# Machine/Deep Learning Applications Using the V93000 and Nvidia Jetson TX2

Keith Schaub, Ira Leventhal and Brian Buras - Advantest

Gerard John - Amkor
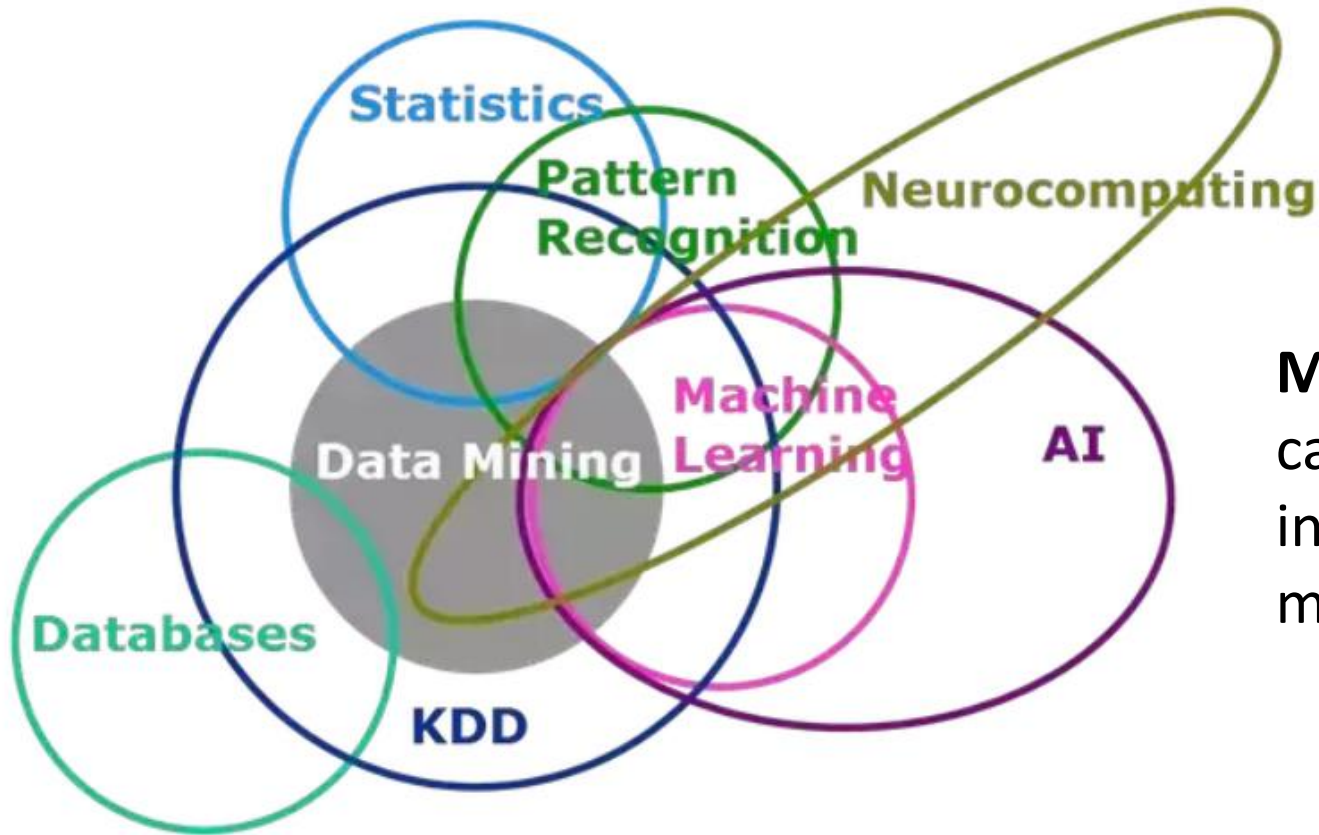
Yiorgos Makris - UTD

**VOICE** 2018

**ADVANTEST**®

# Outline

- **AI Machine Learning / Deep Learning Overview**
- **Problem Statement**
- **Test Compaction:** Hypothesis 1 – Machine learning algorithms analyze test data to optimize the test list.
- **Dynamic Spatial Testing:** Hypothesis 2 – Machine learning algorithms learn wafer spatial correlations to dynamically optimize test coverage
- **Test Compaction**
  - Process / Data Analysis
  - Results
  - Conclusions
- **Dynamic Spatial Testing**
  - Process / Data Analysis
  - Results
  - Conclusions
- **Summary**
- **Next Steps**
- **Machine Learning Image Classifier Integrated into the V93000 Environment (Kiosk)**
- **Future Considerations**

# Machine Learning Overview



**Machine Learning** is an AI sub-category focused on finding patterns in data and using those patterns to make predictions

# Machine Learning Training

**Input, feed a lot of data**

**Machine Learns patterns in the data**

**MODEL**

**"OK, I see the patterns and understand the data now"**

# Machine Learning Training Example

Input a bunch of Chihuahuas

Machine Learns to recognize Chihuahua patterns

**MODEL**



"hmm, ok I learned what Chihuahuas look like"
- Pointed ears
- Small typically dark nose
- Little beady eyes
- …

Disclaimer: No dogs were harmed as part of this presentation

# Chihuahua or Muffin?

Input Chihuahuas and "non Chihuahuas"

Algorithm applies Chihuahua model to classify

**MODEL**

**Classification Result**

"You didn't train me what a muffin looks like?!"

