



香港中文大學  
The Chinese University of Hong Kong

# Large Scale VLSI Mask Optimization

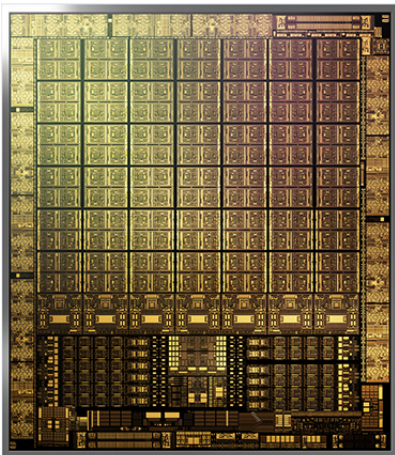
Bei Yu  
Chinese University of Hong Kong  
[byu@cse.cuhk.edu.hk](mailto:byu@cse.cuhk.edu.hk)

July 30, 2023

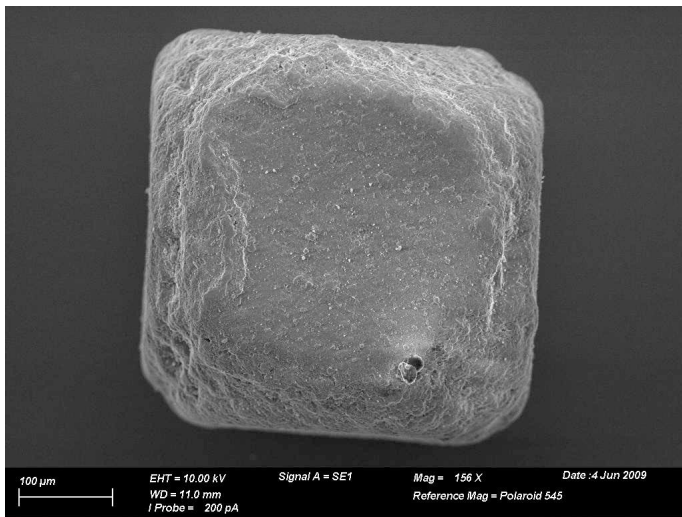


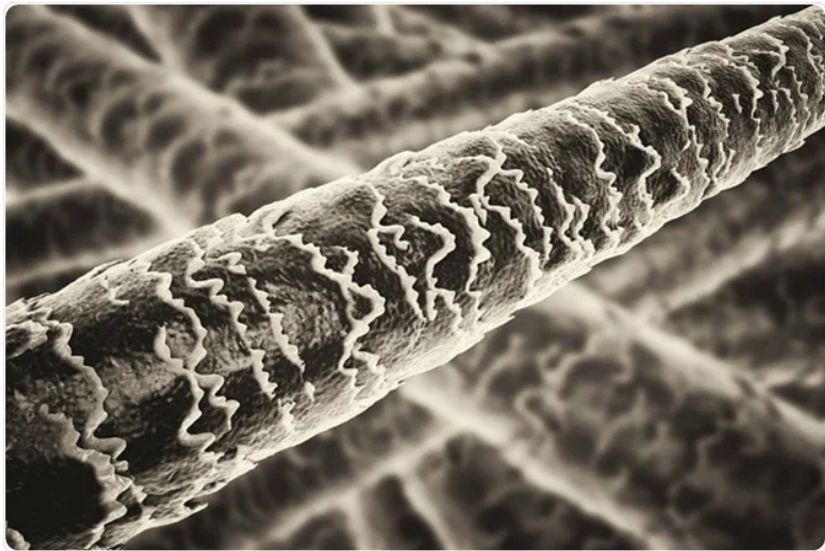
# 1x: Your Laptop



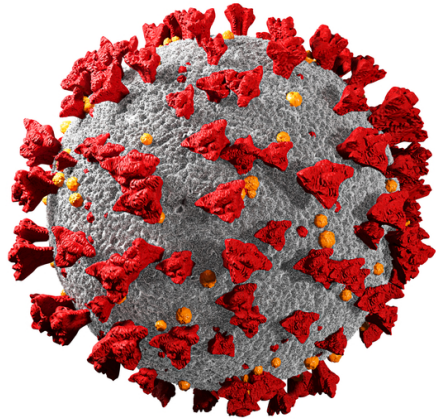






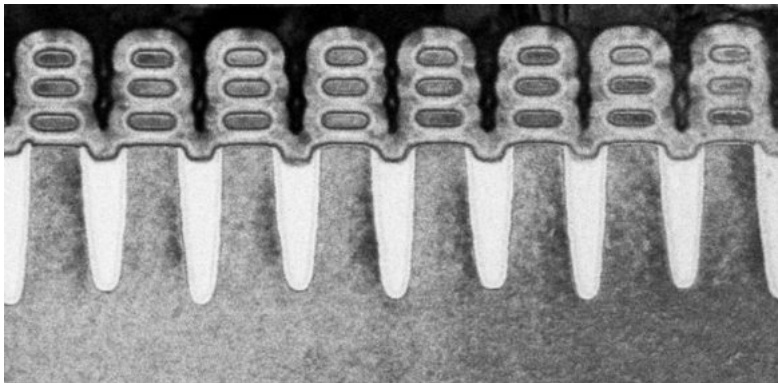


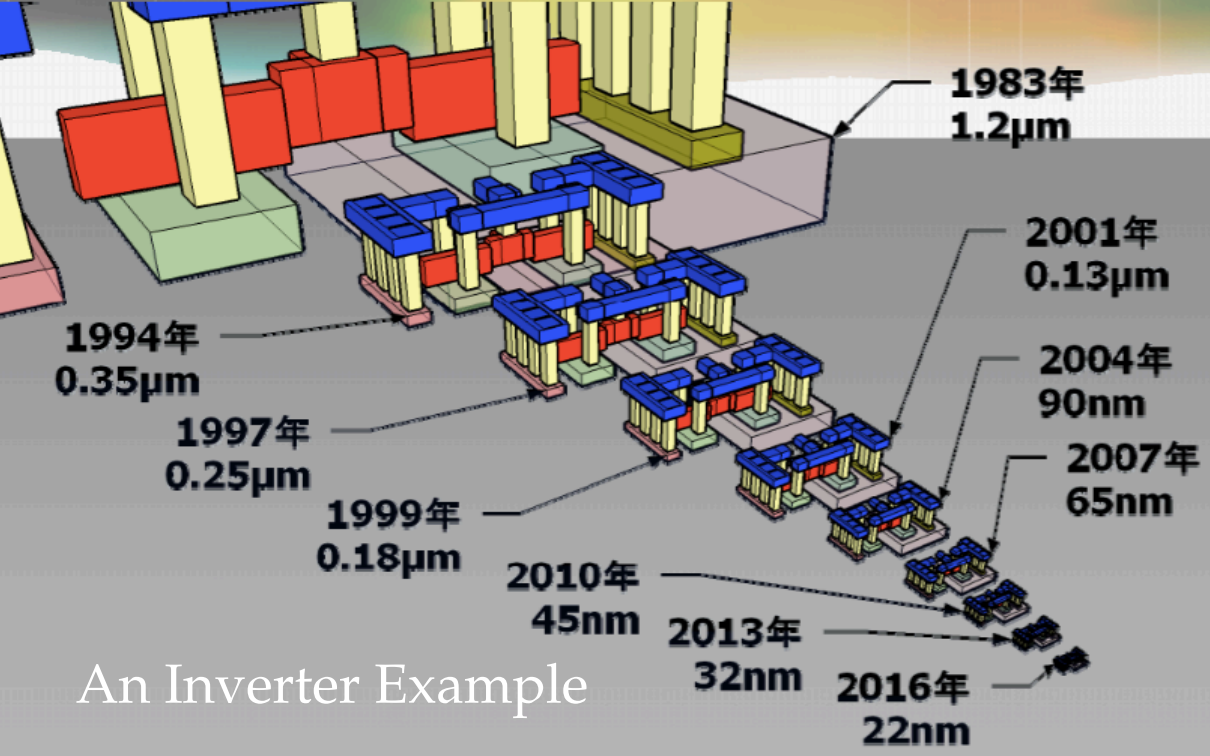






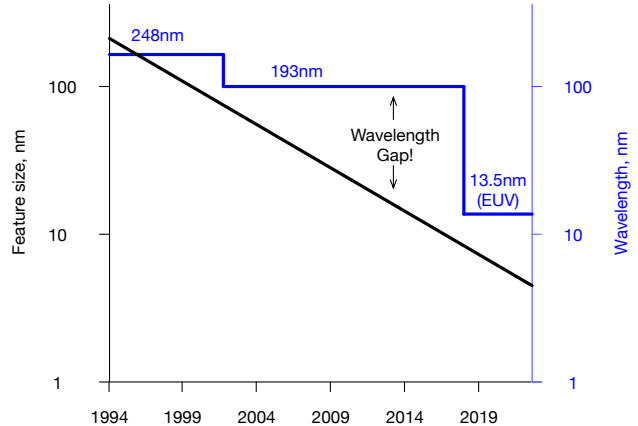
