

# EDA with ML, Rule-Based, or Both?

Youngsoo Shin

KAIST, Korea

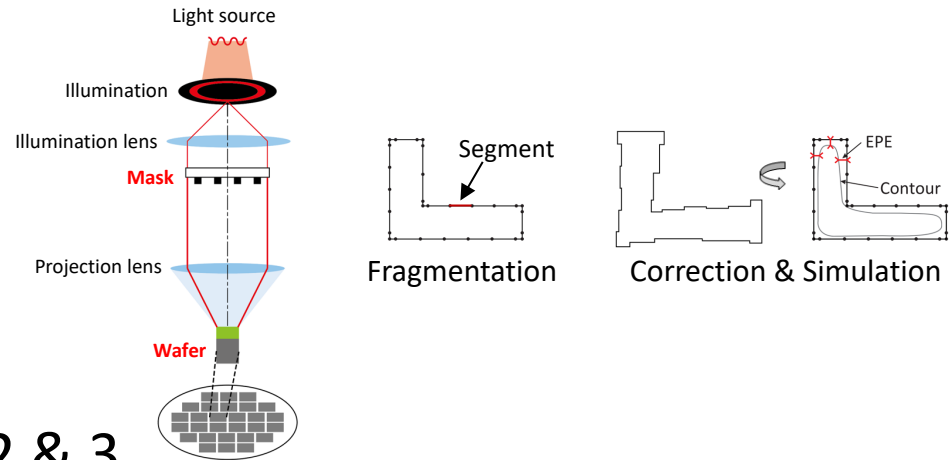
# Background

- ML vs Rule-based
  - ML is applicable for “large volume data”, which is often not the case in semiconductor industry
  - Rule-based is suited for smaller data volume
- Questions addressed in this talk
  - **Compare ML and Rule-based** in large- and small-volume data
  - **Combine ML and Rule-based**, so that Rule-based can be an efficient option in smaller data volume

# Example 1: Re-Fragmentation

- OPC process

1. Fragmentation
2. Correction of segments
3. Lithography simulation to check EPE → Iterate 2 & 3

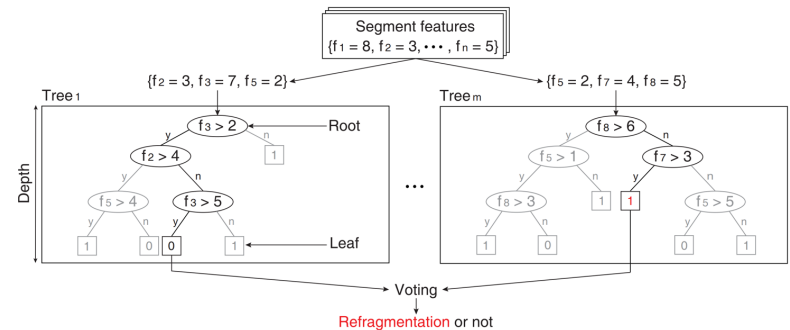


- Fragmentation

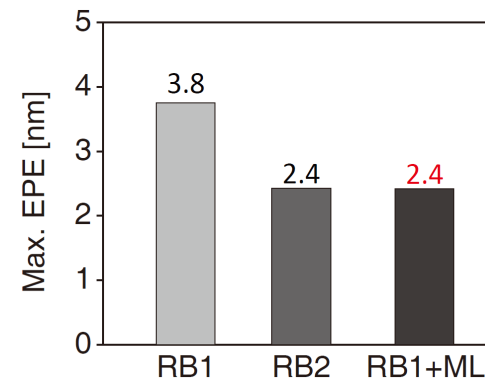
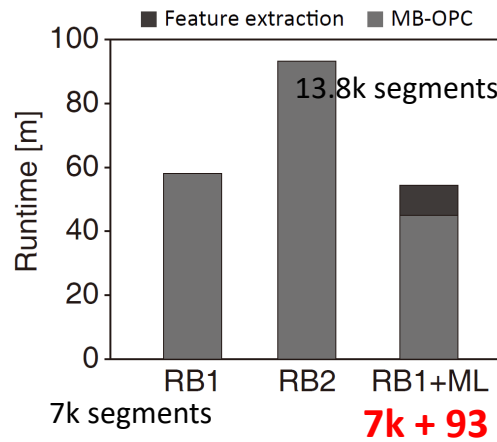
- Based on simple rules (e.g. nominal segment length)
- Segments which are not short enough are trouble
- **Re-Fragmentation: further divide a few segments** (so that OPC can complete faster with smaller EPE)