# Input training data is important!

**Puppy or Bagel?**



**Sheepdog or Mop?**



**Labradoodle or fried chicken?**

# Problem Statement

- **Testing complexity and test cost continues to increase**
  - Quality is the new Cost
  - More testing
  - Multiple domain types and insertions needed
  - Need to avoid longer test times
  - Need to minimize test costs
- **Process variations are not static, yet testing methodologies typically are static**
  - Same tests applied throughout device life cycle
  - Engineers manually adjust
    - Laborious, tedious, "after the fact; late"
      - Negatively impacts quality
      - TONS of data, but humans are not efficient at analyzing it

# Test Compaction – Hypothesis 1

Can a machine learning algorithm learn measurement correlations to automatically optimize testing metrics?



Machine Learning Algorithm Trains

Determine if some of these are highly correlated!

Use the newly organized subset

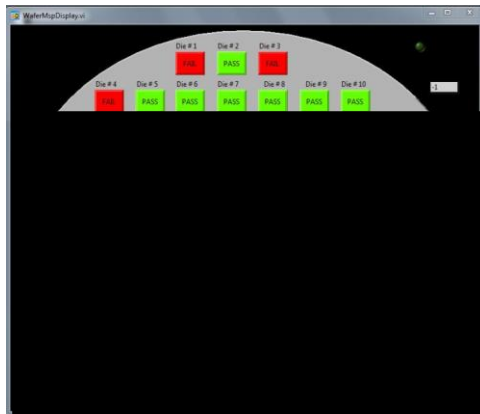Show test results just as accurately
Show same quality
Show reduced cost impact

# Dynamic Spatial Machine Learning of Wafer Testing – Hypothesis 2

Can a machine learning algorithm learn spatial correlations to automatically optimize testing per die?



Machine Learning Algorithm Trains

Trained Model

Apply Model to predict result

Predicted test results

# Specification Test Compaction Concept

T: Total set of n tests

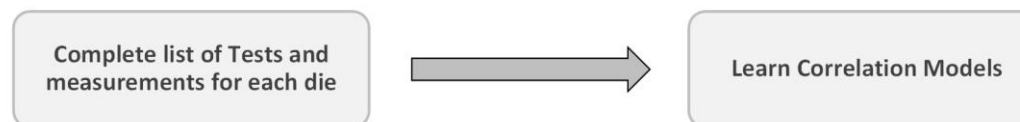S: $S \subset T$ (Subset of k tests)

Et: Number of test escapes for test t

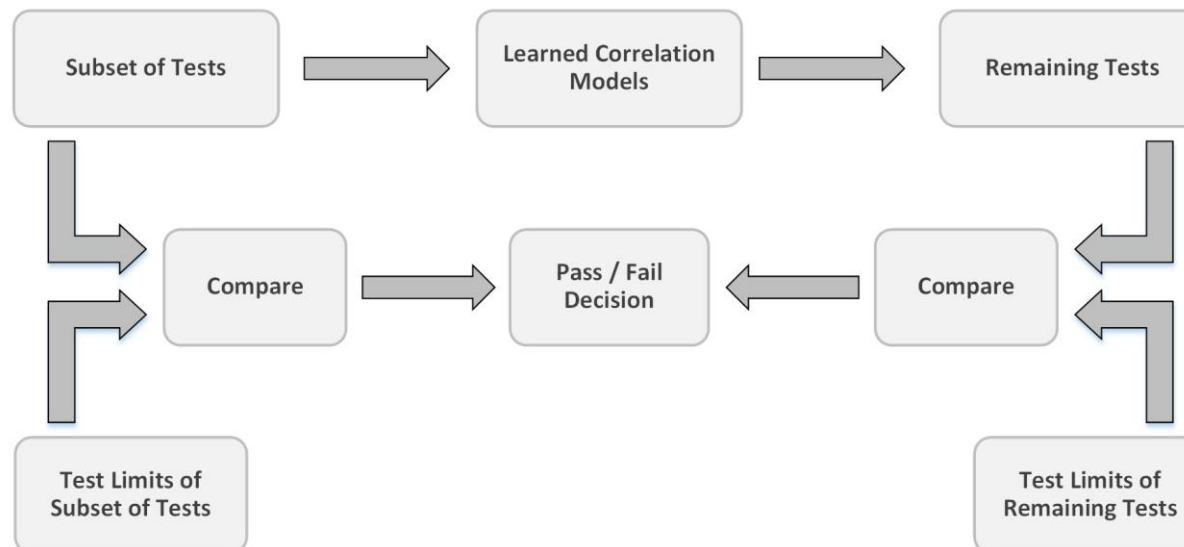The objective is to minimize the size of S while maintaining a low

$$\sum_{i=0}^{k} E_i$$

where Ei is the test escape for $i^{th}$ test in S.

Different sizes of S can be produced depending on what the acceptable escape rate is.

