

# Data Efficient Lithography Modeling with Residual Neural Networks and Transfer Learning

Yibo Lin

The University of Texas at Austin  
yibolin@cerc.utexas.edu

Yuki Watanabe

Toshiba Memory Corporation  
yuki9.watanabe@toshiba.co.jp

Taiki Kimura

Toshiba Memory Corporation  
taiki2.kimura@toshiba.co.jp

Tetsuaki Matsunawa

Toshiba Memory Corporation  
tetsuaki.matsunawa@toshiba.co.jp

Shigeki Nojima

Toshiba Memory Corporation  
shigeki.nojima@toshiba.co.jp

Meng Li

The University of Texas at Austin  
alfred@cerc.utexas.edu

David Z. Pan

The University of Texas at Austin  
dpan@cerc.utexas.edu

## ABSTRACT

Lithography simulation is one of the key steps in physical verification, enabled by the substantial optical and resist models. A resist model bridges the aerial image simulation to printed patterns. While the effectiveness of learning-based solutions for resist modeling has been demonstrated, they are considerably data-demanding. Meanwhile, a set of manufactured data for a specific lithography configuration is only valid for the training of one single model, indicating low data efficiency. Due to the complexity of the manufacturing process, obtaining enough data for acceptable accuracy becomes very expensive in terms of both time and cost, especially during the evolution of technology generations when the design space is intensively explored. In this work, we propose a new resist modeling framework for contact layers that utilizes existing data from old technology nodes to reduce the amount of data required from a target lithography configuration. Our framework based on residual neural networks and transfer learning techniques is effective within a competitive range of accuracy, i.e., 2-10X reduction on the amount of training data with comparable accuracy to the state-of-the-art learning approach.

## ACM Reference Format:

Yibo Lin, Yuki Watanabe, Taiki Kimura, Tetsuaki Matsunawa, Shigeki Nojima, Meng Li, and David Z. Pan. 2018. Data Efficient Lithography Modeling with Residual Neural Networks and Transfer Learning. In *Proceedings of 2018 International Symposium on Physical Design (ISPD'18)*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3177540.3178242>

## 1 INTRODUCTION

Due to the continuous semiconductor scaling from 10nm technology node (N10) to 7nm node (N7) [10, 11], the prediction of printed pattern sizes is becoming increasingly difficult and complicated due to the complexity of manufacturing process and variations. However, complex designs demand accurate simulations to guarantee functionality and yield. Resist modeling, as a key component in

lithography simulation, is critical to bridge the aerial image simulation to manufactured wafer data. Rigorous simulations that perform physics-level modeling suffer from large computational overhead, which are not suitable when used extensively. Thus compact resist models are widely used in practice.

Figure 1(a) shows the process of lithography simulations where the optical model computes the aerial image from the input mask patterns and the resist model determines the output patterns from this. As the aerial image contains the light intensity map, the resist model needs to determine the slicing thresholds for the output patterns as shown in Figure 1(b). With the thresholds, the critical dimensions (CDs) of printed patterns can be computed, which need to match CDs measured from manufactured patterns. In practice, various factors may impact a resist model such as the physical properties of photoresist, design rules of patterns, process variations.

Accurate lithography simulation like rigorous physics-based simulation is notorious for its long computational time, while simulation with compact models suffers from accuracy issues [21, 25]. On the other hand, machine learning techniques are able to construct accurate models and then make efficient predictions. These approaches first take training data to calibrate a model and then use this model to make predictions on testing data for validation. The effectiveness of learning-based solutions has been studied in various lithography related areas including aerial image simulation [15], hotspot detection [13, 16, 22, 26, 28, 29], optical proximity correction (OPC) [5, 8, 14, 17], sub-resolution assist features (SRAF) [24, 27], resist modeling [21, 25], etc. In resist modeling, a convolutional neural network (CNN) that predicts slicing thresholds in aerial images is proposed [25]. The neural network consists of three convolution layers and two fully connected layers. Since the slicing threshold is a continuous value, learning a resist model is a regression task rather than a classification task. Around 70% improvement in accuracy is reported compared with calibrated compact models from Mentor Calibre [18]. Shim et al. [21] propose an artificial neural network (ANN) with five hidden layers to predict the height of resist after exposure. Significant speedup is reported with high accuracy compared with a rigorous simulation.

Although the learning-based approaches are able to achieve high accuracy, they are generally data-demanding in model training. In other words, big data is assumed to guarantee accuracy and generality. Furthermore, one data sample can only be used to train the

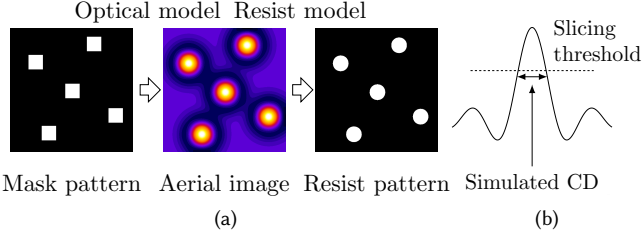
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ISPD'18, March 25–28, 2018, Monterey, CA, USA

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5626-8/18/03...\$15.00

<https://doi.org/10.1145/3177540.3178242>



**Figure 1: (a) Process of lithography simulation with optical and resist models. (b) Thresholds for aerial image determine simulated CD, which should match manufactured CD.**

corresponding model under the same lithography configuration, indicating a low data efficiency. Here data efficiency evaluates the accuracy a model can achieve given a specific amount of data, or the amount of data samples are required to achieve target accuracy. Nevertheless, obtaining a large amount of data is often expensive and time-consuming, especially when the technology node switches from one to another and the design space is under active exploration, e.g., from N10 to N7. The lithography configurations including optical sources, resist materials, etc., are frequently changed for experiments. Therefore, a fast preparation of models with high accuracy is urgently desired.

Different from previous approaches, in this work, we assume the availability of large amounts of data from the previous technology generation with old lithography configurations and small amounts of data from a target lithography configuration. We focus on increasing the data efficiency by reusing those from other lithography configurations and transfer the knowledge between different configurations. The objective is to achieve accurate resist models with significantly fewer data to a target configuration. The major contributions are summarized as follows.

- We propose a high performance resist modeling technique based on the residual neural network (ResNet).
- We propose a transfer learning scheme for ResNet that can reduce the amount of data with a target accuracy by utilizing the data from other configurations.
- We explore the impacts from various lithography configurations on the knowledge transfer.
- The experimental results demonstrate 2-10X reduction in the amount of training data to achieve accuracy comparable to the state-of-the-art learning approach [25].

