or more generally which queries in Q were most influential in a decision made by a black-box model trained on input features supplied by all the query-answers. However, the unreliability of these methods in being faithful to the model they try to explain brings their utility into question [19], [20], [22]. We conjecture that this is because current attribution methods are not designed to generate explanations that are sufficient statistics of the model’s prediction. We illustrate this with a simple experiment using our binary image classification datasets.

For each input image x , we compute the corresponding attribution map $e(x)$ for the model’s predicted class using two popular attribution methods, Integrated gradients (IG) [95] and DeepSHAP [45]. We then compute the L most important 3×3 patches, where L is the number of patches queried by IP for that particular input image. For computing the attribution/importance of a patch we average the attributions of all the pixels in that patch (following [20]). We proceed as follows: (i) Given $e(x)$, compute the patch with maximum attribution and add these pixels to our explanation, (ii) Zero-out the attributions of all the pixels in the previously selected patch and repeat step (i) until L patches are selected. The final explanation consists of L possibly overlapping patches. Now, we evaluate the sufficiency of the generated explanation for the model’s prediction by estimating the MAP accuracy of the posterior over labels given the intensities in the patches included in this explanation. This is done via a VAE trained to learn the joint distribution over image pixels and class labels. We experiment with both the raw attribution scores returned by IG and DeepSHAP and also the absolute values of the attribution scores for $e(x)$. The results are reported in Table 2. In almost all cases (with the exception of DeepSHAP on FashionMNIST), IP generates explanations that are more predictive of the class label than popular attribution methods.

6 CONCLUSION

We have presented a step towards building trustworthy interpretable machine learning models that respect the domain- and user-dependent nature of interpretability. We address this by composing user-defined, inter-

TABLE 2

MAP accuracy of explanations generated by Information Pursuit (IP) v/s other attribution methods. IP explanations (in almost all cases) achieve a higher classification accuracy than explanations of the same length generated using baseline attribution methods. The (absolute) method refers to explanations generated using absolute values of the attribution map scores. On MNIST and KMNIST, IP explanations achieve a 10% and 2.38% relative improvement respectively over the best performing baseline method. On FashionMNIST, IP explanations are second best with a relative decrease of about 3.12% from the best performing baseline.

Explanation Method	MNIST	KMNIST	Fashion-MNIST
INFORMATION PURSUIT	96.78%	91.02%	85.60%
IG	78.48%	84.87%	78.49%
IG (ABSOLUTE)	70.39%	84.72%	64.95%
DEEPSHAP	87.98%	88.90%	88.36%
DEEPSHAP (ABSOLUTE)	84.80%	84.56%	84.35%

interpretable queries into concise explanations. Furthermore, unlike many contemporary attempts at explainability, our method is not post-hoc, but is *interpretable by design* and guaranteed to produce faithful explanations. We formulate a tractable approach to implement this framework through deep generative models, MCMC algorithms, and the information pursuit algorithm. Finally, we demonstrate the effectiveness of our method across various vision and language tasks at generating concise explanations describing the underlying reasoning process behind the prediction. Future work will be aimed at extending the proposed framework to more complex tasks beyond classification such as scene parsing, image captioning, and sentiment analysis.

ACKNOWLEDGMENTS

The authors thank María Pérez Ortiz and John Shawe-Taylor for their contributions to the design of the experiments on document classification presented in Section 5.0.3. This research was supported by the Army Research Office under the Multidisciplinary University Research Initiative contract W911NF-17-1-0304 and by the NSF grant 2031985.

REFERENCES

- [1] D. Gunning, M. Stefik, J. Choi, T. Miller, S. Stumpf, and G.-Z. Yang, “Xai—explainable artificial intelligence,” *Science Robotics*, vol. 4, no. 37, p. eaay7120, 2019.
- [2] C. Rudin, “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead,” *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.
- [3] W. J. Murdoch, C. Singh, K. Kumbier, R. Abbasi-Asl, and B. Yu, “Interpretable machine learning: definitions, methods, and applications,” *arXiv preprint arXiv:1901.04592*, 2019.
- [4] European Commission, “Building trust in human-centric artificial intelligence,” *Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions*, vol. 168, 2019.
- [5] United States Food and Drug Administration, “Virtual public workshop - transparency of artificial intelligence/machine learning-enabled medical devices,” Transcript: <https://www.fda.gov/media/154423/download>, Oct. 14, 2021.
- [6] U. Johansson, C. Sönström, U. Norinder, and H. Boström, “Trade-off between accuracy and interpretability for predictive in silico modeling,” *Future medicinal chemistry*, vol. 3, no. 6, pp. 647–663, 2011.
- [7] J. Wanner, L.-V. Herm, K. Heinrich, and C. Janiesch, “Stop ordering machine learning algorithms by their explainability! an empirical investigation of the tradeoff between performance and explainability,” in *Conference on e-Business, e-Services and e-Society*. Springer, 2021, pp. 245–258.
- [8] F. K. Došilović, M. Brčić, and N. Hlupić, “Explainable artificial intelligence: A survey,” in *2018 41st International convention on information and communication technology, electronics and microelectronics (MIPRO)*. IEEE, 2018, pp. 0210–0215.
- [9] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins *et al.*, “Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai,” *Information fusion*, vol. 58, pp. 82–115, 2020.
- [10] D. Gunning and D. Aha, “Darpa’s explainable artificial intelligence (xai) program,” *AI magazine*, vol. 40, no. 2, pp. 44–58, 2019.
- [11] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Müller, “How to explain individual classification decisions,” *The Journal of Machine Learning Research*, vol. 11, pp. 1803–1831, 2010.
- [12] K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv preprint arXiv:1312.6034*, 2013.
- [13] S. Kolek, D. A. Nguyen, R. Levie, J. Bruna, and G. Kutyniok, “A rate-distortion framework for explaining black-box model decisions,” in *International Workshop on Extending Explainable AI Beyond Deep Models and Classifiers*. Springer, 2022, pp. 91–115.
- [14] A. Shrikumar, P. Greenside, and A. Kundaje, “Learning important features through propagating activation differences,” in *International conference on machine learning*. PMLR, 2017, pp. 3145–3153.
- [15] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks,” in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [16] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [17] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, “Smoothgrad: removing noise by adding noise,” *arXiv preprint arXiv:1706.03825*, 2017.
- [18] A. Subramanya, V. Pillai, and H. Pirsiavash, “Fooling network interpretation in image classification,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2020–2029.
- [19] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim, “Sanity checks for saliency maps,” *Advances in neural information processing systems*, vol. 31, 2018.
- [20] M. Yang and B. Kim, “Benchmarking attribution methods with relative feature importance,” *arXiv preprint arXiv:1907.09701*, 2019.
- [21] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim, “The (un) reliability of saliency methods,” in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer, 2019, pp. 267–280.
- [22] H. Shah, P. Jain, and P. Netrapalli, “Do input gradients highlight discriminative features?” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [23] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju, “Fooling lime and shap: Adversarial attacks on post hoc explanation methods,” in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2020, pp. 180–186.
- [24] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang, “Concept bottleneck models,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 5338–5348.
- [25] C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. K. Su, “This looks like that: deep learning for interpretable image recognition,” *Advances in neural information processing systems*, vol. 32, 2019.
- [26] C. Rudin, C. Chen, Z. Chen, H. Huang, L. Semenova, and C. Zhong, “Interpretable machine learning: Fundamental principles and 10 grand challenges,” *Statistics Surveys*, vol. 16, pp. 1–85, 2022.
- [27] T. M. Janssen and B. H. Partee, “Compositionality,” in *Handbook of logic and language*. Elsevier, 1997, pp. 417–473.
- [28] H. Lakkaraju, S. H. Bach, and J. Leskovec, “Interpretable decision sets: A joint framework for description and prediction,” in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1675–1684.
- [29] A. Wan, L. Dunlap, D. Ho, J. Yin, S. Lee, S. Petryk, S. A. Bargal, and J. E. Gonzalez, “{NBDT}: Neural-backed decision tree,” in