## Learning Phase



## Testing Phase

# Test Compaction – Data Description & Idiosyncrasies



- Dataset contains 6 wafers, with 20 test measurements

- There are 30 failing die out of 402

- Small number of die locations per wafer
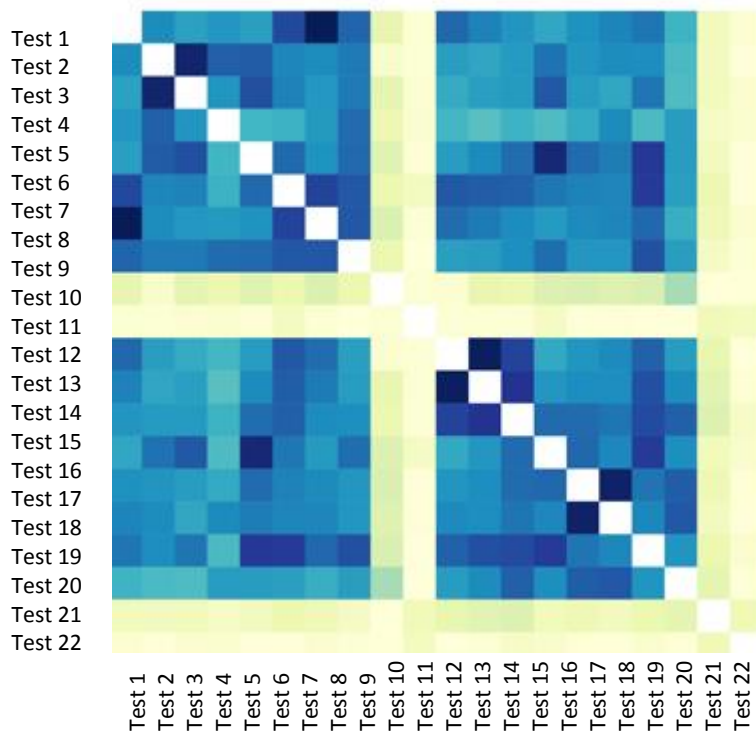
- No test groups or test times

| Wafer # | Pass count | Fail count |
|---------|-----------|-----------|
| 1 | 58 | 9 |
| 2 | 60 | 7 |
| 3 | 57 | 10 |
| 4 | 65 | 2 |
| 5 | 67 | 0 |
| 6 | 65 | 2 |

# Test Compaction – Pre-filter



**Clean the Data (labradoodle or fried chicken?)**
An important pre-step to training the model is to clean up the data before it is fed to the training algorithm.

We removed outliers before training the algorithm

# Test Correlation



**Correlation Results**

**Yellow** = low correlation

**Blue** = high correlation

Many measurements are highly correlated!

High correlation

Low correlation

- Bi-variate correlation of all test pairs using absolute values of Pearson Correlation Coefficients (PCC).
  - This shows the degree by which two variables co-vary
- Multi-variate non-linear regression modeling is a more suitable technique for discovering correlations between tests.

# Test Correlations

- Multi-variate Adaptive Regression Splines (MARS)[1] is a non-linear regression analysis methodology

- Training consists of two phases that aim to select the optimal number of features:
  - **Forward pass**: Starting with the intercept term and progressively adds a basis function that minimizes the prediction error. This usually generates an overfit model
  - **Backward pass**: This stage prunes the basis functions using a metric that penalizes the model based on the number of features

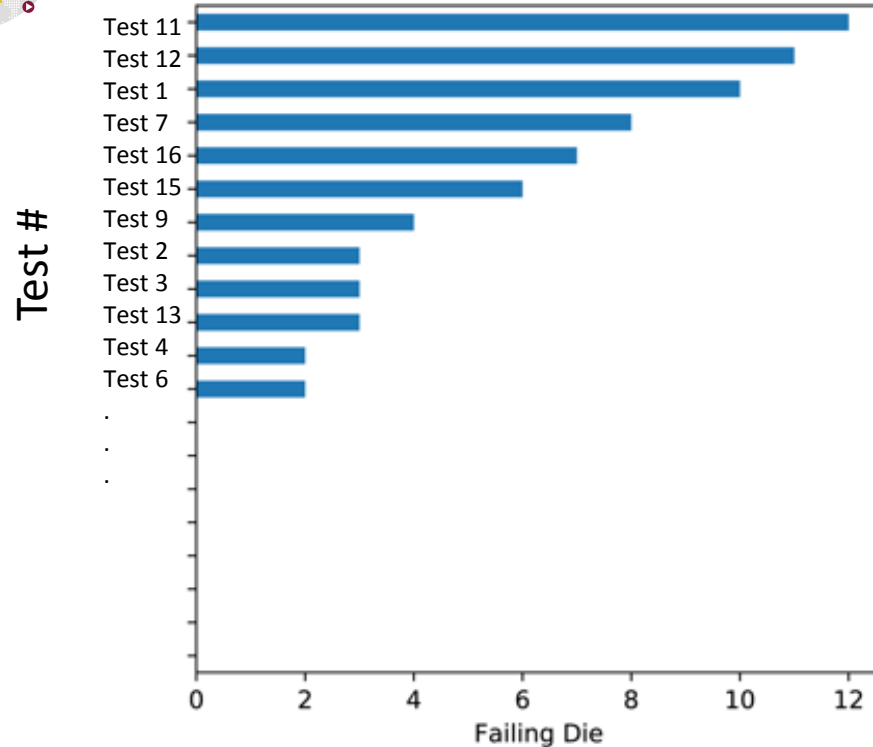[1] Friedman, J. H. (1991). "Multivariate Adaptive Regression Splines". *The Annals of Statistics*.

