

---

# Using ensembles to quantify stability

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

The abstract paragraph should be indented  $\frac{1}{2}$  inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

## 1 Introduction

Quantifying prediction uncertainty is an active area of research in machine learning and elsewhere. A useful way to reason about uncertainty is to divide it into *aleatoric* vs *epistemic* uncertainty [1, 12, 8, 4, 3, 9]. Aleatoric uncertainty refers to uncertainty that is *inherent* to the system under observation, while epistemic uncertainty refers to uncertainty in the learning system *about* the system.

Consider predicting the outcome of flipping a coin whose (unknown) probability of heads is  $p = 1/4$ . If we know  $p$  we should always predict “tails” and will be incorrect with probability  $1/4$ . This uncertainty is irreducible, or aleatory, because no amount of additional data will change it. On the other hand, if we flip the coin  $n$  times and find that  $m$  of the outcomes were “heads” we can estimate the true value of  $p$  to be approximately  $\frac{m}{n} \pm \frac{1}{\sqrt{n}}$ . The range  $\pm \frac{1}{\sqrt{n}}$  is the epistemological or reducible uncertainty, because it relates to our *knowledge* regarding  $p$  which can be improved by increasing  $n$ .

For another example, consider our model to be a DNN for binary classification. We consider a single input  $x$  our goal is to predict the associated label  $y \in \{-1, +1\}$ .

Denote by  $f : X \rightarrow \mathbb{R}$  a specific DNN (with a specific set of weights). Assume first that this function is fixed. The aleatoric, or intrinsic uncertainty is captured by the conditional probability  $p_a = P(y = +1 | f(x) > a)$ . This so-called “calibration function” is “intrinsic”, it depends only on the true distribution and on the function  $f$ .

On the other hand the epistemic uncertainty corresponds to *model uncertainty* i.e. the uncertainty that we have about the function  $f$ .

Measuring aleatory using ensembles

Stable and unstable examples. An interesting and useful property of measuring  $x$  we create different

## 2 Related Work

[14] introduced a framework to deal with multiclass classification task in a one-vs-rest(OVR) way, which is similar to ours that treat each 10-class classification task as 10 binary classification tasks. It also trains multiple classifiers. However, its classifiers are not independent and they work in a stacked order and each classifier is responsible for one class. This paper also pay attention to easy classes and confusing classes and it uses the classes’ easiness to arrange the order of the classifiers. The key differences to our work are(Need further verify, I haven’t carefully read the details yet): (1)It need labels to define the confusing classes. (2)Whether an example is confusing only matters in the training stage, it no longer cares if it’s confusing in the inferencing stage.

[16] focus on building hash for images to make search and retrieve of the images efficient. And it pay attention to the easy and hard examples. However, it needs labels to tell which one is easy or hard. And it only focus on the easiness on training process.

[13] This paper mentioned a concept: Focal Loss(Not originated in this paper). It also focus on the uncertainty of the classifier's prediction. But simply using the difference between the predicted score and ground truth as the measure of uncertainty which is very different from ours.

I will look for more paper focusing on easy-confusing examples. Till now, I find most works that pay attention to the easiness and confusion of examples are only trying to make use of it in the training stage to improve the training process. After they get the trained framework, they no longer care about whether an example is easy or confusing on the inference stage, but to simply predict a label in the classical way. But our framework can distinguishing the hardness of the examples on the inference stage, and is able to tell the uncertainty quantitatively.