- International Conference on Learning Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=mCLVeEppINE>
- [30] J. Mu and J. Andreas, "Compositional explanations of neurons," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 153–17 163, 2020.
- [31] E. Jahangiri, E. Yoruk, R. Vidal, L. Younes, and D. Geman, "Information pursuit: A bayesian framework for sequential scene parsing," *arXiv preprint arXiv:1701.02343*, 2017.
- [32] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, "Network dissection: Quantifying interpretability of deep visual representations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6541–6549.
- [33] E. Hernandez, S. Schettmann, D. Bau, T. Bagashvili, A. Torralba, and J. Andreas, "Natural language descriptions of deep visual features," *arXiv preprint arXiv:2201.11114*, 2022.
- [34] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and regression trees*. Routledge, 2017.
- [35] J. R. Quinlan, "Induction of decision trees," *Machine learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [36] Y. Amit and D. Geman, "Shape quantization and recognition with randomized trees," *Neural computation*, vol. 9, no. 7, pp. 1545–1588, 1997.
- [37] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [38] R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 161–168.
- [39] M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do we need hundreds of classifiers to solve real world classification problems?" *The journal of machine learning research*, vol. 15, no. 1, pp. 3133–3181, 2014.
- [40] H. Xu, K. A. Kinfu, W. LeVine, S. Panda, J. Dey, M. Ainsworth, Y.-C. Peng, M. Kusmanov, F. Engert, C. M. White *et al.*, "When are deep networks really better than decision forests at small sample sizes, and how?" *arXiv preprint arXiv:2108.13637*, 2021.
- [41] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Buló, "Deep neural decision forests," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1467–1475.
- [42] D. Geman and B. Jedynek, "An active testing model for tracking roads in satellite images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 1, pp. 1–14, 1996.
- [43] A. Chattopadhyay, A. Sarkar, P. Howlader, and V. N. Balasubramanian, "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks," in *2018 IEEE winter conference on applications of computer vision (WACV)*. IEEE, 2018, pp. 839–847.
- [44] M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [45] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in neural information processing systems*, vol. 30, 2017.
- [46] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas *et al.*, "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)," in *International conference on machine learning*. PMLR, 2018, pp. 2668–2677.
- [47] B. Zhou, Y. Sun, D. Bau, and A. Torralba, "Interpretable basis decomposition for visual explanation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 119–134.
- [48] C.-K. Yeh, B. Kim, S. Arik, C.-L. Li, T. Pfister, and P. Ravikumar, "On completeness-aware concept-based explanations in deep neural networks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20 554–20 565, 2020.
- [49] D. Alvarez Melis and T. Jaakkola, "Towards robust interpretability with self-explaining neural networks," *Advances in neural information processing systems*, vol. 31, 2018.
- [50] M. Bohle, M. Fritz, and B. Schiele, "Convolutional dynamic alignment networks for interpretable classifications," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 029–10 038.
- [51] M. Wu, S. Parbhoo, M. C. Hughes, V. Roth, and F. Doshi-Velez, "Optimizing for interpretability in deep neural networks with tree regularization," *Journal of Artificial Intelligence Research*, vol. 72, pp. 1–37, 2021.
- [52] V. Pillai and H. Pirsiavash, "Explainable models with consistent interpretations," *UMBC Student Collection*, 2021.
- [53] Z. Chen, Y. Bei, and C. Rudin, "Concept whitening for interpretable image recognition," *Nature Machine Intelligence*, vol. 2, no. 12, pp. 772–782, 2020.
- [54] M. De-Arteaga, A. Romanov, H. Wallach, J. Chayes, C. Borgs, A. Chouldechova, S. Geyik, K. Kenthapadi, and A. T. Kalai, "Bias in bios: A case study of semantic representation bias in a high-stakes setting," in *proceedings of the Conference on Fairness, Accountability, and Transparency*, 2019, pp. 120–128.
- [55] A. Galassi, M. Lippi, and P. Torrioni, "Attention in natural language processing," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 10, pp. 4291–4308, 2020.
- [56] D. Pruthi, M. Gupta, B. Dhingra, G. Neubig, and Z. C. Lipton, "Learning to deceive with attention-based explanations," *arXiv preprint arXiv:1909.07913*, 2019.
- [57] M. Craven and J. Shavlik, "Extracting tree-structured representations of trained networks," *Advances in neural information processing systems*, vol. 8, 1995.
- [58] D. Dancey, D. A. McLean, and Z. A. Bandar, "Decision tree extraction from trained neural networks," in *Proceedings of the Nineteenth Conference on Artificial Intelligence*. American Association for Artificial Intelligence, 2004.
- [59] N. Frosst and G. Hinton, "Distilling a neural network into a soft decision tree," *arXiv preprint arXiv:1711.09784*, 2017.
- [60] A. Roy and S. Todorovic, "Monocular depth estimation using neural regression forest," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5506–5514.
- [61] U. C. Biçici, C. Keskin, and L. Akarun, "Conditional information gain networks," in *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018, pp. 1390–1395.
- [62] V. N. Murthy, V. Singh, T. Chen, R. Manmatha, and D. Comaniciu, "Deep decision network for multi-class image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2240–2248.
- [63] R. Vilalta, G. Blix, and L. Rendell, "Global data analysis and the fragmentation problem in decision tree induction," in *European Conference on Machine Learning*. Springer, 1997, pp. 312–326.
- [64] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," *arXiv preprint physics/0004057*, 2000.
- [65] K. Chaloner and I. Verdine, "Bayesian experimental design: A review," *Statistical Science*, pp. 273–304, 1995.
- [66] R. Sznitman and B. Jedynek, "Active testing for face detection and localization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 10, pp. 1914–1920, 2010.
- [67] M. Cuturi, O. Teboul, Q. Berthet, A. Doucet, and J.-P. Vert, "Noisy adaptive group testing using bayesian sequential experimental design," *arXiv preprint arXiv:2004.12508*, 2020.
- [68] S. Branson, G. Van Horn, C. Wah, P. Perona, and S. Belongie, "The ignorant led by the blind: A hybrid human-machine vision system for fine-grained categorization," *International Journal of Computer Vision*, vol. 108, no. 1, pp. 3–29, 2014.
- [69] V. Mnih, N. Heess, A. Graves *et al.*, "Recurrent models of visual attention," in *Advances in neural information processing systems*, 2014, pp. 2204–2212.
- [70] G. Elsayed, S. Kornblith, and Q. V. Le, "Saccader: improving accuracy of hard attention models for vision," in *Advances in Neural Information Processing Systems*, 2019, pp. 702–714.
- [71] M. Li, S. S. Ge, and T. H. Lee, "Glance and glimpse network: A stochastic attention model driven by class saliency," in *Asian Conference on Computer Vision*. Springer, 2016, pp. 572–587.
- [72] H. Li, P. Wang, C. Shen, and A. v. d. Hengel, "Visual question answering as reading comprehension," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6319–6328.
- [73] M. Malinowski, M. Rohrbach, and M. Fritz, "Ask your neurons: A deep learning approach to visual question answering," *International Journal of Computer Vision*, vol. 125, no. 1, pp. 110–135, 2017.
- [74] J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum, and J. Wu, "The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision," *arXiv preprint arXiv:1904.12584*, 2019.
- [75] K. J. Shih, S. Singh, and D. Hoiem, "Where to look: Focus regions for visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4613–4621.
- [76] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," *Advances in neural information processing systems*, vol. 29, 2016.

- [77] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein, "Neural module networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 39–48.
- [78] H. Laurent and R. L. Rivest, "Constructing optimal binary decision trees is np-complete," *Information processing letters*, vol. 5, no. 1, pp. 15–17, 1976.
- [79] M. I. Belghazi, A. Baratin, S. Rajeswar, S. Ozair, Y. Bengio, A. Courville, and R. D. Hjelm, "Mine: mutual information neural estimation," *arXiv preprint arXiv:1801.04062*, 2018.
- [80] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [81] A. Doucet, A. M. Johansen *et al.*, "A tutorial on particle filtering and smoothing: Fifteen years later," *Handbook of nonlinear filtering*, vol. 12, no. 656-704, p. 3, 2009.
- [82] A. Jalal, S. Karmalkar, A. Dimakis, and E. Price, "Instance-optimal compressed sensing via posterior sampling," *Proceedings of Machine Learning Research*, vol. 139, 2021.
- [83] E. Nijkamp, M. Hill, T. Han, S.-C. Zhu, and Y. N. Wu, "On the anatomy of mcmc-based maximum likelihood learning of energy-based models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5272–5280.
- [84] A. Durmus and E. Moulines, "High-dimensional bayesian inference via the unadjusted langevin algorithm," *Bernoulli*, vol. 25, no. 4A, pp. 2854–2882, 2019.
- [85] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *Proceedings of the 28th international conference on machine learning (ICML-11)*. Citeseer, 2011, pp. 681–688.
- [86] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [87] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [88] T. Clanuwat, M. Bober-Irizar, A. Kitamoto, A. Lamb, K. Yamamoto, and D. Ha, "Deep learning for classical japanese literature," *arXiv preprint arXiv:1812.01718*, 2018.
- [89] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011.
- [90] R. Misra, "News category dataset," 06 2018.
- [91] M. Lavin, "Analyzing documents with tf-idf," 2019.
- [92] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [93] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [94] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: A mnist-like fashion product database," in *GitHub*, 2017.
- [95] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International conference on machine learning*. PMLR, 2017, pp. 3319–3328.



Aditya Chattopadhyay is a PhD student in the Computer Science Department, Johns Hopkins University. He received the Bachelor of Technology degree in Computer Science and Master of Science by Research degree in Computational Natural Sciences from the International Institute of Information Technology, Hyderabad in 2016 and 2018 respectively. His research interests include explainable AI, probabilistic graphical models and Bayesian inference.



Stewart Slocum received his BS in computer science and applied mathematics from Johns Hopkins University in 2021. His research interests center on principled deep learning methods with performance and robustness guarantees.



Benjamin Haeffele is an Associate Research Scientist in the Mathematical Institute for Data Science at Johns Hopkins University. His research interests involve developing theory and algorithms for processing high-dimensional data at the intersection of machine learning, optimization, and computer vision. In addition to basic research in data science he also works on a variety of applications in medicine, microscopy, and computational imaging. He received his Ph.D. in Biomedical Engineering at Johns Hopkins University in 2015 and his B.S. in Electrical Engineering from the Georgia Institute of Technology in 2006.



René Vidal received his B.S. degree in Electrical Engineering (valedictorian) from the Pontificia Universidad Católica de Chile in 1997 and his M.S. and Ph.D. degrees in Electrical Engineering and Computer Science from the University of California at Berkeley in 2000 and 2003, respectively. He is currently the Director of the Mathematical Institute for Data Science (MINDS) and the Hershel L. Seder Professor of Department of Biomedical Engineering at The Johns Hopkins University, where he has been since 2004. He is co-author of the book "Generalized Principal Component Analysis" (Springer 2016), co-editor of the book "Dynamical Vision" (Springer 2006) and co-author of over 300 articles in machine learning, computer vision, signal and image processing, biomedical image analysis, hybrid systems, robotics and control. He is or has been Associate Editor in Chief of the IEEE Transactions on Pattern Analysis and Machine Intelligence and Computer Vision and Image Understanding, Associate Editor or Guest Editor of Medical Image Analysis, the IEEE Transactions on Pattern Analysis and Machine Intelligence, the SIAM Journal on Imaging Sciences, Computer Vision and Image Understanding, the Journal of Mathematical Imaging and Vision, the International Journal on Computer Vision and Signal Processing Magazine. He has received numerous awards for his work, including the 2021 Edward J. McCluskey Technical Achievement Award, the 2016 D'Alembert Faculty Fellowship, the 2012 IAPR J.K. Aggarwal Prize, the 2009 ONR Young Investigator Award, the 2009 Sloan Research Fellowship and the 2005 NSF CAREER Award. He is a Fellow of the IEEE, Fellow of IAPR, Fellow of AIMBE, and a member of the ACM and SIAM.



Donald Geman (Life Senior Member, IEEE) received the B.A. degree in literature from the University of Illinois and the Ph.D. degree in mathematics from Northwestern University. He was a Distinguished Professor with the University of Massachusetts until 2001, when he joined the Department of Applied Mathematics and Statistics, Johns Hopkins University, where he is currently a member of the Center for Imaging Science and the Institute for Computational Medicine. His current research interests include statistical learning, computer vision, and computational biology. He is a member of the National Academy of Sciences and a fellow of the IMS and SIAM.

APPENDIX A

In proofs of propositions and lemmas we rewrite the statement (un-numbered) for convenience.

A.1 Characterizing the optimal strategy π^*

In this subsection we characterize the optimal strategy for any such arbitrary query set chosen by the user.

Proposition. *Assuming Q is finite and when d is taken to be the KL-divergence then objective (3) can be rewritten as,*

$$\begin{aligned} H_Q^\epsilon(X; Y) &:= \min_{\pi} \mathbb{E}_X [|\text{expl}_Q^\pi(X)|] \\ \text{s.t. } \sum_{k=1}^{\tau^\pi} I(Y; S_k^\pi(X) | S_{k-1}^\pi(X)) &\geq I(X; Y) - \epsilon \end{aligned} \quad (21)$$

where, $\tau^\pi = \max\{t^\pi(x) : x \in \mathcal{X}\}$ and $t^\pi(X)$ is defined as the number of queries selected by π for input X until q_{STOP} . We define $S_k^\pi(X)$ as a random variable where any realization $S_k^\pi(x^{obs})$, $x^{obs} \in \mathcal{X}$, denotes the event

$$S_k^\pi(x^{obs}) := \{x' \in \mathcal{X} \mid \{q_i, q_i(x^{obs})\}_{1:k} = \{q_i, q_i(x')\}_{1:k}\},$$

where q_i is the i^{th} query selected by π for input x^{obs} . Here we use the convention that $S_0^\pi(X) = \Omega$ (the entire sample space) and $S_l^\pi(X) = S_{t^\pi(X)}^\pi(X) \quad \forall l > t^\pi(X)$.

