

---

# MIXED-TYPE WAFER CLASSIFICATION FOR LOW MEMORY DEVICES USING KNOWLEDGE DISTILLATION

---

**Nitish Shukla**  
*Independent researcher*  
nitishshukla86@gmail.com

**Anurima Dey**  
*Independent researcher*  
anurima10@gmail.com

**Srivatsan K**  
*Independent researcher*  
srivatsanraman3@gmail.com

## ABSTRACT

Manufacturing wafers is an intricate task involving thousands of steps. Defect Pattern Recognition (DPR) of wafer maps is crucial for determining the root cause of production defects, which may further provide insight for yield improvement in wafer foundry. During manufacturing, various defects may appear standalone in the wafer or may appear as different combinations. Identifying multiple defects in a wafer is generally harder compared to identifying a single defect. Recently, deep learning methods have gained significant traction in mixed-type DPR. However, the complexity of defects requires complex and large models making them very difficult to operate on low-memory embedded devices typically used in fabrication labs. Another common issue is the unavailability of labeled data to train complex networks. In this work, we propose an unsupervised training routine to distill the knowledge of complex pre-trained models to lightweight deployment-ready models. We empirically show that this type of training compresses the model without sacrificing accuracy despite being up to 10 times smaller than the teacher model. The compressed model also manages to outperform contemporary state-of-the-art models.

**Keywords** Knowledge Distillation, · Wafer Defect Classification, · Embedded devices

## 1 Introduction

Wafers are silicon bases on which the fabrication of the Integrated Circuits (ICs) takes place. Wafer fabrication and testing are the two most crucial steps of the IC manufacturing process. Wafer fabrication consists of a series of complex steps laid down in the design of the particular type of IC. The layers of fabrication are compassed at the nanoscale by highly sophisticated machines. Defects incorporated in a wafer during this fabrication process, if not amended at an early stage, may cause severe damage to the production line. The root causes of such defects variation can be due to any of the 6Ms[1] namely Man, Method, Machine, Material, Milieu, and Measurement which may have blended during the fabrication process. After wafer testing is done, these defects are observed in the form of failed dies on a wafer, known as a wafer map. Studying these patterns on the wafer map enables the engineers to determine the probable cause of the die failures which in turn allows them to improve the product yield [2]. The study of such defects on wafer maps is commonly known as Defect Pattern Recognition (DPR).

Historically, such defects were manually analyzed using various statistical methods by experienced engineers and subject matter experts. However, with the advent of Deep Neural Networks and CNNs[3], detecting such defects is slowly becoming more efficient. Various supervised and unsupervised models are currently in use for detecting the type of defects on the wafer map. A recent study on semi-supervised pattern recognition[4] on wafer maps, has found a total of 14 defect types [5]. Some of the common ones are “Centre”, “Donut”, “Scratch”, “Edge-ring” etc. However, adding to the problem, these defects do not always stand alone. Due to the complexity of the entire manufacturing process, each wafer passes through a series of chemical processes occurring in layers. The compactness of ICs has increased over time according to Moore’s law [6], which resulted in more layers. It so happens that due to one of the 6M cause and effect [7], a particular defect type may appear in one such layer and another may appear in another layer on a single wafer. When this wafer is manufactured, there can be two or more such defects appearing on a single wafer map. Such patterns are called mixed-type wafer defects and have gained a lot of attention in recent times.

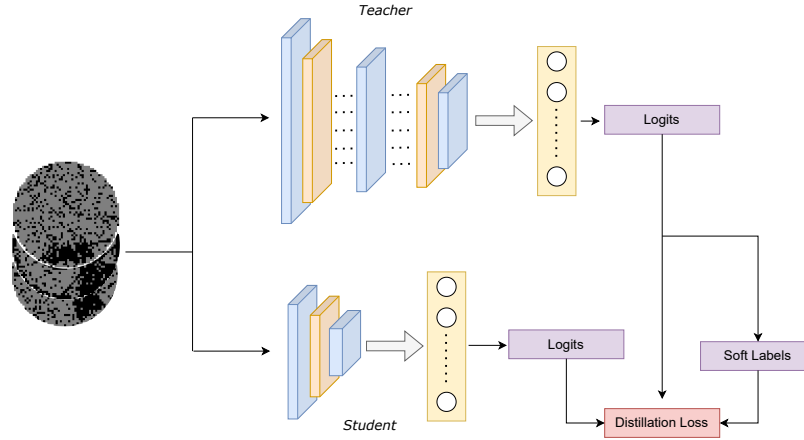


Figure 1: Architecture of the knowledge distillation routine. The distillation loss depends on the logits produced by  $T$ ,  $S$  and the soft labels obtained from the predictions made by  $T$ .

Deep neural networks have been very successful in fields ranging from computer vision [8] [9] [10], reinforcement learning [11] [12] to natural language processing [13] and speech recognition [14]. However, deep models are generally very huge making deployment to low-memory devices with limited computation complexity very challenging. To address this issue, Bucilua et al[15] first proposed model compression to transfer learned information from a large model or set of models to smaller models without a significant drop in accuracy. A semi-supervised version of knowledge transfer was also introduced by Urner et al., 2011[16] allowing the student model to learn without using ground truth labels. The central idea is that the student model mimics the teacher model in order to obtain a competitive or even superior performance. The key problem is how to transfer the knowledge from a large teacher model to a small student model.

## 2 Motivation

Accurately classifying the defect patterns on the wafer is a computer vision problem. In the initial phases of this research, mostly shallow models were used like SVMs, RBFN, and decision tree [17]. However, the efficiency of these methods is limited to very good input features [18] and image dimensionality. In recent times deep learning has received a lot of popularity in this field since they extract feature from raw data automatically. Particularly, CNNs received great admiration when Nakazawa and Kulkarni [19] could identify 22 defects mixed by six basic types (random, ring, edge, scratch, cluster, gross). Kyeong and Kim [20] applied the CNN model to obtain 16 defects combining circle, ring, scratch, and zone. Wang et al. [21] have proposed a Deformable CNN (DC-Net) capable of identifying up to 38 mixed-type defects comprising 8 basic defect types.