- [4]: Add dropout in the inference stage. This can generate an 'ensemble' without really training multiple predictors. However, this paper doesn't propose novel ways to make use of the ensemble outputs. Actually the mc-dropout method can be combined with ours.
- [11] It does three things: (1) Directly predict uncertainty as part of the predictor output(i.e, the predictor's output should be a normal distribution with mean value and uncertainty rather than just a single value) and properly add this uncertainty into the loss function. (2)Use adversarial training (3)Train an ensemble to estimate the predicted mean value and uncertainty. Difference to ours: (1)The way to make use of the ensemble output is different from us: treat the ensemble as a uniformly-weighted gaussian mixture. (2)Doesn't mention predicting a SET.
- [5] Constructs a method to find a proper threshold for confidence-rate(Can be almost any function) so that only the confident predictions will be made(Say don't know for unconfident). Difference to ours: (1)doesn't use ensemble(The author mentioned ensemble saying using ensemble costs too much). (2) Doesn't mention predicting a SET. Seems its way to choose threshold can be applied to our method. We can compare our results with this by regarding the Set size  $\neq 1$  as 'don't know'.
- [7] Following the work in [5], this paper tries to find a good method to measure the confidence-rate. It also uses an ensemble, but the members in the ensemble are not trained independently but rather predictors obtained in different epochs in the training process.
- [6] Designed a way to train predictors that output uncertainty directly. Neither use ensemble nor mention prediction set. We can compare our results with this by regarding the Set size  $\neq 1$  as 'don't know'.

### 3 Theory

Many learning algorithms are based on finding the single best rule. Ensemble methods such as Bagging and Random Forests are based on combining many rules, each of which is close to optimal. Averaging/voting across the rules results in a more stable and more accurate rule. We use a similar logic for classification with one important twist: when the ratio between +1 and -1 predictions is close to 50%/50% our combination rule outputs 0, which can be interpreted as "I don't know" (IDK).

Figure 1 is a sketch of ensemble based classification. The central assumption is that models that are in the close the best model, in terms of the probability of disagreement, are almost as accurate as the best model. models. This implies that under slightly different training conditions that sub-optimal model can become the best. We call this set of models the "support set" and divide the instance space into the "easy" instances, on which all models in the support set agree. The "hard" instances are those that are not "easy". The "hardness" of hard examples is measured by the the difference between the probability <sup>1</sup> of the two parts to which the instance divides the version space.

An ensemble is a sample from the support set. Procedurally, it is generated by perturbing the training process. In this paper we consider three ways of perturbing the training process, a few more possibilities are listed in the conclusions.

---

<sup>1</sup>Assuming some fixed prior distribution. Yoav: ?

1. **Bootstrap** perturb the training process by selecting from the size  $n$  training set a new training set of size  $n$  by sampling from the  $n$  times with replacement.
2. **Starting point** Choose the starting point for gradient descent multiple times independently at random.
3. **Architecture** Choose a different architecture for each ensemble member.
4. **human labels**

We divide our discussion to the We use the binary classification problem as a building block of our solution for multi-label prediction.