Proof. We begin by reformulating our sufficiency constraint in (3) in terms of entropy by taking the ‘‘distance’’ between probability distributions as the KL divergence.⁹ The ϵ -Sufficiency constraint can then be rewritten as,

$$\begin{aligned} \epsilon &\geq \mathbb{E}_X [KL(p(Y | X), p(Y | \text{expl}_Q^\pi(X)))] \\ &= \mathbb{E}_X \left[\mathbb{E}_Y \left[\log \frac{p(Y | X)}{p(Y | \text{expl}_Q^\pi(X))} \mid X \right] \right] \\ &= \mathbb{E}_X \left[\mathbb{E}_Y \left[\log \frac{p(Y | X)}{P(Y)} \mid X \right] \right] + \mathbb{E}_X \left[\mathbb{E}_Y \left[\log \frac{p(Y)}{P(Y | \text{expl}_Q^\pi(X))} \mid X \right] \right] \\ &= I(X; Y) - I(\text{expl}_Q^\pi(X); Y) \\ &= H(Y | \text{expl}_Q^\pi(X)) - H(Y | X) \end{aligned}$$

In the third equality we multiplied the term inside the log by the identity $\frac{P(Y)}{P(Y)} = 1$. The fifth inequality is by definition of mutual information. Thus, we can rewrite (3) as,

$$\begin{aligned} H_Q^\epsilon(X; Y) &:= \min_{\pi} \mathbb{E}_X [|\text{expl}_Q^\pi(X)|] \\ \text{s.t. } H(Y | \text{expl}_Q^\pi(X)) - H(Y | X) &\leq \epsilon \quad (\epsilon\text{-Sufficiency}) \end{aligned} \quad (22)$$

Let’s define $t^\pi(X)$ as the number of queries selected by π for input X until q_{STOP} . Define $\tau^\pi = \max\{t^\pi(X) : X \in \mathcal{X}\}$. For the purpose of analysis we can vacuously modify π such that for any given $x^{obs} \in \mathcal{X}$, π asks a fixed τ^π number of queries by filling the remaining $\tau^\pi - t^\pi(x^{obs})$ queries with q_{STOP} . An immediate consequence of this modification is that $S_l^\pi(X) = S_{t^\pi(X)}^\pi(X) \quad \forall l > t^\pi(X)$.

We will now show that the sufficiency criteria $H(Y | \text{expl}_Q^\pi(X)) - H(Y | X)$ can be rewritten as a sum of successive mutual information terms.

$$\begin{aligned} H(Y | \text{expl}_Q^\pi(X)) - H(Y | X) &= (H(Y) - H(Y | X)) - (H(Y) - H(Y | \text{expl}_Q^\pi(X))) \\ &= I(X; Y) - (H(Y) - H(Y | S_\tau^\pi(X))) \end{aligned} \quad (23)$$

The third equality uses the fact that $H(Y | \text{expl}_Q^\pi(X)) = H(Y | S_\tau^\pi(X))$ since the $\forall x \in \mathcal{X}$, $S_\tau^\pi(x) = S_{t^\pi(x)}^\pi(x)$. We can now write $H(Y) - H(Y | S_\tau^\pi(X))$ as a telescoping series,

$$\begin{aligned} H(Y) - H(Y | S_\tau^\pi(X)) &= H(Y) - H(Y | S_1^\pi(X)) + H(Y | S_1^\pi(X)) - H(Y | S_1^\pi(X), S_2^\pi(X)) \\ &\quad + H(Y | S_1^\pi(X), S_2^\pi(X)) \dots + H(Y | S_1^\pi(X), S_2^\pi(X), \dots, S_{\tau-1}^\pi(X)) - H(Y | S_\tau^\pi(X)) \\ &= I(Y; S_1^\pi(X)) + I(Y; S_2^\pi(X) | S_1^\pi(X)) + \dots + I(Y; S_\tau^\pi(X) | S_{\tau-1}^\pi(X)) \end{aligned} \quad (24)$$

The last equality is obtained by noticing that

9. In favour of a clearer exposition, we abuse notation here and use $p(Y | \text{expl}_Q^\pi(X))$ to denote $P(Y | S_{t^\pi(X)}^\pi(X))$. In reality, $\text{expl}_Q^\pi(x)$ refers to the sequence of query-answer pairs chosen for x as defined in (2) whereas $S_{t^\pi(X=x)}^\pi(X=x)$ refers to the event which is the set of all possible data-points that agree on the first $t^\pi(x)$ query-answers observed for x .

- 1) $H(Y | S_1^\pi(X), \dots, S_{k-1}^\pi(X)) - H(Y | S_1^\pi(X), \dots, S_k^\pi(X)) = I(Y; S_k^\pi(X) | S_{k-1}^\pi(X))$,
- 2) $H(Y | S_\tau^\pi(X)) = H(Y | S_1^\pi(X), S_2^\pi(X), \dots, S_\tau^\pi(X))$ since the events $\{S_k^\pi(x)\}_{1:\tau}$ are nested $\forall x \in \mathcal{X}$.

Putting it all together we can rewrite (22) as,

$$\begin{aligned} H_Q^\epsilon(X; Y) &:= \min_{\pi} \mathbb{E}_X [|\text{expl}_Q^\pi(X)|] \\ \text{s.t. } I(X; Y) - \sum_{k=1}^{\tau^\pi} I(Y; S_k^\pi(X) | S_{k-1}^\pi(X)) &\leq \epsilon \end{aligned} \quad (25)$$

□

A.2 Approximation guarantees for IP

In Proposition 2 we prove that the IP strategy comes within 1 bit of $H(Y)$ (the entropy of Y) under the assumption that one has access to all possible binary functions of X , that are also binary functions of Y , as queries. Given such a query set, it is well-known that the optimal strategy π^* is given by the Huffman Code for Y which is also within 1 bit of $H(Y)$. Thus, the result is immediate that IP asks at most 1 query more than π^* on average. We restate Proposition 2 below for ease.

Proposition. *Let Y be discrete. Let $\tilde{H}_Q(X; Y)$ be the expected description length obtained by the IP strategy. If $H(Y|X) = 0$ and Q is the set of all possible binary functions of X such that $H(q(X) | Y) = 0 \forall q \in Q$, then $\tilde{H}_Q(X; Y) \leq H(Y) + 1$.*

We make two remarks before turning to the proof.

Remark 1: We have observed data X , a categorical discrete r.v. Y , and binary queries $\{q(X), q \in Q\}$ from which Y can be estimated. In fact, let's suppose that Y is determined by X , say $Y = f(X)$; equivalently, $H(Y|X) = 0$. We assume that Y represents some high-level, possibly semantic, interpretation of X which can only be seen through the eyes of the queries $q \in Q$. Ideally, we would like to be able to query the membership of Y in any subset of \mathcal{Y} (the set of possible values of Y); this is the information we would have in coding Y . We will say that the query $I_{Y \in D}$ is *realizable* for some $D \subset \mathcal{Y}$ if the query $I_{X \in f^{-1}(D)}$ is in Q , i.e., we can test for $Y \in D$ by one of our observable data queries. In general, not all subsets of \mathcal{Y} can be associated with attributes or features of X , e.g., "Napoleon" or "Dead" in "20 Questions", or "black beak" in bird species classification. If there was a query $q(X)$ for every subset D of values of Y , then our theorem says that mean number of queries needed to determine the state of Y with IP is bounded below by $H(Y)$ and above by $H(Y) + 1$.

Remark 2: The sequence of queries q_1, q_2, \dots generated by the IP algorithm for a particular data point can be seen as one branch, root to leaf, of a decision tree constructed by the standard machine learning strategy based on successive reduction of uncertainty about Y as measured by mutual information: $q_1 = \arg \max_{q \in Q} I(q(X); Y)$, $q_{k+1} = \arg \max_{q \in Q} I(q(X); Y | S_k^{\text{IP}}(x^0))$ where the $S_k^{\text{IP}}(x^0)$ is the event that for the first k questions the answers agree with those for x^0 . We stop as soon as Y is determined. Whereas a decision tree accommodates all x simultaneously, the questions along the branch depends on having a particular, fixed data point. But the learning problem in the branch version ("active testing") is exponentially simpler.

Proof. The lower bound $H(Y) \leq \tilde{H}_Q(X; Y)$ comes from Shannon's source coding theorem for stochastic source Y .

Now for the upper bound, since $I(q(X); Y | S_k^{\text{IP}}(x^0)) = H(q(X) | S_k^{\text{IP}}(x^0)) - H(q(X) | Y, S_k^{\text{IP}}(x^0))$ and since Y determines $q(X)$ and hence also $q(X)$, the second entropy term is zero (since given $H(q(X) | Y) = 0$). So our problem is maximize the conditional entropy of the binary random variable $q(X)$ given $S_k^{\text{IP}}(x^0)$. So the IP algorithm is clearly just "divide and conquer":

$$\begin{aligned} q_1 &= \arg \max_{q \in Q} H(q(X)), \\ q_{k+1} &= \arg \max_{q \in Q} H(q(X) | S_k^{\text{IP}}(x^0)). \end{aligned}$$

Equivalently, since entropy of a binary random variable ρ is maximized when $P(\rho) = \frac{1}{2}$,

$$q_{k+1} = \arg \min_{q \in Q} |P(q(X) = 1 | S_k^{\text{IP}}(x^0)) - \frac{1}{2}|.$$

Let \mathcal{Y}_k be the set of "active hypotheses" after k queries (denoted as \mathcal{A}_k), namely those y with positive posterior probability: $P(Y = y | S_k^{\text{IP}}(x^0)) > 0$. Indeed,

$$\begin{aligned} P(Y = y | S_k^{\text{IP}}(x^0)) &= \frac{P(S_k^{\text{IP}}(x^0) | Y = y) p(y)}{\sum_y P(S_k^{\text{IP}}(x^0) | Y = y) p(y)} \\ &= \frac{1_{\mathcal{Y}_k} p(y)}{\sum_{y \in \mathcal{A}_k} p(y)} \end{aligned}$$

since

$$P(S_k^{\text{IP}}(x^0) | Y = y) = \begin{cases} 1, & \text{if } y \in \mathcal{A}_k \\ 0, & \text{if } y \notin \mathcal{A}_k \end{cases}$$

In particular, the classes in the active set have the *same relative weights* as at the outset. In summary:

$$p(y|S_k^{\text{IP}}(x^0)) = \begin{cases} p(y)/\sum_{\mathcal{A}_k} p(l), & y \in \mathcal{A}_k \\ 0, & \text{otherwise} \end{cases}$$

The key observation to prove the theorem is that if a hypothesis y generates the same answers to the first m or more questions as y^0 , and hence is active at step m , then its prior likelihood $p(y)$ is at most $2^{-(m-1)}$, $m = 1, 2, \dots$. This is intuitively clear: if y has the same answer as y^0 on the first question, and $p(y^0) > \frac{1}{2}$, then only one question is needed and the active set is empty at step two; if $q_1(y) = q_1(y^0)$ and $q_2(y) = q_2(y^0)$ and $p(y^0) > \frac{1}{4}$, then only two questions are needed and the active set is empty at step three, etc.

Finally, since C , the code length, takes values in the non-negative integers $\{0, 1, \dots\}$:

$$\begin{aligned} \tilde{H}_Q(X; Y) &:= \mathbb{E}[C] \\ &= \sum_{m=1}^{\infty} P(C \geq m) \\ &\leq \sum_{m=1}^{\infty} P(p(Y) < 2^{-(m-1)}) \\ &= \sum_{m=1}^{\infty} \sum_{y: p(y) < 2^{-(m-1)}} p(y) \\ &= \sum_{y \in \mathcal{Y}} \sum_{m=1}^{\infty} 1_{\{p(y) < 2^{-(m-1)}\}} p(k) \\ &= \sum_{y \in \mathcal{Y}} p(k)(1 - \log p(k)) \\ &= H(Y) + 1 \end{aligned}$$

□

A.3 Termination Criteria for IP

We would first analyze the termination criteria for the exact case, that is, $\epsilon = 0$ in (3), and then move on to the more general case.