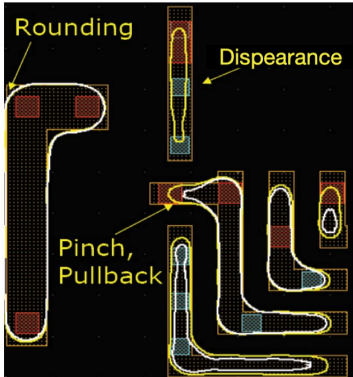
# 50,000,000×: 5nm Transistor

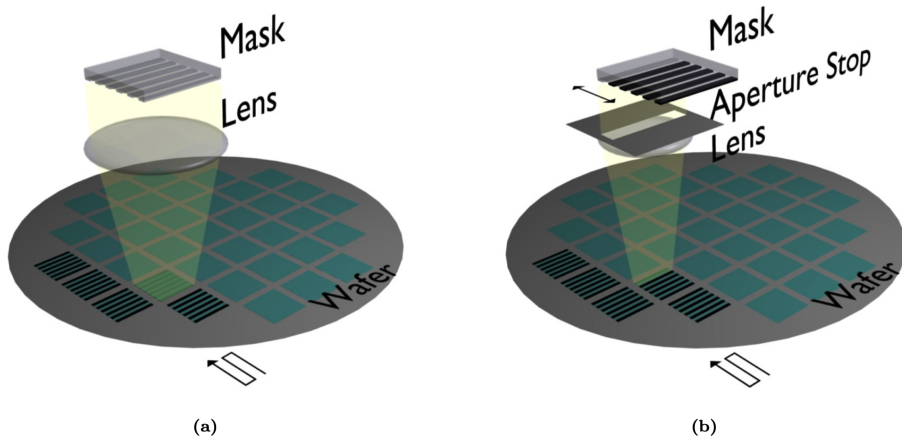




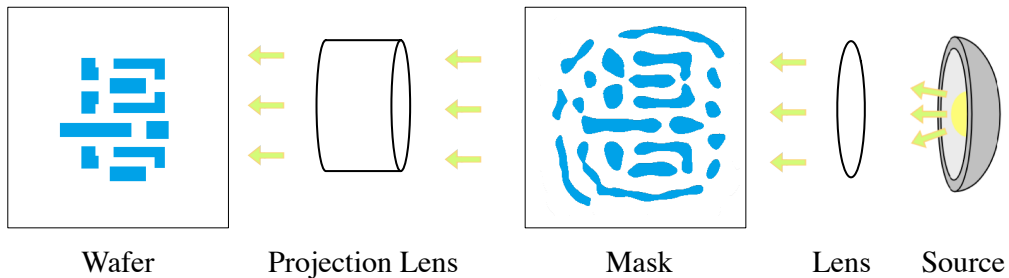
An Inverter Example

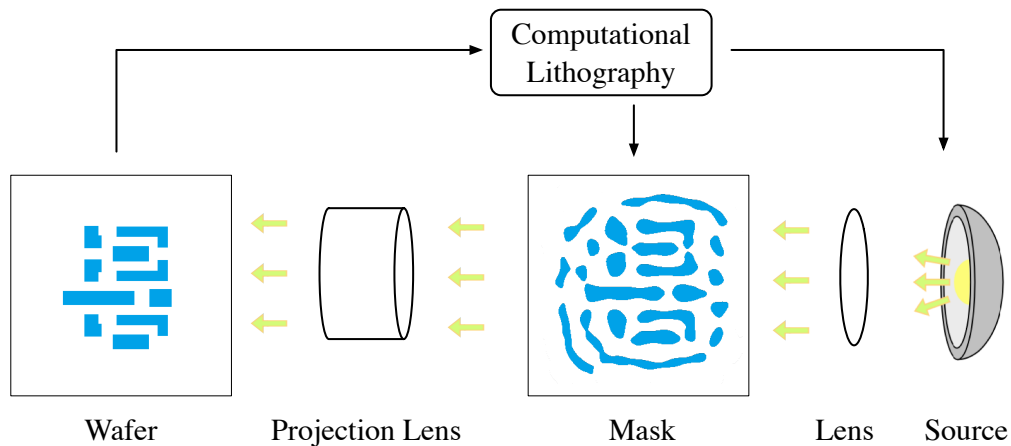
# When feature is small: what you see $\neq$ what you want

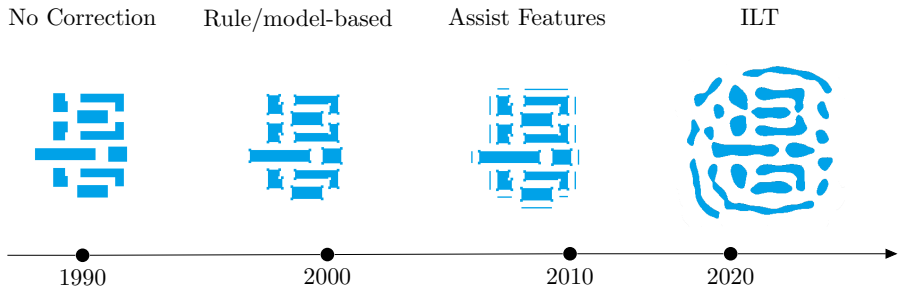




**Figure 2.3:** Exposure procedure of state-of-the-art projection tools: (a) Step-and-repeat: the entire die is exposed in one step; (b) Step-and-scan: the die is gradually exposed by a scanning slit.





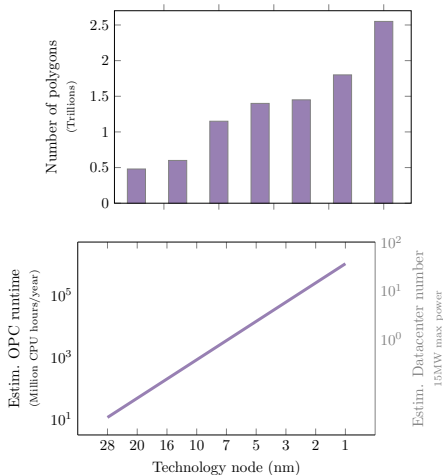
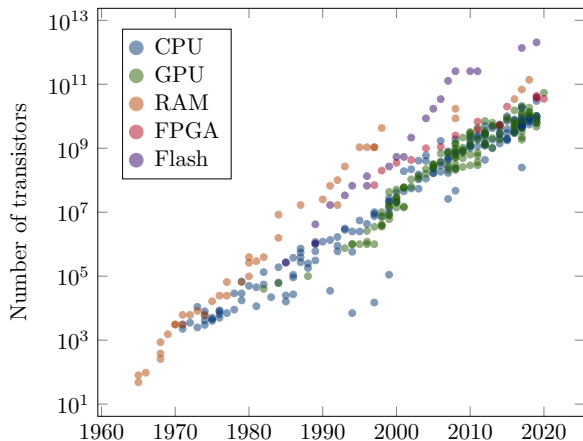


- From simple distortion to complex modification
- Computation becomes more and more complicated

# Moore's Law to Extreme Scaling



- Billions of transistors on a chip → ... Trillions of polygons
- OPC runtime goes to  $10^6$  (million) CPU hours/year



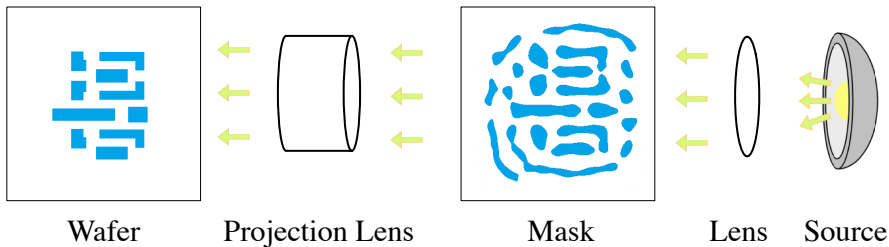
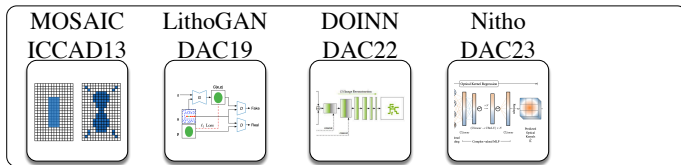




- ① Lithography Modeling
- ② Optical proximity correction (OPC)
- ③ Full-chip OPC

# Lithography Modeling

## Lithography Modeling





## Hopkins Model and Transmission Cross-Coefficient (TCC)

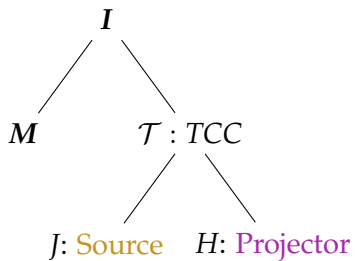
The imaging equation:

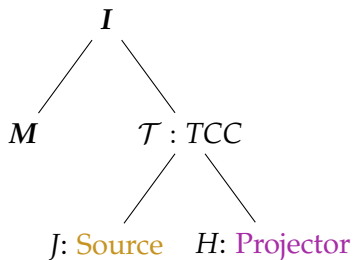
$$\mathcal{F}(\mathbf{I})(f, g) = \iint_{-\infty}^{\infty} \mathcal{T}((f' + f, g' + g), (f', g')) \mathcal{F}(\mathbf{M})(f' + f, g' + g) \mathcal{F}(\mathbf{M})^*(f', g') df' dg', \quad (1)$$

where  $\mathbf{M}$  is the mask,  $(f, g)$  is its frequencies.  $\mathcal{T}$  is TCC given by:

$$\mathcal{T}((f', g'), (f'', g'')) := \iint_{-\infty}^{\infty} \mathcal{F}(\mathbf{J})(f, g) \mathcal{F}(\mathbf{H})(f + f', g + g') \mathcal{F}(\mathbf{H})^*(f + f'', g + g'') df dg, \quad (2)$$

where the weight factor  $J$  solely depends on effective source,  $H$  is projector transfer function.





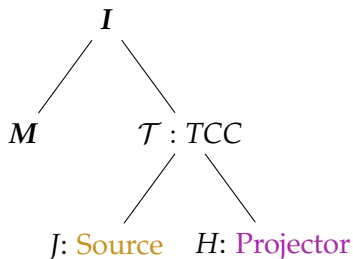
When the projector and source are fixed,

$J$ : constant matrix

$H$ : constant matrix



$\mathcal{T}$ : TCC is a constant matrix



When the projector and source are fixed,

$J$ : constant matrix

$H$ : constant matrix



$\mathcal{T}$ : TCC is a constant matrix

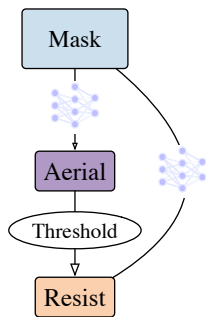


Instead of learning an image-to-image mapping,

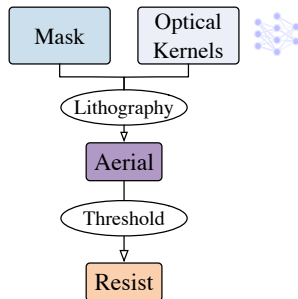
Would it be Possible to learn the **TCC optical kernels**?



- Get rid of negative influence of layer types & dataset distribution.
- Less training data required & smaller model size.



(c) Previous

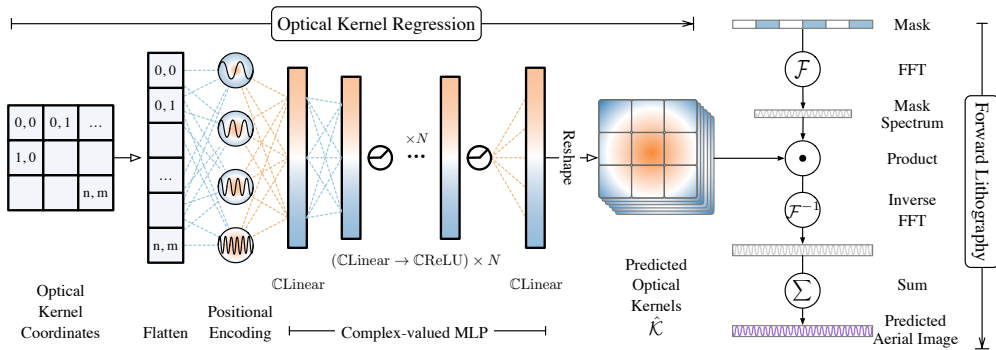


(d) Ours



## Nitho

### Physics-Informed Optical Kernel Regression Using Complex-valued Neural Fields



The overall aerial image prediction pipeline of Nitho framework, which separates mask-related linear operations from optical kernel regression using coordinate-based CMLP.

<sup>1</sup>Guojin Chen, Zehua Pei, et al. (2023). "Physics-Informed Optical Kernel Regression Using Complex-valued Neural Fields". In: *Proc. DAC*.



- SVD Approximation of Partial Coherent System [Cobb,1998]

$$\mathbf{I} = \sum_{k=1}^{N^2} w_k |\mathbf{M} \otimes \mathbf{h}_k|^2. \quad (3)$$

- Reduced Model [Gao+,DAC'14]

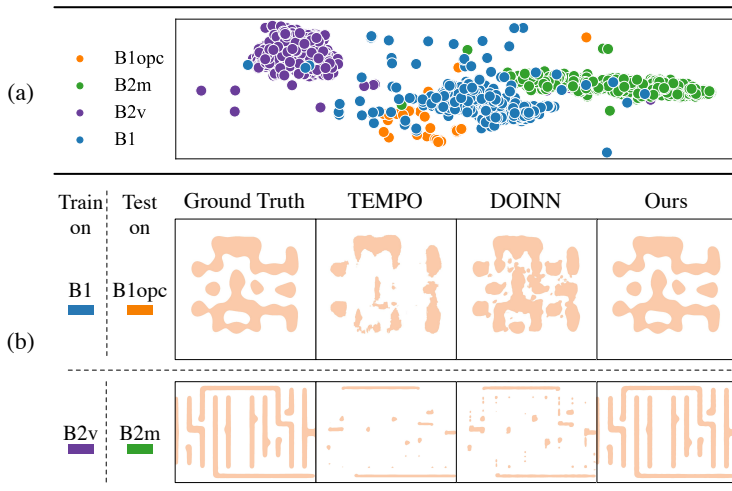
$$\mathbf{I} = \sum_{k=1}^{N_h} w_k |\mathbf{M} \otimes \mathbf{h}_k|^2. \quad (4)$$

- Etch Model

$$\mathbf{Z}(x, y) = \begin{cases} 1, & \text{if } \mathbf{I}(x, y) \geq I_{th}, \\ 0, & \text{if } \mathbf{I}(x, y) < I_{th}. \end{cases} \quad (5)$$

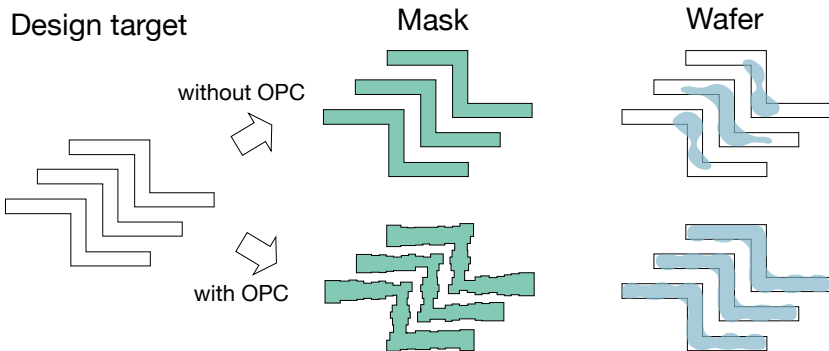
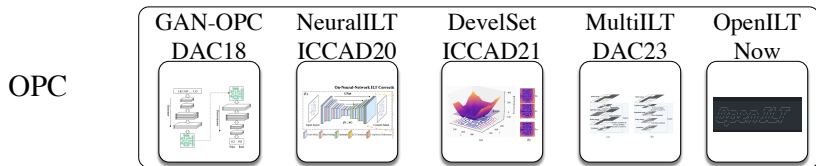
---

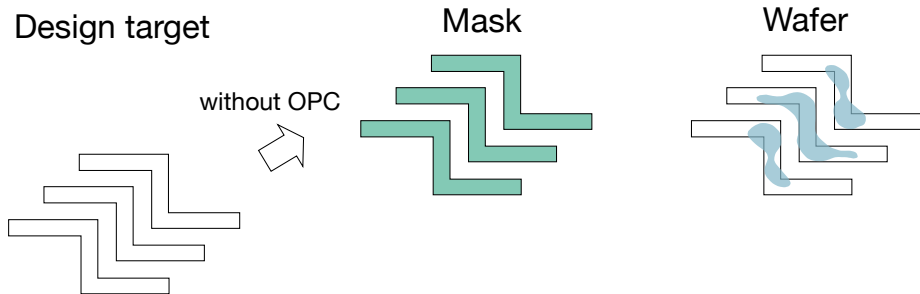
<sup>2</sup>Jih-Rong Gao et al. (2014). "MOSAIC: Mask Optimizing Solution With Process Window Aware Inverse Correction". In: *Proc. DAC*. San Jose, California, 52:1–52:6.

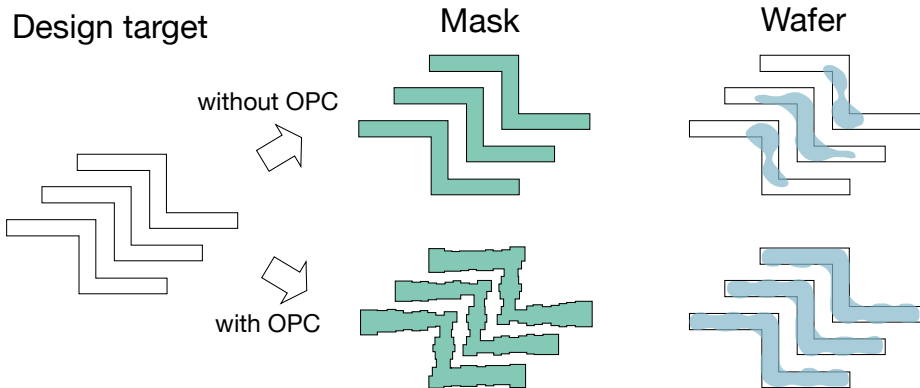


(a) t-SNE distribution of datasets; (b) Comparison of generalization capability on out-of-distribution (OOD) datasets.

