
DeepGD3: Unknown-Aware Deep Generative/Discriminative Hybrid Defect Detector for PCB Soldering Inspection (Supplementary Material)

Ching-Wen Ma¹

Yanwei Liu¹

¹College of Artificial Intelligence, National Yang Ming Chiao Tung University, Tainan, Taiwan

A VISUAL EXPLANATION OF THE PREDICTION CONVERTER

τ_0 controls the variance of Gaussians. If it is small, the Gaussian looks shaper. If it is large, the Gaussian looks wider. Figure. 1 explain this concept. Bayesian optimization then optimizes the shape of Gaussians, making them proper for thresholding.

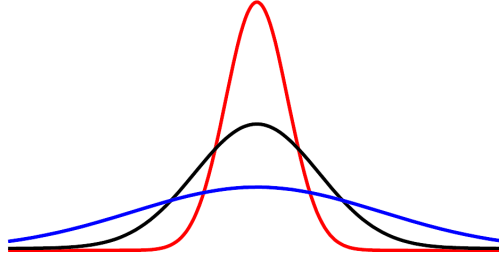


Figure 1: The effect of τ_0 on the shape of probability density function. Black curve: the baseline shape of the probability density function. Red curve: smaller τ_0 makes the shape of the probability density function concentrated. Blue curve: larger τ_0 makes the shape of the probability density function wider.

τ_1 and τ_2 adjust the thresholds $h_{j,n}^1$ and $h_{j,n}^2$ by adjusting the variance of the corresponding Gaussians. A larger τ_1 makes $h_{j,n}^1$ smaller, resulting in more samples being classified as good samples. A larger τ_2 makes $h_{j,n}^2$ smaller, resulting in fewer

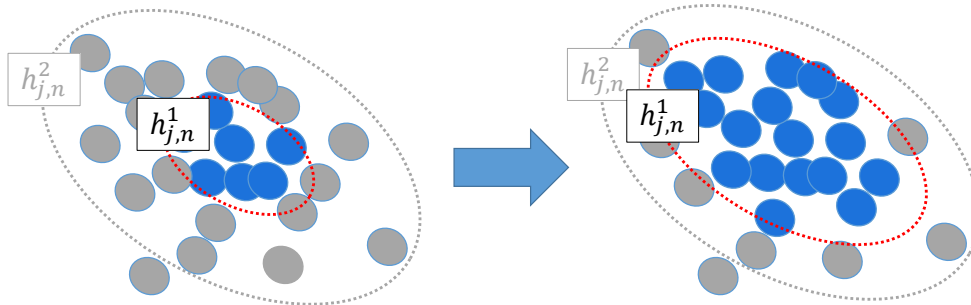


Figure 2: The effect of τ_1 and $h_{j,n}^1$. The larger the τ_1 is, or similarly, the smaller the $h_{j,n}^1$ is, the more test samples will be classified as good samples.

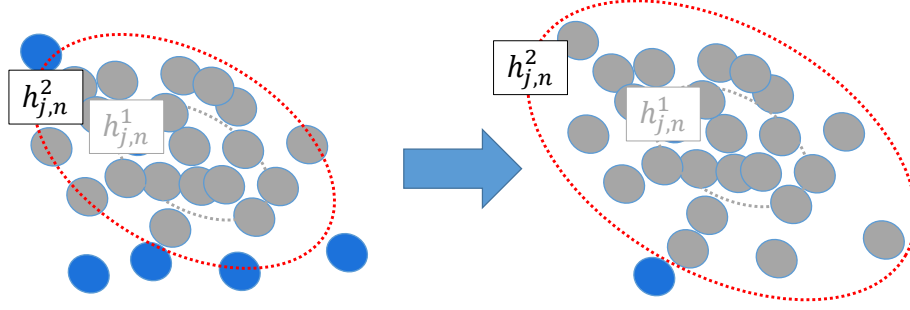


Figure 3: The effect of τ_2 and $h_{j,n}^2$. The larger the τ_2 is, or similarly, the smaller the $h_{j,n}^2$ is, the fewer test samples will be classified as unknown samples.

samples being classified as unknown samples. Figure 2 and Figure 3 explain this concept. Bayesian optimization then optimizes τ_1 and τ_2 for maximizing harmonic score H in main paper’s Equation 5.