Termination Criteria when $\epsilon = 0$ Ideally for a given input x^{obs} , we would like to terminate (IP outputs q_{STOP}) after L steps if

$$p(y | x^{\text{obs}}) = p(y | x') \quad \forall x' \in S_L^{\text{IP}}(x^{\text{obs}}), \quad y \in \mathcal{Y} \quad (26)$$

Recall, $S_L^{\text{IP}}(x^{\text{obs}}) = \{x' \in \mathcal{X} \mid \{q_i, q_i(x')\}_{1:L} = \{q_i, q_i(x^{\text{obs}})\}_{1:L}\}$. In other words, its the event consisting of all $x' \in \mathcal{X}$ which share the first L query-answer pairs with x^{obs} .

If (26) holds for all $x^{\text{obs}} \in \mathcal{X}$, then it is easy to see that this is equivalent to the sufficiency constraint in the case $\epsilon = 0$,

$$p(y | x) = p(y | \text{expl}_Q^{\text{IP}}(x)) \quad \forall (x, y) \in (\mathcal{X} \times \mathcal{Y})$$

where $p(y | \text{expl}_Q^{\text{IP}}(x)) := p(y | S_{t^{\text{IP}}(x)}^{\text{IP}}(x^{\text{obs}})) \quad \forall (x, y) \in (\mathcal{X} \times \mathcal{Y})$ and $t^{\text{IP}}(x)$ is the number of iterations IP takes on input x before termination.

Unfortunately, detecting (26) is difficult in practice. Instead we have the following lemma which justifies our stopping criteria for IP.

Lemma 1. *For a given input x^{obs} if event $S_L^{\text{IP}}(x^{\text{obs}})$ (after asking L queries) satisfies the condition specified by (26) then for all subsequent queries q_m , $m \geq L$, $\max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) = 0$.*

Refer to Appendix A.4 for a proof.

Inspired from Lemma 1 we formulate an optimistic stopping criteria as,

$$L = \inf \{k \in \{1, 2, \dots, |Q|\} : \max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) = 0 \quad \forall m \geq k, m \leq |Q|\} \quad (27)$$

Evaluating (27) would be computationally costly since it would involve processing all the queries for every input x . We employ a more practically amenable criteria

$$q_{L+1} = q_{STOP} \quad \text{if} \quad \max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) = 0 \quad \forall m \in \{L, L+1, \dots, L+T\} \quad (28)$$

$T > 0$ is a hyper-parameter chosen via cross-validation. Note, it is possible that there does not exist any informative query in one iteration, but upon choosing a question there suddenly appears informative queries in the next iteration. For example, consider the XOR problem. $X \in \mathbb{R}^2$ and $Y \in \{0, 1\}$. Let Q be the set to two axis-aligned half-spaces. Both

half-spaces have zero mutual information with Y . However, upon choosing any one as q_L , the other half-space is suddenly informative about Y . Equation (28) ensures that we do not stop prematurely.

Termination Criteria for general ϵ when d is taken as the KL-divergence For a general $\epsilon > 0$ we would like IP to terminate such that on average,

$$\mathbb{E}_X[KL(p(Y | X), p(Y | \text{expl}_Q^{\text{IP}}(X)))] \leq \epsilon \quad (29)$$

Detecting this is difficult in practice since IP is an online algorithm and only computes query-answers for a given input x . So it is not possible to know apriori when to terminate such that in expectation the KL divergence would be less than ϵ . Instead we opt for the stronger requirement that,

$$KL(p(Y | x), p(Y | \text{expl}_Q^{\text{IP}}(x))) \leq \epsilon \quad \forall x \in \mathcal{X}. \quad (30)$$

It is easy to see that (30) implies (29).

As before, $p(y | \text{expl}_Q^{\text{IP}}(x)) := p(y | S_{t^{\text{IP}}(x)}^{\text{IP}}(x^{\text{obs}})) \forall (x, y) \in (\mathcal{X} \times \mathcal{Y})$ and $t^{\text{IP}}(x)$ is the number of iterations IP takes on input x before termination. We have the following lemma (analogous to the $\epsilon = 0$ case).

Lemma 2. *We make the following assumptions:*

- 1) \mathcal{Y} is a countable set (recall $Y \in \mathcal{Y}$).
- 2) For any $x^{\text{obs}} \in \mathcal{X}$ and $x_1, x_2 \in S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$, we have $p(Y | x_1)$ and $p(Y | x_2)$ have the same support.

Then, for given input x^{obs} if event $S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$ (after asking $t^{\text{IP}}(x^{\text{obs}})$ queries) satisfies the condition specified by (30) then for all subsequent queries q_m , $m \geq t^{\text{IP}}(x^{\text{obs}})$, $\max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) \leq \epsilon'$, where $\epsilon' = C\epsilon$ for some constant $C > 0$.

Refer to Appendix A.5 for a proof. The assumption of \mathcal{Y} being countable is typical for supervised learning (the scenario considered in this paper) where the set of labels is often finite. The second assumption intuitively means that $P(y | x_2) = 0 \implies P(y | x_1) = 0$ for any $y \in \mathcal{Y}$. This is a reasonable assumption since we envision practical scenarios in which ϵ is close to 0 and thus different inputs which share the same query-answers until termination by IP are expected to have very ‘‘similar’’ posteriors.¹⁰

Inspired from Lemma 2 we formulate an optimistic stopping criteria $\forall x^{\text{obs}} \in \mathcal{X}$ as,

$$t^{\text{IP}}(x^{\text{obs}}) = \inf\{k \in \{1, 2, \dots, |Q|\} : \max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) \leq \epsilon' \forall m \geq k, m \leq |Q|\} \quad (31)$$

Evaluating (31) would be computationally costly since it would involve processing all the queries for every input x^{obs} . We employ a more practically amenable criteria

$$q_{t^{\text{IP}}(x^{\text{obs}})+1} = q_{\text{STOP}} \quad \text{if} \quad \max_{q \in Q} I(q(X); Y | S_m^{\text{IP}}(x^{\text{obs}})) \leq \epsilon' \forall m \in \{L, L+1, \dots, L+T\} \quad (32)$$

$T > 0$ is a hyper-parameter chosen via cross-validation.

A.4 Proof of Lemma 1

Proof. Recall each query q partitions the set \mathcal{X} and $S_L^{\text{IP}}(x^{\text{obs}})$ is the event $\{x' \in \mathcal{X} \mid \{q_i, q_i(x^{\text{obs}})\}_{1:L} = \{q_i, q_i(x')\}_{1:L}\}$. It is easy to see that if $S_L^{\text{IP}}(x)$ satisfies the condition specified by (26) then

$$P(y | S_m^{\text{IP}}(x^{\text{obs}})) = P(y | x') \forall x' \in S_m^{\text{IP}}(x^{\text{obs}}) \forall m \geq L, \forall q \in Q \quad (33)$$

This is because subsequent query-answers partition a set in which all the data points have the same posterior distributions. Now, $\forall q \in Q, \forall a \in \text{Range}(q), y \in \mathcal{Y}$

$$p(q(X) = a, y | S_m^{\text{IP}}(x^{\text{obs}})) = p(q(X) = a | S_m^{\text{IP}}(x^{\text{obs}}))p(y | q(X) = a, S_m^{\text{IP}}(x^{\text{obs}})) \quad (34)$$

(34) is just an application of the chain rule of probability. The randomness in $q(X)$ is entirely due to the randomness in X . For any $a \in \text{Range}(q), y \in \mathcal{Y}$

$$\begin{aligned} p(y | q(X) = a, S_m^{\text{IP}}(x^{\text{obs}})) &= \sum_{x'} p(y, X = x' | a, S_m^{\text{IP}}(x^{\text{obs}})) \\ &= \sum_{x'} p(y | X = x', a, S_m^{\text{IP}}(x^{\text{obs}}))p(X = x' | a, S_m^{\text{IP}}(x^{\text{obs}})) \\ &= \sum_{x'} p(y | X = x')p(X = x' | a, S_m^{\text{IP}}(x^{\text{obs}})) \\ &= p(y | S_m^{\text{IP}}(x^{\text{obs}})) \sum_{x'} p(X = x' | a, S_m^{\text{IP}}(x^{\text{obs}})) \\ &= p(y | S_m^{\text{IP}}(x^{\text{obs}})) \end{aligned} \quad (35)$$

10. We refer to the distribution $p(y | x)$ for any $x \in \mathcal{X}$ as the posterior distribution of x .

The first equality is an application of the law of total probability, third due to conditional independence of the history and the hypothesis given $X = x'$ (assumption) and the fourth by invoking ((33)).

Substituting (35) in (34) we obtain $Y \perp\!\!\!\perp q(X) \mid S_m^{\text{IP}}(x^{\text{obs}}) \forall m \geq L, q \in Q$. This implies that for all subsequent queries $q_m, m > L, \max_{q \in Q} I(q(X); Y \mid S_m^{\text{IP}}(x^{\text{obs}})) = 0$. Hence, Proved. \square

A.5 Proof of Lemma 2

Proof. **Condition (30) implies bounded KL divergence between inputs on which IP has identical query-answer trajectories** Recall $S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$ is the event $\{x' \in \mathcal{X} \mid \{q_i, q_i(x^{\text{obs}})\}_{1:t^{\text{IP}}(x^{\text{obs}})} = \{q_i, q_i(x')\}_{1:t^{\text{IP}}(x^{\text{obs}})}\}$. If $S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$ satisfies (30) then using Pinsker's inequality we conclude,

$$\delta(p(Y \mid x^{\text{obs}}), p(Y \mid S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}}))) \leq \sqrt{\frac{\epsilon}{2}} \quad (36)$$

Here δ is the total variational distance between the two distributions. Since δ is a metric we conclude for any $x_1, x_2 \in S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$,

$$\begin{aligned} \delta(p(Y \mid x_1), p(Y \mid x_2)) &\leq \delta(p(Y \mid x_1), p(Y \mid S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}}))) + \delta(p(Y \mid x_2), p(Y \mid S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}}))) \\ &\leq \sqrt{2\epsilon} \end{aligned} \quad (37)$$

Since \mathcal{Y} is countable, define $\eta = \min\{p(y \mid \hat{x}) : y \in \mathcal{Y}, p(y \mid \hat{x}) > 0, \hat{x} \in \mathcal{X}\}$. Then, by the reverse Pinsker's inequality we conclude,

$$KL(p(Y \mid x_1), p(Y \mid x_2)) \leq \frac{\epsilon}{\eta} =: \epsilon' \quad \forall x_1, x_2 \in S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}}) \quad (38)$$

Note, the above upper bound holds since by assumption 2, $p(Y \mid x_1)$ and $p(Y \mid x_2)$ have the same support.