VOICE
2018

# Test Correlations

- Description of the MARS-based experiment:
  - Train a MARS model for every test in the dataset and calculate the accuracy of the model using a hold-out set of wafers
  - Identify the most accurately modeled tests based on the prediction error
- Most accurately predicted tests: Test 1, Test 2, Test 3, Test 7 , Test 11 Test 15, Test 16
- For this experiment the python implementation of MARS (pyearth) was used

# Test Compaction & Reordering – Trained algorithm suggests subset of tests



- Greedy Algorithm for test compaction:
  - Start by including the test that captures the most failing devices. Test 11 in our dataset
  - Iteratively add the test that minimizes the test-escapes. This can skip tests based on the overlap
- e.g. tests that capture all 30 failing die are: Test 1, 3, 5, 8, 7
- Algorithm suggests to use these 5 tests
- Algorithm could automatically re-order tests to optimize test flow (i.e. learn and apply most efficient tests and optimize test flow)
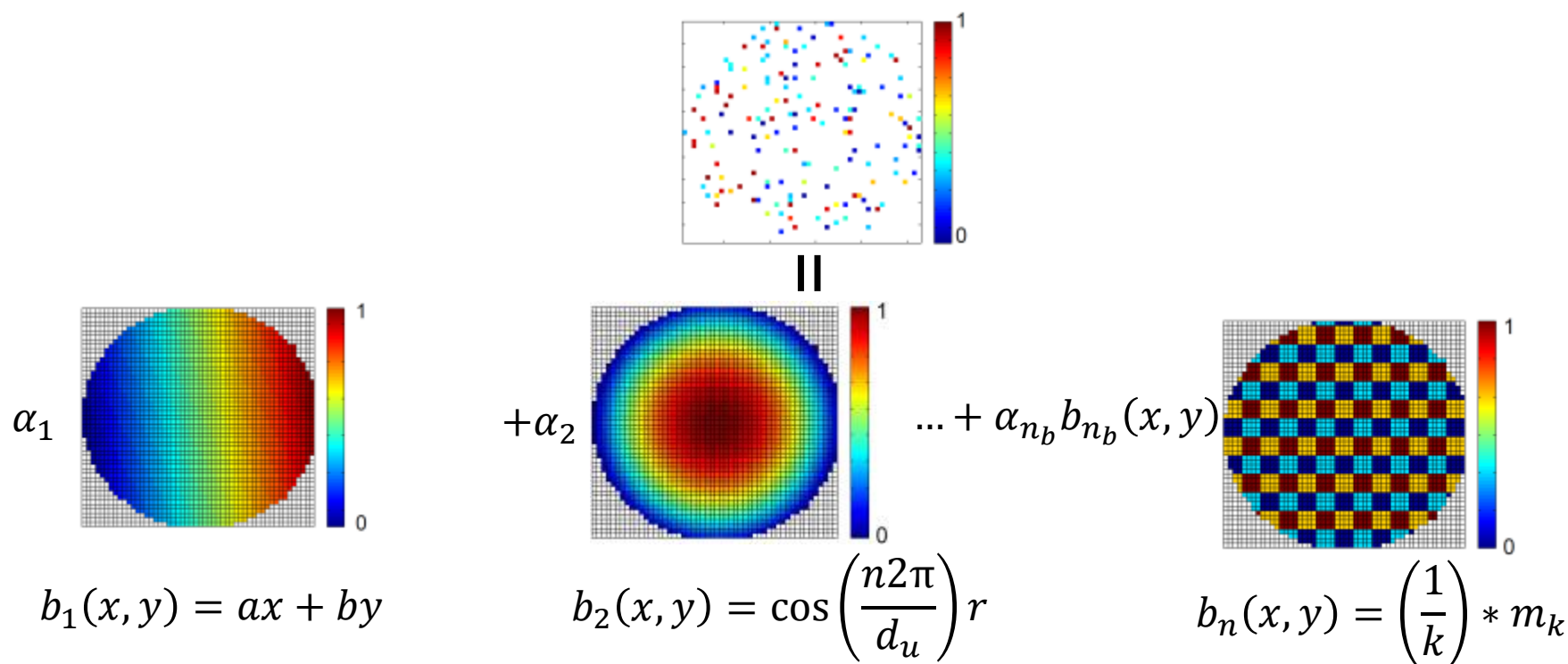- Test time savings reduces cost

Other algorithm examples: Support vector machines, decision trees, neural networks

# Dynamic Spatial Machine Learning of Wafer Testing

Spatial decomposition of wafer measurements

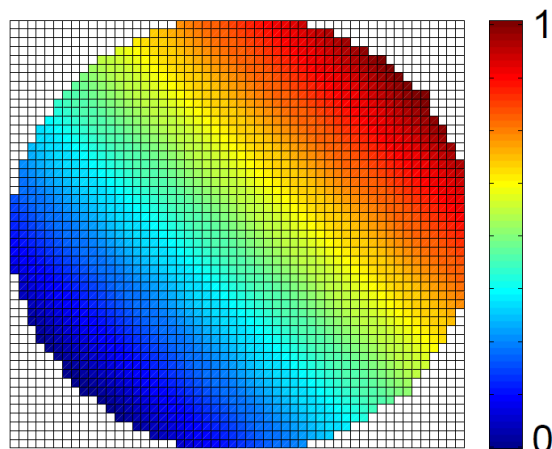$$g(x,y) = \alpha_1 b_1(x,y) + \cdots + \alpha_{n_b} b_{n_b}(x,y)$$



$$=$$

$\alpha_1$  $+\alpha_2$  $\ldots + \alpha_{n_b} b_{n_b}(x,y)$ 

$$b_1(x,y) = ax + by \qquad b_2(x,y) = \cos\left(\frac{n2\pi}{d_u}\right)r \qquad b_n(x,y) = \left(\frac{1}{k}\right) * m_k$$