B ADDITIONAL EXPERIMENTAL DETAILS

We fine-tuned the MobileNetV3 large model pre-trained on ImageNet [Deng et al., 2009] using PyTorch and trained the defect and component classifiers.

The default hyperparameters were as follows: The number of training epochs was 30. The initial learning rate was 0.05. We used linear warm-up for the first 5 epochs and decay the learning rate with the cosine annealing scheduler; For the gradient descent optimization algorithm, stochastic gradient descent (SGD) was used; the momentum and weight decay were set at 0.9 and 0.0001. The input image was 224×224 color images.

In addition, the early stop technique was applied to prevent overfitting. We applied random horizontal flip, random vertical flip, color jitter, and random grayscale to the data loader of the training set. We did not apply any image augmentation to the validation and test set. A single NVIDIA Tesla V100 GPU was used for all experiments.

C ADDITIONAL SIMULATION RESULTS FOR ABLATION STUDY

We investigated the effects of (a) not including the good new component samples in the training data and (b) not using the combiner.

Comparisons between not including and including good new component samples in the training data. Table 1 shows the effects of including and not including good new component samples in the training data for all test samples, i.e. both old and new component test samples. Table 2 and Table 3 show the performances for the old and new component test samples, respectively. The $(\cdot)^-$ symbol indicates training without including good new component samples.

Observing all the three tables, we notice that Expert 1 shows a lower overkill rate with including good new component samples in the training data, suggesting that the good new component samples help improve performance. Expert 2’s leakage rates are greatly reduced when good new component samples are available during training. For the Hybrid Expert, if there are no good new component samples during the training period, the shared fully connected network $f_{\theta 2}$ may not be able to form proper clusters for these new components. As a result, the GMM model will not be able to correctly classify the new component samples, resulting in a higher unknown rate. However, both the overkill rates and the leakage rates still remains low. This implies that even if good new component samples are not available during training, the Hybrid Expert remains trustworthy.

Comparisons between not using and using the prediction combiner Table 4 shows the performance of \hat{y}_{def}^1 , \hat{y}_{def}^2 , and \hat{y}_{def} . The results show that \hat{y}_{def}^1 produces a high leakage rate with a high standard deviation, and \hat{y}_{def}^2 achieves $0.691\% \pm 0.974\%$ in terms of leakage rate, but also produces a high overkill rate. The prediction combiner combines \hat{y}_{def}^1 and \hat{y}_{def}^2 , achieving an overkill rate of $0.108\% \pm 0.033\%$ and a leakage rate of $0.063\% \pm 0.075\%$, with an unknown rate of $3.706\% \pm 2.270\%$.

Table 1: Comparisons between not including $(\cdot)^-$ and including (\cdot) good new component samples as additional training data. The test performances for all test samples, i.e. both old and new component test samples.

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
Expert 1 ⁻	1.773 \pm 1.459	0.327 \pm 0.463	-
Expert 1	0.015 \pm 0.008	1.827 \pm 3.063	-
Expert 2 ⁻	3.650 \pm 2.902	11.419 \pm 19.321	0.000 \pm 0.000
Expert 2	1.954 \pm 0.724	1.942 \pm 1.337	0.000 \pm 0.000
Hybrid Expert ⁻	0.189 \pm 0.266	0.074 \pm 0.114	26.089 \pm 40.843
Hybrid Expert	0.108 \pm 0.033	0.063 \pm 0.075	3.706 \pm 2.270

Table 2: Comparisons between not including $(\cdot)^-$ and including (\cdot) good new component samples as additional training data. The test performances for old component test samples.