While all these models have commendable performance in detecting the patterns, a major disadvantage lies in the fact that these deep neural networks have lots of parameters adding to the inherent complexity. Due to the size of the models, they often become very computationally expensive and time-consuming while dealing with real-life fabrication data sets. This makes it nearly impossible for deploying on smaller machines with limited memory capacity. Even when resources are abundant, a lightweight efficient model does no harm since it would be able to serve more clients at a lower cost. As a result, the serviceability of deep neural network models comes with a lot of constraints. One probable way of mitigating this problem is by using simpler models like decision trees, however as mentioned previously they would require very comprehensive feature engineering for comparable accuracy. Even CNNs with reasonably limited capacity do not promise good performance. Naturally, a high-performing lightweight model is the optimal middle ground considering performance and size. This motivated us to apply Knowledge Distillation(KD) which can help in creating a smaller model, with comparable accuracy to the complicated bigger model. Knowledge Distillation is a way of model compression where the smaller network is trained to approximate the teacher network by trying to replicate its output at every level [22].

In deep learning, Knowledge Distillation is a widely used effective technique that was first defined and generalized by Hilton et al.[2]. Knowledge Distillation methods for knowledge transfer have traditionally used supervised learning and semi-supervised learning. Here, we propose an unsupervised Knowledge Distillation learning routine that competes well against current state-of-the-art methods for identifying mixed-type wafer defects.

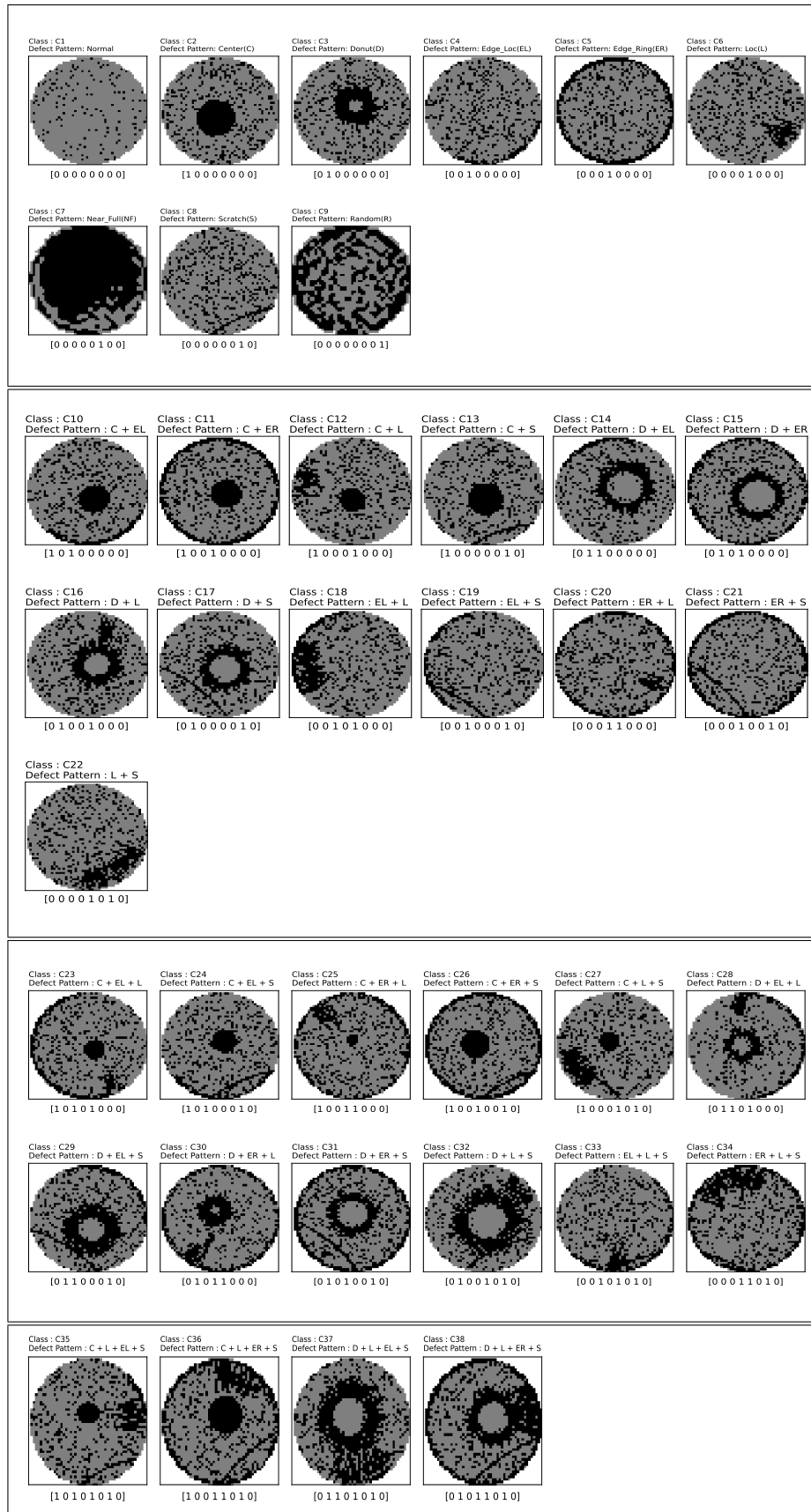


Figure 2: The 38 types of original wafer map defects in the MixedWM38 dataset.

### 3 Knowledge Distillation

Knowledge distillation (KD) is an important technique in the domain of model compression. In knowledge distillation, a small lightweight *student* model is trained to mimic the performance of a compute-intensive high-performing *teacher* model. Surprisingly enough, sometimes the student performs much better than the teacher even though both have the same capacity [23][24][25]. Large-scale deep models have achieved profuse successes in various fields ranging from computer vision [26] to NLP [27], however, the huge computational complexity and massive storage requirements make it a great challenge to deploy them in real-time applications, especially on devices with limited resources, such as video surveillance and autonomous driving cars.