Learn these functions from the data…

*K. Huang, N. Kupp, J. Carulli, and Y. Makris, "Process Monitoring through Wafer-level Spatial Variation Decomposition," ITC 2013
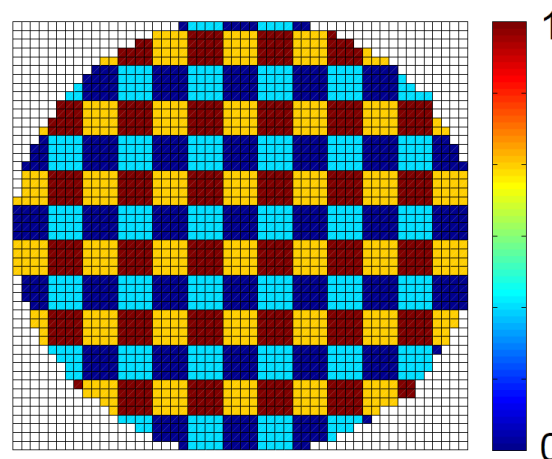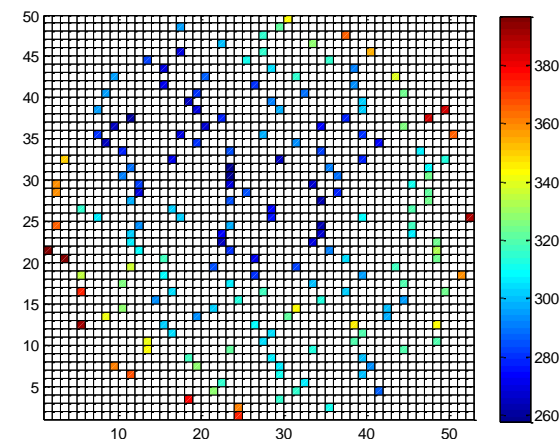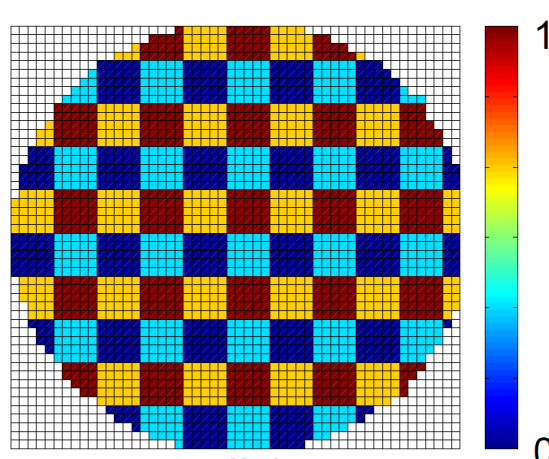
# Examples of spatial basis functions