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
Expert 1 ⁻	0.097 \pm 0.146	0.028 \pm 0.018	-
Expert 1	0.017 \pm 0.007	0.021 \pm 0.011	-
Expert 2 ⁻	3.293 \pm 5.306	9.268 \pm 15.373	0.000 \pm 0.000
Expert 2	1.282 \pm 0.192	2.257 \pm 1.495	0.000 \pm 0.000
Hybrid Expert ⁻	0.007 \pm 0.010	0.042 \pm 0.055	18.695 \pm 31.605
Hybrid Expert	0.129 \pm 0.110	0.019 \pm 0.013	3.529 \pm 3.002

Table 5 and Table 6 show the performance of \hat{y}_{def}^1 , \hat{y}_{def}^2 , and \hat{y}_{def} for the old component test set and new component test set, respectively. For the old component settings, \hat{y}_{def} yields the best overkill and leakage rate performance. For the new component setting, \hat{y}_{def} achieves 0.063% \pm 0.075% in terms of leakage rate, indicating excellent performance for detecting new component defects.

D ALGORITHMS

Training procedure: The detailed training procedure of Figure 3 is described in Algorithm 1.

Inference procedure: The Hybrid Expert inference procedure is described in Algorithm 2.

References

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

Table 3: Comparisons between not including $(\cdot)^-$ and including (\cdot) good new component samples as additional training data. The test performances for new component test samples.

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
Expert 1 ⁻	4.591 \pm 3.687	0.611 \pm 0.824	-
Expert 1	0.010 \pm 0.010	3.380 \pm 5.540	-
Expert 2 ⁻	3.368 \pm 3.112	17.666 \pm 30.599	0.000 \pm 0.000
Expert 2	3.739 \pm 2.459	0.989 \pm 1.713	0.000 \pm 0.000
Hybrid Expert ⁻	0.713 \pm 1.033	0.096 \pm 0.167	35.055 \pm 46.395
Hybrid Expert	0.126 \pm 0.062	0.090 \pm 0.156	3.324 \pm 1.553

Table 4: Performance of \hat{y}_{def}^1 , \hat{y}_{def}^2 , and \hat{y}_{def} for all test samples, i.e. both old and new component test samples

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
\hat{y}_{def}^1	0.390 \pm 0.320	0.681 \pm 1.114	-
\hat{y}_{def}^2	2.287 \pm 1.424	0.691 \pm 0.974	0.000 \pm 0.000
\hat{y}_{def}	0.108 \pm 0.033	0.063 \pm 0.075	3.706 \pm 2.270

Table 5: Performance of \hat{y}_{def}^1 , \hat{y}_{def}^2 , and \hat{y}_{def} for old component test samples.

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
\hat{y}_{def}^1	0.671 \pm 0.794	0.034 \pm 0.032	-
\hat{y}_{def}^2	2.227 \pm 1.768	0.891 \pm 1.343	0.000 \pm 0.000
\hat{y}_{def}	0.129 \pm 0.110	0.019 \pm 0.013	3.529 \pm 3.002

Table 6: Performance of \hat{y}_{def}^1 , \hat{y}_{def}^2 , and \hat{y}_{def} for new component test samples.

	Overkill rate (%)	Leakage rate (%)	Unknown rate (%)
\hat{y}_{def}^1	0.150 \pm 0.073	1.202 \pm 2.064	-
\hat{y}_{def}^2	2.315 \pm 0.592	0.090 \pm 0.156	0.000 \pm 0.000
\hat{y}_{def}	0.126 \pm 0.062	0.090 \pm 0.156	3.324 \pm 1.553