A knowledge distillation system is generally composed of three key components: knowledge, distillation algorithm, and teacher-student architecture. Figure 1 shows a common teacher-student framework for knowledge distillation. Typically, knowledge distillation routines can be classified into the following categories:

**Response-Based Knowledge Distillation:** The main idea in response-based knowledge is to directly mimic the outputs of the teacher. Response-based knowledge distillation is a simple yet effective technique for model compression and it is widely used in different applications. Given a teacher model  $T$  and student model  $S$ , response-based knowledge distillation optimizes the distillation loss between the logits  $z_t = T(x)$  and  $z_s = S(x)$  formulated as

$$\mathcal{L}_{KD} = \mathcal{L}(z_t, z_s) \tag{1}$$

More than often,  $\mathcal{L}_{KD}$  is employed as Kullback-Leibler loss or even simpler mean-square-loss(MSE) loss. Clearly, minimizing  $\mathcal{L}_{KD}$  forces the student logits to match the teacher logits, this essentially results in the student model mimicking the teacher model.

**Feature-Based Knowledge Distillation** Neural Networks are exceptionally well at learning representations of input concepts in increasing levels of abstraction. Feature-based knowledge distillation employs the distillation loss between each feature layer in the teacher and student model respectively. Feature-based distillation is an effective extension of response-based learning, especially for thinner and deeper models. The loss functionally generally takes the form

$$\mathcal{L}_{KD} = \mathcal{L}(\Phi_T(T(x)), \Phi_S(S(x))) \tag{2}$$

where  $T(x)$  and  $S(x)$  are the feature maps of the intermediate layers of the teacher and student model. The transformation  $\Phi_T$  and  $\Phi_S$  are applied when the feature map has different shapes.

**Relation-Based Knowledge Distillation:** Relation-based knowledge distillation further explores the relationships between the different layers in the teacher and student model. In general, the distillation loss of relation-based knowledge based on the relations of feature maps can be formulated as

$$\mathcal{L}_{KD} = \mathcal{L}_R(\Psi_T(\hat{T}(x), \check{T}(x)), \Psi_S(\hat{S}(x), \check{S}(x))) \tag{3}$$

where  $T(x)$  and  $S(x)$  are feature maps from different layers in teacher and student models, respectively. Pairs of feature maps are picked from both teacher and student models,  $\hat{T}(x)$  and  $\check{T}(x)$  from the teacher and  $\hat{S}(x)$ ,  $\check{S}(x)$  from the student.  $\Psi_T$  and  $\Psi_S$  are the similarity functions for the pair of feature maps sampled and  $\mathcal{L}_R$  is the correlation function between the teacher and student network.

## 4 Method

### 4.1 Dataset and Preprocessing

We conduct our experiments on MixedWM38 WaferMap[28] dataset. It comprises more than 38000 original wafer maps as well as synthetically generated wafer maps divided into 38 classes based on their defect patterns. Out of the 38 defect pattern classes, one class is the normal class where the wafer does not contain any defects, 8 classes are of the single defect patterns and the rest are various combinations of single defect patterns. ‘Center(C)’ indicates a defect pattern concentrated at the center of the wafer, ‘Donut(D)’ refers to a defect pattern in the shape of an annular disc/donut, ‘Edge-Loc (EL)’ is a defect pattern where the defect is concentrated around the edge of the wafer, ‘Edge-Ring (ER)’ is a global defect pattern occurring around the entire circumference of the wafer, ‘Loc(L)’ refers to a defect pattern that is in a concentrated region other than the center. ‘Near-Full (NF)’ is a seldom occurring defect that is present throughout the entirety of the wafer, ‘Scratch(S)’ indicates a defect pattern that is in the form of a narrow line along the wafer. Finally, ‘Random(R)’ occurs when dies are failed at random locations on the wafer.

Other than these single defect patterns, there are various defect patterns that comprise two or more single defect patterns in the same wafer, which are known as mixed defect patterns. Based on the number of single defects, the mixed defect patterns have been named as 2 mixed-type, 3 mixed-type, and 4 mixed-type.

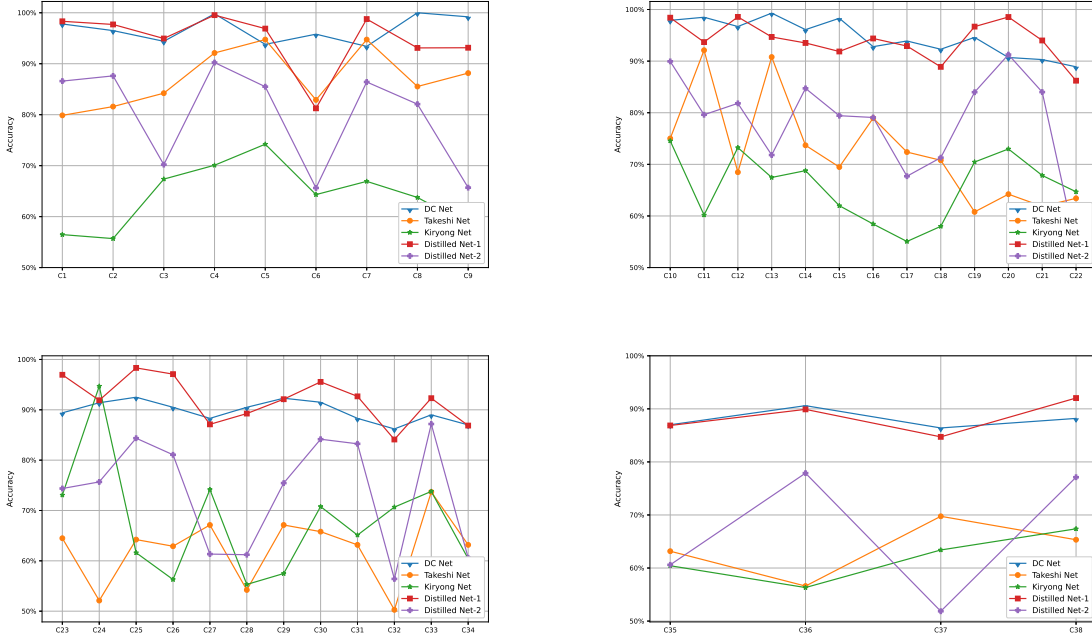


Table 1: Accuracy of compared models. (Top, Left to Right) The accuracy of compared models on single-type and 2-type mixed defects. The bigger distilled network (in red) competes well with contemporary models. (Bottom, Left to Right) Accuracy of compared models on 3 and 4-type mixed defects.

