

Challenges in Adopting ML in Manufacturing

Jacob George

Sr. Al Engineer

Analysis group

+	+		+	+	+	+	+		+	+	+
+	+			+	+		+		+	+	
+	+		+	+	+	+	28	Sept	emb	er 20	21



Agenda

- Wafer Inspection
- AI/ML applications
- Challenges in Adopting AI/ML
 - Limitations from optical physics
 - Data challenges
 - Throughput challenges

Wafer Inspection

Semiconductor Manufacturing Process



We Must Find and Classify Really Small Defects



KLA

defect size

4 KLA Non-Confidential | Unrestricted

Defect Types



protrusion



SPIE Photomask Technology, 104510L

break



SPIE Advanced Lithography, 113231J ,2020



Proc. SPIE 10809, International Conference on Extreme Ultraviolet Lithography 2018



Electron-beam Review System Scanning Electron Microscope (SEM) Technology



Broadband Plasma Optical Patterned Wafer Inspectors



Optical-Based Inspection (photons)

Wavelength Optical Physics Throughput : DUV - Visible: Diffraction of light: 1000x SEM Inspection



Broadband Plasma Optical Inspection

Electron Beam

Ebeam Wavelength: 2.5 pm

Defect size < 20 nm



Image: KLA



Photons
Wavelength >> Defect size

Image: KLA



KLA







Capture more Defects of Interest at lower nuisance rate



Capture them all at the lower time rate

Its not only how good you classify, but how fast you can do that!



Classification prior to ML

Defect Attributes (several hundreds)





Classification with ML







Challenges in Adopting AI/ML

Limitations from Optical Physics



Image source: https://en.wikipedia.org/wiki/Front_end_of_line



Interferences from underneath layers





Data Challenges



Obtaining Labelled Data

Process Variation

Model Maintenance

13 KLA Non-Confidential | Unrestricted





Obtaining 1st defects to train













Settles, Burr. "Active learning literature survey." (2009).







Human Annotator









Human Annotator

KLA

System

Process Variation

1. Cleaning

2. Film Deposition

3. Resist Coating

4. Exposure

5. Development

- **Material Changes**
 - **Design variations**
 - Focus/ Exposure changes in Photolithography

Adaptive Models

Software solutions and infrastructure to track model performance

- Reproducible
- Easy to Maintain

Reproducible: Need to store Training Set

Nuisance

- Bookkeeping required
- Need explainable ML

Need to have an ecosystem that would keep track of models

Throughput Challenges

 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +

Achieving Throughput Expectation

Need to have cheaper compute

Optimize best on the hardware

Hardware Aware Algorithms

DL Inference

Overheads

Custom Implementations

Simple convolution (15x15)

Number of	Tensorflow	TensorRT (FP32)	CuDNN
Images	Inference (seconds)	Inference(seconds)	(Seconds)
1M	17.74	12.6	10.9

These frameworks are made to optimize large DL networks.

DL Inference

Can we create an optimized framework that could hyper boost DL inference tailored for inspection use case?

Layers	Parameters	Size	Cache	Max Capacity
Convolution 2D	32 x 5 x 5 x 3	4.6 KB	Constant Cache	8 KB per SM
Fully Connected	32 x 64	4 KB	Shared Memory	64 KB per Block
Fully Connected	64 x 2	0.2 KB	Shared Memory	64 KB per Block

All Parameters are stored as float16 (2 Bytes)

No of Images	TensorFlow (in seconds)	Custom cuDNN (in seconds)	Custom CUDA using Tensor core (in seconds)
50k	1.02	0.45	0.20
100k	1.81	0.76	0.33
500k	7.91	4.19	2.08
1M	16.15	8.92	4.54

Number of Images(32 x 32 x 3)

DL Inference

How can we train models that improves the inference speed?

BaselinePruning
(80%)Post Training
Quantization1284kb115kb39kb (tflite)

No Accuracy degradation after quantization

Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding Song Han, Huizi Mao, William J. Dally 2016

A model that fits into cache can be ninja optimized for inference

Need of the Hour...

Expert in AI/ML

- Explainable models
- Few shot learning
- Statistical ML
- Model optimization for compute
- Embrace Modern C++
- Heterogenous Compute

Thank You

 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +