

453 A Missing Details

454 A.1 Motivations for working with model latent space

455 In Section 3, we introduced the confusion density matrix that allows us to categorize suspicious
 456 examples at testing time. Crucially, this density matrix relies on kernel density estimations in the
 457 latent space \mathcal{H} associated to the model f through Assumption 1. Why are we performing a kernel
 458 density estimation in latent space rather than in input space \mathcal{X} ? The answer is fairly straightforward:
 459 we want our density estimation to be coupled to the model and its predictions.

460 Let us now make this point more rigorous. Consider two input examples $x_1, x_2 \in \mathcal{X}$. The model
 461 assigns a representations $g(x_1), g(x_2) \in \mathcal{H}$ and class probabilities $f(x_1), f(x_2) \in \mathcal{Y}$. If we define
 462 our kernel κ in latent space \mathcal{H} , this often means¹ that $\kappa[g(x_1), g(x_2)]$ grows as $\|g(x_1) - g(x_2)\|_{\mathcal{H}}$
 463 decreases. Hence, examples that are assigned a similar latent representation by the model f are
 464 related by the kernel. Since our whole discussion revolves around model *predictions*, we would
 465 like to guarantee that two examples related by the kernel are given similar predictions by the model
 466 f . In this way, we would be able to interpret a large kernel density $\kappa[g(x_1), g(x_2)]$ as a hint that
 467 the predictions $f(x_1)$ and $f(x_2)$ are similar. We will now show that, under Assumption 1, such a
 468 guarantee exists. Similar to [36], we start by noting that

$$\begin{aligned} \|(l \circ g)(x_1) - (l \circ g)(x_2)\|_{\mathbb{R}^C} &= \|l[g(x_1) - g(x_2)]\|_{\mathbb{R}^C} \\ &\leq \|l\|_{\text{op}} \|g(x_1) - g(x_2)\|_{\mathcal{H}}, \end{aligned}$$

469 where $\|\cdot\|_{\mathbb{R}^C}$ is a norm on \mathbb{R}^C and $\|l\|_{\text{op}}$ is the operator norm of the linear map l . In order to extend
 470 this inequality to black-box predictions, we note that the normalizing map in Assumption 1 is often
 471 a Lipschitz function with Lipschitz constant $\lambda \in \mathbb{R}$. For instance, a Softmax function with inverse
 472 temperature constant λ^{-1} is λ -Lipschitz [37]. We use this fact to extend our inequality to predicted
 473 class probabilities:

$$\begin{aligned} \|f(x_1) - f(x_2)\|_{\mathcal{Y}} &= \|(\varphi \circ l \circ g)(x_1) - (\varphi \circ l \circ g)(x_2)\|_{\mathcal{Y}} \\ &\leq \lambda \|(l \circ g)(x_1) - (l \circ g)(x_2)\|_{\mathbb{R}^C} \\ &\leq \lambda \|l\|_{\text{op}} \|g(x_1) - g(x_2)\|_{\mathcal{H}}. \end{aligned}$$

474 This crucial inequality guarantees that examples $x_1, x_2 \in \mathcal{X}$ that are given a similar latent representa-
 475 tion $g(x_1) \approx g(x_2)$ will also be given a similar prediction $f(x_1) \approx f(x_2)$. In short: two examples
 476 that are related according to a kernel density defined in the model latent space \mathcal{H} are guaranteed to
 477 have similar predictions. This is the motivation we wanted to support the definition of the kernel κ in
 478 latent space.

479 An interesting question remains: is it possible to have similar guarantees if we define the kernel in
 480 input space? When we deal with deep models, the existence of adversarial examples indicates the
 481 opposite [38]. Indeed, if x_2 is an adversarial example with respect to x_1 , we have $x_1 \approx x_2$ (and hence
 482 $\|x_1 - x_2\|_{\mathcal{X}}$ small) with two predictions $f(x_1)$ and $f(x_2)$ that are significantly different. Therefore,
 483 defining the kernel κ in input space might result in relating examples that are given a significantly
 484 different prediction by the model. For this reason, we believe that the latent space is more appropriate
 485 in our setting.