The 38 classes identified as C1 through C38, have C1 as the normal no-defect class while C2-C9 form the single defect patterns. The 13 mixed defect patterns from C10-C22 belong to wafers with 2 single-type defects on them, while the 12 mixed defect patterns identified as C23 through C34 belong to wafers with 3 single-type defects on them. Finally, the wafers with 4 single-type defects on them are identified as C35-C38.

The pattern name for mixed defect patterns indicates which single type pattern it consists of. For example, class C24, named C + EL + S, consists of three single defect patterns: Center, Edge-Loc, and Scratch. define the dataset

**Label encoding.** The classes are represented in the form of a one-hot encoded vector. Since there are 8 primary defects, an 8-length binary vector is used as the class label. The presence of ‘1’ at a position indicates a certain type of defect being present. For example, class C24 is encoded as [1 0 1 0 0 1 0] having 1’s at positions 1, 3, and 7, which corresponds to Center, Edge-Loc, and Scratch patterns being together simultaneously.

## 4.2 Architecture

We follow the response-based teacher-student model of knowledge distillation. Typically, the teacher model is a large neural network capable of performing complex tasks. The student model, on the other hand, is generally, a small network that we wish to train for specific tasks with the aim that the student model mimics closely the performance of the teacher model for the chosen task.

Let  $S$  be the student model and  $T$  be the teacher model parameterized by  $\theta$  and  $\theta'$  respectively. The learning for the student model aims to minimize the loss

where  $\sigma$  is the *sigmoid* function defined as  $\sigma(x) = \frac{1}{1+e^{-x}}$ ,  $\mathcal{L}_{mse}$  is the *mean-square-error*,  $\mathcal{L}_{bce}$  is the *binary cross entropy* error defined as  $\mathcal{L}_{bce}(x, y) = -y \log(x) - (1 - y) \log(1 - x)$  and  $\mathbb{1}_{x \geq \beta}$  is an indicator variable defined as

$$\mathbb{1}_{x \geq \beta} = \begin{cases} 1 & x \geq \beta \\ 0 & \text{otherwise} \end{cases}$$

The first part of the loss function encourages  $S$  to learn logits that are similar to what  $T$  has produced by minimizing the  $L_2$  distance between both sets of logits obtained on training images. The second part, on the other hand, penalizes misclassification by  $S$  on the predictions made by  $T$ . The parameters  $\alpha$  and  $\beta$  are hyperparameters that control 1)

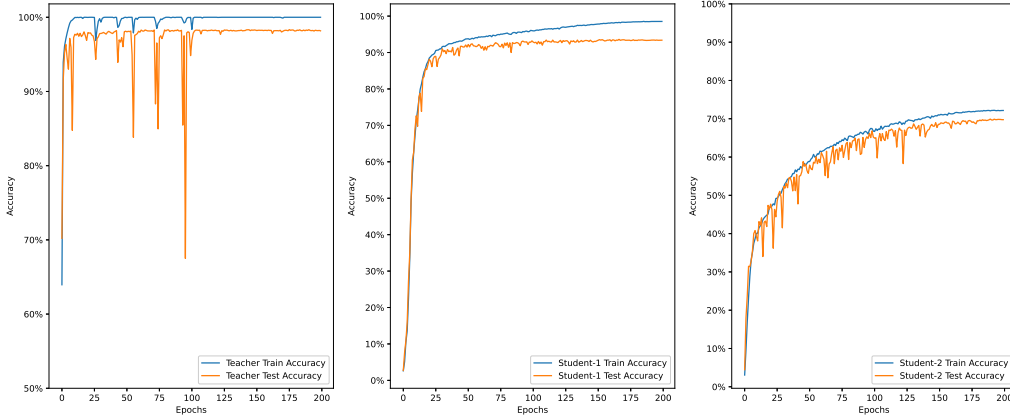


Figure 3: The accuracy curve of teacher model (left), bigger distilled network (middle) and smaller distilled network(right).

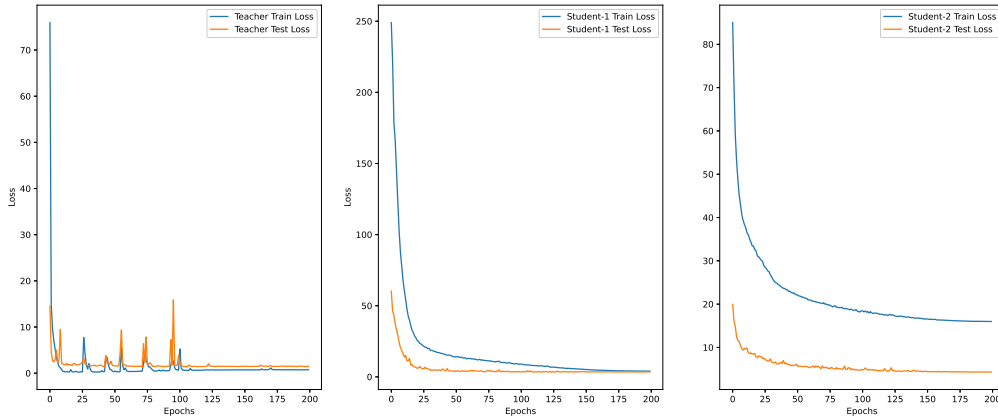


Figure 4: The loss curve of teacher model (left), bigger distilled network (middle) and smaller distilled network(right).

emphasis given to logits or final classification and 2) sensitivity of predictions made by  $T$ .  $\mathcal{L}_{KD}$  only depends on  $\theta$  as the weights of  $T$  are frozen during training and is free of hard labels  $\mathbf{y}$  which are approximated by  $T$ 's predictions on  $\mathbf{x}$ . For our experiments, we set both  $\alpha$  and  $\beta$  as 0.5.