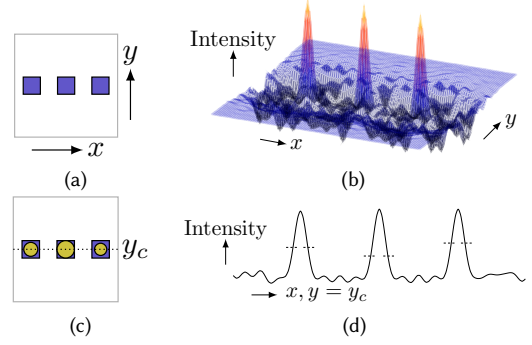
The rest of the paper is organized as follows. Section 2 illustrates the problem formulation. Section 3 explains the details of our approach. The effectiveness of our approach is verified in Section 4 and the conclusion is drawn in Section 5.

## 2 PRELIMINARIES

In this section, we will briefly introduce the background knowledge on lithography simulation and resist modeling. Then the problem formulation is explained. We mainly focus on contact layers in this work, but our methodology shall be applicable to other layers.

### 2.1 Lithography Simulation

Lithography simulation is generally composed of two stages, i.e., optical simulation and resist simulation, where optical and resist models are required, respectively. In the optical simulation, an optical model, characterized by the illumination tool, takes mask patterns to



**Figure 2: (a) Design target of 3 contacts and (b) the light intensity plot of aerial image. Assume that RETs such as SRAF and OPC have been already applied to the contacts before optical simulation. (c) A dotted line horizontally crosses the centers at  $y = y_c$  and the circles denote the contours of printed patterns. (d) Light intensity profiling along the dotted line at  $y = y_c$  extracted from the aerial image and different slicing thresholds for each contact.**

compute aerial images, i.e., light intensity maps. Then in the resist simulation, a resist model finalizes the resist patterns with the aerial images from the optical simulation. Generally, there are two types of resist models. One is a variable threshold resist (VTR) model in which the thresholds vary according to aerial images, and the other is a constant threshold resist (CTR) model in which the light intensity is modulated in an aerial image. We adopt the former since it is suitable to learning-based approaches [25].

Figure 2 shows an example of lithography simulation for a clip with three contacts. We assume that proper resolution enhancement techniques (RETs) such as OPC and SRAF have been applied before the computation of the aerial image [12]. The optical simulation generates the aerial image, as shown in Figure 2(b). Resist simulation then computes the thresholds in the aerial image to predict printed patterns. If we consider the horizontal sizes of contacts along the dotted line in Figure 2(c), the light intensity profiling can be extracted from the aerial image along the line and calculates the CDs for each contact with the thresholds.

### 2.2 Historical Data and Transfer Learning

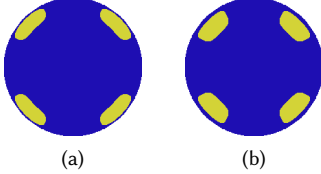
Since the lithography configurations evolve from one generation to another with the advancement of technology nodes, there are plenty of historical data available for the old generation. As mentioned in Section 1, accurate models require a large amount of data for training or calibration, which are expensive to obtain during the exploration of a new generation. If the lithography configurations have no fundamental changes, the knowledge learned from the historical data may still be applicable to the new configuration, which can eventually help to reduce the amount of new data required.

Transfer learning represents a set of techniques to transfer the knowledge from one or multiple source domains to a target domain, utilizing the underlying similarity between the data from these domains. Various studies have explored the effectiveness of knowledge transfer in image recognition and robotics [6, 19, 20], while it is not clear whether the knowledge between different resist models is transferable or not.

In this work, we consider the evolution of the contact layer from the cutting edge technology node N10 to N7 [10, 11]. A large amount

**Table 1: Lithography Configurations for N10 and N7**

	N10	N7	
		N7 <sub>a</sub>	N7 <sub>b</sub>
Design Rule	A	B	B
Optical Source	A	B	B
Resist Material	A	A	B

**Figure 3: Optical sources (yellow) for (a) N10 and (b) N7.**

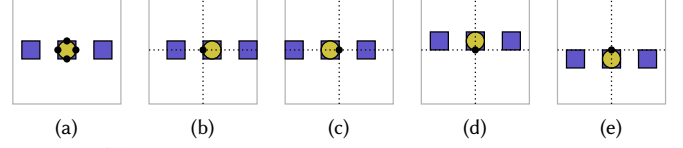
of available N10 data are assumed. During the evolution to N7, different design rules for mask patterns, optical sources and resist materials for lithography are explored. Table 1 shows the lithography configurations considered for N10 and N7. Differences in letters *A*, *B* represent different configurations of design rules, optical sources, or resist materials. One configuration for N10 is considered, while two configurations are considered for N7, i.e., N7<sub>a</sub>, N7<sub>b</sub>, with two kinds of resist materials (about 20% difference in the slopes of dissolution curves). From N10 to N7, both the design rules and optical sources are changed. For N10, we consider a pitch of 64nm with double patterning lithography, while for N7, the pitch is set to 45nm with triple patterning lithography [10]. The width of each contact is set to half pitch. The lithography target of each contact is set to 60nm for both N10 and N7. Optical sources calibrated with industrial strength for N10 and N7 are shown in Figure 3, with the same type of illumination shapes.

Various combinations of knowledge transfer can be explored from Table 1, such as N10→N7, N7<sub>i</sub>→N7<sub>j</sub>, and N10+N7<sub>i</sub>→N7<sub>j</sub>, where  $i \neq j, i, j \in \{a, b\}$ .

### 2.3 Learning-based Resist Modeling

The thresholds of positions near the contacts are of significant importance since they usually determine the boundaries of printed contacts. Hence we consider the middle of the left, right, bottom and top edges for each contact, as shown in Figure 4(a), where the positions for prediction are highlighted with black dots. In addition, the threshold is mainly influenced by the surrounding mask patterns. Therefore, resist models typically compute the threshold using a clip of mask patterns centered by a target position. To measure the thresholds in Figure 4(a), we select a clip where the target position lies in its center, as shown in Figure 4(b) to Figure 4(e). The task of a resist model is to compute the thresholds for these positions of each contact [25].

