

# 칩 간 공간적 유사성과 검사항목의 상관관계를 고려한 반도체 웨이퍼 테스트 데이터의 결측치 대체 방법 개발

김주영 · 배영목 · 최승현 · 김광재<sup>†</sup>

포항공과대학교 산업경영공학과

## Development of a Missing Value Imputation Method for Semiconductor Wafer Test Data Considering Spatial Similarity among Chips and Correlation between Test Items

Juyeong Kim · Youngmok Bae · Seunghyun Choi · Kwangjae Kim

Department of Industrial and Management Engineering, Pohang University of Science and Technology

In semiconductor manufacturing, each wafer comprises multiple chips, and each chip is tested before the packaging process. Wafer test data on electrical characteristics of chips are collected during the wafer test process. However, missing values often occur due to various manufacturing environments. In this study, a new missing value imputation method based on Generative Adversarial Imputation Nets (GAIN) is proposed. The proposed method takes into account the two characteristics of wafer test data, namely, spatial similarity among chips and test item correlation. Spatial similarity refers to the property of having similar test item values between chips in adjacent or symmetrical positions. Test item correlation refers to the positive correlation between test items with similar physical properties. Spatial similarity and test item correlation are reflected by the addition of locational information of chips and modification of the loss function in GAIN, respectively. The performance of the proposed method is validated with a real wafer test dataset by a comparison with those of existing methods in various circumstances.

**Keywords:** Semiconductor manufacturing, Wafer test, Missing Value Imputation, Artificial intelligence

### 1. 서 론

반도체 시장은 웨어러블 디바이스, PC, 스마트폰, 자동차용 반도체의 수요 확대와 인공지능, 빅데이터 기술 활용의 증가로 꾸준히 성장하고 있으며(Kim *et al.*, 2022), 첨단 메모리 기술과 설계 소프트웨어 등의 3세대 반도체 개발이 본격화됨에 따라 기술 패권 경쟁이 치열해지고 있다(Yeon, 2021; Cho, 2020). 반도체 제조 공정은 복잡한 구조로 이루어져 있기 때문에 반도체 제조 공정의 생산 효율성을 증대하여 고품질의 제품을 양산하는 것은 반도체 시장의 경쟁우위 확보에 매우 중요한 요

소이다. 이를 위해, 반도체 제조 기업들은 반도체의 불량 여부를 조기에 탐지하고, 그 원인을 파악하여 반도체 수율을 개선하고 있다(Kang and Baek, 2020).

반도체의 불량 여부는 생산 공정이 끝난 후, 웨이퍼 테스트 공정을 통해 판별된다. 웨이퍼 테스트 공정은 반도체의 일반적인 사용 환경보다 강한 외부 자극을 웨이퍼에 가하여 칩 단위의 전기적 검사항목을 측정한다(Huang *et al.*, 2013). 예를 들어, Wafer Burn In 항목은 웨이퍼에 일정 온도의 열을 전도한 다음 교류 및 직류 전압을 가하여 스트레스를 주는 방식으로 측정한다. 웨이퍼 테스트 데이터는 웨이퍼 테스트 공정에서

이 논문은 2022년도 정부의 재원으로 한국연구재단의 지원을 받아 수행된 연구임(No.2022R1A2C1004614, No.2022K2A9A2A11097291).

<sup>†</sup> 연락저자 : 김광재 교수, 경상북도 포항시 남구 청암로 77, Tel : 054-279-8249, Fax : 054-279-2870, E-mail : kjk@postech.ac.kr

2023년 1월 26일 접수; 2023년 4월 12일 수정본 접수; 2023년 4월 17일 게재 확정.

측정된 다양한 검사항목 계측치를 의미하며, 이를 분석하여 반도체의 품질 수준이 고객의 기대 수준에 부합하는지 확인할 수 있다. 예를 들어, 불량률 확률이 높은 칩을 패키징 과정 전에 선별하여 시장 판매 이후 발생하는 품질 이슈를 최소화하거나(Park and Kim, 2014), 검사항목 간 연계 분석을 통해 불량과 양품의 경계에 있는 칩을 수선할 수 있는 근거를 도출한다(Han et al., 2005).

그러나 웨이퍼 테스트 데이터는 테스트 장비 센서의 결합, 선행 검사항목의 불량 판정, 데이터의 통합 과정에서의 누락, 이상치 발생 등의 이유로 결측치가 발생할 수 있다(Schranner et al., 2019). 결측치로 인한 정보 손실은 반도체 불량 여부 판단의 정확성을 저해하고(Zhang et al., 2022), 불량 칩의 발생원인 파악에 악영향을 끼친다. 따라서 반도체 불량 여부 탐지의 신뢰성을 확보하기 위해서는 웨이퍼 테스트 데이터의 결측치 대체가 필수적이다.

결측치를 효과적으로 대체하기 위해서는 결측치 대체 대상의 특성을 고려해야 한다(Aittokallio et al., 2010). 반도체의 불량 분석은 웨이퍼 단위에서 이루어지므로(Jung and Jung, 2022), 웨이퍼 단위에서 발견되는 특성을 고려한다면 효과적인 결측치 대체가 가능할 것이다. 웨이퍼 단위에서 발견되는 데이터의 특성에는 칩 간 공간적 유사성(이하 공간적 유사성)과 검사항목 간 상관관계가 있으며, 두 특성은 반도체 수율에 영향을 미친다고 알려져 있다(Hsu et al., 2013). 공간적 유사성은 공간적 군집성과 공간적 대칭성으로 구분된다. <Figure 1>은 공간적 유사성의 개념적 예시를 나타내며, 칩의 특정 검사항목 계측치(Test item values)가 클수록 진한 색으로 표현하였다. <Figure 1(a)>는 공간적으로 인접한 칩들이 유사한 계측치를 가지는 공간적 군집성(Spatial Cluster)을, <Figure 1(b)>는 웨이퍼의 중심을 기준으로 대칭 위치의 칩들이 유사한 계측치를 가지는 공간적 대칭성(Spatial Symmetry)을 나타낸다. 검사항목 간 상관관계는 물리적 특성이 유사한 검사항목 간 나타나는 강한 상관관계를 의미하며, 동일한 전기적 환경에서 검사항목들을 반복적으로 측정하는 과정에서 발생한다(Schranner et al., 2019).