# Optical proximity correction (OPC)









The main objective in ILT is minimizing the lithography error through gradient descent.

$$E = \|\mathbf{Z}_t - \mathbf{Z}\|_2^2, \quad (6)$$

where  $\mathbf{Z}_t$  is the target and  $\mathbf{Z}$  is the wafer image of a given mask.

Apply translated sigmoid functions to make the pixel values close to either 0 or 1.

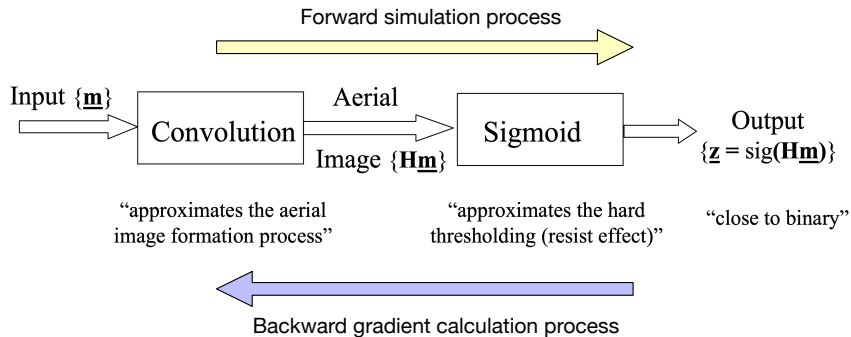
$$\mathbf{Z} = \frac{1}{1 + \exp[-\alpha \times (\mathbf{I} - \mathbf{I}_{th})]}, \quad (7)$$

$$\mathbf{M}_b = \frac{1}{1 + \exp(-\beta \times \mathbf{M})}. \quad (8)$$

Combine Equations (3)–(8) and the analysis in [Poonawala,TIP'07],

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{M}} = & 2\alpha\beta \times \mathbf{M}_b \odot (1 - \mathbf{M}_b) \odot \\ & (((\mathbf{Z} - \mathbf{Z}_t) \odot \mathbf{Z} \odot (1 - \mathbf{Z}) \odot (\mathbf{M}_b \otimes \mathbf{H}^*)) \otimes \mathbf{H} + \\ & ((\mathbf{Z} - \mathbf{Z}_t) \odot \mathbf{Z} \odot (1 - \mathbf{Z}) \odot (\mathbf{M}_b \otimes \mathbf{H})) \otimes \mathbf{H}^*). \end{aligned} \quad (9)$$

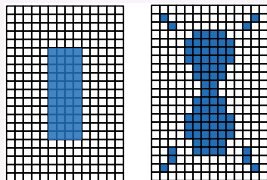






## Typical ILT

- Mask  $\rightarrow$  Image  $\rightarrow$  Matrix
- Calculate gradient on each pixel.



## Level-set method

- Boundary-based update
- Implicit representation; focus on boundaries

$$\begin{cases} \phi(t, \mathbf{x}) < 0 & \text{if } \mathbf{x} \in \Omega(t) \\ \phi(t, \mathbf{x}) = 0 & \text{if } \mathbf{x} \in \Gamma(t) \\ \phi(t, \mathbf{x}) > 0 & \text{if } \mathbf{x} \in \overline{\Omega(t)} \end{cases}$$

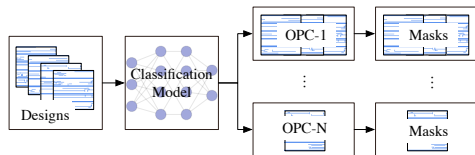
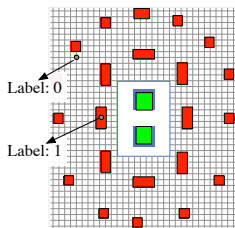
<sup>1</sup>Jhjh-Rong Gao et al. (2014). "MOSAIC: Mask Optimizing Solution With Process Window Aware Inverse Correction". In: *Proc. DAC*. San Jose, California, 52:1–52:6.

<sup>2</sup>Yuzhe Ma et al. (2017). "A Unified Framework for Simultaneous Layout Decomposition and Mask Optimization". In: *Proc. ICCAD*, pp. 81–88.

<sup>3</sup>Ziyang Yu et al. (2021). "A GPU-enabled Level Set Method for Mask Optimization". In: *Proc. DATE*.

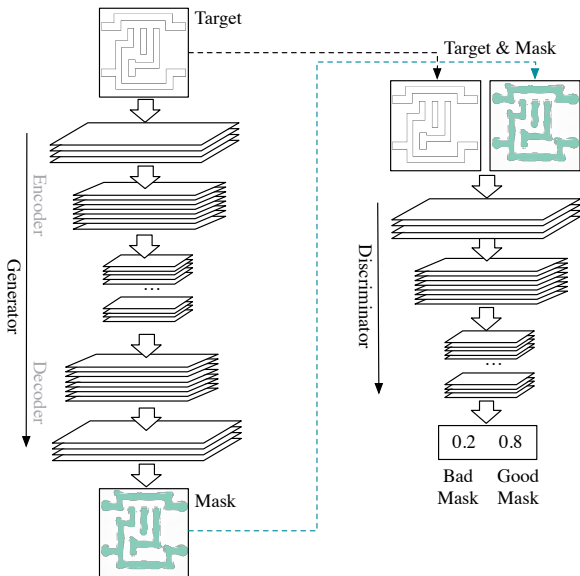
## Discriminative models [TCAD'20]<sup>4</sup> [ASPdac'20]<sup>5</sup>

- Pixel-wise classification
- Printed image estimation/quality estimation



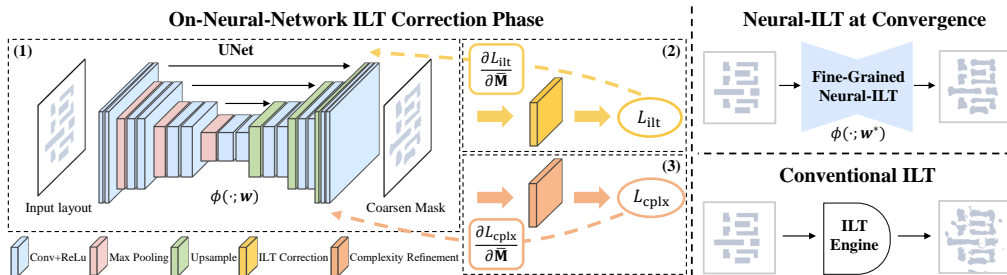
<sup>4</sup>Hao Geng et al. (2020). "SRAF Insertion via Supervised Dictionary Learning". In: *IEEE TCAD*.

<sup>5</sup>Haoyu Yang, Wei Zhong, et al. (2020). "VLSI Mask Optimization: From Shallow To Deep Learning". In: *Proc. ASPDAC*, pp. 434–439.