Learning-based resist modeling consists of two phases, training and testing. In the training phase, training dataset with both aerial images and thresholds are used to calibrate the model, while in the testing phase, the model predicts thresholds for the aerial images from the testing dataset and compares with the golden thresholds to validate the model.

**Figure 4: (a) The thresholds for the middle of the 4 edges of the center contact are predicted. (b) (c) (d) (e) The clip window is shifted such that the target position lies in the center of the clip.**

### 2.4 Problem Formulation

The accuracy<sup>1</sup> of a model is evaluated with root mean square (RMS) error defined as follows,

$$\epsilon = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y} - y)^2}, \quad (1)$$

where  $N$  denotes the amount of samples,  $y$  denotes the golden values and  $\hat{y}$  denotes the predicted values. We further define relative RMS error,

$$\epsilon_r = \sqrt{\frac{1}{N} \sum_{i=1}^N \left( \frac{\hat{y} - y}{y} \right)^2}, \quad (2)$$

where a relative ratio of error from the golden values can be represented. Both metrics can refer to errors in either CD or threshold. Although during model training, the RMS error of threshold is generally minimized due to easier computation, the eventual model is often evaluated with the RMS error of CD for its physical meaning to the patterns. The RMS errors in threshold and CD essentially have almost the same fidelity, and usually yield consistent comparison. For convenience, we report relative RMS error in threshold ( $\epsilon_r^{th}$ ) for comparison of different models since it removes the dependency to the scale of thresholds, and use RMS error in CD ( $\epsilon^{CD}$ ) for data efficiency related comparison.

**Definition 1** (Data Efficiency). *The amount of target domain data required to learn a model with a given accuracy.*

Given a specific amount of data from a target domain, if one can learn a model with a higher accuracy than another, it also indicates higher data efficiency. Thus improving model accuracy benefits data efficiency as well.

The resist modeling problem is defined as follows.

**Problem 1** (Learning-based Resist Modeling). *Given a dataset containing information of aerial images and thresholds at their centers, train a resist model that can maximize the accuracy for the prediction of thresholds.*

In practice, accuracy is not the only objective. The amount of training data should be minimized as well due to the high cost of data preparation. Therefore, we propose the problem of data efficient resist modeling as follows.

**Problem 2** (Data Efficient Resist Modeling). *Given datasets from N10 and N7 containing information of aerial images and thresholds, train a resist model for target dataset N7<sub>i</sub> that can achieve high accuracy and meanwhile minimize the amount of data required for N7<sub>i</sub>, where  $i \in \{a, b\}$ .*

<sup>1</sup> Note that the accuracy we talk about in this paper refers to the accuracy at end of lithography flow including all RETs.

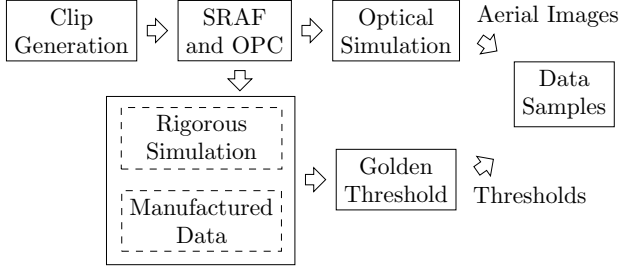
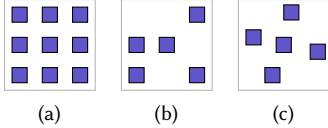


Figure 5: Flow of data preparation.

Figure 6: (a) A clip of  $3 \times 3$  contact array. (b) A clip of  $3 \times 3$  randomized contact array. (c) A clip of contacts with random positions.

### 3 ALGORITHMS

In this section, we will explain the structure of our models and then the details regarding the transfer learning scheme.

#### 3.1 Data Preparation

Figure 5 gives the flow of data preparation. We first generate clips and perform SRAF insertion and OPC. The aerial images are then computed from the optical simulation, and at the same time, the golden thresholds need to be computed from either the rigorous simulation or the manufactured data. Each data sample consists of an aerial image and the threshold at its center.

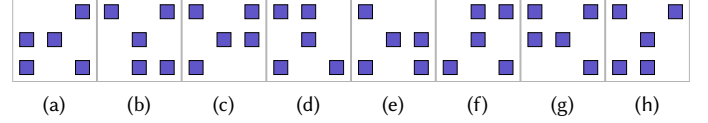
**3.1.1 Clip Generation.** Following the design rules such as minimum pitch of contacts, we generate three types of  $2 \times 2 \mu\text{m}$  clips. It is necessary to ensure that there is a contact in the center of each clip since that is the target contact for threshold computation.

**Contact Array.** All possible  $m \times n$  arrays of contacts within the dimensions of clips are enumerated. The steps of the arrays can be multiple times of the minimum pitch  $p$ , i.e.,  $p, 2p, 3p, \dots$ , in horizontal or vertical directions. An example of  $3 \times 3$  contact array with a certain pitch is shown in Figure 6(a). It needs to mention that the same  $3 \times 3$  contact array with different steps should be regarded as different clips due to discrepant spacing.

**Randomized Contact Array.** The aforementioned contact arrays essentially distribute contacts on grids and fill all the slots in the grid maps. The randomization of contact arrays is implemented by a random distribution of contacts in those grid maps. Fig 6(b) shows an example of randomized contact array from the  $3 \times 3$  contact array in Figure 6(a). Various distribution of contacts can be generated even from the same grid maps.

**Contacts with Random Positions.** Contacts in this type of clips do not necessarily align to any grid map, as their positions are randomly generated, while the design rules are still guaranteed. An example is shown in Figure 6(c). No matter how the surrounding contacts change, the contact in the center of the clip should remain the same.

**3.1.2 Data Augmentation.** Due to the symmetry of optical sources in Figure 3, data can be augmented with rotation and flipping, improving the data efficiency [4]. Eight combinations of rotation and flipping are shown in Figure 7, where new data samples are obtained

Figure 7: Combinations of rotation and flipping. (a) Original. (b) Rotate  $90^\circ$ . (c) Rotate  $180^\circ$ . (d) Rotate  $270^\circ$ . (e) Flip. (f) Flip and rotate  $90^\circ$ . (g) Flip and rotate  $180^\circ$ . (h) Flip and rotate  $270^\circ$ .

without new thresholds. Data augmentation inflates datasets to obtain models with better generalization.