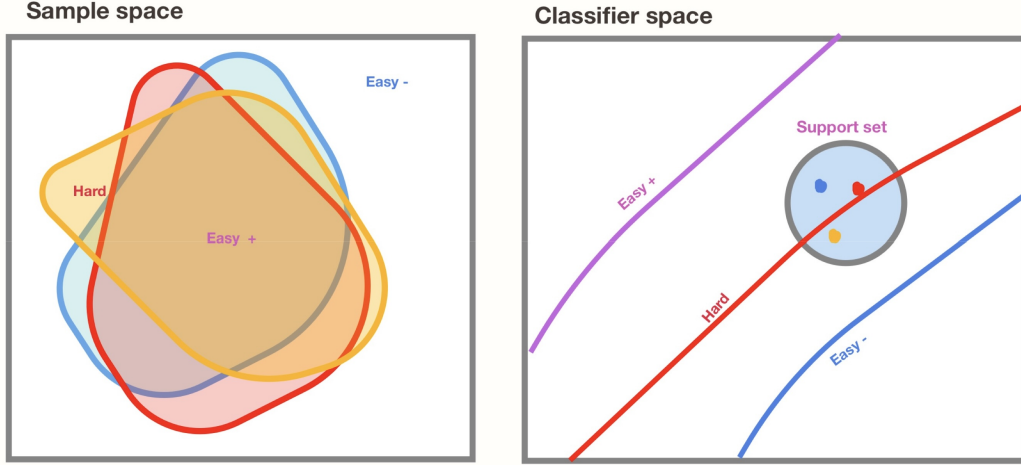


Figure 1: **Duality between classifiers and examples** Each point in the sample space figure corresponds to a sample, the three shaded sets corresponds to an ensemble of 3 classification rules with similar training error. “Easy -” is an instance on which the ensemble predicts unanimously “-”, while “Easy+” is an instance on which the unanimous prediction is “+”. It is this unanimity that makes the instance “easy”. The hard instance is one on which classifiers disagree. The classifier space (parameter space), demonstrates the same relationships between instances and classifiers, but here each classifier is a point and the instances define boundaries between classifiers that predict “+” or “-” on the instance. The *support set* is the set of classifiers whose performance is close to optimal. The three dots correspond to the ensemble of three classifiers. Finally, easy instances are ones on which the rules in the support set are unanimous, while the hard example splits the support set in two.

### 3.1 Binary Case

In a standard binary classification task, given an input  $\mathbf{x}$ , the predictor will output either  $+1$  or  $-1$ . In our framework the model can predict 0 in addition to  $-1$  and  $+1$ , where 0 stands for “hard instance” or “I don’t know the label”.

use the binary classification or the logit value generate from the score, instad, we use the score by a DNN.

Usually in a binary classification task, a neural networks predictor calculates a real value. In the training stage this value can be further transfered by a sigmoid funtion to a real value ranging from 0 to 1 and then plugged into the loss function; while in the inferencing stage, this value can be compared with 0 to decide whether this example belongs to the positive class or the negative class. Denote this value calculated by predictor  $f(\cdot, \mathcal{P})$  as  $O_f(\cdot, \mathcal{P})$ . If  $O_f(x, \mathcal{P})$  is larger than 0, the example will be classified into the positive class, otherwise the negative class.

**Assumption** As we change the perturbation  $\mathcal{P}$  to get  $O_f(x, \mathcal{P}_1), O_f(x, \mathcal{P}_2), \dots, O_f(x, \mathcal{P}_i), \dots, O_f(x, \mathcal{P}_E)$ , the values  $O_f(x, \mathcal{P}_i)$  obeys a normal distribution  $\mathcal{N}(\mu, \sigma^2)$ .

To better fit the assumption of normality we use the score generated by the DNN rather than the logit computed from the score,.

### Yoav to simplify

The mean value  $\mu$  should be positive if this example belongs to the positive class and negative if it belongs to the negative class. After training an ensemble  $f(\cdot, \mathcal{P}_1), f(\cdot, \mathcal{P}_2), \dots, f(\cdot, \mathcal{P}_i), \dots, f(\cdot, \mathcal{P}_E)$ , we can estimate how confident we are to say  $\mu$  is positive or negative based on  $O_f(x, \mathcal{P}_1), O_f(x, \mathcal{P}_2), \dots, O_f(x, \mathcal{P}_i), \dots, O_f(x, \mathcal{P}_E)$ . This can be done by *t-test*. The null hypothesis is  $\mu = 0$ , and the *t*-statistic is:

$$t = \frac{\overline{O}_f(x, \mathcal{P})}{S/\sqrt{E}} \quad (1)$$

where

$$\overline{O}_f(x, \mathcal{P}) = \frac{1}{E} \sum_{i=1}^E O_f(x, \mathcal{P}_i) \quad (2)$$

$$S^2 = \frac{1}{E-1} \sum_{i=1}^E (O_f(x, \mathcal{P}_i) - \overline{O}_f(x, \mathcal{P}))^2 \quad (3)$$

denote the *p*-value of the test as *p*, we define our confidence as:

$$confidence = -sign(\overline{O}_f(x, \mathcal{P})) * \log(p) \quad (4)$$

If the predictors are very confident that the example belongs to the positive class, the confidence should be a very large positive value; if the predictors are confident of the negative class, it should be a very small negative value. If the predictors are unsure, the confidence value would be close to 0. We can set a threshold to discriminate class+ and class-. Further, with this confidence value, we can set one threshold for confidence class+ and another threshold for confident class-. In the regime between the two thresholds, the predictors will say "I don't know".