Bounded KL divergence between inputs implies subsequent queries have mutual information bounded by ϵ For any subsequent query $q \in Q$ that IP asks about input x^{obs} we have $\forall x \in S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$,

$$KL(p(Y \mid x), p(Y \mid S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}}))) = \sum_Y p(Y \mid x) \log p(Y \mid x) - \sum_Y p(Y \mid x) \log p(Y \mid S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})) \quad (39)$$

where, $S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}}) := S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}}) \cap \{x' \in \mathcal{X} : q(x') = q(x)\}$. For brevity, we denote $S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})$ as B ,

$$\begin{aligned} \sum_Y p(Y \mid x) \log p(Y \mid B) &= \sum_Y p(Y \mid x) \log \left[\sum_{x' \in B} p(Y \mid x', B) P(x' \mid B) \right] \\ &= \sum_Y p(Y \mid x) \log \left[\sum_{x' \in B} p(Y \mid x') P(x' \mid B) \right] \\ &\geq \sum_Y p(Y \mid x) \sum_{x' \in B} p(x' \mid B) \log p(Y \mid x') \\ &= \sum_{x' \in B} p(x' \mid B) \sum_Y p(Y \mid x) \log p(Y \mid x') \end{aligned} \quad (40)$$

In the third inequality Jensen's inequality was used. Substituting (40) in (39),

$$\begin{aligned} KL(p(Y \mid x) \parallel p(Y \mid B)) &\leq \sum_{x' \in B} p(x' \mid B) \sum_Y p(Y \mid x) \log \frac{p(Y \mid x)}{p(Y \mid x')} \\ &\leq \epsilon' \sum_{x'} p(x' \mid B) \\ &= \epsilon' \end{aligned} \quad (41)$$

In the second inequality we substituted from (38) since $x, x' \in B \subseteq S_{t^{\text{IP}}(x^{\text{obs}})}^{\text{IP}}(x^{\text{obs}})$. In the third equality we used the identity $\sum_{x' \in B} p(x' \mid B) = 1$,

It is easy to see that (41) holds for all $x' \in S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})$ and thus $I(X; Y \mid S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})) \leq \epsilon$

Since $Y \rightarrow X \rightarrow q(X)$ we can apply the data-processing inequality to obtain,

$$I(q(X); Y \mid S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})) \leq I(X; Y \mid S_{t^{\text{IP}}(x^{\text{obs}})+1}^{\text{IP}}(x^{\text{obs}})) \leq \epsilon' \quad \forall q \in Q.$$

This implies that for all subsequent queries $q_m, m > t^{\text{IP}}(x^{\text{obs}}), \max_{q \in Q} I(q(X); Y \mid S_m^{\text{IP}}(x^{\text{obs}})) \leq \epsilon$. Hence, Proved. \square

A.6 Complexity of the Information Pursuit Algorithm

For any given input x , the per-iteration cost of the IP algorithm is $\mathcal{O}(N + |Q|m)$, where $|Q|$ is the total number of queries, N is the number of ULA iterations, and m is cardinality of the product sample space $q(X) \times Y$. For simplicity we assume that the output hypothesis Y and query-answers $q(X)$ are finite-valued and also that the number of values query answers can take is the same but our framework can handle more general cases.

More specifically, to compute $q_{k+1} = \arg \max_{q \in Q} I(q(X); Y \mid S_k^{\text{IP}}(x^{\text{obs}}))$. We first run ULA for N iterations to get samples from $p(z \mid y, S_k^{\text{IP}}(x))$ which are then used to estimate the distribution $p(q(x), y \mid S_k^{\text{IP}}(x))$ (using (15)) for every query $q \in Q$ and every possible query-answer hypothesis, $(q(x), y)$, pair. This incurs a cost of $\mathcal{O}(N + |Q|m)$. We then numerically compute the mutual information between query-answer $q(X)$ and Y given history for every $q \in Q$ as described in (18). This has a computational complexity $\mathcal{O}(|Q|m)$. Finally, we search over all $q \in Q$ to find the query with maximum mutual information (refer (6)).

It is possible to reduce N by using advanced MCMC sampling methods which converge faster to the required distribution $p(z \mid y, S_k^{\text{IP}}(x))$. The linear cost of searching of all queries can also be reduced by making further assumptions about the structure of the query set. For example, we conjecture that this cost can be reduced to $\log |Q|$ using hierarchical query sets where answers to queries would depend upon to answers to queries higher up in the hierarchy. Note that if the query answers $q(X)$ and Y are continuous random variables then we would need to resort to sampling to construct stochastic estimates of the mutual information between $q(X)$ and Y instead of carrying our explicit numerical computations. We would explore these directions in future work.

A.7 Network Architectures and Training Procedure

Here we describe the architectures and training procedures for the β -VAEs used to calculate mutual information as described in Section 4.

A.7.1 Architectures

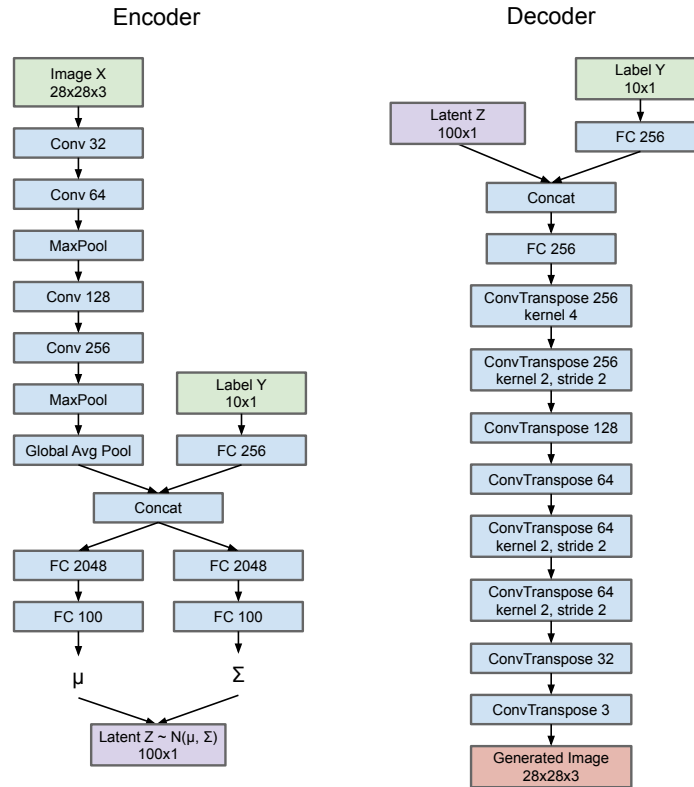


Fig. 7. Binary Image VAE

Binary Image Classification In the encoder, all convolutional layers use kernel size 3 and stride 1, and max pool layers use a pooling window of size 2. In the decoder all transposed convolutions use kernel size 3 and stride 1 unless otherwise noted. In both the encoder and decoder, all non-pooling layers are followed by a BatchNorm layer and LeakyReLU activation (with slope -0.3) except for the final encoder layer (no nonlinearities) and the final decoder layer (sigmoid activation).

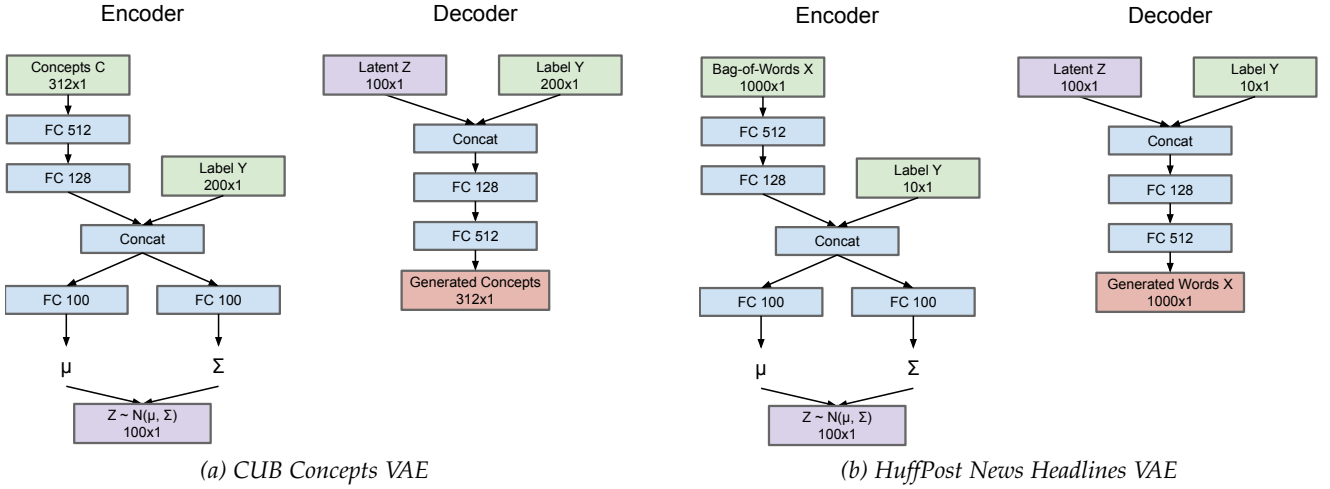


Fig. 8. CUB Concepts and HuffPost News Headlines VAEs

CUB Bird Species Classification and HuffPost News Headline Classification For CUB and HuffPost News, we use essentially the same VAE architecture, only designed to handle different-sized inputs X and one-hot labels Y . All layers are followed by a BatchNorm layer and ReLU activation except for the final encoder layer (no nonlinearities) and the final decoder layer (sigmoid activation).

A.7.2 Training

The β -VAE was trained by optimizing the Evidence Lower Bound (ELBO) objective

$$\max_{\omega, \phi} \text{ELBO}(\omega, \phi) = \sum_{i=1}^n \left[\mathbb{E}_{p_{\phi}(z|y^{(i)}, x^{(i)})} [\log p_{\omega}(x^{(i)} | z, y^{(i)})] - \beta D_{KL}(p_{\phi}(z | y^{(i)}, x^{(i)}) || p(z)) \right] \quad (42)$$

where $p_{\phi}(z | y^{(i)}, x^{(i)})$ denotes the encoder and $p_{\omega}(x^{(i)} | z, y^{(i)})$ the decoder. The prior over latents $p(z)$ is taken to be standard Gaussian.

For all binary image datasets, we trained our VAE for 200 epochs using Adam with learning rate 0.001 and the β -VAE parameter $\beta = 5.0$. We also trained the CUB VAE for 200 epochs using Adam with learning rate 0.001 but with $\beta = 1.0$. Finally, we trained the HuffPost News VAE for 100 epochs using Adam with learning rate 0.0002 and $\beta = 1.0$.

APPENDIX B

B.1 HuffPost News Task and Query Set

The Huffington Post News Category dataset consists of 200,853 articles published by the Huffington Post between 2012 and 2018. Each datapoint contains the article “headline”, a “short description” (a one to two-sentence-long continuation of the headline), and a label for the category/section it was published under in the newspaper. We concatenate the article headline and short description to form one extended headline. Additionally, many of the 41 category labels are redundant (due to changes in how newspaper sections were named over the years), are semantically ambiguous and HuffPost-specific (e.g. “Impact”, “Fifty”, “Worldpost”), or have very few articles. Therefore, we combine category labels for sections with equivalent names (e.g. “Arts & Culture” with “Culture & Arts”), remove ambiguous HuffPost-specific categories, and then keep only the 10 most frequent categories to ensure that each category has an adequate number of samples. This leaves us with a final dataset of size 132,508 and category labels “Entertainment”, “Politics”, “Queer Voices”, “Business”, “Travel”, “Parenting”, “Style & Beauty”, “Food & Drink”, “Home & Living”, and “Wellness”.

B.2 Comparison Models

B.2.1 MAP using Q

For IP, we use ULA to sample from $p(z | y, S_k^{IP}(x))$, but now that we have access to the all the query answers, $Q(x)$, we can improve performance by making use of the VAE’s encoder instead. Following equation 16, we draw many samples from the encoder $p(z | Q(x), y)$ and then decode these samples to estimate the VAE’s posterior distribution $p(y|Q(x))$. For each problem, we set the number of samples to be the same as what we draw during each iteration of IP (12,000 for binary image tasks, 12,000 for CUB, 10,000 for HuffPost News). We expect the accuracy of this model to serve as an upper bound for what we can achieve with IP given that it uses all queries. Our task-specific VAE architectures can be found in Appendix A.7.

