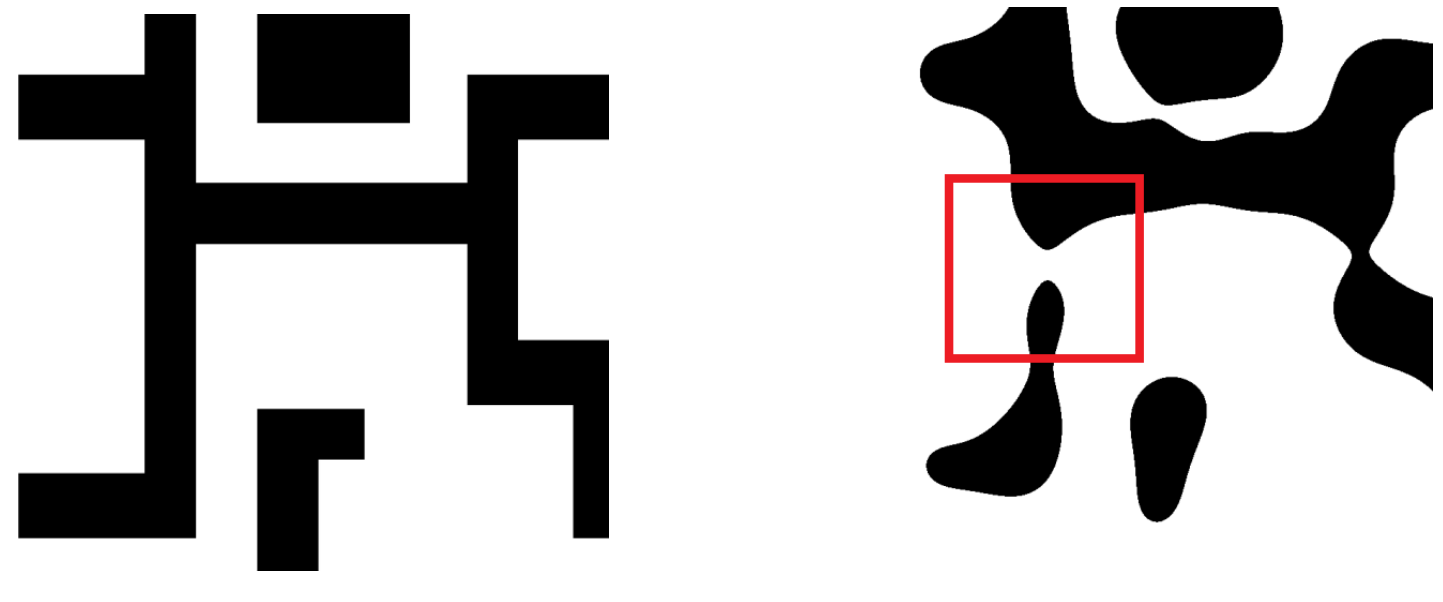


## Background



- What you see  $\neq$  what you get
- RETs: OPC, SRAF, MPL
- Still exists hotspots: low fidelity patterns
- Lithography simulation: time consuming

### Pattern Matching-based Hotspot Detection

- Characterize the hotspots as explicit patterns and identify the hotspots by matching these patterns
- [Yu+, ICCAD'14] [Nosato+, JM3'14] [Kahng+, SPIE'06] [Su+, TCAD'15] [Wen+, TCAD'14] [Yang+, TCAD'17]
- Fast but hard to detect unseen patterns

### Machine Learning based Hotspot Detection

- Build implicit models by learning from existing training data
  - SVM, Bayesian, Decision-tree, Boosting, NN, ...
- [Ding+, ASPDAC'11] [Yu+, DAC'13] [Matsunawa+, SPIE'15] [Zhang+, ICCAD'16] [Wen+, TCAD'14]
- Possible to detect the unseen hotspots but may cause false alarm issues

### Deep Learning based Hotspot Detection

- Belongs to ML-based hotspot detection but different from conventional ML models:
  - Feature Crafting v.s. Feature Learning
  - Stronger scalability
- [Yang+, DAC'17]
- Drawback: not storage and computational efficient

## Preliminaries

### Problem Formulation

#### Definition: Accuracy

The ratio of correctly predicted hotspots among the set of actual hotspots.

$$Accuracy = \frac{\#TP}{\#TP + \#FN}$$

#### Definition: False Alarm

The number of incorrectly predicted non-hotspots.

$$False\ Alarm = \#FP$$

#### Definition: ODST

Abbreviation of Overall Detection and Simulation Time. The sum of the lithography simulation time for layout patterns detected as hotspots (including real hotspots and false alarms) and the learning model evaluation time.

$$ODST = (\#FP + \#TP) t_{ls} + (\#TN + \#FN + \#FP + \#TP) t_{ev}$$

where  $t_{ls}$  is lithography simulation time per instance and  $t_{ev}$  is the model evaluation time per instance.

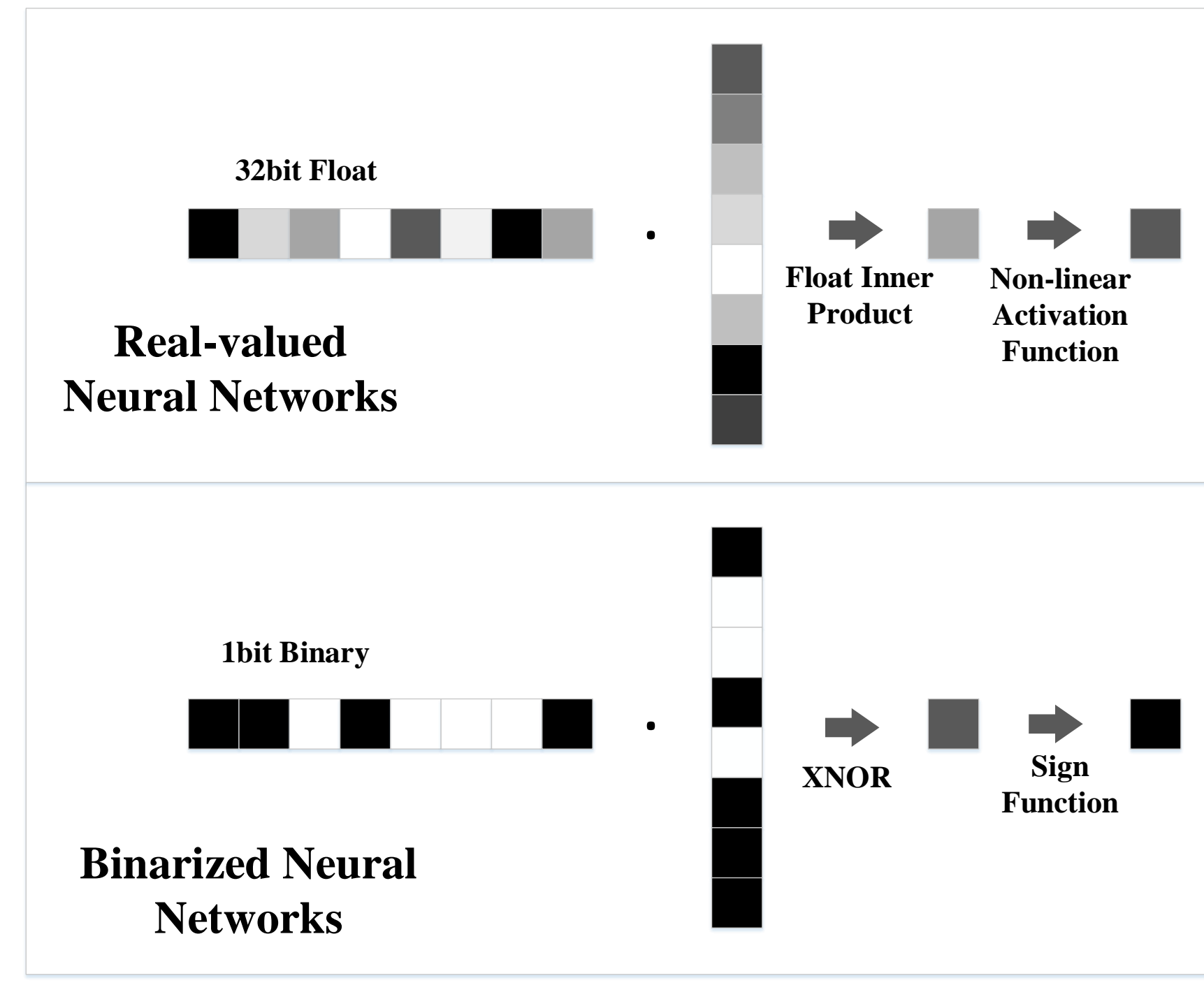
### Problem1: Hotspot Detection

Given a dataset that contains hotspot and non-hotspot instances, train a classifier that can maximize the accuracy and minimize the false alarm.