웨이퍼 테스트 데이터의 결측치 대체에 다양한 방법이 활용되고 있으나(Lee et al., 2022), 웨이퍼 단위에서 발생하는 공간적 유사성과 검사항목 간 상관관계를 고려한 결측치 대체 방

법에 대한 연구는 아직 미비한 실정이다. 본 연구는 반도체 테스트 데이터의 특성인 공간적 유사성과 검사항목 간 상관관계를 반영한 결측치 대체 방법을 제안한다. 제안 방법은 Generative Adversarial Imputation Nets(GAIN, Yoon et al., 2018)의 구조를 차용하였으며, 웨이퍼 테스트 데이터에 위치 데이터를 추가하여 공간적 유사성을, GAIN의 손실 함수에 결측치 대체 전과 후의 상관관계 분포의 차이를 추가하여 검사항목 간 상관관계가 결측치 대체 과정에서 유지될 수 있도록 하였다. 실제 반도체 회사의 웨이퍼 테스트 데이터를 활용하여 다양한 결측 환경에서 제안 방법의 결측치 대체 성능을 검증하였다. 결측률 10%에서 50%까지 10%p 단위에서 검증한 결과, 대체를 수행한 검사항목 총 22개 중 18개 이상에서 제안 방법의 성능이 대조 방법 대비 우수하였다. 제안 방법은 웨이퍼 테스트 데이터 분석의 효과를 제고하여 반도체 수율 향상에 기여할 것으로 기대된다.

본 논문의 구조는 다음과 같다. 2절에서는 기존 결측치 대체 방법을 공간적 유사성과 검사항목 간 상관관계 측면에서 분석한다. 3절에서는 제안 방법의 구조를 설명한다. 4절에서는 다양한 실험 환경에서 제안 방법의 성능을 검증한다. 5절에서는 제안 방법의 구성 요소에 따라 대체 성능을 분석하였다. 마지막으로 6절에서는 본 연구의 기대효과를 설명하고 향후 연구를 제안한다.

## 2. 문헌 리뷰

본 절은 기존의 결측치 대체 방법을 웨이퍼 테스트 데이터의 두 가지 특성의 관점에서 리뷰하였다. 기존의 방법은 전통적 대체 방법과 다중대체(Ko et al., 2014), 그리고 복잡한 데이터 특성을 고려한 머신러닝 기반 대체 방법으로 분류할 수 있다(Lee et al., 2022). 전통적 대체 방법은 결측치가 존재하는 관측치를 일괄 제거하는 완전 제거법과 각 변수에 대한 대평균(예: 평균, 중앙값) 또는 회귀모형의 예측값으로 결측치를 대체하는 단일 대체법 등으로 구성된다(Ko et al., 2014). 전통적 대체 방법은 대체 방법의 원리가 간단하지만, 결측치를 변수의 대평균으로 대체하거나 분석 대상에서 제외하므로 공간적 유사성과 검사항목 간 상관관계를 모두 반영할 수 없어 웨이퍼 테스트 데이터의 결측치 대체에 활용하기 부적합하다.

다중대체 방법은 전통적 대체 방법이 가지는 대체 편향을 개선하기 위해 결측치를 처리한 복수의 데이터를 통합하여 결측치를 대체하는 방법이다. 대표적으로 Multiple Imputation by Chained Equations(MICE, Van et al., 2011)가 있다. MICE는 대평균으로 결측치를 대체하고 예측 모형을 통해 반복적으로 결측치를 대체하는 방법이다. 반복 대체 과정에서, 이전에 대체된 값이 다음에 대체될 값의 초기값이 되는 연쇄 방정식을 이용하며, 예측 모형의 모수가 충분히 수렴하는 것을 목표로 한다. 다중대체 방법은 복수의 데이터를 통합하는 과정에서 검

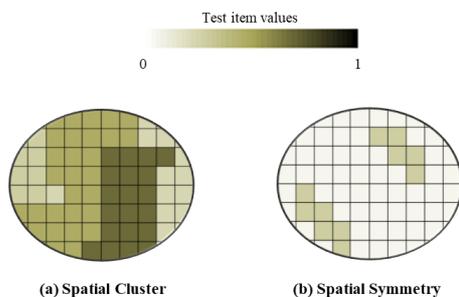


Figure 1. Spatial Similarity of Wafer Test Data

사항목 간 상관관계를 반영할 수 있으나, 예측 모형을 학습하는 과정에서 웨이퍼 별로 발생하는 공간적 유사성을 반영할 수 없다는 한계점을 가진다.

머신러닝 기반 대체 방법은 데이터 중 결측치가 부재한 데이터를 활용하여 학습된 알고리즘으로 결측치를 대체하는 방법으로, 대량의 데이터를 활용하고 복잡한 연산을 수행한다는 점에서 다중대체와 차이가 있다. 머신러닝 기반 대체 방법은 Missforest(Stekhoven *et al.*, 2012), kNN(Troyanskaya *et al.*, 2001) 등이 있다. Missforest는 각 변수를 단일 대체한 후, 랜덤 포레스트 모형을 학습시켜 결측치를 대체한다. kNN은 결측치를 포함하는 관측치와 인접한  $k$ 개 관측치에 대하여, 결측치가 발생한 검사항목의 평균값으로 대체한다. 그러나 머신러닝 기반 대체 방법은 다중대체 방법과 동일한 이유로 검사항목 간 상관관계는 반영할 수 있으나, 공간적 유사성을 반영하기 어렵다는 특징이 있다.

