

# Deep Neural Networks' Applications in EDA and Their Acceleration Techniques

Yiran Chen

Electrical and Computing Engineering Department

Duke Center for Evolutionary Intelligence (CEI)

NSF IUCRC For Alternative Sustainable and Intelligent Computing (ASIC)

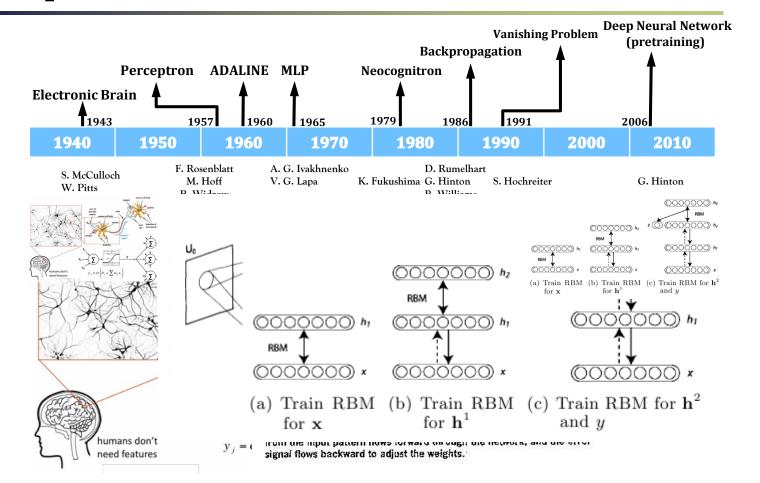
#### **Outline**

- Introduction
- Learning Structured Sparsity in Deep Neural Networks
  - NIPS 2016
- TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning
  - NIPS 2017 (oral)
- ML Application Example in EDA: RouteNet
- Our prospective

# Introduction

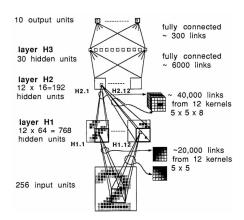


### **Development of Neural Networks**



#### Rise and Decline of Neural Network

## Convolutional Network (1980s)



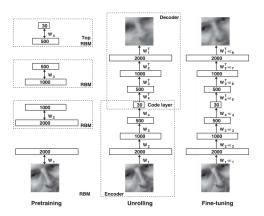


# Dark period (1990s)

- Serious problem: Vanishing gradient
- No benefits observed by adding more layers
- No high performance computing devices



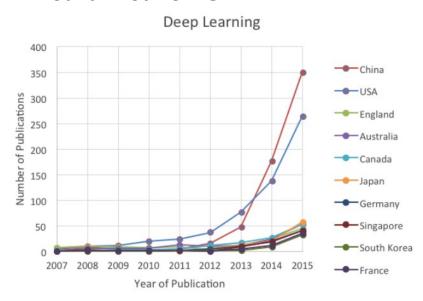
## Renaissance (2006 ~ Present)





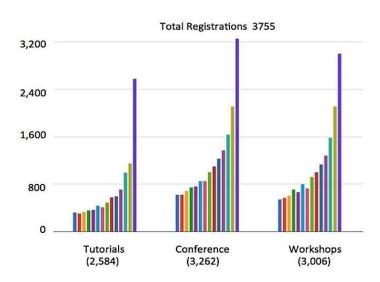
### **Machine Learning in Academia**

# Journal articles mentioning "deep learning" or "deep neural networks"



Source: Office of Science and Technology Policy/The White House

# NIPS registrations growth 2015: 3755, 2016: 6000+



NIPS does not organize sponsor presentations at the Conference or Workshops an sponsors. Satellite meetings should be organized by sponsors immediatley before

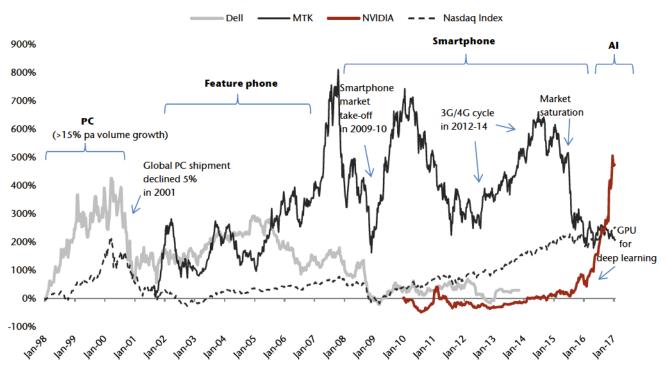
#### 2017 Levels of Sponsorship:

All Levels of Sponsorships are SOLD OUT.