Linear
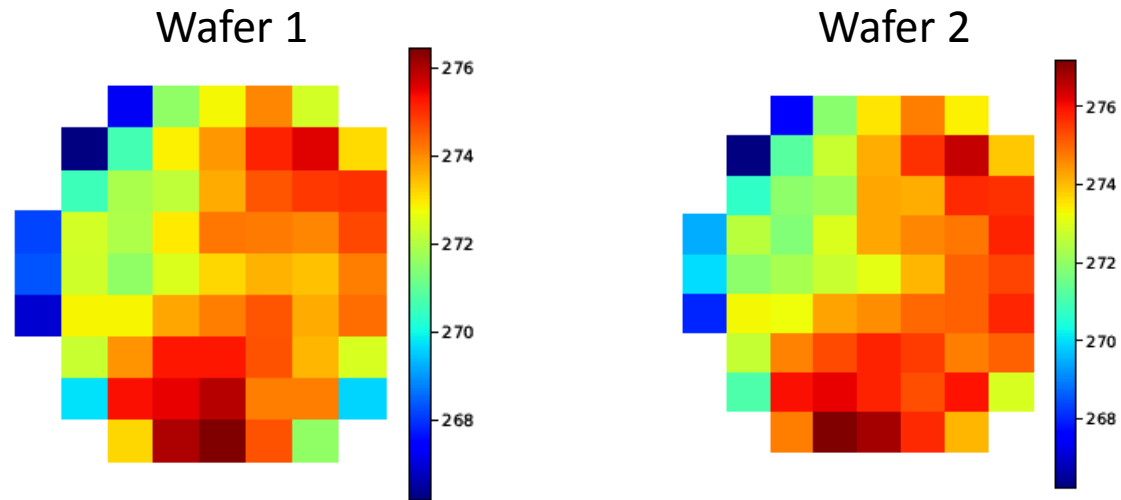


Radial





Checkerboard #1



Checkerboard #2

$$A = [\alpha_1, \alpha_2, \alpha_3, \alpha_4] \ ?$$

**Basis function learned using domain-specific knowledge**
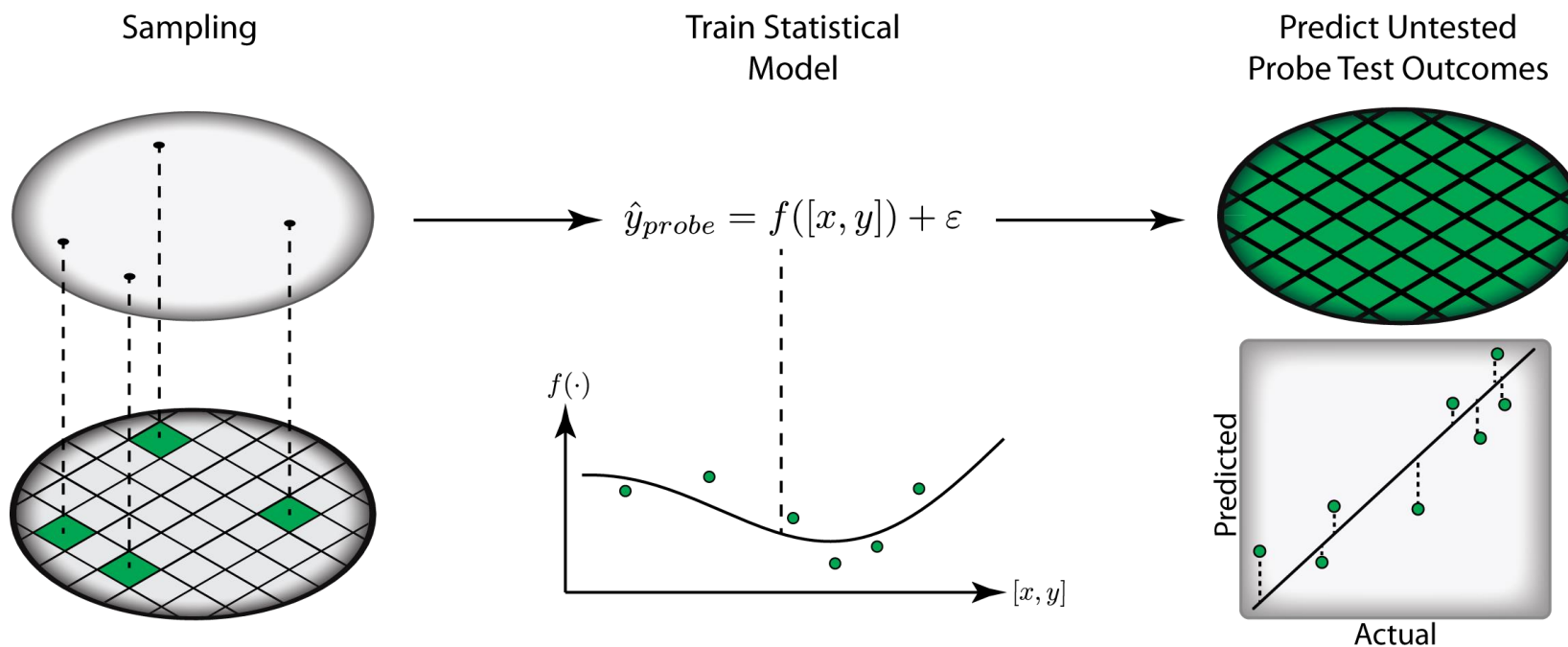
# Algorithm Learns Spatial Correlation Pattern

- Spatial correlation refers to the relationship that certain test measurements have as a function of the die locations

- One way to identify such wafer-level spatial correlations is to perform visual inspection on the wafer maps of each test.



Wafer 1



Wafer 2

# Spatial Correlation Modeling

- In our experiments we performed spatial-correlation modeling using Gaussian processes[2]



Sampling     Train Statistical Model     Predict Untested Probe Test Outcomes
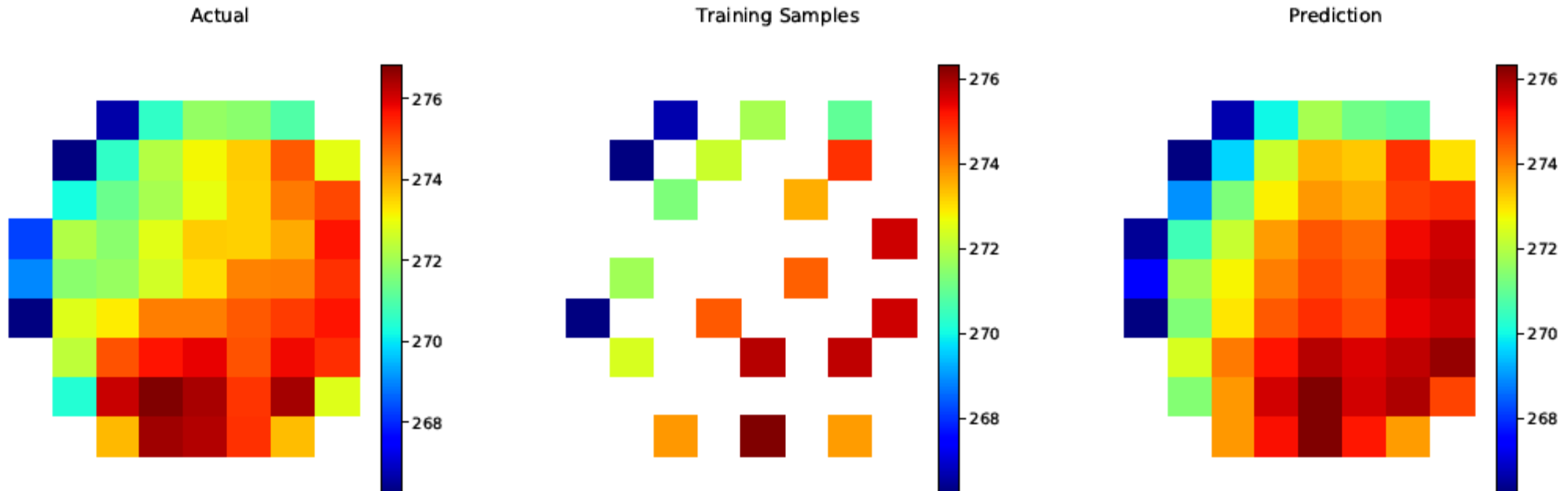
$$\hat{y}_{probe} = f([x, y]) + \varepsilon$$