Algorithm 1 Learning algorithm of the proposed method

```
1: Input: A mini batch of defect training samples  $(x_{def}, y_{def})$ , and component training samples  $(x_{com}, y_{com}) \in D_{train}$ ,  
   A mini batch of validation samples  $(x_{val}, y'_{def}, y'_{com}) \in D_{val}$ , Max number of training epochs  $E$   
2: Output: Feature extractor  $f_{\theta 1}$ , Encoder network  $f_{\theta 2}$ , Defect classifier  $\psi$ , Projection network  $\phi$ , Gaussian mixture model  
    $\gamma_j$  where  $j \in \{1, 2, \dots, 23\}$ .  
3: Initialize  $f_{\theta 1}$ ,  $f_{\theta 2}$ ,  $\psi$ ,  $\phi$ , model selection loss  $\ell_{\omega} = \infty$   
4: for  $e \leftarrow 1$  to  $E$  do  
5:    $\ell_{def} \leftarrow \text{CrossEntropy}(\psi(f_{\theta 2}(f_{\theta 1}(x_{def}))), y_{def})$  // Train the upper branch, i.e.  $f_{\theta 1}$ ,  $f_{\theta 2}$ , and Classifier  $\psi$ .  
6:   Update  $\psi$ ,  $f_{\theta 2}$ ,  $f_{\theta 1}$  using  $\ell_{def}$   
7:    $\ell_{com} \leftarrow \text{MultiSimilarity}(\phi(f_{\theta 2}(f_{\theta 1}(x_{com}))), y_{com})$  // Train the lower branch, i.e.  $f_{\theta 1}$ ,  $f_{\theta 2}$ , and Projector  $\phi$ .  
8:   Update  $\phi$ ,  $f_{\theta 2}$ ,  $f_{\theta 1}$  using  $\ell_{com}$   
9:    $\ell_{def'} \leftarrow \text{CrossEntropy}(\psi(f_{\theta 2}(f_{\theta 1}(x_{val}))), y'_{def})$  // Evaluate the upper branch  
10:   $\ell_{com'} \leftarrow \text{MultiSimilarity}(\phi(f_{\theta 2}(f_{\theta 1}(x_{val}))), y'_{com})$  // Evaluate the lower branch  
11:   $\ell_{eval} \leftarrow \ell_{def'} + \ell_{com'}$  // Procedure for model selection  
12:  if  $\ell_{\omega} > \ell_{eval}$  then  
13:     $\ell_{\omega} \leftarrow \ell_{eval}$   
14:    saveModel( $f_{\theta 1}$ ,  $f_{\theta 2}$ ,  $\psi$ ,  $\phi$ )  
15:  end if;  
16: end for;  
17: for  $j \leftarrow 1$  to 23 do  
18:    $\gamma_j \leftarrow \gamma_j(f_{\theta 2}(f_{\theta 1}(x_{com,j})), y_{com})$  // Gaussian mixture model fitting on  $f_{\theta 2}(f_{\theta 1}(x_{com,j}))$  guided by  $y_{com}$ .  
19: end for
```

Algorithm 2 Inference Procedure of the Proposed Method

```
1: Input: A mini batch of test samples  $x_{test} \in D_{test}$   
2: Output: Defect prediction  $\hat{y}_{def}$   
3: Require Feature extractor  $f_{\theta 1}$ , Shared fully connected network  $f_{\theta 2}$ , Defect classifier  $\psi$ , Gaussian mixture model  $\gamma_j$   
   where  $j \in \{1, 2, \dots, 23\}$ , and prediction converter  $\Lambda$ ;  
4:  $\hat{y}_{def}^1 \leftarrow \psi(f_{\theta 2}(f_{\theta 1}(x_{test})))$  // Defect classification.  
5: for  $j \leftarrow 1$  to 23 do  
6:    $P(\cdot | y_{com} = j) \leftarrow \gamma_j(f_{\theta 2}(f_{\theta 1}(x_{test})))$  // Probabilistic component type prediction.  
7: end for  
8:  $\hat{y}_{def}^2 \leftarrow \Lambda(\hat{y}^{j=1,2,\dots,23})$  // Prediction conversion.  
9: if  $\hat{y}_{def}^1 == \hat{y}_{def}^2$  then  
10:   $\hat{y}_{def} \leftarrow \hat{y}_{def}^1, \hat{y}_{def}^2$  // If the predictions match,  $\hat{y}_{def}$  will be the same as  $\hat{y}_{def}^1$  and  $\hat{y}_{def}^2$ .  
11: else  
12:   $\hat{y}_{def} \leftarrow \text{unknown}$  // If the predictions do not match,  $\hat{y}_{def}$  is unknown.  
13: end if  
14: return  $\hat{y}_{def}$ 
```