### 3.2 Multiclass Case

In a multiclass classification task,  $\mathbf{O}_f(x, \mathcal{P})$  is no longer a real value but a  $R^C$  vector,  $C$  is the number of classes. In the training stage, a softmax can be operated on  $\mathbf{O}_f$  and further used in the loss function. In the inferencing stage,  $\max_j \mathbf{O}_{f,i}$  will be taken as the final prediction result, where  $\mathbf{O}_{f,j}$  refers to the *j*th dimension of the  $\mathbf{O}_f$  vector.

To apply our *t-test*-based method, we take each multiclass classification task as multiple binary classification tasks. Based on  $\mathbf{O}_{f,j}(x, \mathcal{P}_1), \dots, \mathbf{O}_{f,j}(x, \mathcal{P}_i), \dots, \mathbf{O}_{f,j}(x, \mathcal{P}_E)$ , we can calculate the *t*-statistic  $t_j$  and get the *p*-value  $p_j$  and finally the confidence value  $confidence_j$  for class *j*.

**Yoav:** Define set predictions and associated error, size, define performance measure (relatively to random), roughly:

- associate a set of labels with each labeler: if confidence rated, use threshold on the *p*-value. If human: sort labels in decreasing order of count and take prefix that corresponds to 95% of the counts.
- measure jaccard distance between the two sets.

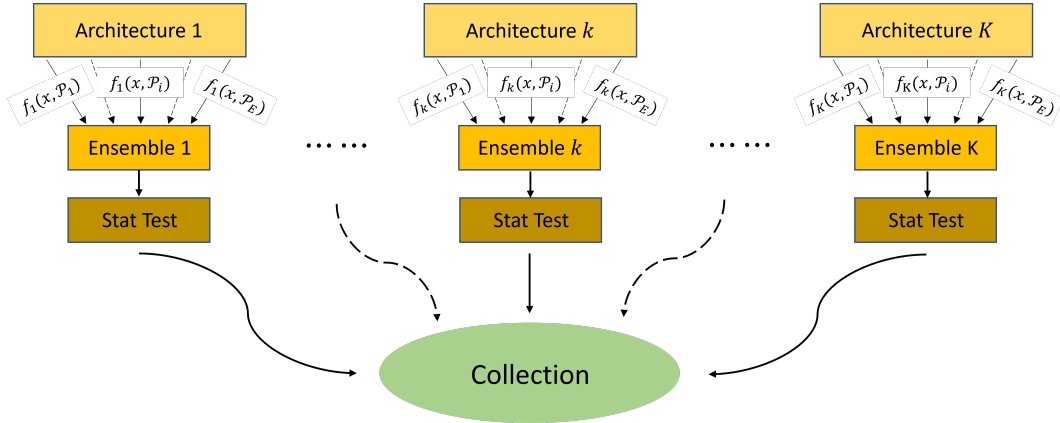
## 4 Performance on multi-label

- compute different ensembles (bootstrap of weak, random starting point of week, bootstrap of strong, random starting point of strong, combining many architectures) and create table with row per size of prediction set, column for range of Jaccard distance, in each box of table: fraction of examples that fall there.

We can have different ways to make use of the confidence values, for example:

1. The class with the largest confident value will be the predicted class.
2. Instead of predicting one class, the class with the largest confident value and classes whose confidence is no smaller than the largest by a certain amount will form a prediction set.
3. We can set a threshold as in the binary case. If  $confidence_j$  is smaller than the threshold, then the example doesn't belong to class *j*, otherwise the predictors think the example belongs to class *j*. Similarly, we can set two thresholds for confident of class *j*, confident not class *j* and "I don't know".

The second and the third ways bring a problem: the example may belong to multiple classes or none of the classes. When this happens, it suggests that this example is a "hard" example, because it makes the predictors confused among these classes. The experiment section shows that this confusion indeed happens for human. Therefore, it can be benefit to allow the predictors to say that they are confused and further careful actions are needed to make a decision(For example, call stronger predictors to classify the example), rather than simply give a prediction without confidence.