# Re-Fragmentation with RFC

- RFC process
  - Each decision tree receives a random subset of segment features & predicts 1 (split) or 0 (no-split)
  - Voting is collected & segment is divided in half if #votes > threshold



- RFC has been “trained” with 28k segments

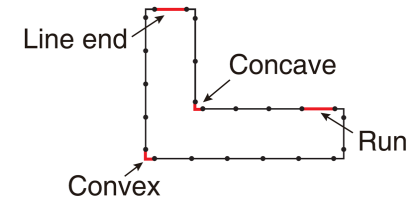


# Rule-Based Re-Fragmentation

- Same amount of data (28k segments) is used to set up a few rules

- $\sigma$  for length and  $2\sigma$  for |initial EPE|

Segment type	Length	Initial EPE
Line-end	>25nm	>31.9nm
Convex	>42nm	>7.4nm
Concave	>40nm	>11.8nm
Run (adjacent to corner)	>38nm	>8.8nm
Run (not adjacent to corner)	>44nm	>5.6nm



- Rule-based is worse (in max EPE) than RFC, as expected, when data volume is enough

Refragmentation	Max. EPE [nm]	#Segments
No	3.83	7,000
RFC (big data)	2.42	7,096
Rules (big data)	3.07	7,163

# RFC vs Rule-Based in Small Data Volume

- Sample data is reduced from 28k segments to 1.4k segments
- RFC model is re-trained; rules are also set up again
- Rule-based is better than RFC, this time
  - RFC is over fitted
  - Rules are less sensitive to the amount of data

Refragmentation	Max. EPE [nm]	#Segments
No	3.83	7,000
RFC (big data)	2.42	7,096
Rules (big data)	3.07	7,163
RFC (small data)	3.41	7,165
Rules (small data)	3.13	7,177

# Revising Rules through RFC

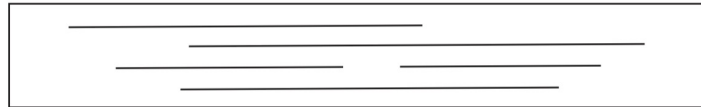
- Intuitions from RFC (trained with small data volume!)
  - All the tree roots carry “ $\phi_1 > 0.73?$ ” (top decision maker) if  $\phi_1$  is a feature
    - First optical signal ( $\phi_1$ ) is a main component in light intensity calculation
  - Such trees carry “ $|\text{initial EPE}| > x?$ ” in leaves (final decision maker)
    - $x$  values are collected and average is calculated
- Key observation
  - **Rule-based approach can be made very efficient, with intuitions extracted from ML model**

Segment type	Length	Initial EPE  if $\phi_1 \leq 0.73$	Initial EPE  if $\phi_1 > 0.73$
Line-end	>25nm	>29.4nm	>26.5nm
Convex	>42nm	>7.8nm	>7.0nm
Concave	>40nm	>10.2nm	>9.2nm
Run (adjacent to corner)	>37nm	>9.7nm	>8.7nm
Run (not adjacent to corner)	>44nm	>5.9nm	>5.3nm

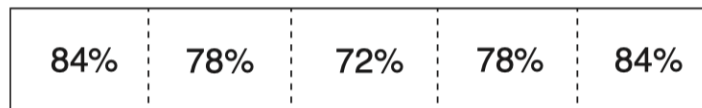
Refragmentation	Max. EPE [nm]	#Segments
No	3.83	7,000
RFC (big data)	2.42	7,096
Rules (big data)	3.07	7,163
RFC (small data)	3.41	7,165
Rules (small data)	3.13	7,177
<b>Revised rules (small data)</b>	<b>2.59</b>	7,168

# Example 2: Placement Utilization

- Very low aspect ratio design
  - Insufficient horizontal routing resources → very low placement utilization (with lots of whitespace)



- Different utilization for different sub-regions
  - Higher utilization towards left- and right-ends; Lower utilization in the center
  - CNN has been used to identify utilization distribution





# Placement Utilization

- Rule-based approach
  - Utilization distribution is assumed to be linear from center to left- or right-end
    - #sub-regions: proportional to % of GRCs with overflow
    - Rules are set up to identify “center utilization” and “linear slope”
- CNN vs Rule-based
  - Large data volume
    - CNN: ~0% overflows ↔ Rule-based: 3% overflows
  - Small data volume
    - CNN: 5.5% overflows ↔ Rule-based: 3.3% overflows

# Summary

- ML is not a silver bullet
  - Lack of training samples with high coverage, in semiconductor industry
- ML model may be a foundation for highly efficient rule-based method
  - Even when training samples are not enough!