486 A.2 Details: Flagging IDM and Bnd Examples with Thresholds

In order to understand uncertainty, it will be clearer to map those scores into binary classes with
 thresholds. In our experiments, we use empirical quantiles as thresholds. e.g., to label an example as
 IDM, we specify an empirical quantile number q , and calculate the corresponding threshold based on
 the order statistics of IDM Scores for test examples: $S_{\text{IDM}}^{(1)}, \dots, S_{\text{IDM}}^{(|\mathcal{D}_{\text{test}}|)}$, where $S_{\text{IDM}}^{(n)}$ denotes the
 n -th smallest IDM score out of $|\mathcal{D}_{\text{test}}|$ testing-time examples. Then, the threshold given quantile
 number q is

$$\tau_{\text{IDM}}(q) \equiv S_{\text{IDM}}^{(\lfloor |\mathcal{D}_{\text{test}}| \cdot q \rfloor)}.$$

¹This is the case for all the kernels that rely on a distance (e.g. the Radial Basis Function Kernel, the Matern kernel or even Polynomial kernels [32]).

Similarly, we can define quantile-based threshold in flagging Bnd examples based on the order statistics of Bnd Scores for test examples, such that for given quantile q ,

$$\tau_{\text{Bnd}}(q) \equiv S_{\text{Bnd}}^{(\lfloor |\mathcal{D}_{\text{test}}| \cdot q \rfloor)}.$$

487 Practically, a natural choice of q is to use the validation accuracy: when there are $1 - q$ examples
 488 misclassified in the validation set, we also expect the testing-time in distribution examples with the
 489 highest $1 - q$ to be marked as Bnd or IDM examples.

490 B Improving Predicting Performance of Uncertain Examples

491 Knowing the category that a suspicious example belongs to, can we improve its prediction? For ease
 492 of exposition, we focus on improving predictions for $\mathcal{S}_{\text{B\&I}}$.

493 Let $p(x | \mathcal{S}_{\text{B\&I}})$ be the latent density be defined as in Definition 1. We can improve the prediction
 494 performance of the model on $\mathcal{S}_{\text{B\&I}}$ examples by focusing on the part of examples in the training
 495 set that are closely related to those suspicious examples. We propose to refine the training dataset
 496 $\mathcal{D}_{\text{train}}$ by only keeping the examples that resembles the latent representations for the specific type of
 497 test-time suspicious examples, and train another model on this subset of the training data:

$$\tilde{\mathcal{D}}_{\text{train}} \equiv \{x \in \mathcal{D}_{\text{train}} | p(x | \mathcal{S}_{\text{B\&I}}) \geq \tau_{\text{test}}\}, \quad (6)$$

498 where τ_{test} is a threshold that can be adjusted to keep a prespecified proportion q of the related
 499 training data. Subsequently, new prediction model $f_{\text{B\&I}}$ is trained on $\tilde{\mathcal{D}}_{\text{train}}$.

500 Orthogonal to ensemble methods that require multiple models trained independently, and improve
 501 overall prediction accuracy by bagging or boosting, our method is targeted at improving the model’s
 502 performance on a *specified* subclass of test examples by finding the most relevant training examples.
 503 Our method is therefore more transparent and can be used in parallel with ensemble methods if
 504 needed.

Threshold $\tau_{\text{test}}(q)$ For every training example $x \in \mathcal{D}_{\text{train}}$, we have the latent density $p(x | \mathcal{D}_{\text{B\&I}})$
 over the B&I class of the test set. With their order statistics $p_{(1)}(x | \mathcal{D}_{\text{B\&I}}), \dots, p_{(|\mathcal{D}_{\text{train}}|)}(x | \mathcal{D}_{\text{B\&I}})$.
 Given quantile number q , our empirical quantile based threshold τ_{test} is chosen as

$$\tau_{\text{test}}(q) \equiv p_{(\lfloor q \cdot |\mathcal{D}_{\text{train}}| \rfloor)}(x | \mathcal{D}_{\text{B\&I}}).$$

505 During the inverse training time, we train our model to predict those B&I class of test examples only
 506 with the training data with higher density than $\tau_{\text{test}}(q)$. We experiment with different choices of q in
 507 the experiment (Figure 6 in Sec. 4.4).

508 C Additional Experiments

509 C.1 Categorization of Uncertainty under Different Thresholds

510 In the main text, we provide results with $\tau_{\text{Bnd}} = \tau_{\text{IDM}} = 0.8$, which approximates the accuracy on
 511 validation set—as a natural choice. In this section, we vary these thresholds and show in Figure 7
 512 that changing those thresholds does not significantly alter the conclusions drawn above.

513 Figure 8 looks more closely into the top 25% uncertain examples for each method, and the accuracy
 514 on each of the uncertainty classes. As expected, the accuracy of the B&I examples is always lower
 515 than that of the trusted class, meaning that those examples are most challenging for the classifier.
 516 And the accuracy of flagged classes are always lower than the *other* class, verifying the proposed
 517 categorization of different classes.

518 C.2 Inverse Direction: More Results

519 In the main text, we show the results on improving prediction performance on the B&I class with
 520 training example filtering (On the Covtype, Digits dataset). More results on other classes of examples
 521 are provided in this section.

522 We experiment on three UCI datasets: **Covtype**, **Digits**, and **Spam**. And experiment with three
 523 classes we defined in this work:

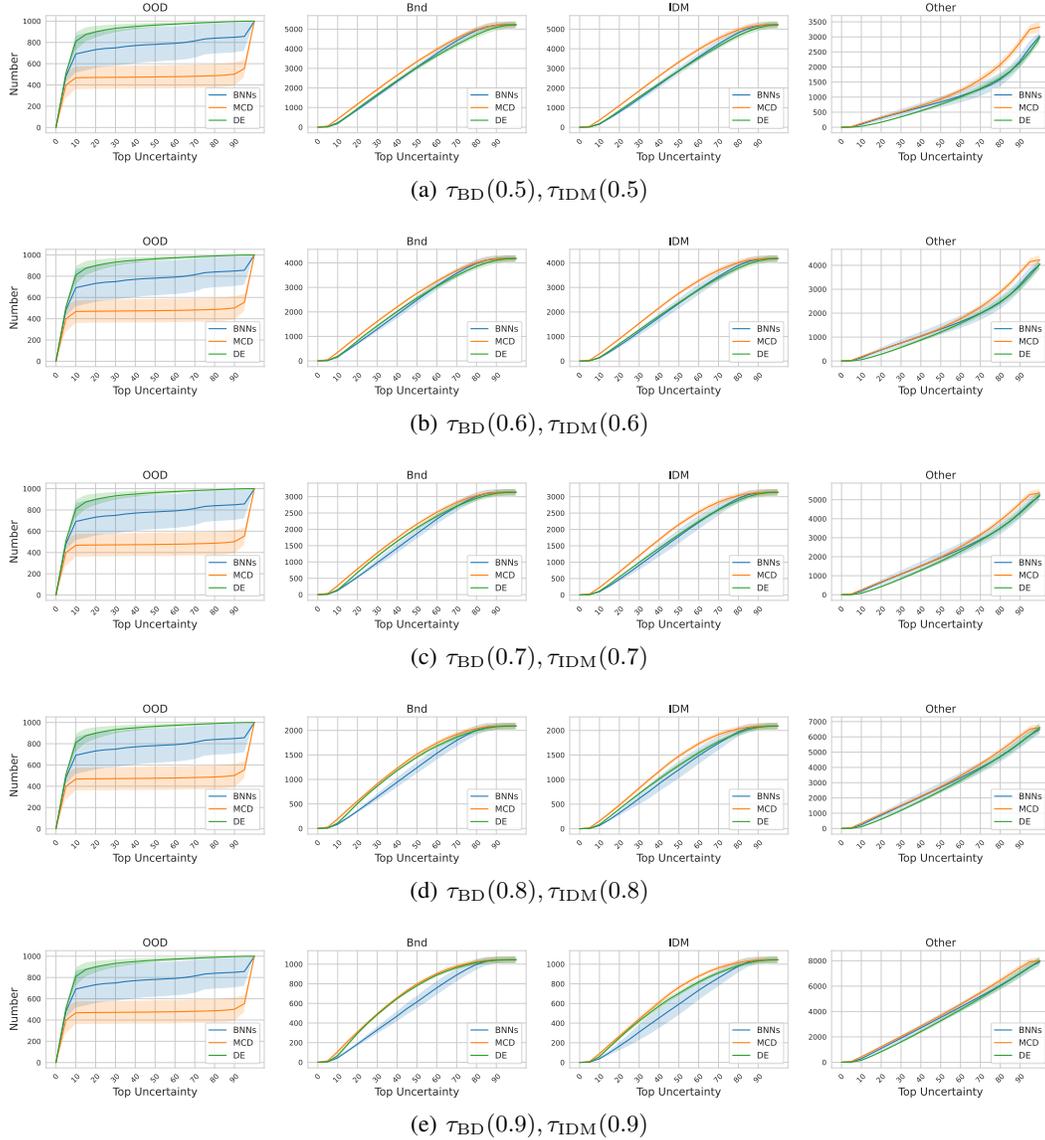


Figure 7: Experiments on different choices of thresholds.

- 524 1. **B&I** class (Figure 9). As we have discussed in our main text, the prediction accuracy on the
525 B&I class are always the lowest among all classes. By training with filtered examples in
526 $\mathcal{D}_{\text{train}}$ rather than the entire training set, the B&I class of examples can be classified with a
527 remarkably improved accuracy.
- 528 2. **Bnd** class (Figure 10). This class of examples are located at boundaries in the latent space
529 of validation set, but not necessarily have been misclassified. Therefore, their performance
530 baseline (training with the entire $\mathcal{D}_{\text{train}}$) is relatively high. The improvement is clear but not
531 as much as on the other two classes.
- 532 3. **IDM** class (Figure 11). For this class of examples, similar mistakes have been make in
533 the validation set, yet those examples are not necessarily located in the boundaries—the
534 misclassification may be caused by ambiguity in decision boundary, imperfectness of either
535 the model or the dataset. The primal prediction accuracy on this class of examples is lower
536 than the Bnd class but higher than the B&I class, training with filtered $\mathcal{D}_{\text{train}}$ also clearly
537 improve the performance on this class of examples.

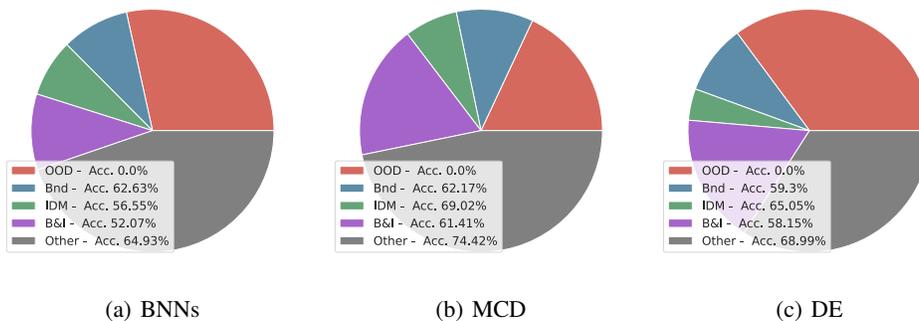


Figure 8: The top 25% uncertain examples identified by different methods. Legend of each figure provide the accuracy and proportion of each class. As the classifier can not make correct predictions on the OOD examples, it’s always better for uncertainty estimators to flag more OOD examples.

Table 4: DAUC is not the only choice in identifying OOD examples. On the Dirty-MNIST dataset, DAUC, Outlier-AE and the IForest can identify most outliers in the test dataset. (Given threshold = 1.0 for those two benchmark methods).

Method	Precision	Recall	F1-Score
DAUC	1.0000 ± 0.0000	1.0000 ± 0.0000	1.0000 ± 0.0000
Outlier-AE	1.0000 ± 0.0000	1.0000 ± 0.0000	1.0000 ± 0.0000
IForest [40]	0.9998 ± 0.0004	1.0000 ± 0.0000	0.9999 ± 0.0002

538 C.3 Alternative Approach in Flagging OOD

539 As we have mentioned in the main text, although DAUC has a unified framework in understanding
 540 all three types of uncertainty the uncertain caused by OOD examples can also be identified by
 541 off-the-shelf algorithms. We compare DAUC to two existing outlier detection methods in Table 4,
 542 where all methods achieve good performance on the Dirty-MNIST dataset. Our implementation is
 543 based on Alibi Detect [39].

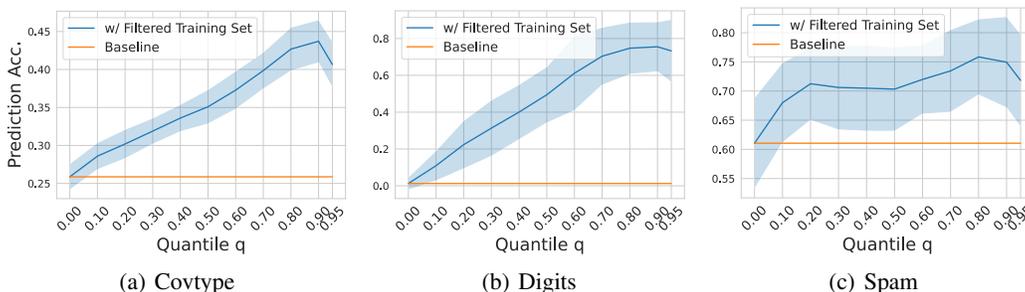


Figure 9: Experiments on the B&I class (reported in the main text).

544 C.4 Experiments on Dirty-CIFAR-10

545 **Dataset Description** In this experiment, we introduce a revised version of the CIFAR-10 dataset
 546 to test DAUC’s scalability. Similar to the Dirty-MNIST dataset [34], we use linear combinations of
 547 the latent representation to construct the “boundary” class. In the original CIFAR-10 Dataset, each
 548 of the 10 classes of objects has 6000 training examples. We split the training set into training set
 549 (40%), validation set (40%) and test set (20%). To verify the performance of DAUC in detecting
 550 OOD examples, we randomly remove one of those 10 classes (denoted with class- i) during training
 551 and manually concatenate OOD examples with the test dataset, with label i . In our experiment, we

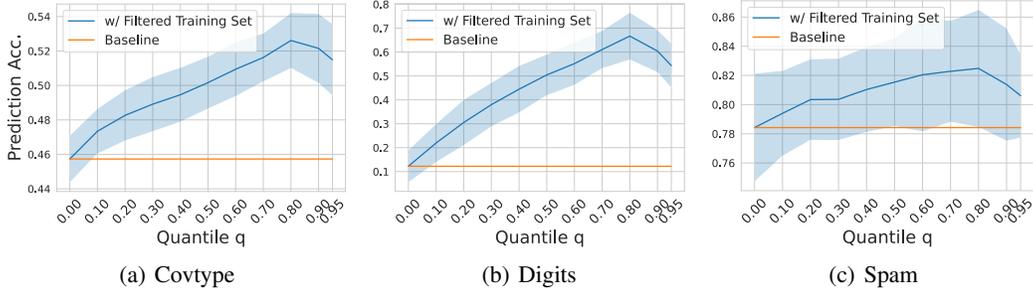


Figure 10: Experiments on the Bnd class.

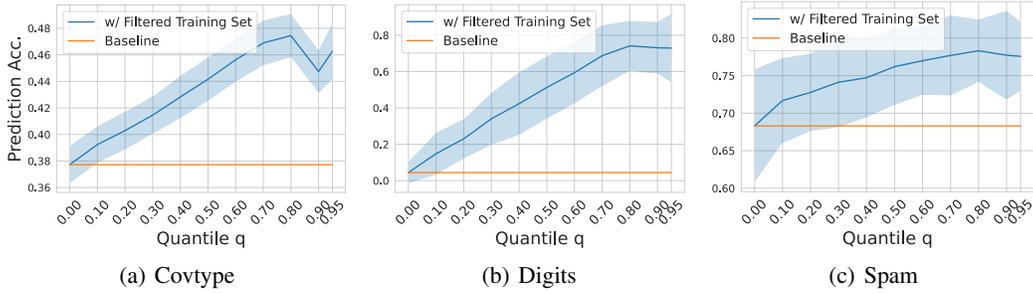


Figure 11: Experiments on the IDM class.