머신러닝 기반 대체 방법 중 딥러닝 기반의 방법에는 대표적으로 GAIN이 있다. GAIN은 계측 데이터에 대한 결합확률분포를 추정하여 결측치를 대체하며, 다양한 분야에서 결측치 대체 성능이 검증된 방법이다(Luo *et al.*, 2022). 그러나 데이터를 일괄 대체하고 칩별 위치 정보가 부재한 상태로 학습이 진행되어 공간적 유사성을 결측치 대체 과정에 반영하기 어려우며, 학습 과정에서 데이터 대체 정확성만이 고려되므로(Han *et al.*, 2019) 검사항목 간 상관관계를 반영하지 못하여 웨이퍼 테스트 데이터의 대체에 적용하기 부적합하다.

한편, GAIN은 데이터의 특성에 따라 결합확률분포를 추정하는 모형을 변형하여 활용할 수 있다. 예를 들어, 시공간 데이터의 공간적 특성을 반영하기 위해 네트워크를 구성하는 레이어를 Convolution 레이어로 변형하거나(Adeli *et al.*, 2021), 명목형 변수를 반영하기 위해 GAIN과 분류 모형을 결합하거나(Wang *et al.*, 2020), 데이터의 시계열적 특성을 반영하기 위해 검증 데이터에 대한 GAIN의 평균 절대 오차를 손실함수에 반영하기도 한다(Lee *et al.*, 2022). 본 연구는 이와 같은 GAIN의 유연성에 착안하여 웨이퍼 테스트 데이터의 결측치 대체에 적합하도록 GAIN을 개선하여 공간적 유사성과 검사항목 간 상관관계를 고려한 웨이퍼 테스트 데이터 결측치 대체 방법을 개발하고자 한다.

### 3. 제안 방법론

본 절에서는 GAIN에 칩 간 공간적 유사성과 검사항목 간 상관관계를 반영한 제안 방법의 구조에 관하여 설명한다. 3.1절에서는 GAIN의 결측치 대체 원리를, 3.2절에서는 GAIN을 기반으로 제안 방법의 구조를 설명한다.

#### 3.1 GAIN

GAIN은 딥러닝 생성모형 GAN 기반의 결측치 대체 방법으로,

계측된 데이터에 임의로 결측치( $x$ )를 발생시켜 대체 값과 계측치의 차이를 최소화하는 방향으로 손실함수를 정의하여 네트워크를 학습한다.<Figure 2>는 GAIN의 구조를 웨이퍼 테스트 데이터에 적용하여 설명한다. 이때, 웨이퍼 테스트 데이터는 하나의 웨이퍼에 대한 칩(Chip) 별 검사항목(Test Item)으로 구성되며, 결측치를 대체하는 생성자(Generator;  $G$ )와 대체 값과 계측치를 예측하는 구별자(Discriminator;  $G$ )로 구성된다.

생성자는 웨이퍼 테스트 데이터에 임의로 결측치를 발생시킨 계측 행렬(Wafer Test Data with Missing Values ;  $X$ )과 계측 행렬의 값이 결측이면 0, 결측이 아니면 1로 표기하는 Mask 행렬(Mask Matrix;  $M$ )을 입력받아 결측치가 대체된 데이터(Imputed Data;  $X_G$ )를 출력한다. 위 과정은  $X$ 의 결측치를 다변량 정규분포로부터 추출된 랜덤 노이즈로 임의 대체한 후 진행된다. 구별자는  $X_G$ 를 입력받아 계측치가 기존에 결측치였는지 생성자로부터 대체된 것인지 분류하는 행렬인 Estimated Mask 행렬(Estimated Mask Matrix;  $M_D$ )을 출력한다. 생성자는 구별자의 분류 성능을 하락시키는 방향으로, 구별자는 생성자의 출력을 정확히 분류하는 방향으로 순차적으로 학습된다. 이 과정에서 Hint Generator는 Mask 행렬의 일부를 제공하는 Hint 행렬(Hint Matrix;  $H$ )을 도출하며,  $H$ 에서 제공되지 않는 계측치는 모두 0.5로 제공한다. 구별자에 대한  $H$ 의 불완전한 정보 제공은 구별자의 학습 속도를 생성자의 학습 속도와 맞추게 함으로써 이들 간의 균형 학습을 지원한다.

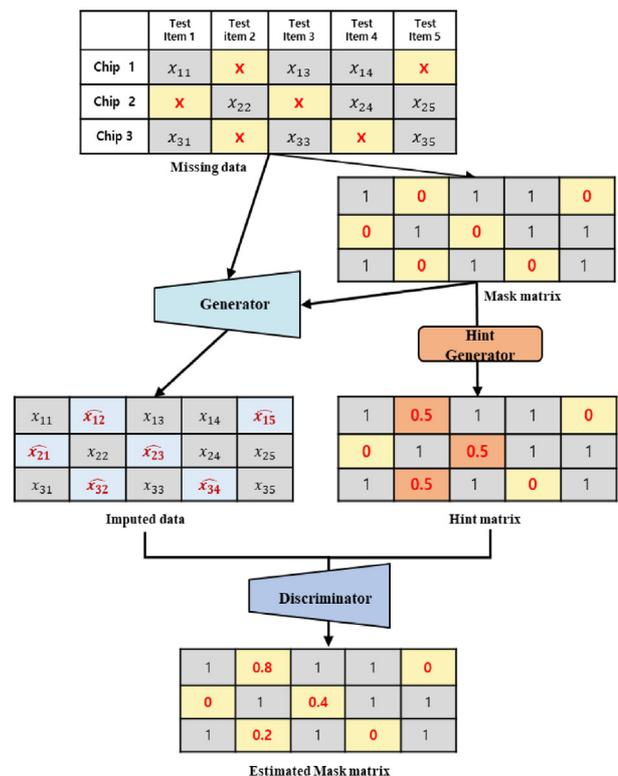


Figure 2. Overview of Generative Adversarial Imputation Nets (GAIN)