B.2.2 Black-Box Using Q

We also compare to non-interpretable supervised models which receive all queries $Q(x)$ as input and try to predict the associated label y . This allows a comparison between the accuracy of the posterior of our generative model and traditional supervised approaches on the chosen interpretable query set.

For the binary image datasets, we use a simple CNN where all convolutional layers use kernel size 3×3 and a stride of 1, all max pooling uses a kernel size of 2. For CUB and HuffPost News we use simple MLPs.

TABLE 3
Binary Image CNN Architecture

Layer	Input Size/Channels	Output Size/Channels	Nonlinearity
Convolution	3	32	BatchNorm + ReLU
Convolution	32	64	BatchNorm + ReLU + MaxPool
Convolution	64	128	BatchNorm + ReLU
Convolution	128	256	BatchNorm + ReLU + MaxPool + Global Avg Pool
Fully-connected	256	2048	BatchNorm + ReLU
Fully-connected	2048	10	Sigmoid

TABLE 4
CUB Attributes MLP Architecture

Layer	Input Size/Channels	Output Size/Channels	Nonlinearity
Fully-connected	312	100	ReLU
Fully-connected	100	25	ReLU
Fully-connected	25	200	Sigmoid

TABLE 5
HuffPost Bag-of-Words MLP Architecture

Layer	Input Size/Channels	Output Size/Channels	Nonlinearity
Fully-connected	1000	100	ReLU
Fully-connected	100	25	ReLU
Fully-connected	25	10	Sigmoid

B.2.3 Black-Box

Since we ourselves pre-processed the cleaned 10-class version of HuffPost News, there are no reported accuracies in the literature to compare IP with. Therefore, as a strong black-box baseline, we fine-tune a pre-trained Bert Large Uncased Transformer model [92] with an additional dropout layer (with dropout probability 0.3) and randomly initialized fully-connected layer. Our implementation is publicly available at <https://www.kaggle.com/code/stewyslocum/news-classification-using-bert>.

B.3 Additional Example Runs

B.3.1 IP with Various Patch Scales

For binary image classification, we also experimented with patch queries of sizes other than 3×3 , from single pixel 1×1 queries up to 4×4 patches.

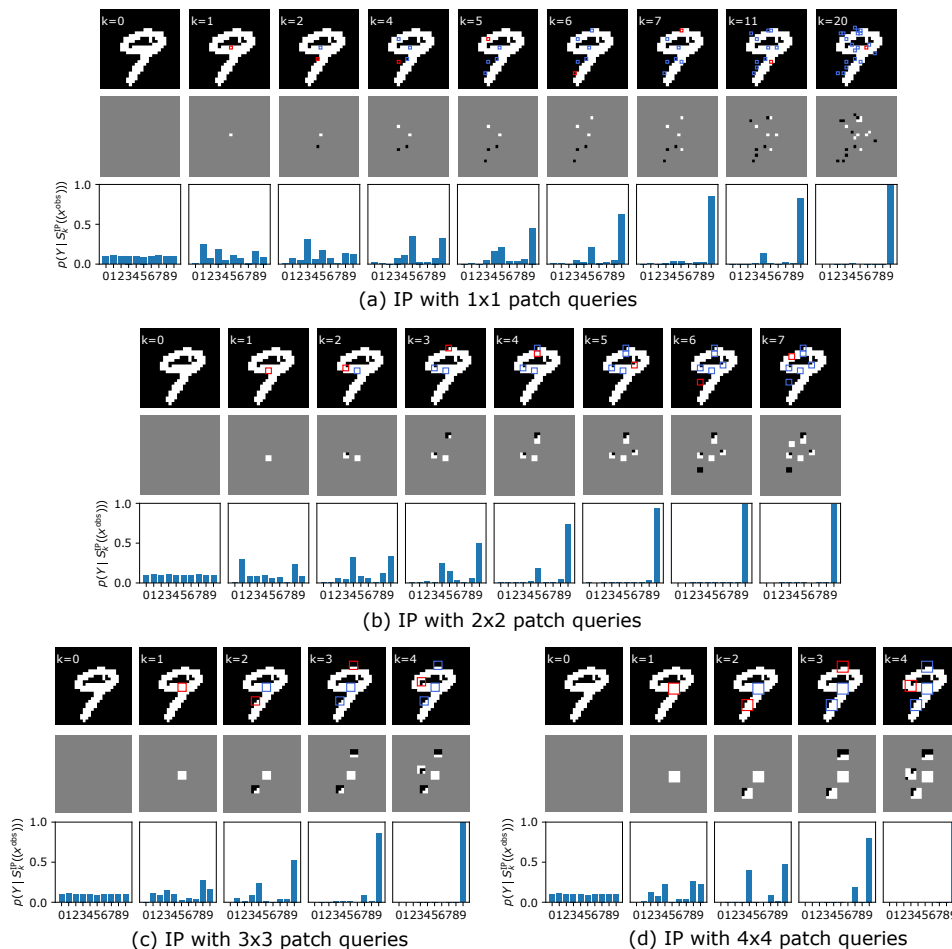


Fig. 9. IP on MNIST with different sized patch queries. In each subfigure, the top row displays the test image with red boxes denoting the current queried patch and blue boxes denoting previous patches. The second row shows the revealed portion of the image that IP gets to use at each query. The final row shows the model's estimated posteriors at each query. For conciseness in (a), we only display the 8 iterations of IP with highest KL divergence between successive posteriors (i.e. the most influential iterations). Observe that at all patch sizes, the queries chosen by IP cover roughly the same parts of the image, illustrating the importance of this region during classification. Reaching the stopping criteria of 99% posterior confidence takes 20 queries with 1×1 patches, 7 queries with 2×2 patches, 4 queries with 3×3 patches, and 4 queries with 4×4 patches.

TABLE 6
Number of queries and pixels gleaned by IP (until termination) using query sets of different patch scales on MNIST.

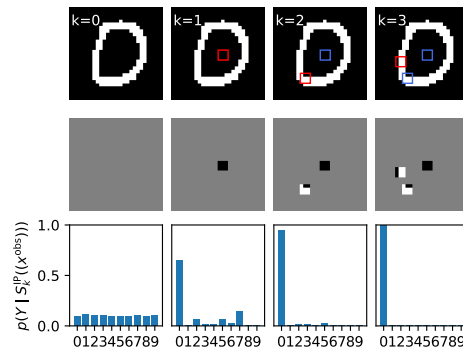
Patch Size	1×1	2×2	3×3	4×4
# Queries	21.1	9.6	5.2	4.6
# Pixels	21.1	32.8	44.7	54.7
% Pixels	2.7	4.2	5.7	6.9

While 1×1 patch queries use the smallest total number of pixels and most similarly resemble existing post-hoc attribution maps, they lead to a large number of scattered queries that are hard to interpret individually. On the other extreme, as the patch size grows larger, the number of total queries decreases, but queries becomes harder to interpret since each patch would contain many image features. On the MNIST dataset, we found 3×3 patches to be a sweet spot where explanations tended to be very short, but were also at a level of granularity where each patch could be individually interpreted as a single edge or stroke. Remarkably, at all chosen patch scales, only a very small fraction (2-7%) of the image needs to be revealed to classify the images with high confidence.

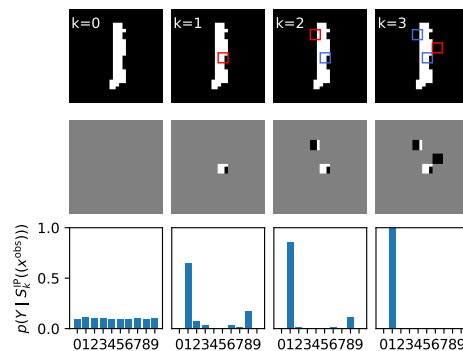
B.3.2 IP for Binary Image Classification

Now we provide additional example runs of IP for each of the three binary image classification datasets.

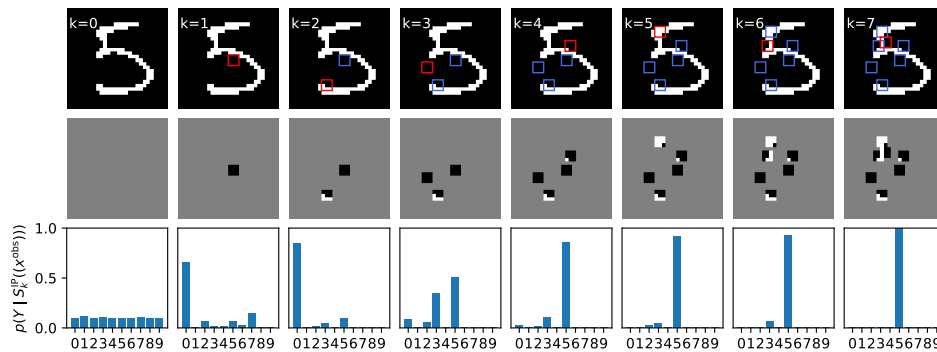
As in Figure 4 in the main paper, in each plot, the top row displays the test image with red boxes denoting the current queried patch and blue boxes denoting previous patches. The second row shows the revealed portion of the image that IP gets to use at each query. The final row shows the model’s estimated posteriors at each query.



(a) Each query reveals the patch with maximum mutual information with Y , conditioned on query history. This is initially independent of the particular image and asks for the pixel intensities in the center patch (see $k = 1$ in row 1). After the first query reveals that the center patch is all black, the posterior concentrates on “0” and “7”. After observing a white corner in the bottom left (which would be black for a “7”), the model becomes confident that the image is a “0”, and even more so after one final query when it reaches termination.



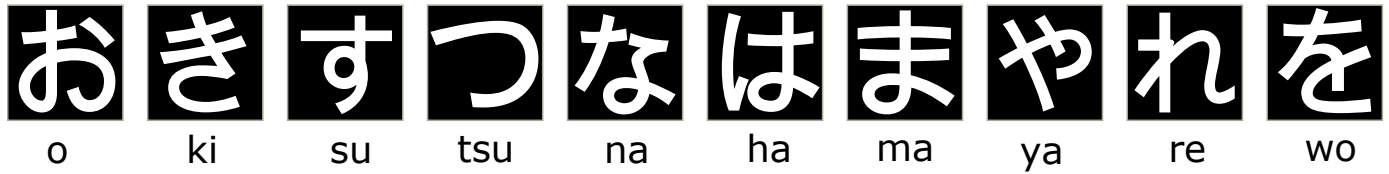
(b) The first query reveals a vertical white stroke in the center of the image, leading to a concentration of the posterior on “1”. In the next two queries, IP determines that there is a single long vertical stroke center stroke taking up the entire height of the image, and so it reaches a 99% confidence of the image being a “1”.