552 use 1000 MNIST digits as the OOD examples, with zero-padding to make those digits share the same
 553 input shape as the CIFAR-10 images. Combining those boundary examples, OOD examples and the
 554 vanilla CIFAR-10 examples, we get a new benchmark, dubbed as Dirty-CIFAR-10, for quantitative
 555 evaluation of DAUC.

556 **Quantify the performance of DAUC on Dirty-CIFAR-10** Quantitatively, we evaluate the per-
 557 formance of DAUC in categorizing all three classes of uncertain examples. Results of averaged
 performance and standard deviations based on 8 repeated runs are provided in Table 5.

Table 5: Quantitative results on the Dirty-CIFAR-10 dataset. DAUC scales well and is able to categorize all three classes of uncertain examples.

Category	Precision	Recall	F1-Score
OOD	0.986 ± 0.003	0.959 ± 0.052	0.972 ± 0.027
Bnd	0.813 ± 0.002	0.975 ± 0.000	0.887 ± 0.001
IDM	0.688 ± 0.041	0.724 ± 0.017	0.705 ± 0.027

558

559 **Categorize Uncertain Predictions on Dirty-CIFAR-10** Similar to Sec. 4.3 and Figure 5, we can
 560 categorize uncertain examples flagged by BNNs, MCD and DE using DAUC—see Figure 12. We find
 561 that in the experiment with CIFAR-10, DE tends to discover more OOD examples as top uncertain
 562 examples. Differently, although BNNs flags less OOD examples as top-uncertain, it continuously
 563 discover those OOD examples and is able to find most of them for the top 50% uncertainty. On the
 564 contrary, MCD performs the worst among all three methods, similar to the result drawn from the
 565 DMNIST experiment. On the other hand, while BNN is good at identifying OOD examples, it flags
 566 less uncertain examples in the Bnd and IDM classes. DE is the most apt at flagging both Bnd and
 567 IDM examples, and categorizes far less examples into the *Other* class. **These observations are well
 568 aligned with the experiment results we had with DMNIST in Sec. 4.3, showing the scalability of
 569 DAUC to large-scale image dataset.**

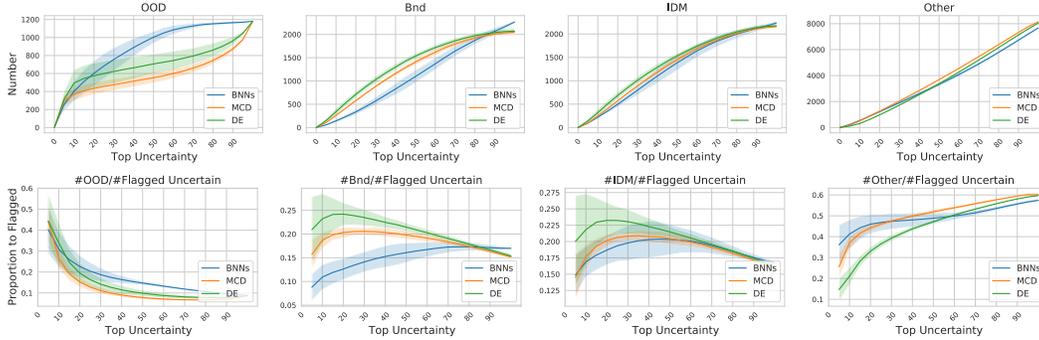


Figure 12: Experiments on the CIFAR-10 dataset. Results of applying DAUC in categorizing different uncertainty estimation methods. First row: comparisons on the numbers in different classes of examples. Second row: comparisons on the proportion of different classes of flagged examples to the total number of identified uncertain examples. Different methods tend to identify different certain type of uncertain examples.