For the teacher network, we use a ResNet-18 [9] modified to output 8 scores which is the number of unique defects. The network is trained to minimize the loss

$$\mathcal{L}_{\mathcal{T}}(\theta', \mathbf{x}, \mathbf{y}) = \mathcal{L}_{bce}(\sigma(T(\mathbf{x})), \mathbf{y}) \quad (4)$$

where  $\sigma$  is the *sigmoid* function described above.  $\mathbf{x}$  is the set of input wafer maps and  $\mathbf{y}$  is the corresponding labels.

The first student model contains three normal convolution layers, where the input wafer maps are convoluted with 6, 16, and 32 convolution kernels of size  $5 \times 5$  pixels. Each of the convolution layers is interspersed with a max pool layer of size  $2 \times 2$  and ReLU nonlinearity operation. Finally, a fully-connected classification head is attached. Similarly, the second student contains two normal convolution layers, where the input wafer maps are convoluted with 6, 16 convolution kernels of size  $5 \times 5$  pixels. The number of learnable parameters in both models is roughly  $60K$  and  $4K$  respectively contrasting to a total of around 11M learnable parameters in the teacher model.

Class	Distilled Net-1			Distilled Net-2		
	<i>Precision</i>	<i>Recall</i>	<i>F1</i>	<i>Precision</i>	<i>Recall</i>	<i>F1</i>
C1	92.76%	99.51%	96.02%	40.79%	100.0%	57.95%
C2	94.62%	98.32%	96.44%	73.81%	86.59%	79.69%
C3	93.01%	97.71%	95.3%	76.1%	87.61%	81.45%
C4	94.47%	94.95%	94.71%	90.26%	70.2%	78.98%
C5	95.11%	99.53%	97.27%	89.81%	90.23%	90.02%
C6	92.12%	96.89%	94.44%	68.18%	85.49%	75.86%
C7	100.0%	81.25%	89.66%	95.45%	65.62%	77.78%
C8	91.95%	98.77%	95.24%	68.63%	86.42%	76.5%
C9	99.26%	93.1%	96.09%	100.0%	82.07%	90.15%
C10	97.94%	93.14%	95.48%	82.21%	65.69%	73.02%
C11	92.54%	98.41%	95.38%	88.08%	89.95%	89.01%
C12	92.79%	93.69%	93.24%	71.0%	79.61%	75.06%
C13	95.81%	98.56%	97.17%	68.67%	81.82%	74.67%
C14	93.19%	94.68%	93.93%	80.84%	71.81%	76.06%
C15	93.95%	93.52%	93.74%	85.92%	84.72%	85.31%
C16	90.57%	91.87%	91.21%	71.24%	79.43%	75.11%
C17	96.35%	94.39%	95.36%	72.09%	79.08%	75.43%
C18	92.92%	92.92%	92.92%	79.69%	67.7%	73.21%
C19	97.96%	88.89%	93.2%	83.7%	71.3%	77.0%
C20	93.58%	96.69%	95.11%	88.37%	83.98%	86.12%
C21	96.67%	98.54%	97.6%	89.95%	91.26%	90.6%
C22	93.53%	94.0%	93.77%	66.93%	84.0%	74.5%
C23	94.09%	86.21%	89.97%	79.14%	54.19%	64.33%
C24	97.7%	96.95%	97.32%	90.15%	74.37%	81.5%
C25	89.47%	91.89%	90.67%	89.17%	75.68%	81.87%
C26	93.62%	98.32%	95.91%	84.36%	84.36%	84.36%
C27	93.02%	97.09%	95.01%	72.29%	81.07%	76.43%
C28	92.86%	87.11%	89.89%	75.32%	61.34%	67.61%
C29	95.02%	89.25%	92.05%	81.88%	61.21%	70.05%
C30	93.75%	92.11%	92.92%	87.76%	75.44%	81.13%
C31	86.94%	95.54%	91.04%	84.16%	84.16%	84.16%
C32	87.62%	92.67%	90.08%	62.85%	83.25%	71.62%
C33	94.8%	84.1%	89.13%	75.34%	56.41%	64.52%
C34	91.84%	92.31%	92.07%	85.0%	87.18%	86.08%
C35	94.51%	86.87%	90.53%	76.92%	60.61%	67.8%
C36	91.67%	89.9%	90.78%	88.52%	77.88%	82.86%
C37	93.85%	84.72%	89.05%	77.78%	51.85%	62.22%
C38	90.24%	92.04%	91.13%	86.59%	77.11%	81.58%
Average	<b>93.74%</b>	<b>93.33%</b>	<b>93.44%</b>	<b>79.71%</b>	<b>77.12%</b>	<b>77.41%</b>

Table 2: Precision and recall of Distilled Nets

## 5 Experiments and Results

We Experiment on the MixedWM38 dataset containing 38K wafer images as described in 4.1. We split the dataset into train-test using 80% of the images as the training set and the remaining as a test set. All the networks are trained on possible different train and test set which is sampled before each experiment independently. The experiments are conducted in Python compiler, Pytorch 1.12, and CUDA 11.3, the computer with option: Linux system, Intel(R) Xeon(R) CPU @ 2.20GHz, and Tesla T4 GPU.

We train all the networks with the setting mentioned in Table 3. We use ResNet-18 as the teacher network which is trained for 200 epochs gaining an overall accuracy of 98.34%. Next, we train two much smaller student networks that try to mimic the learning of the teacher model from unlabelled data. The size of the student model weights are 248KB and 21KB contrasting to 43MB of the teacher model. We compare our distilled models with state-of-the-art methods.