### 4.1 Three types of Perturbations: BootStrap, Architecture Starting Points and Architecture

**Bootstrap** We can construct a training set  $\mathcal{T}$  by randomly draw  $M$  examples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^M$  from  $\mathcal{D}$  with replacement. We can sample a set  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_i, \dots, \mathcal{T}_E$  independently, and train a  $f$  on each of the training set. This will generate a set of predictors  $f(\cdot, \mathcal{T}_1), f(\cdot, \mathcal{T}_2), \dots, f(\cdot, \mathcal{T}_i), \dots, f(\cdot, \mathcal{T}_E)$  with the same architecture but different parameters.

## 5 Experiment Result

Table 1: **Basic Properties of the Architectures:** Parameter numbers refers to the number of trainable parameters in each architecture. Forward Time refers to the averaged forward propagation time for one CIFAR image. We tested these values on 3 different devices: Nvidia RTX 3080Ti, Nvidia RTX 3060 and 11th Gen Intel(R) Core(TM) i5-11400F CPU. The last two columns show the accuracy of single predictors trained with different perturbation methods. Bootstrap samples 30000 examples from the 50000 training data.

|                | Parameter Numbers | Forward Time(ms) |          |        | Single Predictor Accuracy |              |
|----------------|-------------------|------------------|----------|--------|---------------------------|--------------|
|                |                   | RTX 3080Ti       | RTX 3060 | CPU    | Bootstrap                 | Random Start |
| shufflenet     | 925,618           | 2.016            | 7.815    | 13.995 | 0.88802                   | 0.895520     |
| shufflenetv2   | 1,263,854         | 1.867            | 7.341    | 5.931  | 0.88535                   | 0.890430     |
| mobilenetv2    | 2,296,922         | 1.717            | 6.694    | 7.592  | 0.89370                   | 0.889671     |
| regnetx        | 2,321,946         | 2.115            | 7.400    | 8.260  | 0.93330                   | 0.938220     |
| mobilenet      | 3,217,226         | 0.777            | 2.981    | 3.424  | 0.87632                   | 0.873603     |
| efficientnetb0 | 3,599,686         | 2.474            | 10.068   | 10.929 | 0.88585                   | 0.896610     |
| googlenet      | 6,166,250         | 2.990            | 10.722   | 30.060 | 0.93784                   | 0.947831     |
| densenet121    | 6,956,298         | 4.905            | 18.762   | 22.279 | 0.93735                   | 0.945434     |
| resnext29      | 9,128,778         | 1.335            | 4.723    | 24.190 | 0.93952                   | 0.948639     |
| resnet18       | 11,173,962        | 0.930            | 3.375    | 7.130  | 0.93688                   | 0.948173     |
| senet18        | 11,260,354        | 1.417            | 5.282    | 8.088  | 0.93591                   | 0.945290     |
| simplerla      | 15,142,970        | 1.769            | 6.415    | 12.596 | 0.93045                   | 0.943860     |
| vgg19          | 20,040,522        | 0.855            | 3.154    | 6.660  | 0.91733                   | 0.931955     |
| dpn92          | 34,236,634        | 4.712            | 17.596   | 58.489 | 0.94122                   | 0.948508     |

### 5.1 Stability of easy examples

As explained in section 3.2, each multiclass example can be transformed to multiple binary case. In our settings, each image can be divided to 10 image-class pairs. Each pair will have a label + if the image belongs to this class, or - if it doesn't. To compare the difference of the confidence values by weak and strong architectures, we assign a  $[S, W]$  coordinates for each image-class pair. The  $S$  value represents the confidence value for a strong architecture, while the  $W$  represents a weak. In this way, the 10000 test examples will be 100000 points on the  $S - W$  plane. Figure 3 shows the density distribution of these points with KDE plots.