[2] N. Kupp, K. Huang, J. Carulli, Y. Makris, "Spatial Estimation of Wafer Measurement Parameters Using Gaussian Process Models", *Proceedings of the IEEE International Test Conference (ITC)*

# Spatial Correlation Accuracy Results

- Spatial correlation modeling example on Test 9
- Relative prediction error = 0.4%

$$Relative\ Error = |\frac{(acutal\ - predicted)}{actual}|$$



Actual                Training Samples                Prediction

[2] N. Kupp, K. Huang, J. Carulli, Y. Makris, "Spatial Estimation of Wafer Measurement Parameters Using Gaussian Process Models", *Proceedings of the IEEE International Test Conference (ITC)*

# Summary

- Both hypothesis were shown to be true
  - Machine learning algorithms can automatically learn test optimization techniques by analyzing the data
    - They can learn which tests are most important
    - They can automatically generate the relevant/sub-set test list
    - They can automatically optimize the test flow by re-organizing the test list
  - Machine learning algorithms
    - Can find correlations and dependencies in the data
    - Use that information to optimize testing and lower test cost
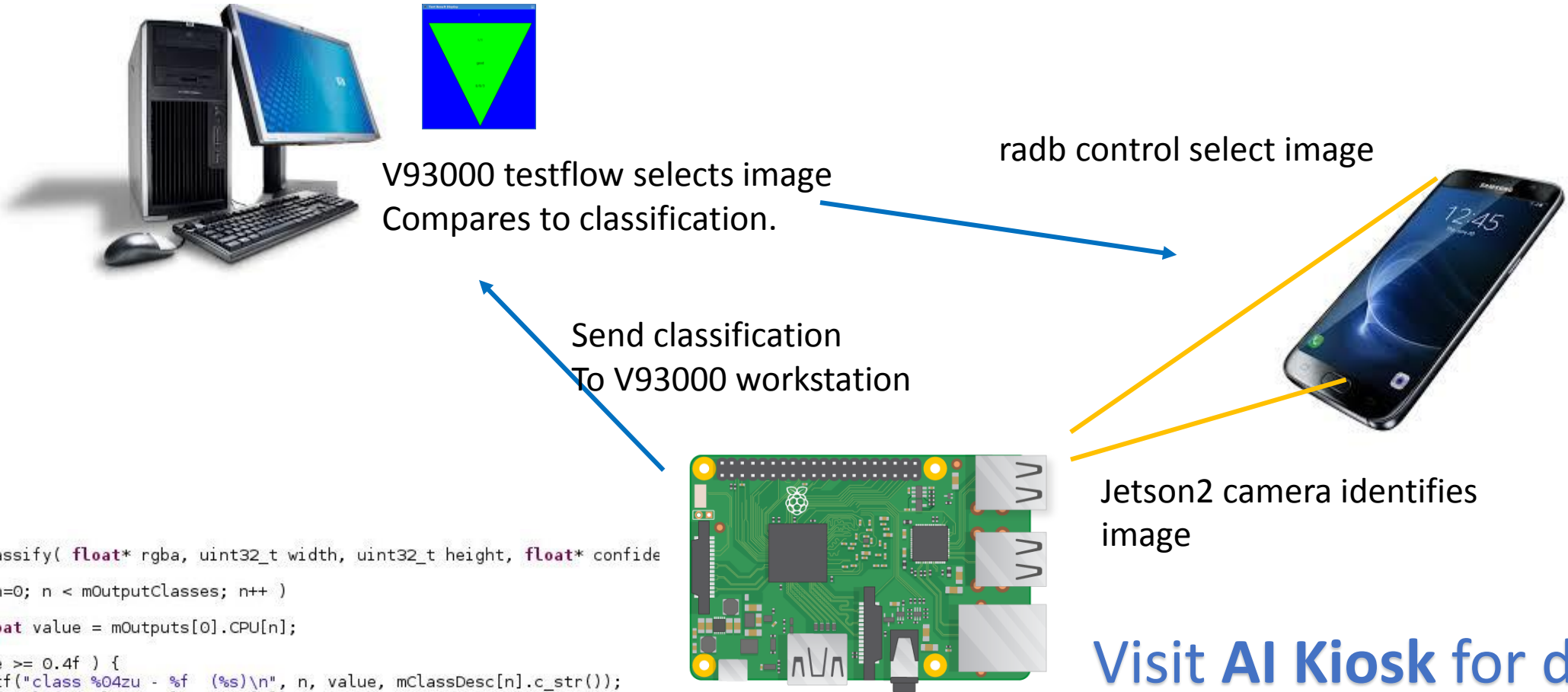    - Example: the foreknowledge could be used to eliminate re-testing