각 네트워크에 대한 손실함수는 식 (1)과 (2)이며,  $N$ 개의 칩과  $d$ 개의 검사항목에 대한 계측 행렬  $X = \{x^1, x^2, \dots, x^N\} \in \mathbb{R}^{N \times d}$ 과 Mask 행렬  $M = \{m^1, m^2, \dots, m^N\} \in \mathbb{R}^{N \times d}$ 로 정의한다. 이때,  $x^k$ 는  $X$ 의  $k$ 번째 칩의 벡터로, 각 벡터는  $x^k = \{x_{11}^k, x_{21}^k, \dots, x_{d1}^k\} \in \mathbb{R}^{1 \times d}$ 로 구성된다.  $x_{R1}^k$ 는 결측치가 대체된  $x^k$ 로,  $x_{R1}^k = \{x_{R,1}^k, x_{R,2}^k, \dots, x_{R,d}^k\} \in \mathbb{R}^{1 \times d}$ 로 구성된다.  $m^k$ 는  $M$ 의  $k$ 번째 칩의 벡터,  $m_D^k$ 는  $M_D$ 의  $k$ 번째 칩의 벡터로 정의한다.

$$\min_G \frac{1}{N} \sum_{k=1}^N L_G(m^k, m_D^k) + \alpha L_R(x^k, x_{R1}^k)$$

$$\text{where } L_G(m, m_D) = -(1-m) \log(m_D), L_R(x_i^k, x_{R,i}^k) = (x_i^k - x_{R,i}^k)^2$$

$$\min_D \frac{1}{N} \sum_{k=1}^N L_D(m^k, m_D^k),$$

$$\text{where } L_D(m, m_D) = -m \log(m_D) - (1-m) \log(1-m_D)$$

$L_G$ 와  $L_D$ 는 각각  $M$ 과  $M_D$ 의 차이를 최소화하는 생성자와 구별자의 손실함수로서, 생성자에서는 구별자의 분류 성능이 하락하도록, 구별자에서는 생성자의 결측치에 대한 예측 성능이 하락하도록 설정한다. 생성자의 손실함수인 식 (1)은 구별자가 추정한 계측치의 결측 유무에 대한 정확도인  $L_G$ 와 생성자가 추정한 대체 값의 정확도인  $L_R$ 의 합으로 구성되며,  $L_R$ 은 하이퍼 파라미터  $\alpha$ 의 가중치를 통해 생성자 학습 과정에 관여함으로써 생성자가 계측치와 유사한 값을 도출하도록 네트워크를 학습한다. 구별자의 손실함수인 식 (2)는 구별자가 예측한 계측치의 결측 유무인  $M_D$ 를 실제 계측치의 결측 유무인  $M$ 에 가깝도록 네트워크를 학습한다.

### 3.2. 제안 방법

제안 방법은 공간적 유사성과 검사항목 간 상관관계를 고려하여 GAIN을 개선하였으며, <Figure 3>은 제안 방법의 구조를 웨이퍼 테스트 데이터에 적용하여 설명한다.

공간적 유사성을 반영하기 위해 웨이퍼 테스트 데이터에 칩의 위치 정보를 추가한다. 위치 정보는 각 칩의 직교 좌표(Orthogonal), 극좌표(Polar) 그리고 대칭 좌표(Symmetrical)로 구성되며, <Figure 3>의 A에서 확인할 수 있다. 이때, 대칭 좌표는 직교 좌표를 -1에서 1 사이로 정규화한 좌표를 의미한다. 위치 좌표들은 결측치 대체 과정에서 데이터의 변수로 추가되며, 생성자가 웨이퍼 별 결측치를 대체할 때 다른 칩의 공간 정보와 계측치를 함께 반영할 수 있도록 한다. 이를 통해, 위치 좌표들은 웨이퍼 내 칩 간 발생하는 계측치의 유사성 추정을 위한 생성자 학습 과정에 관여한다. 직교 좌표를 통해 공간적 근접성을, 극좌표와 대칭 좌표를 통해 공간적 대칭성을 반영한다. 또한, 웨이퍼 단위로 결측치를 대체하여 웨이퍼 내에서 관측될 수 있는 공간적 유사성을 반영한다.

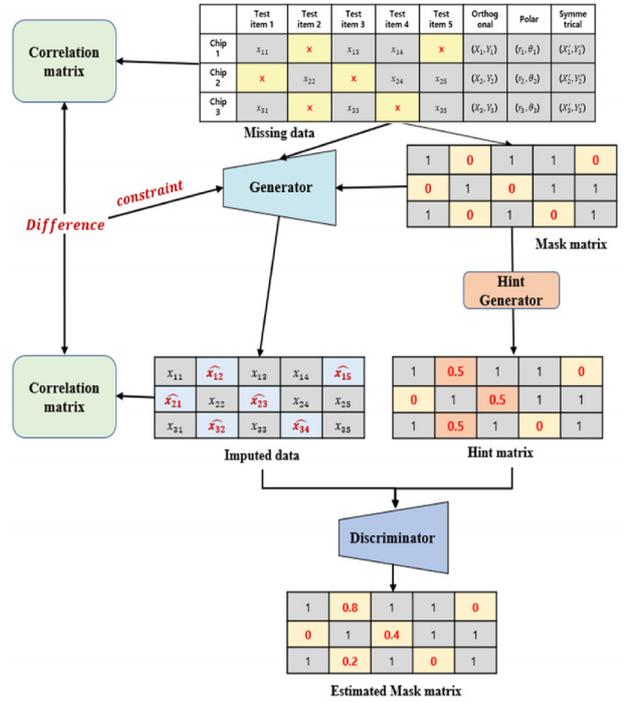


Figure 3. Overview of the Proposed Method

검사항목의 상관관계를 반영하기 위해 생성자에 결측치 대체 전후의 상관관계가 보존될 수 있도록 손실함수를 수정한다. 손실함수에 상관관계의 항을 추가하여 변수 간 상관관계를 네트워크 학습 과정에서 강화할 수 있으므로(Chander *et al.*, 2016), 결측치 대체 전의 상관관계( $Cor_0$ )와 대체 후의 상관관계( $Cor_1$ )를 각각 계산한 뒤, 상관관계의 거리를 Jensen-Shannon Divergence( $JSD$ )를 통해 계산하여 생성자의 손실함수에 추가한다.  $Cor_0$ 와  $Cor_1$ 의 정의는 식 (3),  $JSD$ 의 정의는 식 (4)와 같다.  $JSD$ 는 근사된 분포의 정보 손실량을 측정하는 Kullback-Leibler Divergence를 기반으로 두 확률 분포 간 차이를 측정하는 척도이다.