Description	Value
Initial Learning rate	0.01
Optimizer	SGD
Scheduler	Cosine Annealing
Batch Size	64
Decay	$5e^{-4}$
Momentum	0.9

Table 3: Experimental Parameter Settings used while Training

During training, the accuracy and loss curves of all three models are presented in Figure 3 and Figure 4. We observe that the first student model very quickly learns to mimic the teacher model and converges nicely over the course of 200 epochs while the second lags behind limited by its capacity. Despite its size and complexity, the second student model still outperforms Takeshi CNN and Kiryong CNN as presented in Table 5.

As reported in Figure 1, the test accuracy on the validation set of both the distilled networks is competitive in their class. The average accuracy of a bigger distilled network is 93.33% which is 20.81% higher than Takeshi-CNN and 27.48% higher than Kiryong-CNN and compares very well to DC-Net with 0.10% higher test accuracy whereas the smaller network achieves an accuracy of 77.11% which is 4.6% higher than Takeshi CNN and 11.27% higher than Kiryong CNN. We also observe the accuracy based on a number of defects on the wafer and report in Table 4. The bigger distilled network achieves competitive performance among the evaluated models especially when the number of defects on the wafer is greater than 2. The results demonstrate that despite being very small and by only using a tiny number of convolutional layers, the distilled networks were effective for detecting mixed-type patterns when compared to their *hard* trained counterparts. From a manufacturing point of view, the detection of wafers without any defects *i.e* normal wafers is crucial since any misclassification would result in passing a faulty wafer forward, which leads to downstream reliability issues. The accuracy of the proposed method with “C1” is 99.51%, which demonstrates the effectiveness of knowledge distillation in industrial practice.

We also evaluate the performance of our student networks in the false recognition and missing recognition of wafer defects. We use the statistical metrics of Precision, Recall, and F1 score which is defined as:

True Positive (TP): predicting positive, the actual is positive.

False Positive (FP): predicting positive, the actual is negative.

False Negative (FN): predicting negative, the actual is positive.

True Negative (TN): predicting negative, the actual is negative.

$$Recall = \frac{TP}{TP + FN} \quad (5)$$

$$Precision = \frac{TP}{TP + FP} \quad (6)$$

$$F1 = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}} \quad (7)$$

We report the performance of both the distilled networks in Table 2. Higher values of *Precision* indicate less false recognition and higher value of *Recall* indicate less missing recognition. As observed, the average precision and recall of our distilled net are very high which demonstrates that the distilled networks make a very less false prediction.

Class	Takeshi CNN	Kiryong CNN	DC Net	Distilled Net-1	Distilled Net-2
C1-C9	87.37%	64.8%	97.04%	95.32%	81.99%
C10-C22	72.46%	65.67%	94.64%	94.03%	78.37%
C23-C34	62.36%	67.78%	89.74%	92.02%	73.76%
C35-C38	66.05%	63.24%	89.61%	90.0%	70.99%

Table 4: Accuracy of the comparable models in single type defect and mixed type defects.



Class	Takeshi CNN	Kiryong CNN	DC Net	Distilled Net-1	Distilled Net-2
C1	89.89%	70.76%	99.7%	99.51%	100.0%
C2	79.87%	56.47%	97.8%	98.32%	86.59%
C3	81.59%	55.69%	96.5%	97.71%	87.61%
C4	84.22%	67.36%	94.4%	94.95%	70.2%
C5	92.11%	70.08%	99.8%	99.53%	90.23%
C6	94.74%	74.21%	93.8%	96.89%	85.49%
C7	82.9%	64.32%	95.8%	81.25%	65.62%
C8	94.74%	66.92%	93.4%	98.77%	86.42%
C9	85.53%	63.76%	100.0%	93.1%	82.07%
C10	88.16%	58.41%	99.2%	93.14%	65.69%
C11	75.01%	74.54%	97.9%	98.41%	89.95%
C12	92.11%	60.18%	98.5%	93.69%	79.61%
C13	68.48%	73.24%	96.7%	98.56%	81.82%
C14	90.79%	67.46%	99.3%	94.68%	71.81%
C15	73.69%	68.78%	96.1%	93.52%	84.72%
C16	69.48%	61.99%	98.3%	91.87%	79.43%
C17	78.96%	58.44%	92.8%	94.39%	79.08%
C18	72.38%	55.06%	93.9%	92.92%	67.7%
C19	70.79%	57.99%	92.3%	88.89%	71.3%
C20	60.79%	70.45%	94.6%	96.69%	83.98%
C21	64.22%	72.98%	90.7%	98.54%	91.26%
C22	61.86%	67.85%	90.3%	94.0%	84.0%
C23	63.42%	64.69%	88.9%	86.21%	54.19%
C24	64.49%	73.06%	89.4%	96.95%	74.37%
C25	52.11%	94.69%	91.4%	91.89%	75.68%
C26	64.22%	61.61%	92.5%	98.32%	84.36%
C27	62.9%	56.29%	90.5%	97.09%	81.07%
C28	67.12%	74.16%	88.3%	87.11%	61.34%
C29	54.22%	55.31%	90.5%	89.25%	61.21%
C30	67.12%	57.48%	92.3%	92.11%	75.44%
C31	65.8%	70.78%	91.5%	95.54%	84.16%
C32	63.17%	65.11%	88.3%	92.67%	83.25%
C33	50.27%	70.67%	86.2%	84.1%	56.41%
C34	73.69%	73.79%	89.0%	92.31%	87.18%
C35	63.17%	60.43%	87.0%	86.87%	60.61%
C36	56.6%	56.32%	90.6%	89.9%	77.88%
C37	69.75%	63.4%	86.4%	84.72%	51.85%
C38	65.33%	67.39%	88.2%	92.04%	77.11%
Average	<b>72.52%</b>	<b>65.85%</b>	<b>93.23%</b>	<b>93.33%</b>	<b>77.12%</b>

Table 5: Accuracy of the distilled models compared with state-of-the-art models.