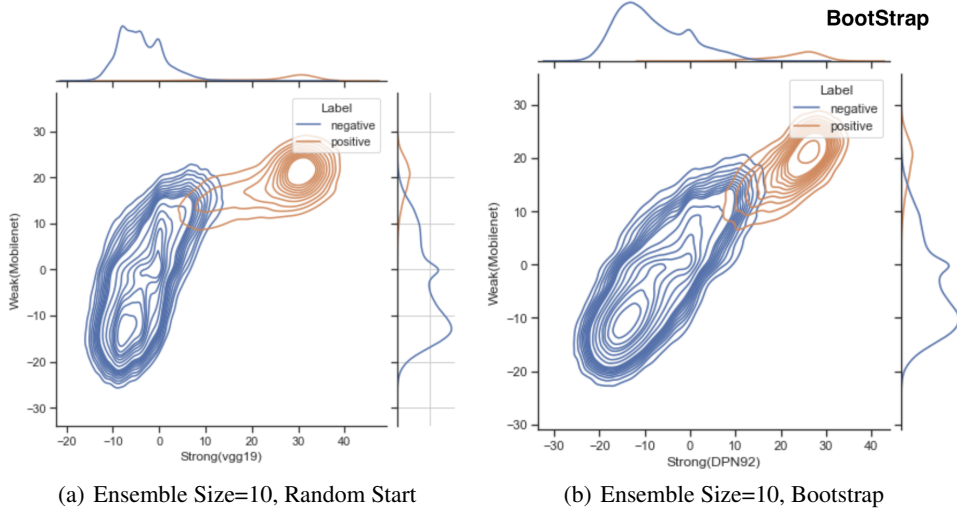


Figure 3: **KDE Plot of Confidence of Strong and Weak Architectures:** The strong architecture is Vgg19, the weak architecture is MobileNet  
Add to this figure: compare bootstrap to random starting point on weak, same on strong. Mark with circles the easy positive and the easy negative.

## 5.2 Performance on multi-label

## 5.3 Comparison with Human labels

With our t-test based method, we can have different ways to combine the outputs of different architectures. Here we illustrate two ways and compare their results with human prediction.

Following the framework shown in 4, we use the 3rd way in section 3.2 and only set a single threshold for each architecture and allow prediction sets rather than only a single class as the output. Then for each image, each architecture will give a prediction set. We regard every class in the set as getting one vote from the architecture. Combine the sets by all the architectures, we get a collection. In this collection, each class among the 10 will get zero or one or multiple votes from the architectures. Based on the frequencies of each class's votes, we can measure how likely it is for each class to be chosen as a final prediction. With CIFAR10H, we can compare the results of human labeler and our machine predictors. Figure 4 illustrates some examples that human and machine have the same confusion.