$$Cor_0 = \begin{pmatrix} cor(x_1, x_1) & \dots & cor(x_1, x_d) \\ \dots & \ddots & \dots \\ cor(x_d, x_1) & \dots & cor(x_d, x_d) \end{pmatrix},$$

$$Cor_1 = \begin{pmatrix} cor(x_{R,1}, x_{R,1}) & \dots & cor(x_{R,1}, x_{R,d}) \\ \dots & \ddots & \dots \\ cor(x_{R,d}, x_{R,1}) & \dots & cor(x_{R,d}, x_{R,d}) \end{pmatrix}$$

$$JSD(P_0, P_1) = \frac{1}{2} D_{KL}(P_0 \parallel \frac{P_0 + P_1}{2}) + \frac{1}{2} D_{KL}(P_1 \parallel \frac{P_0 + P_1}{2}),$$

$$\text{where } D_{KL}(X \parallel Y) = \sum_{i=1}^I x_i \log \left( \frac{x_i}{y_i} \right),$$

$$X = (x_1, x_2, \dots, x_l), Y = (y_1, y_2, \dots, y_l),$$

$$P_0 = (cor(x_1, x_1), \dots, cor(x_d, x_d))$$

$$P_1 = \begin{matrix} \in \mathbb{R}^{1 \times \binom{d(d+3)}{2}}, \\ (cor(x_{R,1}, x_{R,1}), \dots, cor(x_{R,d}, x_{R,d})) \\ \in \mathbb{R}^{1 \times \binom{d(d+3)}{2}} \end{matrix}$$

최종적으로 수정된 손실함수는 식 (5)와 같으며, <Figure 3>의 B에서 확인할 수 있다. 생성자의 결측치 대체 과정에서 JSD로 추정된 상관관계 행렬 분포 간 거리가 최소화될 수 있도록 하였으며, 하이퍼 파라미터  $\beta$ 로 구별자에 대한 상관관계 항의 영향력을 조절한다.

$$\min_G \frac{1}{N} \sum_{k=1}^N L_D(m^k, m_D^k) + \alpha L_R(x^k, x_R^k) + \beta JSD(Cor_0, Cor_1)$$

### 4. 실험

본 절에서는 제안 방법의 성능을 다양한 대체 방법과 비교하여 검증한다. 4.1절에서는 실험 데이터, 실험 절차 그리고 성능 지표에 대해 설명한다. 4.2절에서는 실험 결과를 설명한다. 4.3절에서는 결측치가 대체된 웨이퍼 테스트 데이터를 활용하여 반도체 칩의 불량 여부를 예측하는 실험을 통해 제안 방법의 실용성을 검증한다.

#### 4.1 실험 환경 설정

실험에는 국내 반도체 생산 기업에서 수집된 웨이퍼 테스트 데이터를 활용하였으며, 60장의 웨이퍼와 22개의 검사항목으로 구성된다. 실험 절차는 다음과 같다. 먼저, 웨이퍼 테스트 데이터를 7대 3의 비율로 학습과 검증 데이터로 분리하였다. 다음으로, 학습과 검증 데이터에 결측치를 각각 생성하였다. 웨이퍼 테스트 과정에서 특정 검사항목의 결측은 다른 검사항목의 결측과 연관이 있는 것으로 알려져 있다. 따라서 결측치의 발생이 변수와는 관련이 있으나, 데이터 분석을 통해 얻고자 하는 예측값과는 관련이 없는 무작위 결측 유형인 MAR (Missing At Random) 유형으로 결측치를 생성하였다. 실제 웨이퍼 테스트 데이터에 나타나는 결측률을 참고하여 10%에서 50%까지 10%p 단위로 임의 결측치를 생성하였다. 실험 결과의 신뢰성을 확보하기 위해 결측률마다 30개의 임의 결측 데이터를 생성하였다. 결측 데이터에 제안 방법과 대조 방법을 적용하여 웨이퍼 단위로 결측치를 대체하였으며, 0.001에서 1000 사이의 수들을 무작위 샘플링하여  $\alpha$ 와  $\beta$ 의 최적 조합을 모색하였다. 대조 방법으로는 GAIN, Missforest, MICE, kNN, 그리고 평균 대체법(Mean)을 활용하였다. 이후, 결측치 대체 방법별로 대체된 데이터의 성능을 Normalized Root Mean Square Error(NRMSE)지표를 통해 측정하여 대체 방법 간 성능을 비교하였다. NRMSE의 정의는 식 (6)과 같다.

$$NRMSE_k = \frac{\|x_{true}^k - x_{impute}^k\|_2}{\sigma_{true}^k}$$

$x_{true}^k$ 와  $x_{impute}^k$ 는 각각 결측치 대체 전의 k번째 검사항목 벡터와 결측치 대체 후의 k번째 검사항목 벡터를,  $\sigma_{true}^k$ 는 결측치 대체 전의 k번째 검사항목의 표준편차를 의미한다. 즉,  $NRMSE_k$ 는 k번째 검사항목의 계측치의 표준편차 대비 결측치의 Root Mean Square Error의 비율이며 0에 가까울수록 좋은 성능을 보인다고 해석할 수 있다.