570 D Implementation Details

571 D.1 Code

572 Our code is anonymously available at <https://anonymous.4open.science/r/DAUC-B234/>

573 D.2 Hyperparameters

574 D.2.1 Bandwidth

575 In our experiments, we use (z-score) normalized latent representations and bandwidth 1.0. In
 576 the inverse direction, as the sample sizes are much smaller, a bandwidth of 0.01 is used as the
 577 recommended setting. There is a vast body of research on selecting a good bandwidth for Kernel
 578 Density Estimation models [41-43] and using these to adjust DAUC’s bandwidth to a more informed
 579 choice may further improve performance.

580 D.3 Inverse Direction: Quantile Threshold q

581 As depicted in Appendix A.2, a natural choice of q is to use the validation accuracy. We use this
 582 heuristic approach in our experiments for the inverse direction.

583 D.4 Model Structure

584 In our experiments, we implement **MCD** and **DE** with 3-layer-CNNs with ReLU activation. Our
 585 experiments on **BNNs** are based on the IBM UQ360 software [44]. More details of the convolutional
 586 network structure are provided in Table 6.

Table 6: Network Structure

Layer	Unit	Activation	Pooling
Conv 1	(1, 32, 3, 1, 1)	ReLU()	MaxPool2d(2)
Conv 2	(32, 64, 3, 1, 1)	ReLU()	MaxPool2d(2)
Conv 3	(64, 64, 3, 1, 1)	ReLU()	MaxPool2d(2)
FC	(64 × 3 × 3, 40)	ReLU()	-
Out	(40, N_{Class})	SoftMax()	-

587 D.5 Implementation of Kernel Density Estimation and Repeat Runs

588 Our implementation of KDE models are based on the sklearn’s KDE package [45]. Gaussian kernels
 589 are used as default settings. In all experiments, we run with 10 random seeds and report the averaged

590 results. In our experiments, we find using different kernels in density estimation provides highly
591 correlated scores. We calculate the Spearman’s ρ correlation between scores DAUC gets over 5
592 runs with Gaussian, Tophat, Exponential kernels under the same bandwidth. Changing the kernel
593 brings highly correlated scores (all above **0.86**) for DAUC and, hence, has minor impact on DAUC’s
594 performance. We preferred KDE since the latent representation is relatively low-dimensional. We
595 found that a low-dim latent space (e.g., 10) works well for all experiments (including CIFAR-100).

596 **D.6 Hardware**

597 All results reported in our paper are conducted with a machine with 8 Tesla K80 GPUs and 32
598 Intel(R) E5-2640 CPUs. The computational cost is mainly in density estimation, and for low-dim
599 representation space, such an estimation can be efficient: running time for DAUC on the Dirty-MNIST
600 dataset with KDE is approximately 2 hours.

601 **Assumptions and Limitations**

602 In this work, we introduced the confusion density matrix that allows us to categorize suspicious
603 examples at testing time. Crucially, this density matrix relies on kernel density estimations in the
604 latent space \mathcal{H} associated to the model f through Assumption [\[1\]](#). We note this assumption generally
605 holds for most modern uncertainty estimation methods.

606 While the core contribution of this work is to introduce the concept of confusion density matrix
607 for uncertainty categorization, the density estimators leveraged in the latent space can be further
608 improved. We leave this to the future work.

609 **Broader Impact**

610 While previous works on uncertainty quantification (UQ) focused on the discovery of uncertain
611 examples, in this work, we propose a practical framework for categorizing uncertain examples that
612 are flagged by UQ methods. We demonstrated that such a categorization can be used for UQ method
613 selection — different UQ methods are good at figuring out different uncertainty sources. Moreover,
614 we show that for the inverse direction, uncertainty categorization can improve model performance.

615 With our proposed framework, many real-world application scenarios can be potentially benefited.
616 e.g., in Healthcare, a patient marked as uncertain that categorized as OOD — preferably identified by
617 Deep Ensemble, as we have shown — should be treated carefully when applying regular medical
618 experience; and an uncertain case marked as IDM — preferably identified by MCD — can be
619 carefully compared with previous failure cases for a tailored and individualized medication.