#### 3.2 Convolutional Neural Networks

Convolutional neural networks (CNN) have demonstrated impressive performance on mask related applications in lithography such as hotspot detection, and resist modeling [25, 29]. The structure of CNN mainly includes convolution layers and fully connected layers. Features are extracted from convolution layers and then classification or regression is performed by fully connected layers. Figure 10(a) illustrates a CNN structure with three convolution layers and two fully connected layers [25]. The first convolution layer has 64 filters with dimensions of  $7 \times 7$ . Although not explicitly shown most of the time, a rectified linear unit (ReLU) layer for activation is applied immediately after the convolution layer, where the ReLU function is defined as,

$$x^l = \begin{cases} x^{l-1}, & \text{if } x^{l-1} \geq 0, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

Then the max-pooling layer performs down-sampling with a factor of 2 to reduce the feature dimensions and improve the invariance to translation [4]. After three convolution layers, two fully connected layers are applied where the first one has 256 hidden units followed with a ReLU layer and a 50% dropout layer, and second one connects to the output.

#### 3.3 Residual Neural Networks

One way to improve the performance of CNN is to increase the depth for a larger capacity of the neural networks. However, the counterintuitive degradation of training accuracy in CNN is observed when stacking more layers, preventing the neural networks from better performance [7]. An example of CNNs with 5 and 10 layers is shown in Figure 8, where the deeper CNN fails to converge to a smaller training error than the shallow one due to gradient vanishing [2, 3], eventually resulting in the failure to achieve a better testing error either. The study from He et al. [7] reveals that the underlying reason comes from the difficulty of identity mapping. In other words, fitting a hypothesis  $\mathcal{H}(x) = x$  is considerably difficult for solvers to find optimal solutions. To overcome this issue, residual neural networks (ResNet), which utilizes shortcut connections, are adopted to assist the convergence of training accuracy.

The building block of ResNet is illustrated in Figure 9, where a shortcut connection is inserted between the input and output of two convolution layers. Let the function  $\mathcal{F}(x)$  be the mapping defined by the two convolution layers. Then the entire function for the building block becomes  $\mathcal{F}(x) + x$ . Suppose the building block targets to fit the hypothesis  $\mathcal{H}(x)$ . The residual networks train  $\mathcal{F}(x) = \mathcal{H}(x) - x$ , while the convolution layers without shortcut connections like that in CNN try to directly fit  $\mathcal{F}(x) = \mathcal{H}(x)$ . Theoretically, if  $\mathcal{H}(x)$  can be approximated with  $\mathcal{F}(x)$ , then it can also be approximated with  $\mathcal{F}(x) + x$ . Despite the same nature, comprehensive experiments have



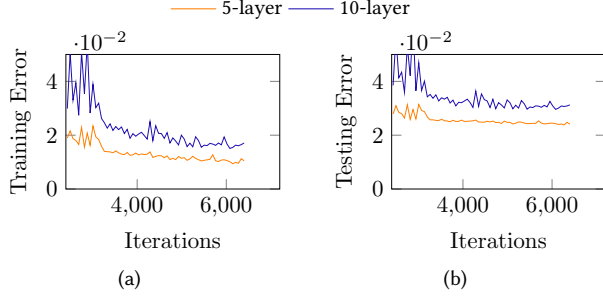


Figure 8: Counterintuitive (a) training and (b) testing errors for different depth of CNN with epochs.

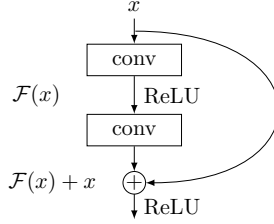


Figure 9: Building block of ResNet.

demonstrated a better convergence of ResNet than that of CNN for deep neural networks [7]. We also observe a better performance of ResNet with the transfer learning schemes than that of CNN in our problem, which has never been explored before.

The ResNet is shown in Figure 10(b) with 8 convolution layers and 2 fully connected layers. Different from the original setting [7], we add a shortcut connection to the first convolution layer by broadcasting the input tensor of  $64 \times 64 \times 1$  to  $64 \times 64 \times 64$ . This minor change enables better empirical results in our problem. For the rest of the networks, 3 building blocks for ResNet are utilized.

### 3.4 Transfer Learning

Transfer learning aims at adapting the knowledge learned from data in source domains to a target domain. The transferred knowledge will benefit the learning in the target domain with a faster convergence and better generalization [4]. Suppose the data in the source domain has a distribution  $P_s$  and that in the target domain has a distribution  $P_t$ . The underlying assumption of transfer learning lies in the common factors that need to be captured for learning the variations of  $P_s$  and  $P_t$ , so that the knowledge for  $P_s$  is also useful for  $P_t$ . An intuitive example is that learning to recognize cats and dogs in the source task helps the recognition of ants and wasps in the target task, especially when the source task has significantly larger dataset than that of the target task. The reason comes from the low-level notions of edges, shapes, etc., shared by many visual categories [4]. In resist modeling, different lithography configurations can be viewed as separate tasks with different distributions.

Typical transfer learning scheme for neural networks fixes the first several layers of the model trained for another domain and finetune the successive layers with data from the target domain. The first several layers usually extract general features, which are considered to be similar between the source and the target domains, while the successive layers are classifiers or regressors that need to be adjusted. Figure 11 shows an example of the transfer learning scheme. We first train a model with source domain data and then use the source domain model as the starting point for the training of the target domain. During the training for the target domain, the first  $k$  layers