#### 4.2 제안 방법의 성능 평가

첫 번째 성능 평가는 결측률 증가에 따른 결측치 대체 성능을 절대적인 수준에서 확인하였다. 다양한 검사항목의 결측치 대체 성능을 대표하기 위해, 검사항목별 NRMSE의 평균값인 T-NRMSE를 평가지표로 설정하였다. T-NRMSE의 정의는 식 (7)과 같으며, 총 22개의 검사항목에 대한 반복 실험의 평균 NRMSE( $E(NRMSE_k)$ )에 대한 평균으로 정의하였다.

$$T-NRMSE = \frac{\sum_{k=1}^{22} E(NRMSE_k)}{22}$$

<Figure 4>는 5가지 결측률에 따른 평가지표의 변화를 나타내는 그래프이며, x축은 결측률을 나타내며 y축은 T-NRMSE를 나타낸다. 실험 결과, 모든 결측률에서 제안 방법의 T-NRMSE가 가장 낮게 측정되었다. 즉, 10%에서 50% 사이의 결측 환경에서 대조 방법 대비 제안 방법이 정확하게 결측치를 대체하고 있음을 알 수 있다.

두 번째 성능 평가는 각 검사항목에서의 제안 방법의 성능을 일록슨 부호 순위 검정을 통해 확인하였다. 본 가설 검정은 단측 검정으로, 귀무가설은 특정 검사항목에 대한 제안 방법의 NRMSE가 대조 방법과 같으며, 대립가설은 특정 검사항목에 대한 제안 방법의 NRMSE가 대조 방법 대비 작다

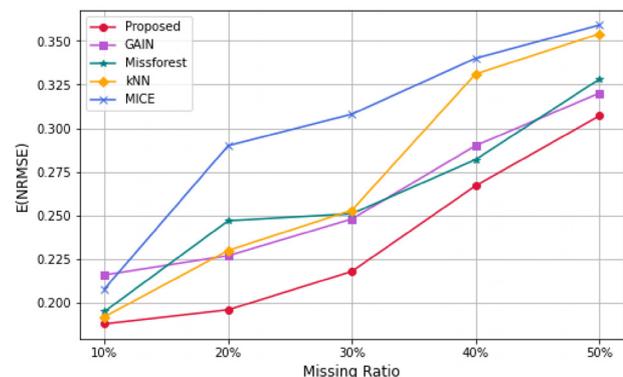


Figure 4. Average NRMSE of Imputation methods by Missing Ratio

**Table 1.** The Number of Significant Test Items in Various Missing Rates (M.R)

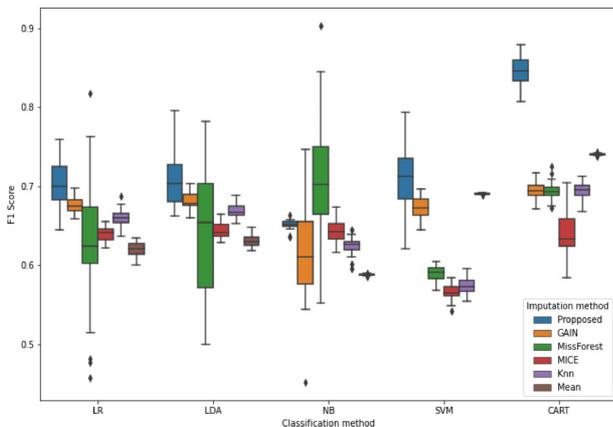
	M.R 10%	M.R 20%	M.R 30%	M.R 40%	M.R 50%
GAIN	20	20	20	20	20
Missforest	22	22	22	22	22
MICE	20	20	18	18	18
kNN	20	18	18	18	18
Mean	22	22	22	22	22

고 설정하였다. 가설 검정의 유의수준은 0.05로 설정하였다. <Table 1>은 전체 22개 검사항목 중 귀무가설이 기각된 검사항목의 수를 나타낸 것이다. 실험 결과, 모든 대조 방법에 대하여 전체 검사항목의 80% 이상에서 귀무가설이 기각되었다. 즉, 이는 전반적인 실험 환경에서 제안 방법이 대조 방법 대비 우수한 대체 성능을 유지하는 것으로 해석할 수 있다.

**4.3 제안 방법의 활용 및 활용 결과**

웨이퍼 테스트 데이터는 주로 반도체 칩의 불량 여부를 예측하기 위해 분석된다. 이에 제안 방법과 대조 방법으로 대체된 데이터로 칩의 불량 여부를 예측하는 이진 분류 성능을 확인하였다. 로지스틱 회귀(Logistic Regression), 선형 판별분석(LDA), 서포트 벡터 머신(SVM) 그리고 의사결정나무(CART)가 이진 분류 모형으로 활용되었다. 분류 모형은 전체 22개 검사항목을 설명 변수로 활용하여 칩이 불량이면 1, 정상이면 0으로 예측한다. 4.2절에서 10%에서 50% 사이의 모든 결측률에 대하여 각 결측치 대체 방법의 성능 순위가 유사하였으므로, 결측률 10% 환경에서 30번의 반복 실험을 통해 성능을 검증하였다. 예측 성능 지표는 F1-score를 활용하여 예측 변수의 불균형으로 인한 결과 왜곡을 최소화하고자 하였다.

<Figure 5>는 제안 방법과 대조 방법의 F1-score를 예측 모



**Figure 5.** Boxplot for the Prediction Results

형에 따라 박스 플롯으로 나타낸 것이다. x축은 예측 모형의 종류를, y축은 F1-score를, 박스 플롯의 색깔은 결측치 대체 방법을 나타낸다. 실험 결과, 모든 대체 방법에 대해 제안 방법으로 대체된 데이터가 우수한 예측성능을 보임을 확인할 수 있다. 제안 방법의 성능은 GAIN 대비 평균 5%, Missforest 대비 평균 24%, MICE 대비 평균 28%, kNN 대비 평균 26%, 그리고 평균 대체법 대비 평균 4% 향상되었다. 즉, 제안 방법으로 대체된 웨이퍼 테스트 데이터를 분석한다면 대조 방법 대비 칩의 불량 여부를 더 정확히 예측할 수 있을 것으로 해석할 수 있다.