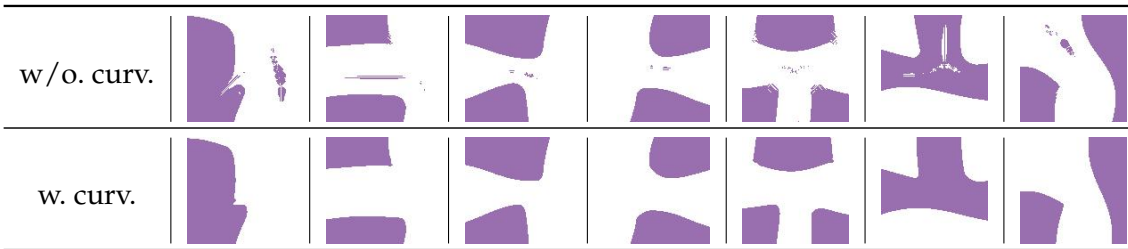
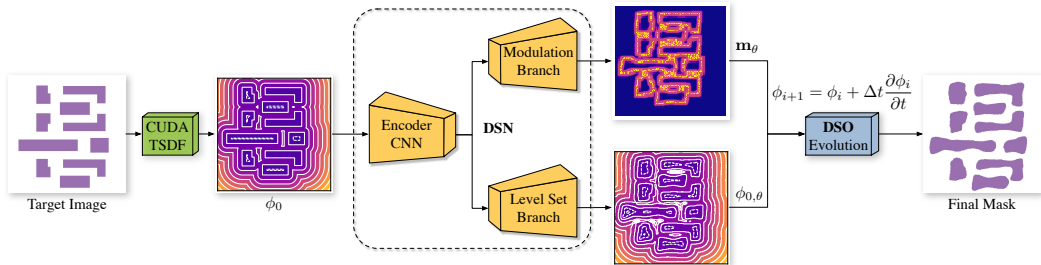


<sup>6</sup>Haoyu Yang, Shuhe Li, et al. (2018). "GAN-OPC: Mask Optimization with Lithography-guided Generative Adversarial Nets". In: *Proc. DAC*, 131:1–131:6.

- A pre-trained UNet for performing layout-to-mask translation.
- An ILT correction layer for minimizing inverse lithography loss.
- A mask complexity refinement layer for removing redundant complex features.

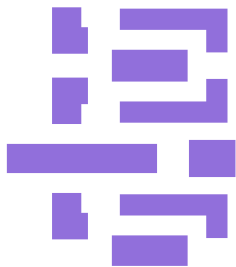


<sup>7</sup>Bentian Jiang et al. (2020). "Neural-ILT: Migrating ILT to Neural Networks for Mask Printability and Complexity Co-optimization". In: *Proc. ICCAD*.

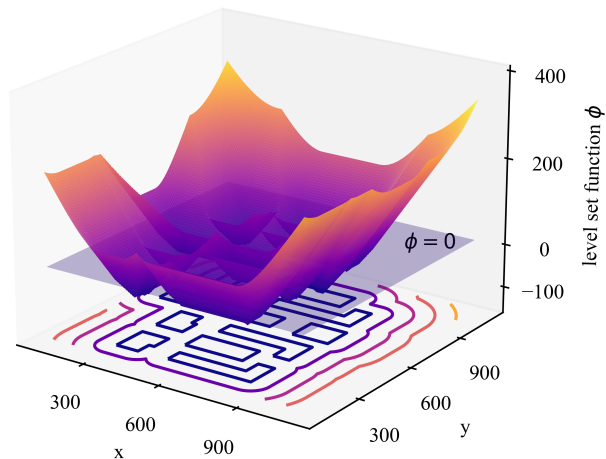


<sup>8</sup>Guojin Chen, Ziyang Yu, et al. (2021). “DevelSet: Deep neural level set for instant mask optimization”. In: *Proc. ICCAD*.

The level set function is defined as: min distance of each point to the boundary.



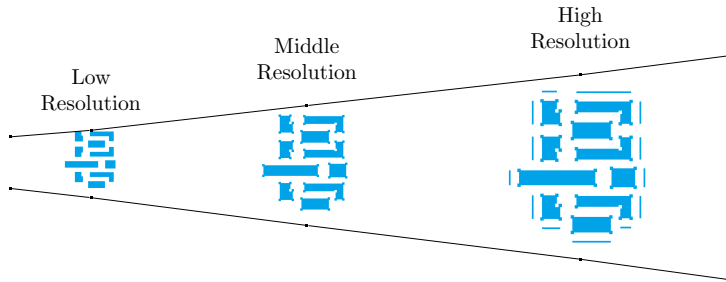
(a) target



(b) levelset function



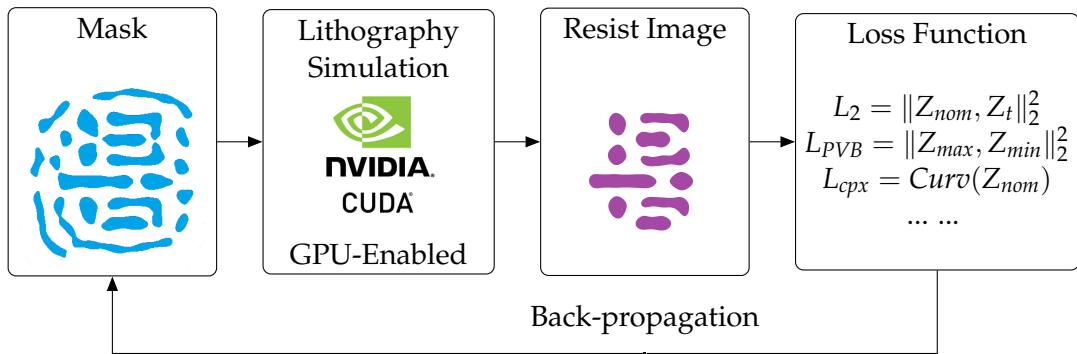
- From low resolution to high resolution
- Save runtime due to less complex computation in low resolution



<sup>9</sup>Shuyuan Sun et al. (2023). "Efficient ILT via Multi-level Lithography Simulation". In: *Proc. DAC.35/49*



- GPU Acceleration → Better Runtime
- PyTorch-Based → Easier Development
- Open-Source → More Accessible



[github.com/OpenOPC/OpenILT/](https://github.com/OpenOPC/OpenILT/)

☰ README.md



## OpenILT: An Open-source Platform for Inverse Lithography Technology Research

OpenILT is an open-source platform for inverse lithography technology (ILT) research. It has a comprehensive and flexible ecosystem of libraries that enable the efficient development and evaluation of ILT algorithms. OpenILT decouples the ILT flow into different components, lithography simulation, initialization, optimization, and evaluation. ILT researchers can implement and evaluate their ideas quickly by replacing a component with a novel method. Moreover, the platform is implemented with *pytorch*, which enables easy GPU acceleration and deep-learning integration.