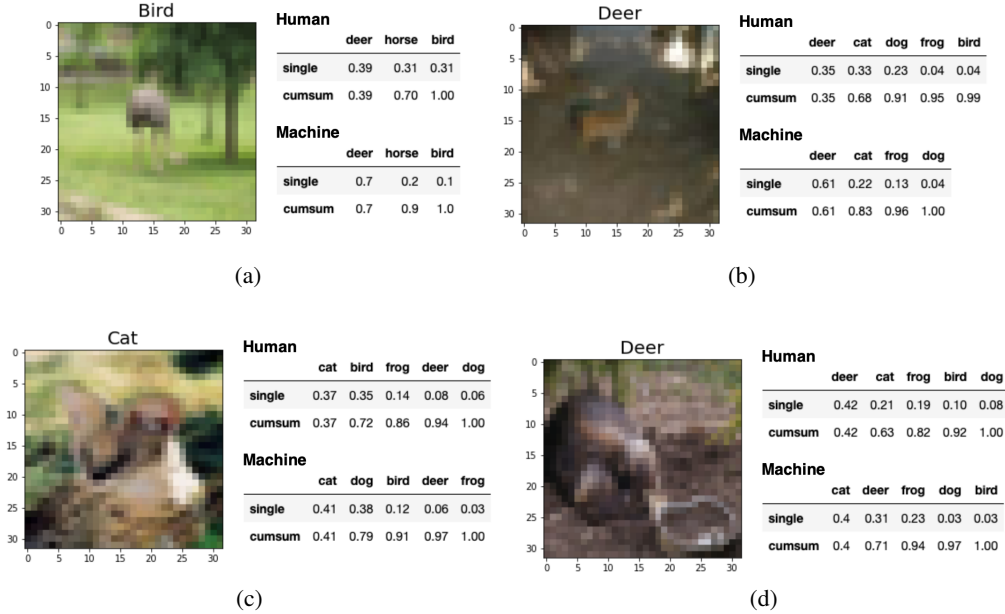


Figure 4: **Examples That Human and Machine have the Same Confusion:** The tables shows the frequencies of votes of human labeler and machine predictors. The rows with index 'single' are the frequencies for each single class, the unshown classes all have zero frequencies. The rows with index 'cumsum' are the culmulative sum of the above row.

Table 2: **Running Time of Weak-Strong Flow:** The weak model is MobileNet, the strong model is DPN92, the ensemble size for both architectures are 10. Batch Size indicates how many images are processed parallelly by the weak-strong flow. The unit for all values is ms.

| Batch Size | Weak Run | Weak T-Test | Strong Run | Strong T-Test | Weak-Strong Flow Total |
|------------|----------|-------------|------------|---------------|------------------------|
| 16         | 0.54016  | 0.04063     | 3.38590    | 0.03598       | 4.00267                |
| 64         | 0.32618  | 0.01138     | 2.32941    | 0.00957       | 2.67653                |
| 256        | 0.42272  | 0.00420     | 2.83185    | 0.00308       | 3.26185                |
| 1024       | 0.33215  | 0.00329     | 2.68724    | 0.00153       | 3.02421                |

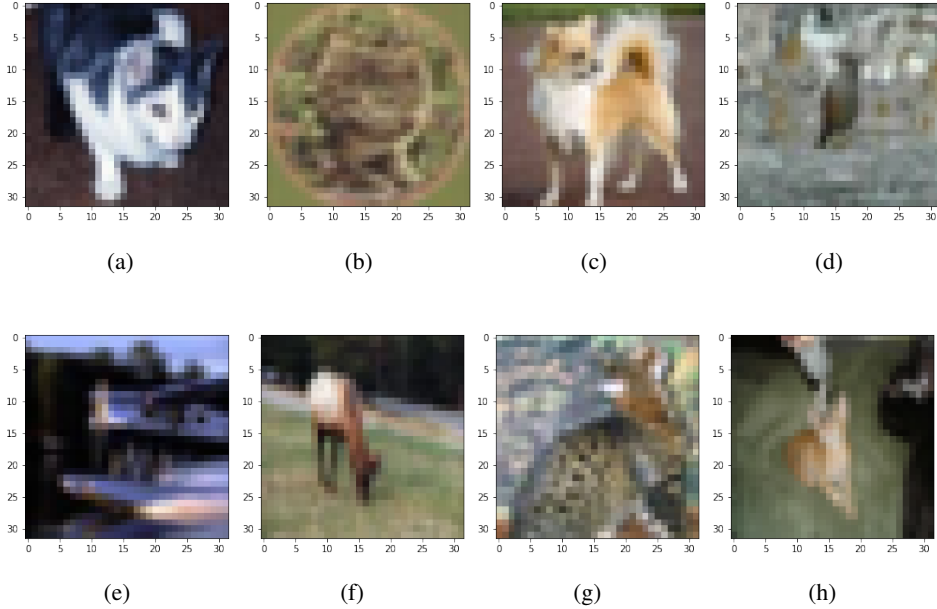


Figure 5: **Examples of Human and Machine Classification** (a)**Human:** Cat(1.0). **Machine:** Dog. (b)**Human:** Frog(0.49), Cat(0.47), Bird(0.04). **Machine:** Frog. (c)**Human:** Dog(1.0). **Machine:** Dog. (d)**Human:** Bird(0.92), Ship(0.02), Frog(0.02), Dog(0.02), Cat(0.02). **Machine:** Cat, Deer. (e)**Human:** Airplane(0.56), Ship(0.44). **Machine:** Truck, Airplane. (f)**Human:** Deer(0.6), Horse(0.4). **Machine:** Deer, Horse. (g)**Human:** Cat(0.72), Frog(0.09), Deer(0.09), Bird(0.08), Airplane(0.02). **Machine:** Cat, Frog, Bird. (h)**Human:** Cat(0.65), Dog(0.21), Deer(0.14). **Machine:** Dog, Deer, Cat