### Parameter Quantization

- Problem with deep neural networks:
  - Enormous computational and storage consumption
- To alleviate this problem:
  - Parameter Quantization
  - 32-bit floating-point weights not necessary: quantized to fixed-point of 8-bit, 3-bit, 1-bit...
  - [Arora+, ICML'14] [Hwang+, SIPS'14] [Soudry+, ANIPS'14] [Rastegari+, ECCV'16]

## Binarized Neural Networks



- Binarized neural network (BNN):
  - Extremely quantized to 1 bit
  - Inherently suitable for hardware implementation
- Layout patterns are binary images
  - BNN might be suitable for that

### Binarization Approach

#### Definition

Let  $W$  be the kernel which is an  $n$ -element vector and  $X$  be the vector of the corresponding block in the input tensor,  $n = w_k \times h_k$ . Let  $W_B, X_B$  be the binarized kernel and input vector and  $\alpha_W, \alpha_X$  be the corresponding scaling factors. Here  $W, X \in \mathbb{R}^n$ ,  $W_B, X_B \in \{-1, +1\}^n$  and  $\alpha_W, \alpha_X \in \mathbb{R}^+$ .

#### Problem2: Binarization

Given the kernel and input vector  $W, X$ , find best  $W_B, X_B, \alpha_W, \alpha_X$  that minimizes the binarization loss  $L_i$ .  $L_i(W_B, X_B, \alpha_W, \alpha_X) = \|W \odot X - \alpha_W W_B \odot \alpha_X X_B\|^2$  where  $\odot$  means inner product.

- Solving the minimization problem:

$$W_B^* = \text{sign}(W), X_B^* = \text{sign}(X)$$

$$\alpha_W^* = \frac{1}{n} \|W\|_{l_1}, \alpha_X^* = \frac{1}{n} \|X\|_{l_1}$$

- The estimated weight and corresponding input vector  $\tilde{W}, \tilde{X}$  are:

$$\tilde{W} = \frac{1}{n} \text{sign}(W) \|W\|_{l_1}$$

$$\tilde{X} = \frac{1}{n} \text{sign}(X) \|X\|_{l_1}$$

## Training BNN

- Gradient for  $\text{sign}$  function [Hubara, 2016]

$$\frac{\partial \text{sign}(x)}{\partial x} = \mathbf{1}_{\|W\| < 1}$$

- Back propagation through the Binarized Layer

$$\frac{\partial l}{\partial W} = \frac{\partial l}{\partial \tilde{W}} \frac{\partial \tilde{W}}{\partial W}$$

$$= \frac{\partial l}{\partial \tilde{W}} \frac{\partial (\frac{1}{n} \|W\|_{l_1} \text{sign}(W))}{\partial W}$$

$$= \frac{\partial l}{\partial \tilde{W}} (\frac{1}{n} + \alpha_W^* \mathbf{1}_{\|W\| < 1})$$

### Algorithm 1 Training a BNN

**Input:**  $(\mathcal{T}_0, Y)$ : a minibatch of input tensors and labels;

- $l(Y, Y_{out})$ : loss function;
- $\mathcal{W}^t$ : current real-valued weight;
- $L$ : number of layers;
- $n$ : kernel size of layers;
- $\eta^t$ : current learning rate;

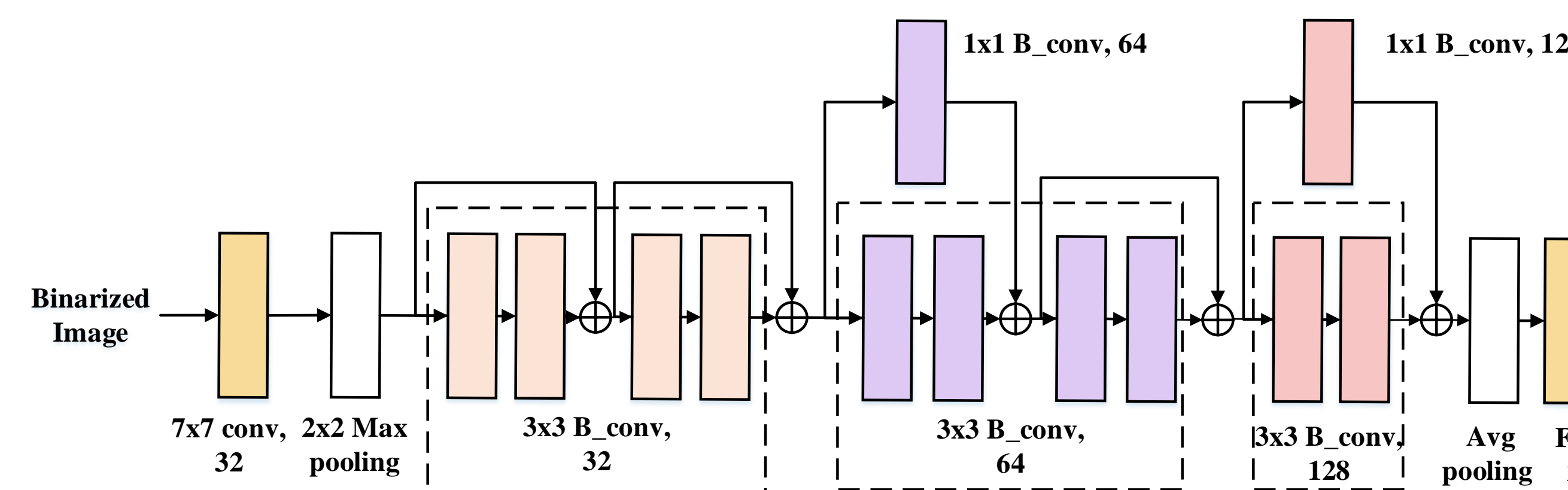
**Output:**  $\mathcal{W}^{t+1}$ : updated real-valued weight;  $\eta^{t+1}$ : updated learning rate.

1. Forward Process:
2. Backward Process:
3. Update Parameters:
4. Return

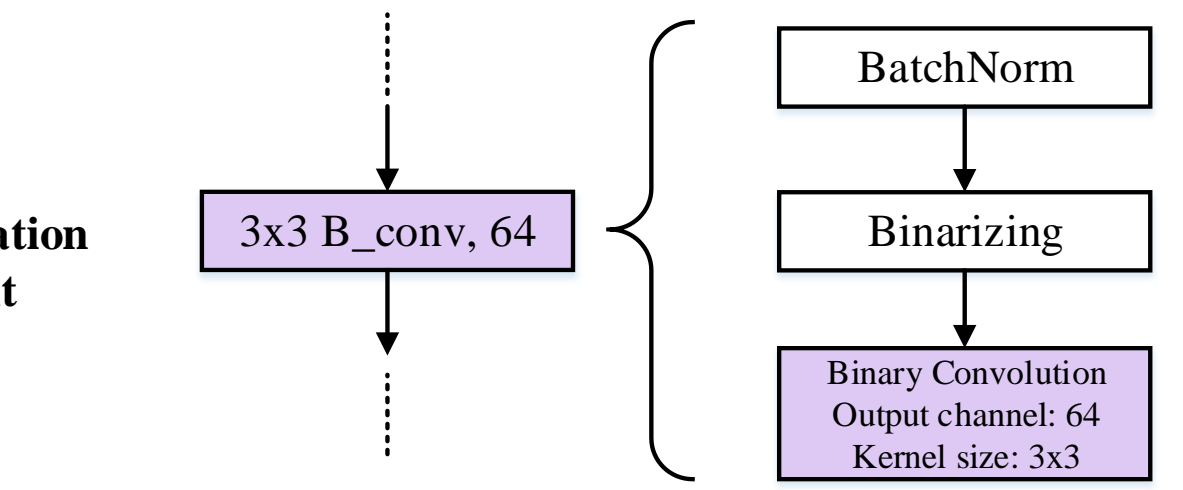
- Data preprocessing
  - Down-sampled to 128x128
- Training hyper-parameters
  - Batch size: 128
  - Learning rate: Initial 0.15, exponentially decay each time loss plateaus
  - Optimizer: NAdam optimizer [Dozat, 2016]
  - Initializer: Xavier initializer [Glorot, 2010]

## Network Architecture

### Residual block-based architecture



### Typical BNN block structure



## Experimental Results

- Benchmark: ICCAD 2012 Contest

Benchmark	# Train HS	# Train NHS	# Test HS	# Test NHS
ICCAD	1204	17096	2524	13503

- Experimental Results on ICCAD 2012 Contest

Method	Accuracy (%)	False Alarm #	Runtime (s)
SPIE'15	84.2	2919	2672
ICCAD'16	97.7	4497	1052
DAC'17	98.2	3413	482
Ours	<b>99.2</b>	<b>2787</b>	<b>60</b>

- Performance Comparisons with Previous Hotspot Detectors

- Accuracy improved from 84.2% to 99.2%
- Least False Alarms: 2787
- Lowest Runtime: 60s, 8x faster

### Conclusions

- Propose a BNN-based architecture to speed up the neural networks in hotspot detection
- Our architecture outperforms previous hotspot detectors and achieves an 8x speedup over the best deep learning-based hotspot detector.