**5. 토의**

제안 방법의 요체는 GAIN 대비 위치 정보의 추가와 손실함수 개선의 두 가지 요소로 구성된다. 본 절에서는 하나의 요소만을 제외하거나, 두 가지 요소 모두 제외한 경우의 성능을 제안 방법과 비교하여 각 요소가 제안 방법의 성능에 미치는 영향을 기여도 관점에서 해석하였다. 기여도 평가지표는 검사항목별 평균 *NRMSE*에 대한 신뢰구간으로 설정하였다. <Table 2>는  $GAIN(M_{00})$ 에 각 요소를 추가했을 때의 성능 향상 정도를 결측률 10% 환경에서 측정한 결과이다. 위치 정보만을 반영한 모델인  $M_{10}$ 은 두 가지 요소를 모두 반영하지 않은  $M_{00}$  대비 표본 평균을 약 11% 감소시켰으며, 손실함수 개선만이 이루어진 모델인  $M_{01}$ 은  $M_{00}$  대비 표본 평균을 약 22% 감소시켰다. 한편, 두 가지 요소를 모두 반영한 모델  $M_{11}$ 은  $M_{00}$ 의 표본 평균을 약 31% 감소시키면서 가장 좋은 결측치 대체 성능을 보였다. 즉, 위치 정보의 추가와 손실함수의 개선을 모두 고려했을 때 결측치 대체 성능을 가장 크게 개선할 수 있었다.

**Table 2.** Confidence Interval of *NRMSE* according to the presence of Location Information and Correlation Loss Term

Model	Confidence Interval of <i>NRMSE</i> (Increase rate compared with GAIN)
GAIN ( $M_{00}$ )	0.366 ± 0.216(Baseline)
$M_{00}$ + Location Information ( $M_{10}$ )	<b>0.329 ± 0.052(11%)</b>
$M_{00}$ + Correlation Loss Term ( $M_{01}$ )	<b>0.298 ± 0.084(22%)</b>
$M_{10}$ + Correlation Loss Term or $M_{01}$ + Location Information ( $M_{11}$ )	<b>0.279 ± 0.044(31%)</b>

## 6. 결론

웨이퍼 테스트 데이터는 분석 과정에서 반도체 불량량의 판별 근거를 제공하여 반도체 공정 전반의 품질 향상에 기여한다. 본 연구는 반도체 웨이퍼 테스트 데이터의 특성을 고려한 결측치 대체 방법을 제안한다. 제안 방법은 기존 결측치 대체 방법인 GAIN에 웨이퍼 내에서 발생하는 공간적 유사성과 검사 항목 간 상관관계를 반영하기 위하여 위치 정보 추가와 GAIN 손실함수 개선의 두 가지 측면을 반영하였다. 제안 방법은 실제 반도체 회사에서 수집된 웨이퍼 테스트 데이터를 활용하여 검증한 결과, 대부분의 검사항목과 다양한 결측률에서 제안 방법의 결측치 대체 성능이 우수함을 확인하였다.

제안 방법은 반도체 칩의 불량 예측성능을 향상하여 웨이퍼 수율 개선에 기여할 수 있을 것으로 기대된다. 또한, 제안 방법은 위치 정보가 존재하고, 수집되는 항목 간의 상관관계가 높은 OLED 패널 품질 데이터, 열간 압연 스트립강 품질 데이터를 비롯한 제조 데이터뿐만 아니라, 지역별 기상 데이터, 거주 지역별 소득 데이터 등의 결측치 대체에도 확장할 수 있을 것으로 기대된다.

제안 방법은 다음과 같은 향후 연구 수행으로 발전될 수 있다. 첫째, 웨이퍼 테스트 데이터의 특성 추가를 통한 결측치 대체 성능 향상이다. 5절에서 웨이퍼 테스트 데이터의 두 가지 특성의 추가를 통한 결측치 대체 성능 향상을 확인하였으므로, 본 연구에서 고려하지 못한 다양한 특성을 고려한다면 결측치 대체 방법의 성능을 개선할 수 있을 것이다. 예를 들어, 웨이퍼 단위에서 웨이퍼 간 관계성과 불량 칩의 분포에 영향을 끼친다고 알려진 웨이퍼가 처리된 공정 이력(Lee et al., 2009) 등이 추가로 반영될 수 있을 것이다. 둘째, 검증 대상의 확장을 통한 제안 방법의 신뢰성 향상이다. 본 연구는 웨이퍼 테스트 데이터 확보의 어려움으로 한정된 검사항목에 대하여 성능을 검증했다는 한계점이 존재하였다. 따라서 검사항목을 추가 확보하여 제안 방법 검증의 신뢰성을 증대할 수 있을 것이다. 마지막으로, 제안 방법의 생성자 손실함수 하이퍼 파라미터인  $\alpha$ 와  $\beta$ 의 관계성을 실험할 수 있을 것이다.  $\alpha$ 가 커지면 생성자로부터 추정된 개별 대체 값의 정확도가 올라갈 것이며,  $\beta$ 가 커지면 검사항목 간 상관관계 분포의 정확도가 올라갈 것이다. 따라서 가중치 간 관계성 파악을 통해 생성자 학습 과정에서의 가중치별 영향도를 확인할 수 있을 것이다. 이를 기반으로 데이터 수집 방법과 결측 유형 등의 기준에 따라 적합한 하이퍼 파라미터 조합을 선정할 수 있을 것으로 기대된다.