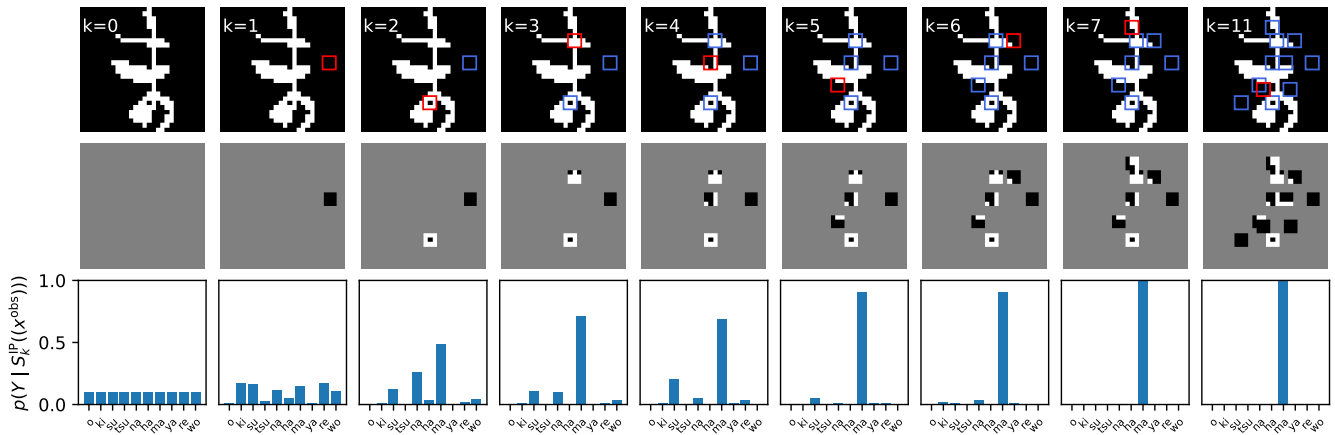
(c) As in the top left example, the first query is all black, and the posterior mass shifts onto “0”. Because the answer to this first query is identical as in the top left example, the second query chosen, in the bottom left of the image, is also the same. The response to this query is also a white corner as in the top left example, and so the posterior continues to concentrate on “0”, and the third query is also in the same left area of the image. However, this third query reveals a black patch, indicating that the image might be a “5” instead of a “0”. In the remaining four queries, IP discovers other portions of the “5” digit and finally arrives at the right answer with high confidence.

Fig. 10. Additional Examples of IP on MNIST

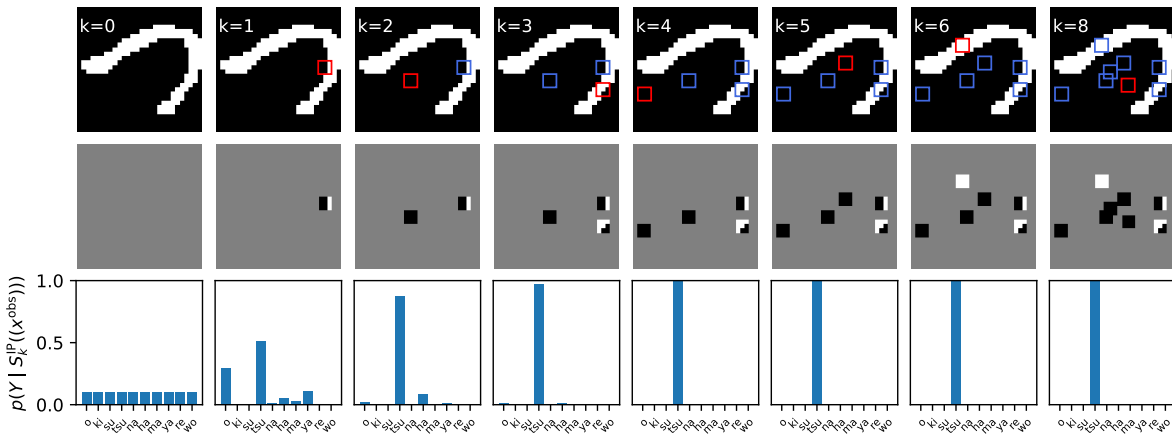
Recall that in order to improve performance on the KMNIST and FashionMNIST datasets (at the expense of asking a few more queries) we modified IP’s termination criteria to include a stability condition: terminate when the original criterion ($\max_Y p(Y|S_k^{IP}(x)) \geq 0.99$) is true for 5 queries in a row.



(a) KMNIST is a 10-class dataset of handwritten Japanese Hiragana characters. To assist the reader in understanding the examples below, we display each typed character along with its romanized name.

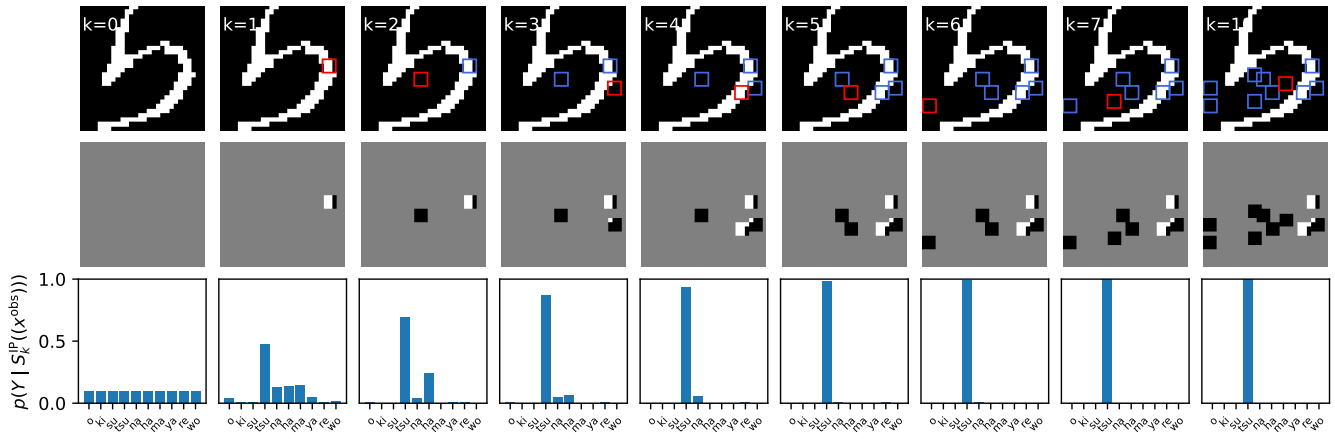


(b) For each image, IP selects the first query to be in the middle right of the image, where several characters are likely to have a stroke. Upon finding none, IP rules out “o”, “tsu”, and “ya” but otherwise distributes probability mass rather equally. On the second query, IP discovers a closed loop in the bottom of the character, a clear sign of “na” and “ma”, which increase the most. The discovery of the double crosses in the remaining queries concentrate the posterior on “ma” until termination. For conciseness, we only display the 8 iterations of IP with highest KL divergence between successive posteriors (i.e. the most influential iterations).



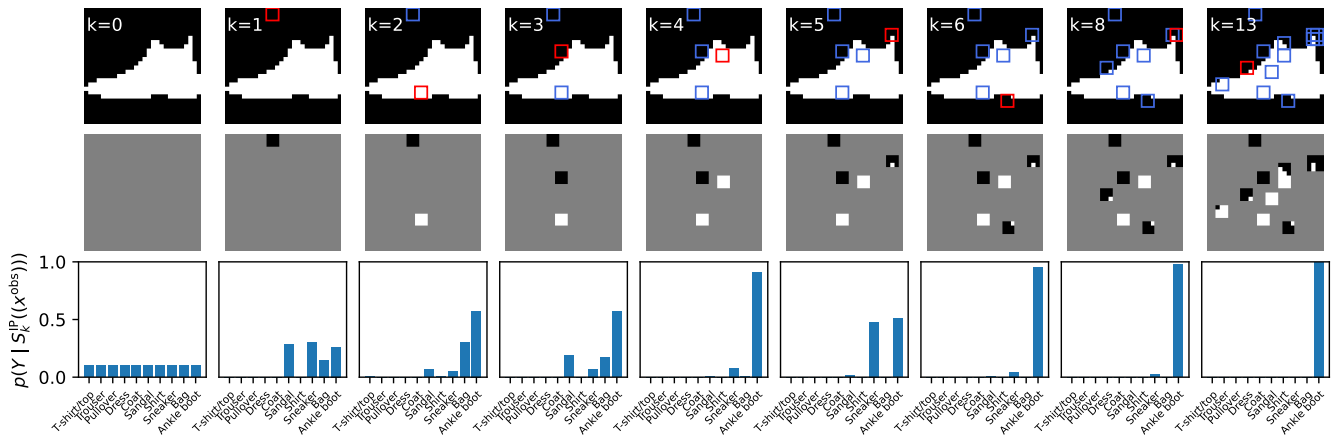
(c) In the first query, IP discovers a left edge, hinting at the presence of a large loop on the right side of the image, features of “o” and “tsu”. The second query is likely intended to disambiguate these two characters as “o” contains white strokes in this region. Upon finding a black patch here, IP is already confident, and reaches 99% confidence in just four queries.

Fig. 11. Additional Examples of IP on KMNIST

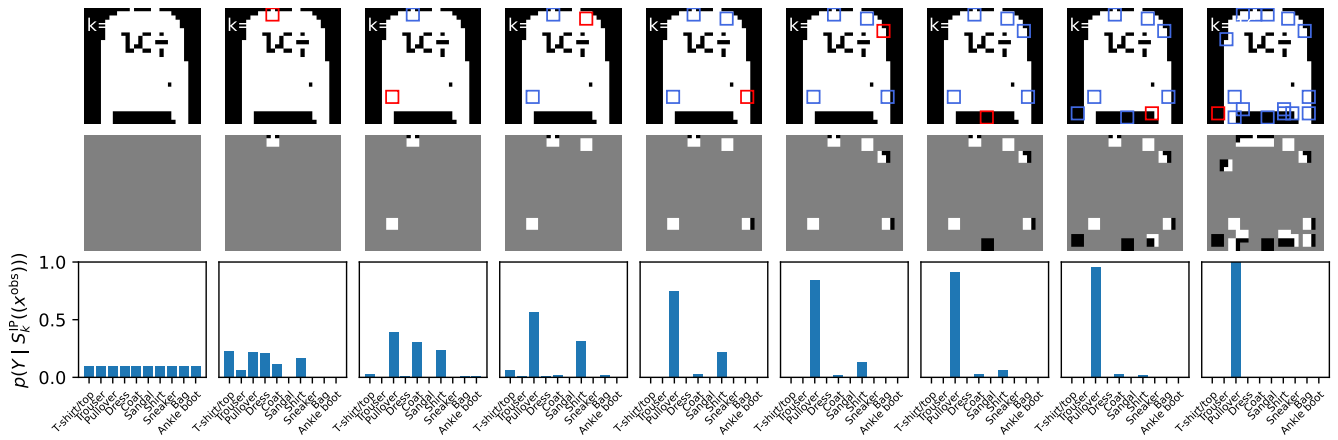


(d) The first query reveals a right edge, suggesting a slightly smaller loop on the right side of the image, a feature shared by several characters. However, most of these queries have a busy center region except for “tsu”, whose probability increases after the second query reveals a black patch. The next two queries outline the shape of the loop which is very large, a distinctive characteristic of “tsu”.