- Easy Implementation of ILT Methods

```
cfg, litho = SimpleCfg(), LithoSim()  
solver = SimpleILT(cfg, litho)  
design = Design("M1_test1.glp")  
Zt,P = PixelInit().run(design)  
l2,pvb,P,M = solver.solve(Zt,P)  
l2,pvb,epe,shot = evaluate(M,Zt,litho)
```



(a) Target Image



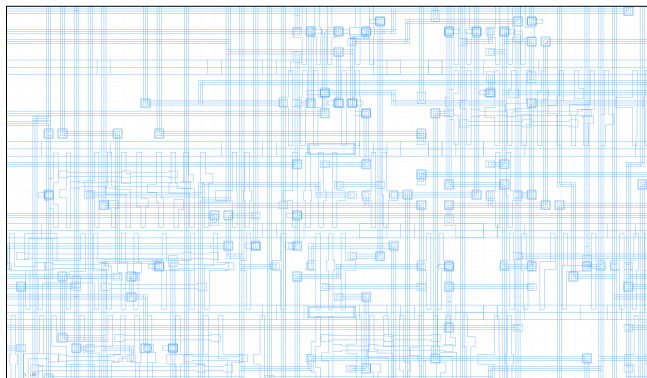
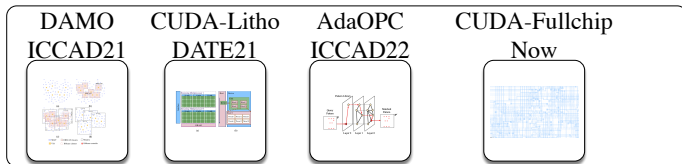
(b) LevelSet

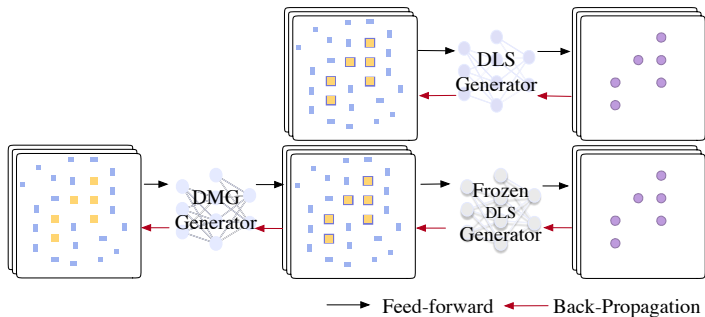


(c) MultiLevel

# Full-chip OPC

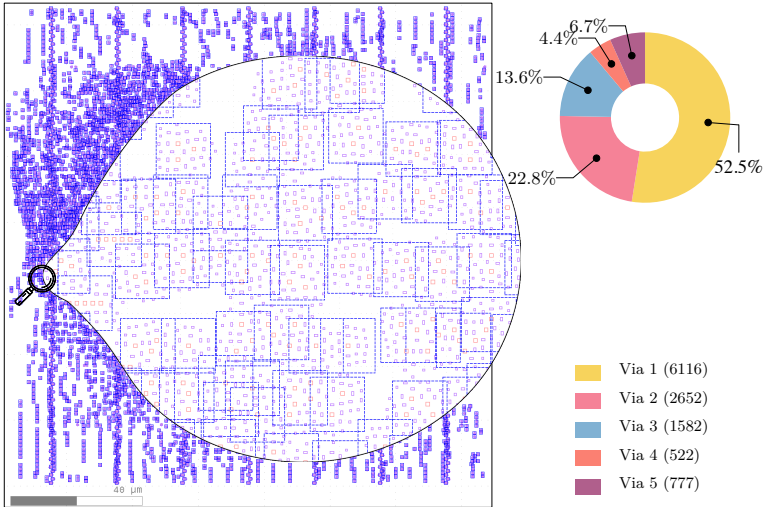
Fullchip  
OPC

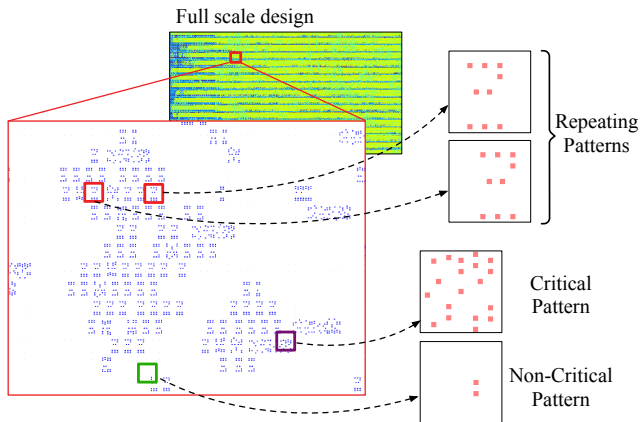




	GAN-OPC			Calibre			DMO		
	$L_2$ (nm)	PV Band ( $nm^2$ )	Runtime (s)	$L_2$ (nm)	PV Band ( $nm^2$ )	Runtime (s)	$L_2$ (nm)	PV Band ( $nm^2$ )	Runtime (s)
case 1	7456	11424	284	5159	11671	1417	<b>4631</b>	<b>11166</b>	352
case 2	7321	11215	281	4987	11463	1406	<b>4432</b>	<b>10955</b>	336
case 3	7102	11265	285	5420	11516	1435	<b>4802</b>	<b>11032</b>	367
case 4	8032	11642	322	5382	11910	1606	<b>4835</b>	<b>11265</b>	399
Average	7478	11386	293	5237	11640	1466	<b>4675</b>	<b>11104</b>	363
Ratio	1.60	1.03	0.80	1.12	1.05	4.04	<b>1.00</b>	<b>1.00</b>	1.00

<sup>8</sup>Guojin Chen, Wanli Chen, et al. (2020). "DAMO: Deep Agile Mask Optimization for Full Chip Scale". In: *Proc. ICCAD*.



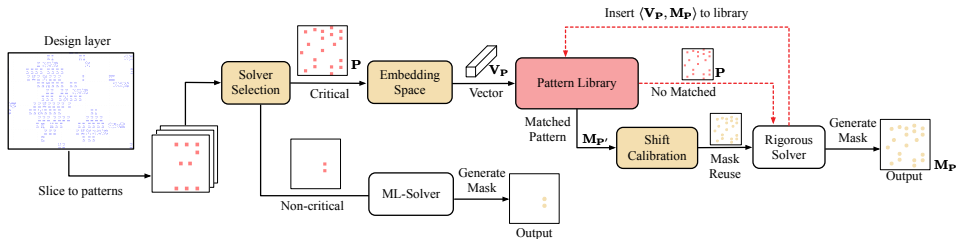


- Uneven scattering. → Solver selection
- Large ratio of repetition. → Mask Reuse

<sup>9</sup>Wenqian Zhao et al. (2022). “AdaOPC: A Self-Adaptive Mask Optimization Framework For Real Design Patterns”. In: *Proc. ICCAD*.