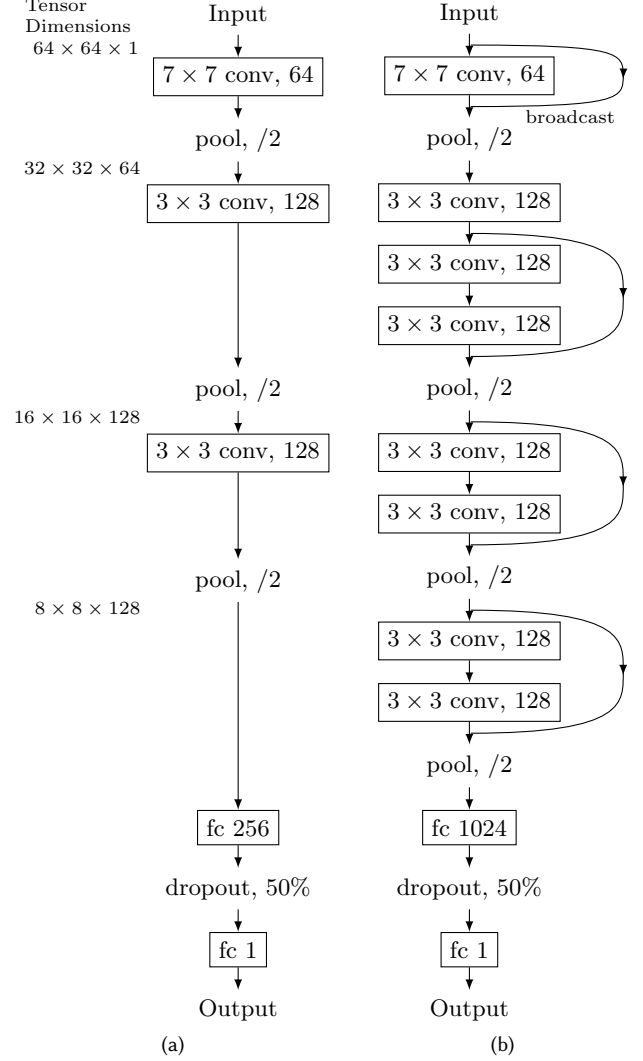


Figure 10: (a) CNN and (b) ResNet structure.

are fixed, while the rest layers are finetuned. We denote this scheme as  $TF_k$ , shortened from “Transfer and Fix”, where  $k$  is the parameter for the number of fixed layers.

## 4 EXPERIMENTAL RESULTS

Our framework is implemented with Tensorflow [1] and validated on a Linux server with 3.4GHz Intel i7 CPU and Nvidia GTX 1080 GPU. Around 980 mask clips are generated according to Section 3.1 for N10 and N7 separately following the design rules in Section 2.2, respectively. N7<sub>a</sub> and N7<sub>b</sub> use the same set of clips, but different lithography configurations. SRAF, OPC and aerial image simulation are performed with Mentor Calibre [18]. The golden CD values are obtained from rigorous simulation using Synopsys Sentaurus Lithography models [23] calibrated from manufactured data for N10, N7<sub>a</sub>, and N7<sub>b</sub> according to Table 1. Then golden thresholds are extracted. Each clip has four thresholds as shown in Figure 4. Hence the N10 dataset contains 3928 samples and each N7 dataset contains 3916 samples, respectively. The data augmentation technique in Section 3.1.2 is applied, so the training set and the testing set will be augmented by a factor of 8 independently. For example, if 50% of the data for N10

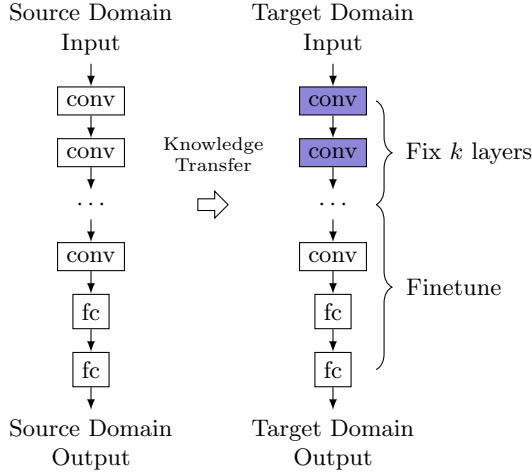


Figure 11: Transfer learning scheme with the first  $k$  layers fixed when training for target domain, denoted as  $TF_k$ .

are used for training, then there are  $3928 \times 50\% \times 8 = 15712$  samples. It needs to mention that always the same 50% portions are used during the validation of a dataset for fair comparison of different techniques. The batch size is set to 32 for training accommodating to the large variability in the sizes of training datasets. Adam [9] is used as the stochastic optimizer and maximum epoch is set to 200 for training.

The training time for one model takes 10 to 40 minutes according to the portions of a dataset used for training, and prediction time for an entire N10 or N7 dataset takes less than 10 seconds, while the rigorous simulation takes more than 15 hours for each N10 or N7 dataset. Thus we no longer report the prediction time which is negligible compared with that of the rigorous simulation. Each experiment runs 10 different random seeds and averages the numbers.

#### 4.1 CNN and ResNet

We first compare CNN and ResNet in Figure 12(a). Column “CNN-5” denotes the network with 5 layers shown in Figure 10(a). Column “CNN-10” denotes the one with 10 layers that has the same structure as that in Figure 10(b) but without shortcut connections. Column “ResNet” denotes the one with 10 layers shown in Figure 10(b). When using 1% to 20% training data, ResNet shows better average relative RMS error  $\epsilon_r^{th}$  than CNN-10, but CNN-5 provides the best error. We will show later that ResNet on the contrary outperforms CNN-5 when transfer learning is incorporated.

The impacts of depth on the performance of ResNet are further explored in Figure 12(b), where we gradually stack more building blocks in Figure 9 before fully connected layers. The x-axis denotes total number of convolution and fully connected layers corresponding to different numbers of building blocks. For instance, 0 building block leads to 4 layers and 3 building blocks result in 10 layers (Figure 10(b)). The testing error decreases to lowest value at 10 layers and then starts to increase, indicating potential overfitting afterwards [4]. Therefore, we use 10 layers for the ResNet in the experiment.

#### 4.2 Knowledge Transfer From N10 to N7

We then compare the testing accuracy between knowledge transfer from N10 to N7 and directly training from N7 datasets in Figure 13(a). In this example, the x-axis represents the percentage of training dataset for the target domain  $N7_a$ , while the percentage of data

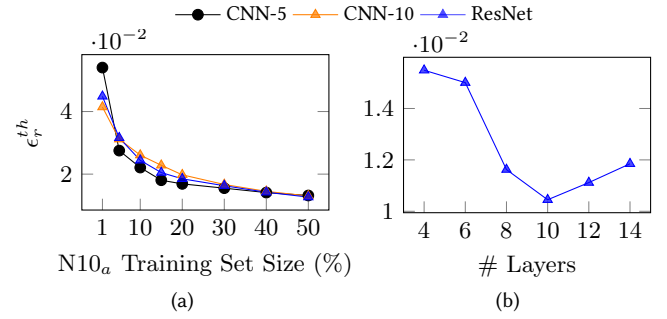


Figure 12: (a) Comparison on testing accuracy of CNN-5, CNN-10, and ResNet on N10. (b) Impact of depth on the testing accuracy of ResNet.

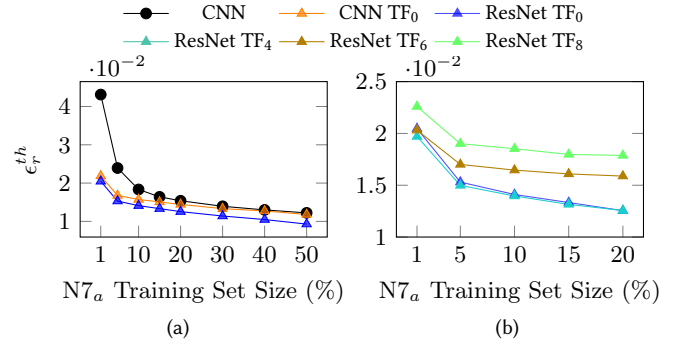


Figure 13: Testing accuracy of transfer learning from N10 to N7. (a) Comparison between CNN and transfer learning. (b) Comparison between transfer learning schemes where different numbers of layers are fixed.

from the source domain N10 is always 50%. Similar trends are also observed for  $N7_b$ . Curve “CNN” denotes training the CNN of 5 layers in Figure 10(a) with data from target domain only, i.e., no transfer learning involved. Curve “CNN  $TF_0$ ” denotes the transfer learning scheme in Section 3.4 for the same CNN with zero layer fixed. Curve “ResNet  $TF_0$ ” denotes applying the same scheme to ResNet. The most significant benefit of transfer learning comes from small training dataset with a range of 1% to 20%, where there are around 52% to 18% improvement in the accuracy from CNN. Meanwhile, ResNet  $TF_0$  can achieve an average of 13% smaller error than CNN  $TF_0$ .

Figure 13(b) further compares the results of fixing different numbers of layers during transfer learning. In this case, ResNet  $TF_0$  and ResNet  $TF_4$  have the best accuracy, while the error increases with more layers fixed. It is indicated that the tasks N10 and N7 are quite different and both feature extraction layers and regression layers need finetuning.

#### 4.3 Knowledge Transfer within N7

The transfer learning between different N7 datasets, e.g., from  $N7_a$  to  $N7_b$ , is also explored in Figure 14. The x-axis represents the percentage of training dataset for the target domain  $N7_b$ , while the percentage of data from the source domain  $N7_a$  is always 50%. Compared with the knowledge transfer from N10 to N7, we achieve even higher accuracy between 1% and 20% training datasets in Figure 14(a). For example, with 1% training dataset, there is around 65% improvement in accuracy from CNN, and with 20% training dataset, the improvement is around 23%. ResNet  $TF_0$  keeps having lower errors than that of CNN  $TF_0$  as well, with an average benefit around 15%.

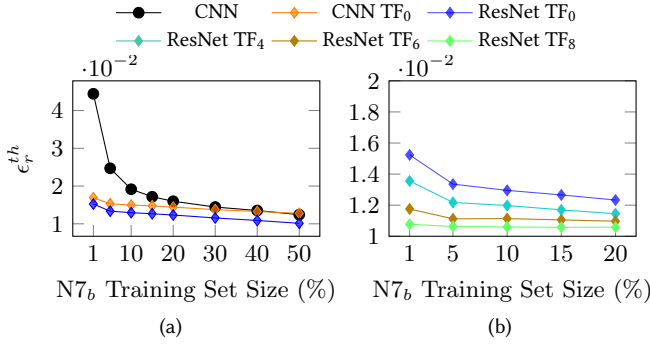


Figure 14: Testing accuracy of transfer learning from N7<sub>a</sub> to N7<sub>b</sub>. (a) Comparison between CNN and transfer learning. (b) Comparison between transfer learning schemes where different numbers of layers are fixed.

The curves in Figure 14(b) show different insights from that of the knowledge transfer from N10 to N7. The accuracy of ResNet TF<sub>0</sub> can be further improved with more layers fixed, e.g., ResNet TF<sub>8</sub>, by around 28% to 14%. This is reasonable since N7<sub>a</sub> and N7<sub>b</sub> have the same design rules and illumination shapes, and the only difference lies in the resist materials. Therefore, the feature extraction layers are supposed to remain almost the same. With the sizes of the training dataset increasing to 15% and 20%, the differences in the accuracy become smaller, because there are enough data to find good configurations for the networks.

#### 4.4 Impact of Various Source Domains

In transfer learning, the correlation between the datasets of source and target domains is critical to the effectiveness of knowledge transfer. Thus, we explore the impacts of source domain datasets on the accuracy of modeling for the target domain. Figure 15 plots the testing errors of learning N7<sub>b</sub> using ResNet TF<sub>0</sub> with various source domain datasets. Curves “N10<sup>50%</sup>” and “N7<sub>a</sub><sup>50%</sup>” indicate that 50% of the N10 or the N7<sub>a</sub> dataset is used to train source domain models, respectively. Curve “N10<sup>50%</sup> + N7<sub>a</sub><sup>1%</sup>” describes the situation where we have 50% of the N10 dataset and 1% of the N7<sub>a</sub> dataset for training. In this case, as shown in Figure 16, we first use the 50% N10 data to train the first source domain model; then train the second source domain model using the first model as the starting point with the 1% N7<sub>a</sub> data; in the end, the target domain model for N7<sub>b</sub> is trained using the second model as the starting point with N7<sub>b</sub> data. Curves “N10<sup>50%</sup> + N7<sub>a</sub><sup>5%</sup>” and “N10<sup>50%</sup> + N7<sub>a</sub><sup>10%</sup>” are similar, simply with different amounts of N7<sub>a</sub> data for training.

The knowledge from N7<sub>a</sub><sup>50%</sup> is the most effective for N7<sub>b</sub> due to the minor difference in resist materials between two datasets. For the rest curves, the accuracy of N10<sup>50%</sup> + N7<sub>a</sub><sup>5%</sup> and N10<sup>50%</sup> + N7<sub>a</sub><sup>10%</sup> is in general better than or at least comparable to that of N10<sup>50%</sup>. This indicates that having more data from closer datasets to the target dataset, e.g., N7<sub>a</sub>, is still helpful.

#### 4.5 Improvement in Data Efficiency

Table 2 presents the accuracy metrics, i.e., relative threshold RMS error ( $\epsilon_r^{th}$ ) and CD RMS error ( $\epsilon_r^{CD}$ ), for learning N7<sub>b</sub> from various source domain datasets. Since we consider the data efficiency of different learning schemes, we focus on the small training dataset for N7<sub>b</sub>, from 1% to 20%. Situations such as no source domain data (0), only source domain data from N10 (N10<sup>50%</sup>), only source domain

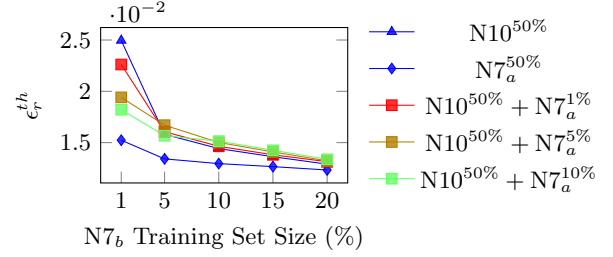


Figure 15: Testing accuracy of ResNet TF<sub>0</sub> for N7<sub>b</sub> from different source domain datasets.

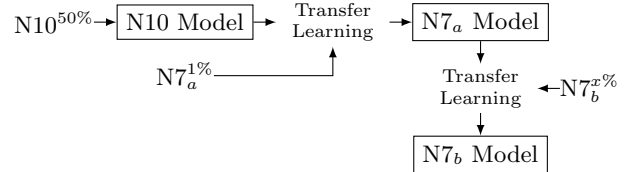


Figure 16: Transfer learning from 50% of N10 dataset and 1% of N7<sub>a</sub> dataset (i.e., N10<sup>50%</sup> + N7<sub>a</sub><sup>1%</sup>) to N7<sub>b</sub> with x% of N7<sub>b</sub> dataset.

data from N7<sub>a</sub> (N7<sub>a</sub><sup>50%</sup>), and combined source domain datasets, are examined. As mentioned in Section 2, the fidelity between relative threshold RMS error and CD RMS error is very consistent, so they share almost the same trends. Transfer learning with any source domain dataset enables an average improvement of 23% to 40% from that without knowledge transfer. In small training datasets of N7<sub>b</sub>, ResNet also achieves around 8% better performance on average than CNN in the transfer learning scheme. At 1% of N7<sub>b</sub>, combined source domain datasets have better performance compared with N10<sup>50%</sup> only, but the benefits vanish with the increase of the N7<sub>b</sub> dataset.

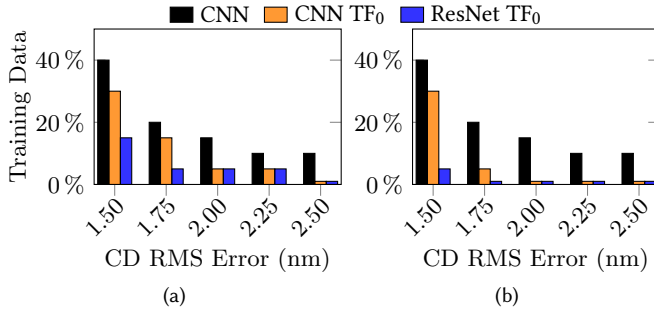
In real manufacturing, models are usually calibrated to satisfy a target accuracy or target CD RMS error. Figure 17 demonstrates the amount of training data required in the target domain for learning the N7<sub>b</sub> model. Curve “CNN” does not involve any knowledge transfer, while curves “CNN TF<sub>0</sub>” and “ResNet TF<sub>0</sub>” utilize transfer learning in CNN and ResNet, respectively. The curves in Fig 17(a) assume the availability of N10 data. Consider the CD RMS error from 1.5nm to 2.5nm, which is around 10% of the half pitch for N7 contacts. This range of accuracy is also comparable to that of the state-of-the-art CNN [25]. ResNet TF<sub>0</sub> requires significantly fewer data than both CNN and CNN TF<sub>0</sub>. For instance, when the target CD error is 1.75nm, ResNet TF<sub>0</sub> demands 5% training data from N7<sub>b</sub>, while CNN requires 20% and CNN TF<sub>0</sub> requires 15%. Figure 17(b) considers the transfer from N7<sub>a</sub> to N7<sub>b</sub>. Both ResNet TF<sub>0</sub> and CNN TF<sub>0</sub> only require 1% training data from N7<sub>b</sub> for most target CD RMS errors, where CNN TF<sub>0</sub> cannot achieve the accuracy unless given 30% data. Overall, ResNet TF<sub>0</sub> can achieve 2-10X reduction of training data within this range compared with CNN. It needs to mention that 1% of dataset only correspond to fewer than 40 samples owing to the data augmentation, indicating only thresholds of 40 clips are required.

## 5 CONCLUSION

A transfer learning framework based on residual neural networks is proposed for resist modeling. The combination of ResNet and transfer learning is able to achieve high accuracy with very few data from the target domains, under various situations for knowledge transfer, indicating high data efficiency. Extensive experiments demonstrate that the proposed techniques can achieve 2-10X reduction according to various requirements of accuracy comparable to the state-of-the-art

**Table 2: Relative Threshold RMS Error and CD RMS Error for  $N7_b$  with Different Source Domain Datasets**

Source Datasets		$\emptyset$		$N10^{50\%}$				$N7^{50\%}_a$				$N10^{50\%} + N7^{5\%}_a$		$N10^{50\%} + N7^{10\%}_a$	
Neural Networks		CNN		CNN TF <sub>0</sub>		ResNet TF <sub>0</sub>		CNN TF <sub>0</sub>		ResNet TF <sub>0</sub>		ResNet TF <sub>0</sub>		ResNet TF <sub>0</sub>	
		$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$	$\epsilon_r^{th}$ ( $10^{-2}$ )	$\epsilon^{CD}$
$N7_b$	1%	4.44	4.76	2.34	2.48	2.29	2.39	1.69	1.79	1.52	1.60	1.94	2.03	1.82	1.91
	5%	2.78	2.96	1.73	1.86	1.60	1.70	1.53	1.64	1.34	1.43	1.67	1.78	1.57	1.67
	10%	1.92	2.04	1.63	1.76	1.47	1.57	1.50	1.60	1.30	1.38	1.50	1.60	1.51	1.61
	15%	1.72	1.84	1.56	1.68	1.39	1.47	1.48	1.55	1.27	1.35	1.41	1.50	1.43	1.52
	20%	1.60	1.71	1.50	1.61	1.31	1.39	1.44	1.55	1.23	1.31	1.32	1.41	1.34	1.43
ratio		1.00	1.00	0.77	0.77	0.70	0.69	0.69	0.69	0.60	0.60	0.69	0.69	0.69	0.68

**Figure 17: Amount of training data required for  $N7_b$  given target CD RMS errors when (a) 50%  $N10$  dataset is available or (b) 50%  $N7_a$  dataset is available.**

learning approach. It is also shown that the performance of transfer learning differs from dataset to dataset and is worth exploring to see the correlation between datasets. Examining the quantitative relation between the correlation of datasets and performance of transfer learning is valuable in the future.