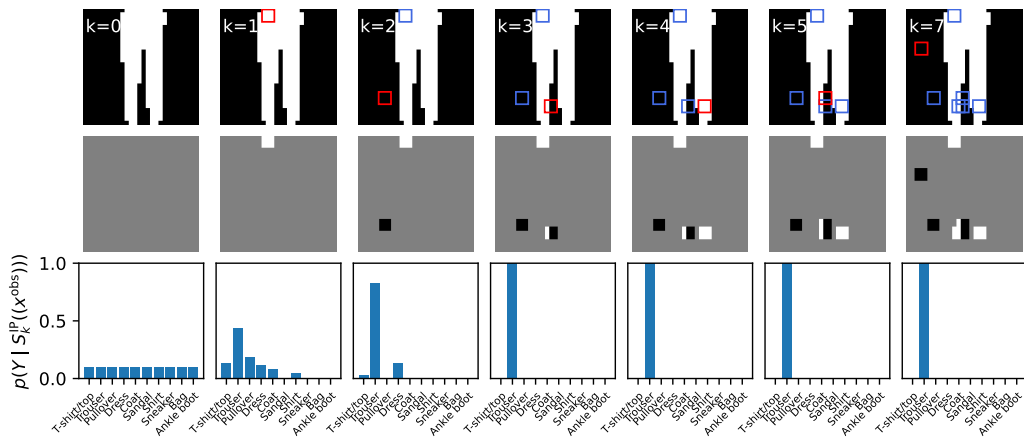
Fig. 11. Additional Examples of IP on KMNIST



(a) On this dataset, IP always selects the first query to be the patch in the top center, which being all black in this case rules out the possibility of the object being a type of pant or upper body garment, which would take up the entire height of the image. Over the next several queries, IP focuses on queries that would allow it to distinguish types of shoes from each other, in particular finding the shoe to have a high top and a small heel, eventually causing the posterior to concentrate on the correct “Ankle Boot” category. For conciseness, we only display the 8 iterations of IP with highest KL divergence between successive posteriors (i.e. the most influential iterations).



(b) The first query detects a white corner in the top center of the image, which hints at the presence of a collar, causing the posterior mass to move to the “Coat” and “Shirt” categories. Determining between these two categories is relatively difficult however, especially with binary images. But in general, coats tend to be bulkier than shirts. Therefore, after finding a white patch in iteration $k = 2$, the probability of “Coat” slightly increases, but as more partially black queries along the outside of the shirt are revealed and the slimmer outline of the shirt comes into view, the posterior converges on the correct category of “Shirt”.

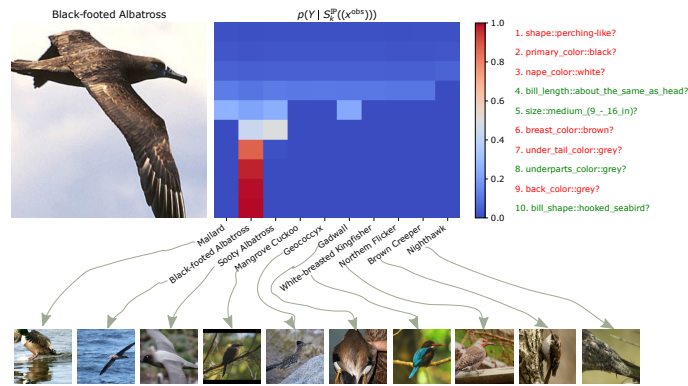


(c) Detecting an all white patch in the top center, IP rules out the possibility of the image being some type of shoe. The second query in the lower left returns a black patch, suggesting that the image is also not an upper body garment, which would take up the width of the image. In query $k = 3$, IP queries the center bottom of the image, discovering the space in between the two legs of the trouser, causing the posterior probability of “Trousers” to jump to nearly 100%, which remains stable over the last few queries.

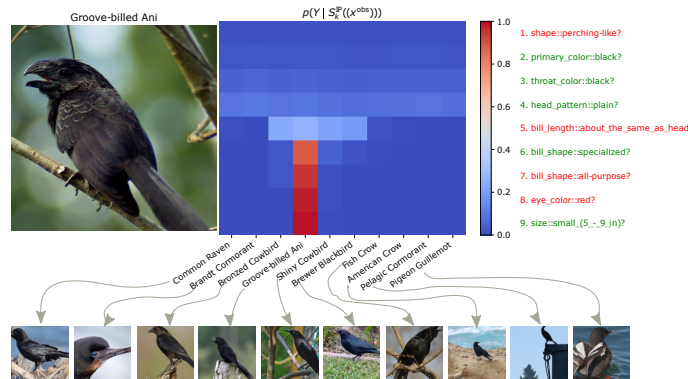
Fig. 12. Additional Examples of IP on FashionMNIST

B.3.3 IP for Bird Species Classification

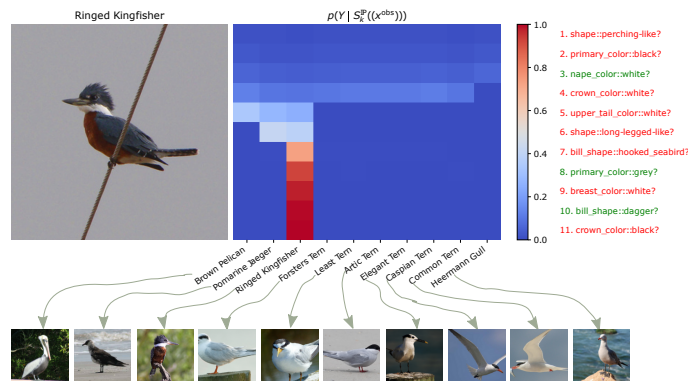
As in the main text, for each plot, on the left, we show the input image and on the right we have a heatmap of the estimated class probabilities per iteration. For readability, we only show the top 10 most probable classes out of the 200. To the right, we display the queries asked at each iteration, with red indicating a “no” response and green a “yes” response.



(a) The first few queries narrow down the potential species from 200 down to a small number. Since we only show the 10 most probable classes, the first few queries increase the probability of all shown classes. The first queries that distinguish among these classes concern bill length (which rules out the short-billed Nighthawk) and size (which rule out the smaller Mangrove Cuckoo, Geococcyx, Kingfisher, Northern Flicker, and Brown Creeper). The remaining birds are quite similar, all medium-sized brown water birds. The next two color queries suggest that the bird in question is a Black-footed Albatross, which is confirmed by the answers to the next few queries, which all match up with the characteristics of that bird.



(b) Again, the first few queries narrow down the likely species into the top 10 displayed classes. In just two queries (bill length and bill shape), IP distinguishes among these similar-looking, black, plain-headed birds that are hard for non-expert humans to differentiate between. Again, the last few queries serve to confirm the posterior prediction that the bird is a Groove-billed Ani.



(c) After establishing the top few most probable classes, IP converges on the class Ringed Kingfisher after just 7 queries. The last four queries simply serve to increase its confidence in its prediction.

Fig. 13. Additional Examples of IP on CUB Bird Species Identification

B.3.4 IP for HuffPost News Headline Category Classification

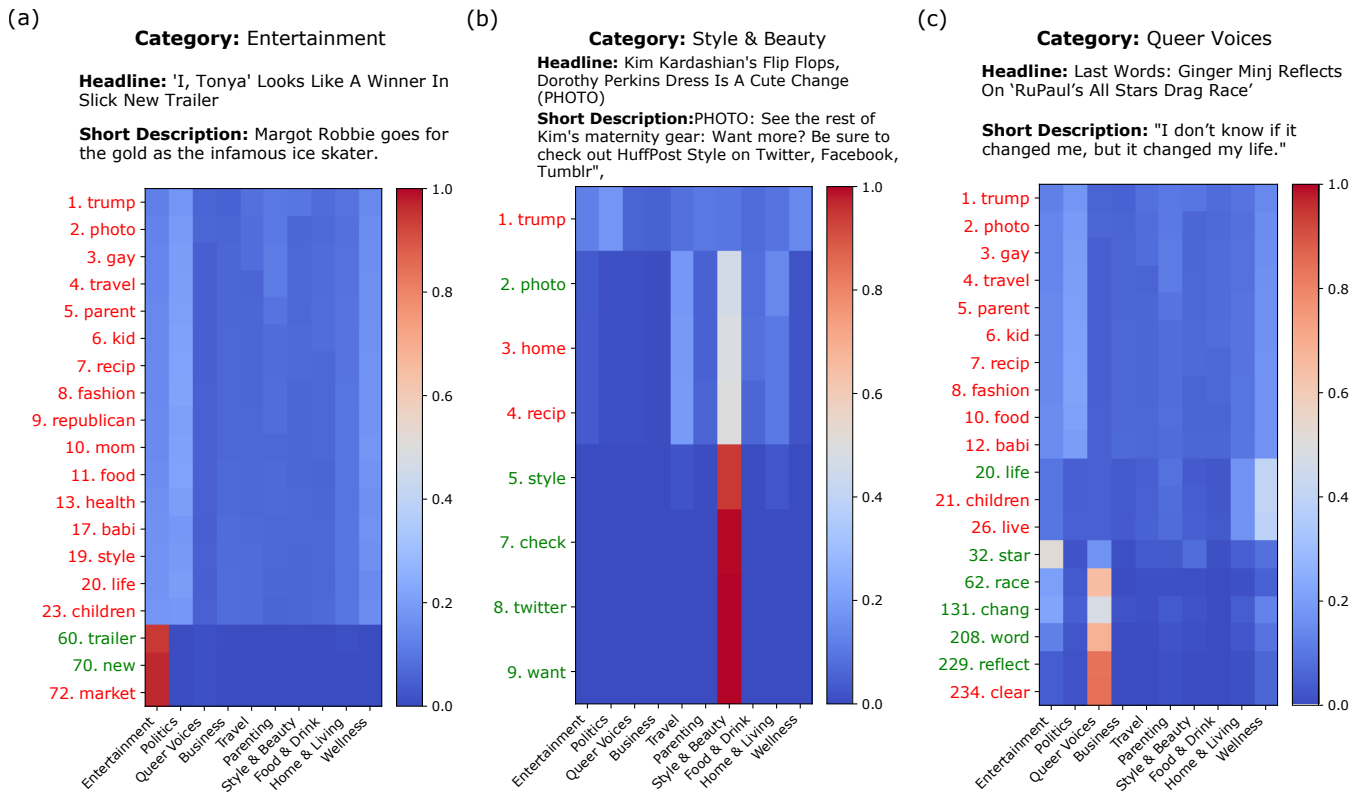


Fig. 14. **Additional Examples of IP on HuffPost News Headline Classification.** As before, when more than 20 queries were asked, we only display the 20 queries that led to greatest KL divergence between successive posteriors. **(a)** Because of the sparse structure of natural language, it typically takes a significant number of queries before the first word that is present in the sentence is found. Until this point, no query is particularly informative, and the posterior distribution remains mostly unchanged from the prior. However, at query 60, IP asks the word “trailer”, which is present in the extended headline. Naturally, the posterior shifts heavily towards “Entertainment”, and a few queries later IP reaches its termination criteria. Analyzing this run, we can say that IP reached its decision primarily because of the presence of the word “trailer”, leading us to say that this is a reasonable and trustworthy prediction. **(b)** This is an example of a relatively short explanation as IP happens to discover words present in the sentence after just two queries. Initially, the presence of “photo” causes the categories “Travel”, “Home & Living”, and “Style & Beauty” to become more probable. Several words later however, the word “style” is found, which is very strongly associated with the “Style & Beauty” category. **(c)** The posterior remains mostly unchanged until IP discovers the word “life”, which reasonably, shifts probability mass onto the “Wellness” category. However, several queries later, “star” is found to be present, which shifts the posterior away from “Wellness” onto “Entertainment” and “Queer Voices”. After discovering several more words that are present, “race”, “change”, “word”, “reflect”, the posterior progressively converges on “Queer Voices”, but still with relatively high uncertainty, likely because it never came across identifying words such as “drag”, which was not present in the vocabulary.