## 참고문헌

Adeli, E., Zhang, J., and Taflanidis, A. A. (2021), Convolutional generative adversarial imputation networks for spatio-temporal missing data in storm surge simulations, arXiv preprint.  
Aittokallio, T. (2010), Dealing with missing values in large-scale

studies: microarray data imputation and beyond, *Briefings in Bioinformatics*, **11**(2), 253-264.  
Baek, S. and Lee, M. (2022), A Study on the Type Classification Model of Defective Semiconductor Wafers Using Deep Learning, *Journal of the Korean Institute of Communications and Information Sciences*, 1158-1159.  
Burkhardt, A., Berryman, S., Brio, A., Ferkau, S., Hubner, G., Lynch, K., ..., and Sonderer, K. (2018), Measuring Manufacturing Test Data Analysis Quality, In *2018 IEEE Autotestcon*, 1-6.  
Chandar, S., Khapra, M. M., Larochelle, H., and Ravindran, B. (2016), Correlational neural networks, *Neural Computation*, **28**(2), 257-285.  
Han, Y. and Lee, C. (2005), Automatic Classification of Failure Patterns in Semiconductor EDS Test for Yield Improvement, *Journal of the Korea Society for Simulation*, **14**(1), 1-8.  
Hsu, C. K., Lin, F., Cheng, K. T., Zhang, W., Li, X., Carulli, J. M., and Butler, K. M. (2013), Test Data Analytics - Exploring Spatial and test-item Correlations in Production Test Data, *Proceedings - International Test Conference*, ITC, 1-4.  
Huang, K., Carulli, J. M., and Makris, Y. (2013), Counterfeit electronics: A rising 43 threat in the semiconductor manufacturing industry, *Proceedings - International Test Conference*, ITC, 1-4.  
Jung, J. and Jung, Y. (2022), Wafer bin map failure pattern recognition using hierarchical clustering, *Journal of Korean International Statistical Society*, **35**(3), 407-419.  
Kang, H. and Baek, J. (2020), Improved Quality Prediction Method by Clustering Data in Semiconductor Manufacturing Process, *Journal of the Korean Institute of Industrial Engineers*, **46**(2), 134-142.  
Kim, D., Park, Y. S., Kim, H. W., Park, K. S., and Moon, I. K. (2022), Inventory policy for postponement strategy in the semiconductor industry with a die bank, *Simulation Modelling Practice and Theory*, **117**, 102498.  
Kim, S. and Kim, J. (2022), A study on the development strategy of the Metrology industry using the modified AHP and IPA, *Journal of the Korean Institute of Plant Engineering*, **27**(2), 49-59  
Ko, G., Tak, H., and Lee, B. (2014), Impact of Missing Values on Survey Research and Relevancy of Multiple Imputation Techniques, *Journal of the Korean Journal of Policy Analysis and Evaluation*, **24**(3), 49-75.  
Lee, Y. H., Ham, M., Yoo, B., & Lee, J. S. (2009), Daily planning and scheduling system for the EDS process in a semiconductor manufacturing facility, *The International Journal of Advanced Manufacturing Technology*, **41**(5), 568-579.  
Lee, S. Y., Connerton, T. P., Lee, Y. -W., Kim, D., and Kim, J. -H. (2022), Semi-GAN: An Improved GAN-Based Missing Data Imputation Method for the Semiconductor Industry, *IEEE Access*, **10**, 72328-72333.  
Lee, Y. H., Ham, M., Yoo, B., and Lee, J. S. (2009), Daily planning and scheduling system for the EDS process in a semiconductor manufacturing facility, *The International Journal of Advanced Manufacturing Technology*, **41**(5), 568-579.  
Luo, M., Wang, S., Wang, C., Chen, W., Zhu, E., and Liu, X. (2022), DICDP: Deep Incomplete Clustering with Distribution Preserving, In *International Conference on Artificial Intelligence and Security* Springer, Cham, 162-175.  
Nuhu, A. A., Zeeshan, Q., Safaei, B., and Shahzad, M. A. (2022), Machine learning-based techniques for fault diagnosis in the semiconductor manufacturing process: A comparative study, *The Journal of Supercomputing*, 1-51.  
Park, J. and Kim, S. (2014), A Prediction Methodology of Package

- Chip Quality using Probe Test Fail bit Count Data, *Journal of the Korean Institute of Industrial Engineers*, 528-536.
- Schranner, S. (2019), Pattern Recognition in Analog Wafer Test Data: A Health Factor for Process Patterns. Ph.D.diss, Graz University of Technology, Austria.
- Tsung, C. K., Hsieh, H. Y., and Yang, C. T. (2019), An implementation of scalable high throughput data platform for logging semiconductor testing results, *IEEE Access*, 7, 26497-26506.
- Van Buuren, S. and Groothuis-Oudshoorn, K. (2011), Mice: Multivariate imputation by chained equations in R, *Journal of Statistical Software*, 45, 1-67.
- Wang, Y., Li, D., Li, X., and Yang, M. (2020), PC-GAIN: Pseudo-label Conditional Generative Adversarial Imputation Networks for Incomplete Data, *Neural Networks*, 141, 395-403.
- Wang, Y., Li, D., Li, X., and Yang, M. (2020), PC-GAIN: Pseudo-label Conditional Generative Adversarial Imputation Networks for Incomplete Data.
- Yeon, W. (2021), Conflict between the US and China and China's strategy and prospects for fostering the semiconductor industry, *KIEP World Economy Focus*, 39(4), 1-19.
- Yoon, J., Jordon, J., and van der Schaar, M. (2018), GAIN: Missing Data Imputation using Generative Adversarial Nets, *International Conference on Machine Learning*, 5689-5698.
- Zhang, Y. and Thorburn, P. J. (2022), Handling missing data in near real-time environmental monitoring: A system and a review of selected methods, *Future Generation Computer Systems*, 128, 63-72.

## 저자소개

**김주영**: 동국대학교 통계학과 학사 졸업 후, 포항공과대학교 산업경영공학과 석사과정에 재학 중이다. 주요 연구분야는 품질공학, 인공지능 등이다.

**배영목**: University of Utah Operations Management 학사, 포항공과대학교 산업경영공학과 석사 졸업 후, SK hynix 재직 중이며, 동 대학 산업경영공학과 박사과정에 재학 중이다. 주요 연구분야는 품질공학, 이상탐지 등이다.

**최승현**: 포항공과대학교 산업경영공학과 학사 졸업 후, 동 대학 산업경영공학과 박사과정에 재학 중이다. 주요 연구분야는 품질공학, 공정최적화 등이다.

**김광재**: 서울대학교 산업공학과 학사, 한국과학기술원 산업공학과 석사, Purdue University 경영과학 박사 학위를 취득하였다. 현재 포항공과대학교 산업경영공학과 교수로 재직 중이다. 주요 연구분야는 품질공학, 서비스공학 등이다.