## ACKNOWLEDGE

This project is supported in part by Toshiba Memory Corporation, NSF, and the University Graduate Continuing Fellowship from the University of Texas at Austin. The authors would like to thank Memory Lithography Group from Toshiba Memory Corporation for helpful discussions and feedback.

## REFERENCES

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, et al. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. (2015). <https://www.tensorflow.org>
- [2] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks* 5, 2 (1994), 157–166.
- [3] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. 249–256.
- [4] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT press.
- [5] Allan Gu and Avideh Zakhori. 2008. Optical proximity correction with linear regression. *IEEE Transactions on Semiconductor Manufacturing (TSM)* 21, 2 (2008), 263–271.
- [6] Josiah P Hanna and Peter Stone. 2017. Grounded Action Transformation for Robot Learning in Simulation. In *AAAI*. 3834–3840.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [8] Ningning Jia and Edmund Y Lam. 2010. Machine learning for inverse lithography: using stochastic gradient descent for robust photomask synthesis. *Journal of Optics* 12, 4 (2010), 045601.
- [9] Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [10] Lars Liebmann, Albert Chu, and Paul Gutwin. 2015. The daunting complexity of scaling to 7nm without EUV: Pushing DTCO to the extreme. In *Proceedings of SPIE*, Vol. 9427.
- [11] Lars Liebmann, Jia Zeng, Xuelian Zhu, Lei Yuan, Guillaume Bouche, and Jongwook Kye. 2016. Overcoming scaling barriers through design technology CoOptimization. In *VLSI Technology, 2016 IEEE Symposium on*. IEEE, 1–2.
- [12] Lars W Liebmann, Scott M Mansfield, Alfred K Wong, Mark A Lavin, William C Leipold, and Timothy G Dunham. 2001. TCAD development for lithography resolution enhancement. *IBM Journal of Research and Development* 45, 5 (2001), 651–665.
- [13] Yibo Lin, Xiaoqing Xu, Jiaojiao Ou, and David Z Pan. 2017. Machine learning for mask/wafer hotspot detection and mask synthesis. In *Photomask Technology*, Vol. 10451. International Society for Optics and Photonics, 104510A.
- [14] Rui Luo. 2013. Optical proximity correction using a multilayer perceptron neural network. *Journal of Optics* 15, 7 (2013), 075708.
- [15] Xu Ma, Xuejiao Zhao, Zhiqiang Wang, Yanqiu Li, Shengjie Zhao, and Lu Zhang. 2017. Fast lithography aerial image calculation method based on machine learning. *Applied Optics* 56, 23 (2017), 6485–6495.
- [16] Tetsuaki Matsunawa, Shigeki Nojima, and Toshiya Kotani. 2016. Automatic Layout Feature Extraction for Lithography Hotspot Detection Based on Deep Neural Network. In *Proceedings of SPIE*.
- [17] Tetsuaki Matsunawa, Bei Yu, and David Z Pan. 2016. Optical proximity correction with hierarchical Bayes model. *Journal of Micro/Nanolithography, MEMS, and MOEMS* 15, 2 (2016), 021009–021009.
- [18] Mentor Graphics. 2008. Calibre Verification User's Manual. (2008).
- [19] Sinno Jialin Pan and Qiang Yang. 2010. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22, 10 (2010), 1345–1359.
- [20] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. 2016. Progressive neural networks. *arXiv preprint arXiv:1606.04671* (2016).
- [21] Seongbo Shim, Suhyeon Choi, and Youngsoo Shin. [n. d.]. Machine Learning-Based Resist 3D Model. In *Proc. of SPIE Vol.*, Vol. 10147. 101471D–1.
- [22] Moojoon Shin and Jee-Hyong Lee. 2016. Accurate Lithography Hotspot Detection Using Deep Convolutional Neural Networks. In *Journal of Micro/Nanolithography, MEMS, and MOEMS (JM3)*.
- [23] Synopsys. 2016. Sentaurus Lithography. <https://www.synopsys.com/silicon/mask-synthesis/sentaurus-lithography.html>. (2016).
- [24] Chin Boon Tan, Kar Kit Koh, Dongqing Zhang, and Yee Mei Foong. 2015. Sub-resolution assist feature (SRAF) printing prediction using logistic regression. In *Proceedings of SPIE*. 94261Y–94261Y.
- [25] Yuki Watanabe, Taiki Kimura, Tetsuaki Matsunawa, and Shigeki Nojima. 2017. Accurate lithography simulation model based on convolutional neural networks. In *SPIE Advanced Lithography*. International Society for Optics and Photonics, 101470K–101470K.
- [26] Jen-Yi Wu, Fedor G Pikus, and Malgorzata Marek-Sadowska. 2011. Efficient approach to early detection of lithographic hotspots using machine learning systems and pattern matching. In *SPIE Advanced Lithography*. International Society for Optics and Photonics, 79740U–79740U.
- [27] Xiaoqing Xu, Yibo Lin, Meng Li, Tetsuaki Matsunawa, Shigeki Nojima, Chikaaki Kodama, Toshiya Kotani, and David Z. Pan. 2017. Sub-Resolution Assist Feature Generation with Supervised Data Learning. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* PP, 99 (2017).
- [28] Haoyu Yang, Yajun Lin, Bei Yu, and F.Y. Evangeline Young. 2017. Lithography Hotspot Detection: From Shallow to Deep Learning. In *IEEE International System-on-Chip Conference (SOCC)*.
- [29] Haoyu Yang, Jing Su, Yi Zou, Bei Yu, and F.Y. Evangeline Young. 2017. Layout Hotspot Detection with Feature Tensor Generation and Deep Biased Learning. In *ACM/IEEE Design Automation Conference (DAC)*.