#### 5.4 Example: Weak-Strong Flow

To see how the above finds can be applied, we establish a weak-strong flow where we consturct an ensemble with a weak yet fast architecture and another ensemble with a strong yet slow architecture. For each coming image, the weak ensemble labels  $+$ ,  $-$ , or  $0$  for each class using the 3rd way in section 3.2 with two thresholds, where  $+$  represents being confident of belonging to this class,  $-$  for confident of not belonging to this class, or  $0$  for "I don't know". If there is only one class with the label  $+$ , then this class will be the final result. Otherwise we call the strong the strong ensemble and make a prediction with the 1st way in section 3.2.

The accuracy of this weak-strong flow is 0000. To compare, the accuracy for purly running an easy ensemble with the 1st way in section 3.2 is 00000. The accuracy of running the strong ensemble only is 0000000

Table 2 and 3 show the running time of the weak-strong flow and only running the strong ensemble.



Table 3: **Running Time of the Strong Ensemble Only:** The strong model is DPN92, the ensemble size is 10. Batch Size indicates how many images are processed parallelly. The unit for all values is ms

| Batch Size | Strong Only Run | Strong Only T-Test | Strong Only Total |
|------------|-----------------|--------------------|-------------------|
| 16         | 5.93767         | 0.03345            | 5.97112           |
| 64         | 4.93924         | 0.00943            | 4.94867           |
| 256        | 5.70755         | 0.00375            | 5.71130           |
| 1024       | 7.01784         | 0.00224            | 7.02008           |

## References

- [1] Leo Breiman. Fitting additive models to regression data: Diagnostics and alternative views. *Computational Statistics & Data Analysis*, 15(1):13–46, 1993.
- [2] Leo Breiman. Bagging predictors. *Machine learning*, 24:123–140, 1996.
- [3] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural safety*, 31(2):105–112, 2009.
- [4] Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [5] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Advances in neural information processing systems*, 30, 2017.
- [6] Yonatan Geifman and Ran El-Yaniv. Selectivenet: A deep neural network with an integrated reject option. In *International conference on machine learning*, pages 2151–2159. PMLR, 2019.
- [7] Yonatan Geifman, Guy Uziel, and Ran El-Yaniv. Bias-reduced uncertainty estimation for deep neural classifiers. *arXiv preprint arXiv:1805.08206*, 2018.
- [8] Stephen C Hora. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety*, 54(2-3):217–223, 1996.
- [9] Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: A tutorial introduction. *arXiv preprint arXiv:1910.09457*, 5, 2019.
- [10] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [11] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [12] Lihong Li, Michael L Littman, and Thomas J Walsh. Knows what it knows: a framework for self-aware learning. In *Proceedings of the 25th international conference on Machine learning*, pages 568–575, 2008.
- [13] Chen Liu, Xiaomeng Dong, Michael Potter, Hsi-Ming Chang, and Ravi Soni. Adversarial focal loss: Asking your discriminator for hard examples. *arXiv preprint arXiv:2207.07739*, 2022.
- [14] Weiwei Liu, Ivor W Tsang, and Klaus-Robert Müller. An easy-to-hard learning paradigm for multiple classes and multiple labels. *The Journal of Machine Learning Research*, 18, 2017.
- [15] Joshua C Peterson, Ruairidh M Battleday, Thomas L Griffiths, and Olga Russakovsky. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9617–9626, 2019.
- [16] Cheng Yan, Guansong Pang, Xiao Bai, Chunhua Shen, Jun Zhou, and Edwin Hancock. Deep hashing by discriminating hard examples. In *Proceedings of the 27th ACM international conference on multimedia*, pages 1535–1542, 2019.