## 6 Discussion

To further evaluate the effectiveness of the training, we visualize the embeddings of test wafers produced by the teacher and both the student models. In each visualization plot, the defect class is represented by a unique color. As is shown in Figure 5, the embeddings produced by bigger distilled net resemble well the separation achieved by the teacher net. This is also apparent in the t-SNE visualization of the pair. The smaller distilled network, though limited greatly by its size, also achieves reasonable separation. This demonstrates graphically the learning capability of the student network to mimic the structure of the teacher network in the latent space. Future work would be to positively apply other distillation schemes to train lightweight DPR models that can be used efficiently in the fabrication process.

## 7 Conclusion

In this study, we discuss the identification of the mixed-type defective pattern produced during the wafer manufacturing process. Correctly identifying and classifying these patterns gives valuable insights into the manufacturing operation, especially during the productivity improvement phase. The contribution of this article can be summarized as follows:

- This study introduces an unsupervised form of knowledge distillation that benefits lightweight models to achieve great accuracy compared to models trained from scratch. Results clearly show that distilled models

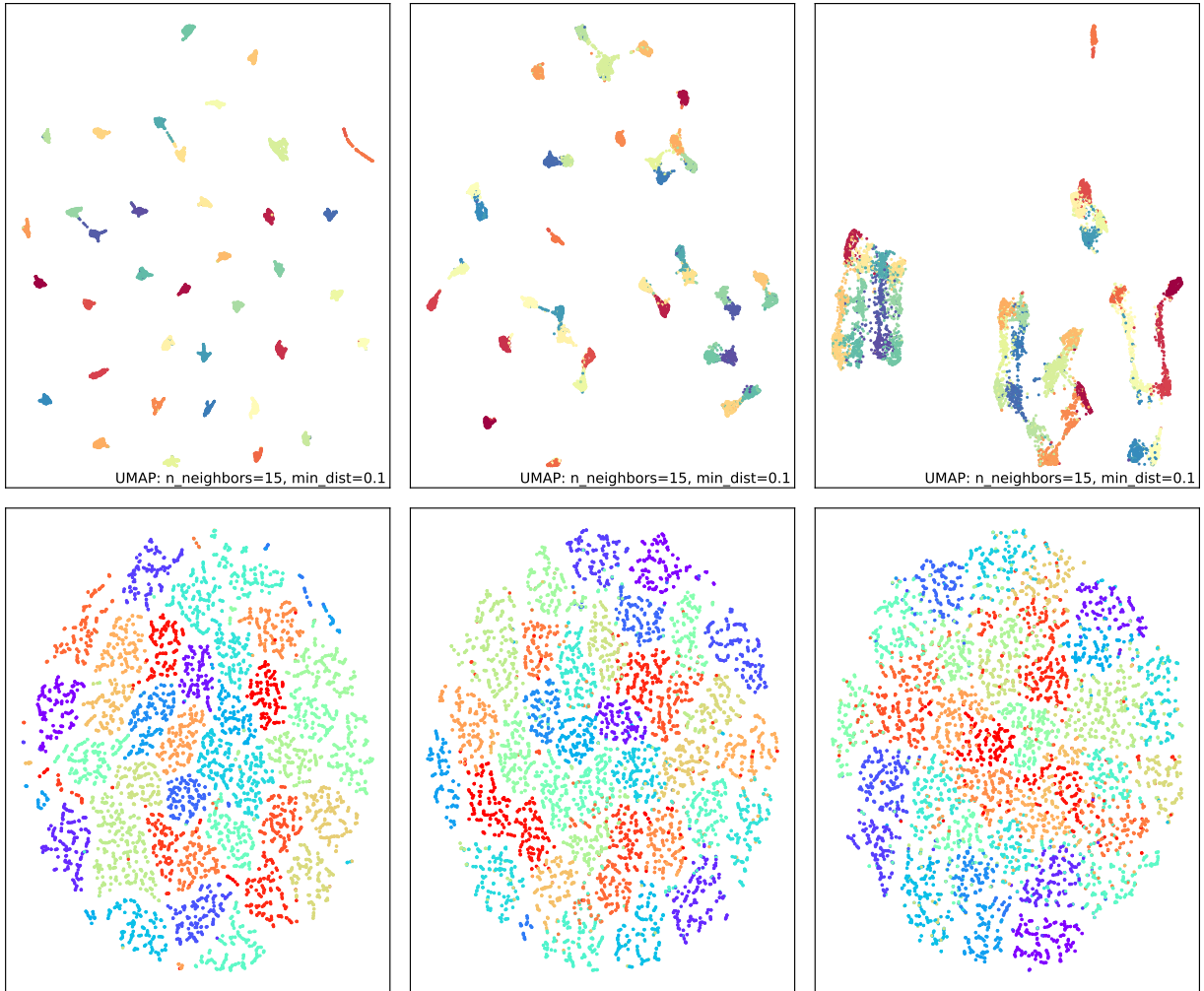


Figure 5: (Top) U-MAP visualization of the embeddings produced by (left to right) Teacher model, Distilled model-1 and Distilled model 2. (Bottom) t-SNE visualization of the embeddings produced by (left to right) Teacher model, Distilled model-1 and Distilled model 2.

benefit the accuracy of mixed-type defect pattern recognition(DPR) of wafer maps without using ground-truth labels.

- Compare with conventional DPR works, the distilled models are lighter, and easily deployable on low-memory devices without compromising the performance. In fact, the distilled model may sometimes perform better than the teacher model it is trained from.

Taking all points into account, future work will be worthwhile to apply the proposed method to detect the root causes of defects. Besides, wafer defect modeling and analysis to improve wafer fabrication quality will be explored in further research.

## References

- [1] Liliana, L. (2016). A new model of Ishikawa diagram for quality assessment. IOP Conference Series: Materials Science and Engineering, 161(1), 012099. <https://doi.org/10.1088/1757-899X/161/1/012099>
- [2] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," NIPS Deep Learning Workshop, 2015.

- [3] A. Krizhevsky, I. Sutskever and G. E. Hinton, "ImageNet classification with deep convolutional neural networks", *Proc. Adv. Neural Inf. Process. Syst.*, pp. 1097-1105, 2012.
- [4] Y. Kong and D. Ni, "A Semi-Supervised and Incremental Modeling Framework for Wafer Map Classification," in *IEEE Transactions on Semiconductor Manufacturing*, vol. 33, no. 1, pp. 62-71, Feb. 2020, doi: 10.1109/TSM.2020.2964581.
- [5] K. S. -M. Li et al., "Wafer Defect Pattern Labeling and Recognition Using Semi-Supervised Learning," in *IEEE Transactions on Semiconductor Manufacturing*, vol. 35, no. 2, pp. 291-299, May 2022, doi: 10.1109/TSM.2022.3159246.
- [6] R. R. Schaller, "Moore's law: past, present and future," in *IEEE Spectrum*, vol. 34, no. 6, pp. 52-59, June 1997, doi: 10.1109/6.591665.
- [7] R. Ufuk Bilsel & Dennis K.J. Lin (2012) Ishikawa Cause and Effect Diagrams Using Capture Recapture Techniques, *Quality Technology & Quantitative Management*, 9:2, 137-152, DOI: 10.1080/16843703.2012.11673282
- [8] Krizhevsky, Alex and Sutskever, Ilya and Hinton, Geoffrey E, "ImageNet Classification with Deep Convolutional Neural Networks", *Advances in Neural Information Processing Systems*, 2012
- [9] He, Kaiming and Zhang, Xiangyu and Ren, Shaoqing and Sun, Jian, "Deep Residual Learning for Image Recognition", 2016 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
- [10] S. Liu and W. Deng, "Very deep convolutional neural network based image classification using small training sample size," 2015 3rd *IAPR Asian Conference on Pattern Recognition (ACPR)*, Kuala Lumpur, Malaysia, 2015.
- [11] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. & Riedmiller, M. (2013). Playing Atari with Deep Reinforcement Learning
- [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charlie Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015.
- [13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, page 6000–6010, 2017
- [14] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2015.
- [15] Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. 2006. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '06)*. Association for Computing Machinery, New York, NY, USA, 535–541. <https://doi.org/10.1145/1150402.1150464>
- [16] Ruth Urner, Shai Shalev-Shwartz, & Shai Ben-David (2011). Access to Unlabeled Data can Speed up Prediction Time. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)* (pp. 641–648). ACM.
- [17] F. Adly et al., "Simplified Subspaced Regression Network for Identification of Defect Patterns in Semiconductor Wafer Maps," in *IEEE Transactions on Industrial Informatics*, vol. 11, no. 6, pp. 1267-1276, Dec. 2015, doi: 10.1109/TII.2015.2481719.
- [18] J. Wang, J. Zhang and X. Wang, "A Data Driven Cycle Time Prediction With Feature Selection in a Semiconductor Wafer Fabrication System," in *IEEE Transactions on Semiconductor Manufacturing*, vol. 31, no. 1, pp. 173-182, Feb. 2018, doi: 10.1109/TSM.2017.2788501.
- [19] T. Nakazawa and D. V. Kulkarni, "Wafer Map Defect Pattern Classification and Image Retrieval Using Convolutional Neural Network," in *IEEE Transactions on Semiconductor Manufacturing*, vol. 31, no. 2, pp. 309-314, May 2018, doi: 10.1109/TSM.2018.2795466.
- [20] K. Kyeong and H. Kim, "Classification of Mixed-Type Defect Patterns in Wafer Bin Maps Using Convolutional Neural Networks," in *IEEE Transactions on Semiconductor Manufacturing*, vol. 31, no. 3, pp. 395-402, Aug. 2018, doi: 10.1109/TSM.2018.2841416.
- [21] X. Wang, K. C. K. Chan, K. Yu, C. Dong and C. C. Loy, "EDVR: Video Restoration With Enhanced Deformable Convolutional Networks," 2019 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, Long Beach, CA, USA, 2019, pp. 1954-1963, doi: 10.1109/CVPRW.2019.00247.
- [22] F. Tung and G. Mori, "Similarity-Preserving Knowledge Distillation," 2019 *IEEE/CVF International Conference on Computer Vision (ICCV)*, Seoul, Korea (South), 2019, pp. 1365-1374, doi: 10.1109/ICCV.2019.00145.
- [23] Bagherinezhad, H., Horton, M., Rastegari, M., Farhadi, A.: Label refinery: Improving imagenet classification through label progression.(2018)

- [24] Dong, B., Hou, J., Lu, Y., Zhang, Z.: Distillation – early stopping? harvesting dark knowledge utilizing anisotropic information retrieval for overparameterized neural network. (2019)
- [25] Furlanello, T., Lipton, Z., Tschannen, M., Itti, L., Anandkumar, A.: Born again neural networks. In: International Conference on Machine Learning. pp. 1607–1616. PMLR (2018)
- [26] N. Aloysius and M. Geetha, "A review on deep convolutional neural networks," 2017 International Conference on Communication and Signal Processing (ICCSP), Chennai, India, 2017, pp. 0588-0592, doi: 10.1109/ICCSP.2017.8286426.
- [27] Z. Li et al., "A Unified Understanding of Deep NLP Models for Text Classification," in IEEE Transactions on Visualization and Computer Graphics, vol. 28, no. 12, pp. 4980-4994, 1 Dec. 2022, doi: 10.1109/TVCG.2022.3184186.
- [28] J. Wang, C. Xu, Z. Yang, J. Zhang and X. Li, "Deformable Convolutional Networks for Efficient Mixed-Type Wafer Defect Pattern Recognition," in IEEE Transactions on Semiconductor Manufacturing, vol. 33, no. 4, pp. 587-596, Nov. 2020, doi: 10.1109/TSM.2020.3020985.