- Superior in both speed& performance



Test Case ID	DAMO-DGS			ILT-GPU			AdaOPC		
	#EPE	PVB ( $nm^2$ )	RT (s)	#EPE	PVB ( $nm^2$ )	RT (s)	#EPE	PVB ( $nm^2$ )	RT (s)
1	22	23323	5.20	23	23329	41.15	22	23232	5.50
2	26	26729	5.26	25	26762	48.5	24	26580	5.41
3	27	26938	5.22	24	26720	55.92	24	26718	5.37
4	36	27975	5.18	29	28127	70.57	25	27934	5.40
5	35	28805	5.32	30	28925	66.89	30	28927	5.44
6	30	26960	5.31	25	26762	55.81	24	26775	5.38
7	33	26382	5.23	28	26453	59.47	28	26281	5.43
8	32	30646	5.38	25	29450	54.88	27	29341	5.42
9	25	24054	5.25	24	24053	70.62	23	24022	5.43
10	24	21939	5.29	23	21701	37.59	22	21644	5.53
Avg.	29.0	26375	<b>5.26</b>	25.6	26228	56.14	<b>24.9</b>	<b>26145</b>	5.43
Ratio	1.165	1.009	<b>0.970</b>	1.028	1.003	10.340	<b>1.000</b>	<b>1.000</b>	1.000

# Roadmap of Full-chip OPC Research



20nm AF  
20nm main feature

SPIE, Leo Pang  
Academic Guidelines

2010s

DATE21  
GLS-ILT  
CUDA-  
Accelerated  
Lithography  
Modeling

(a) (b)

2021

Metal Layer: ML solution for industrial layout

Now



2020

ICCAD20  
DAMO  
First Full-chip  
ML solution

(a) (b) (c) (d)

■ MLP ■ DBSCAN clusters ■ Waiver  
■ VSA ■ KMeans clusters ■ KMeans outliers

2022

ICCAD22  
AdaOPC  
First Adaptive  
Full-chip  
OPC solution

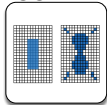
Query Pattern  $P$  Pattern Library Layer 2 Layer 1 Layer 0 Matched Pattern  $P'$

# Conclusion

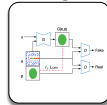


Lithography  
Modeling

MOSAIC  
ICCAD13



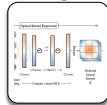
LithoGAN  
DAC19



DOINN  
DAC22



Nitho  
DAC23

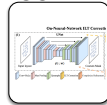


OPC

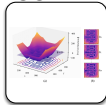
GAN-OPC  
DAC18



NeuralILT  
ICCAD20



DevelSet  
ICCAD21



MultiILT  
DAC23



OpenILT  
Now

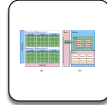


Fullchip  
OPC

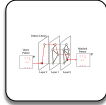
DAMO  
ICCAD21



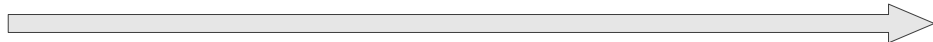
CUDA-Litho  
DATE21

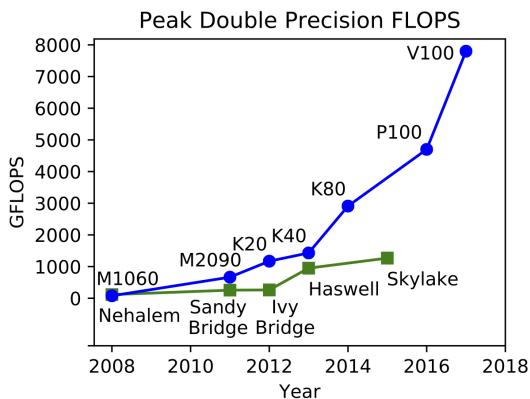


AdaOPC  
ICCAD22

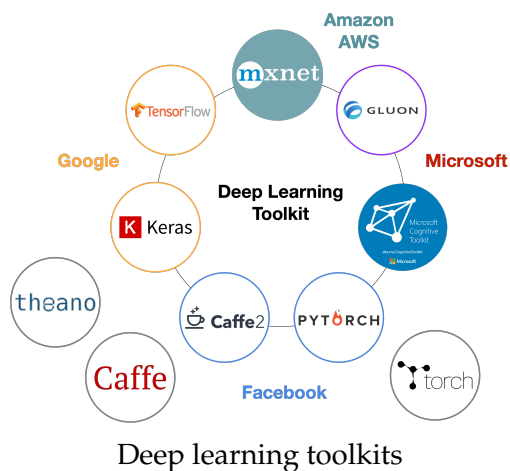


CUDA-Fullchip  
Now

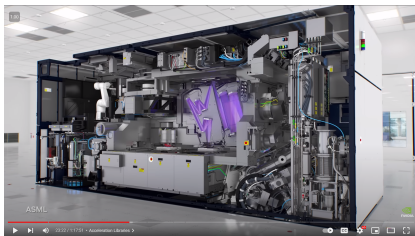




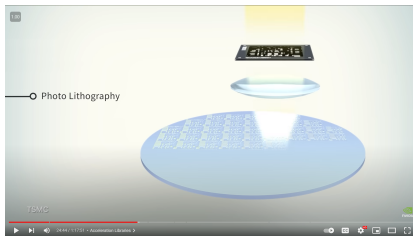
Over **60x** speedup in neural network training since 2013



# Start from NVIDIA GTC CuLitho



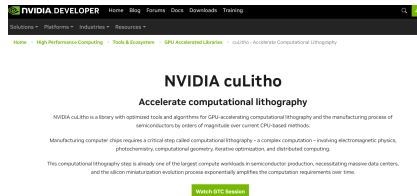
(a)



(b)

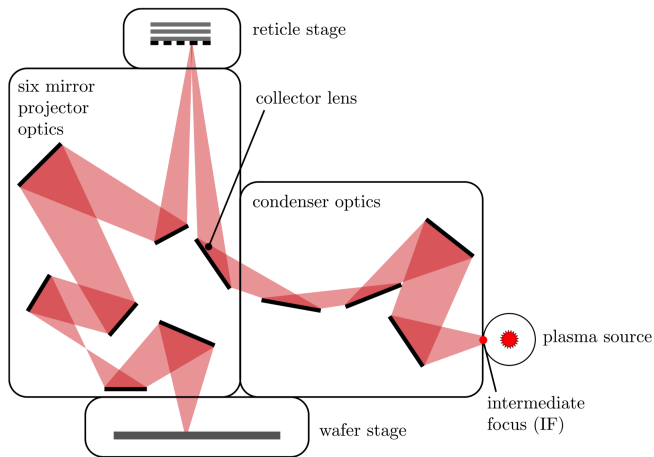


(c)



(d)

(a) EUV lithograph machine (ASML). (b) Photo lithography. (c) and (d) Nvidia-cuLitho has 40× acceleration.



**Figure 2.16:** EUVL projection system employing all-reflective optical components. In this example, a six-mirror projection system suited for NAs of up to 0.35 is depicted. In order to reduce down-times in case of maintenance, condenser, projector, and the mask and wafer stages are placed into individual vacuum chambers.

How will this mask print?