# Next Steps

- Apply same methods to multiple and larger data sets

- Integrate machine learning technique into the SmarTest environment

- Develop an AI V93000 demonstration using Nvidia's Jetson[3] 256 core AI environment

  - Kiosk Demo – AI ML Jetson 2 TX - operating within smartest that classifies smartphone display images

[3] https://developer.nvidia.com/embedded/buy/jetson-tx2

# Machine Learning V93000 Environment

V93000 testflow selects image
Compares to classification.

radb control select image

Send classification
To V93000 workstation

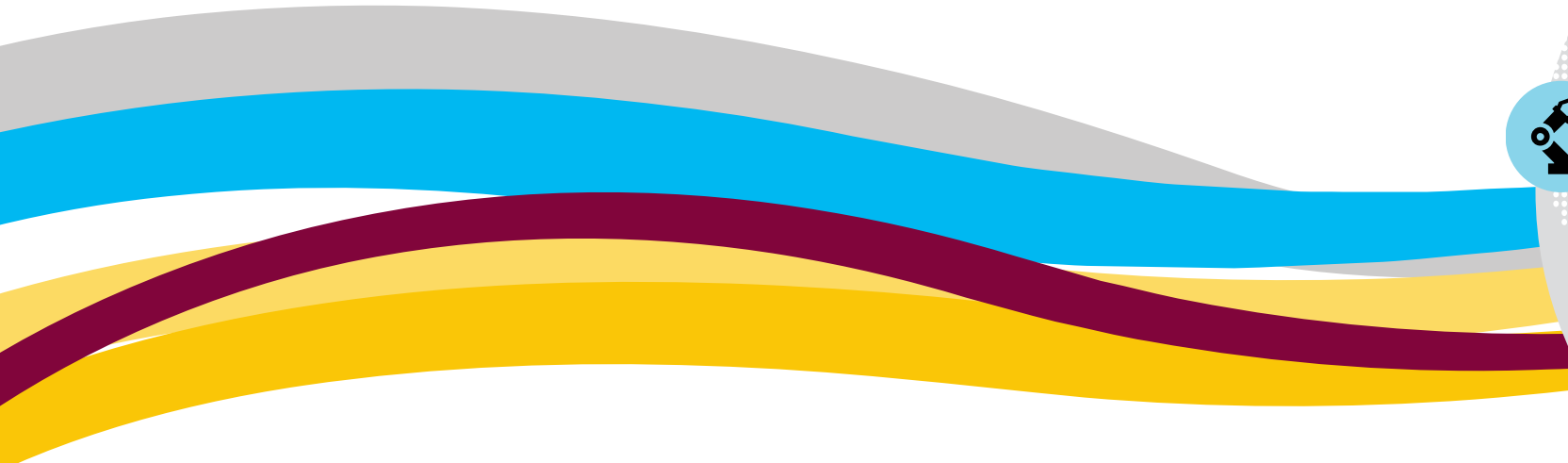Jetson2 camera identifies
image

```
// Classify
int imageNet::Classify( float* rgba, uint32_t width, uint32_t height, float* confide
{
    for( size_t n=0; n < mOutputClasses; n++ )
    {
        const float value = mOutputs[0].CPU[n];

        if( value >= 0.4f ) {
            printf("class %04zu - %f  (%s)\n", n, value, mClassDesc[n].c_str());
            sendResult(n,value, mClassDesc[n]);
```

Visit **AI Kiosk** for demo

# Future Considerations

- Develop Machine Learning APIs for the SmarTest that customers could use from a library

- Develop similar APIs for the Nvidia Jetson II AI environment that could be controlled from SmarTest environment
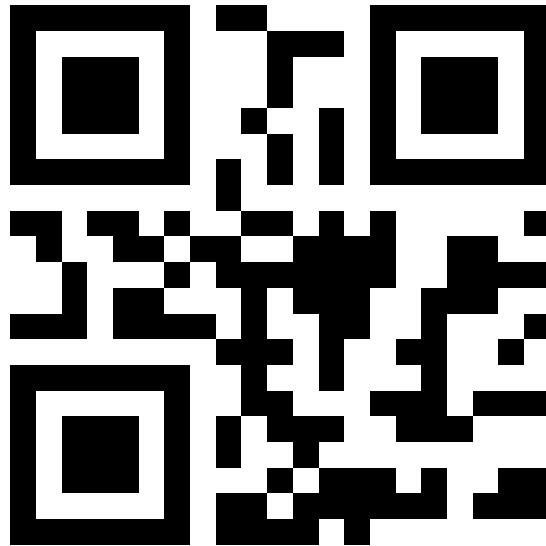  - Customers would have a 256 core AI environment that they can build their own models

# Measure the Connected World
*And Everything in It[SM]*



# Thank You.

# V93000-379-HT - Machine/Deep Learning Applications Using the V93000 and Nvidia Jetson TX2



## San Diego