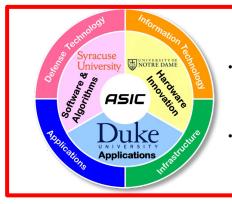
## **Machine Learning in the Market**

**Technology cycle** - from *PC*, to *smartphone*, to *artificial intelligence*?

#### "Pure Play" Share Price Performance



## **NSF IUCRC ASIC Center**



#### What is ASIC?

- The Alternative Sustainable and Intelligent Computing (ASIC) Center is a multi-site, multidisciplinary consortium that explores research frontiers in emerging computing platforms for cognitive applications
- The **ASIC** Center focuses on designing alternative computing platforms for cognitive applications, which perform inefficiently on conventional von Neumann architectures



#### Members include influential faculty across the three research sites:





Yiyu Shi
Site Director

Sharon Hu
Site Co-Director

Michael Niemier
Site Co-Director

Michael Niemier
Site Co-Director

Chaoli Wang

cādence





#### **Facts of ASIC**

**22** 

AFFILIATED FACULTY MEMBERS, INCLUDING 6 FEMALE PROFESSORS 1

MEMBER OF NATIONAL ACADEMY OF ENGINEERING

8

DIFFERENT DEPARTMENTS, INCLUDING ECE, CS, ACMS, BIO, MAE, MATH, PHYS AND SIS 1

ACM FELLOWS AND 3 ACM DISTINGUISHED MEMBERS

13

RECIPIENTS OF NSF CAREER AWARD



**IEEE FELLOWS** 

## Research of ASIC: DNN Acceleration

#### **RESEARCH TOPICS**

Weight quantization

Compact models

**Network compression** 

**Network pruning** 

**Distributed learning** 

#### Representative Industrial Impacts



☐ Our 1-level quantization method (ASP-DAC'17 and DAC'16) is included in the latest SDK of IBM TrueNorth Chip, achieving 6X performance improvement and/or 2/3 hardware cost reduction.



☐ Our structural pruning technique (NIPS'16)



- is supported by the library of Intel Nervana Neural Network Processors.
- · is adopted by Intel NLP accelerator.
- is adopted by SF-Technology, achieving 2X performance improvement.



 Our TenGrad technique (NIPS'17) is supported by Facebook Caffe2 and HP parameter server product for distributed learning.

**Academic Recognitions:** 

ASP-DAC'17 Best Paper Award NIPS'17 Oral Presentation (40 out of 3240 submissions)

DAC'16 Best Paper Nomination

## Research of ASIC: AI Computing Platforms

CPU



Optimizing training and inference on CPU platforms for high cost efficiency. Achieved 3.1-7.3X speedup on Intel Atom, Xeon, and Xeon Phi (ICLR'17).

**GPU** 



Systematic accelerations of deep learning on GPU and heterogeneous platforms. Received 3<sup>rd</sup> Place of 2018 IEEE Low-Power Image Recognition Challenge (track 2).

**FPGA** 



Comprehensive technical portfolio of FPGA-based deep learning, including the first RNN acceleration work published in major FPGA conferences (FCCM 2015).

**ASIC** 



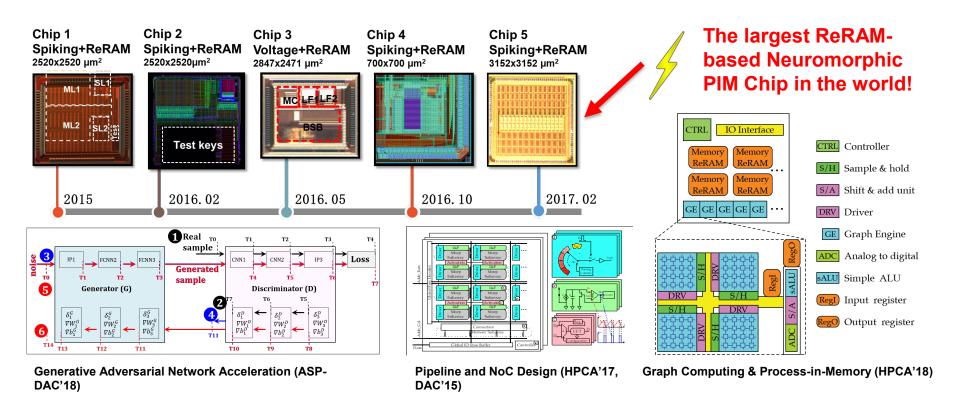
Extensive experience on Al Chip design, including both architecture and circuit. Taping out chips in every 6-9 months.

**Emerging** 

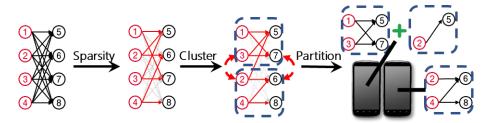


Access to emerging architectures and device technologies, such as IBM TrueNorth, Intel Loihi, and ReRAM, Spin Memory, etc.

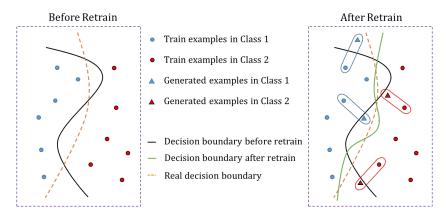
## Research of ASIC: AI Chips



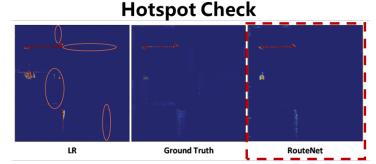
## Research of ASIC: Miscellaneous



#### Mobile DNN computing platforms (DATE'17 Best Paper)



Robustness and safety of DNN models



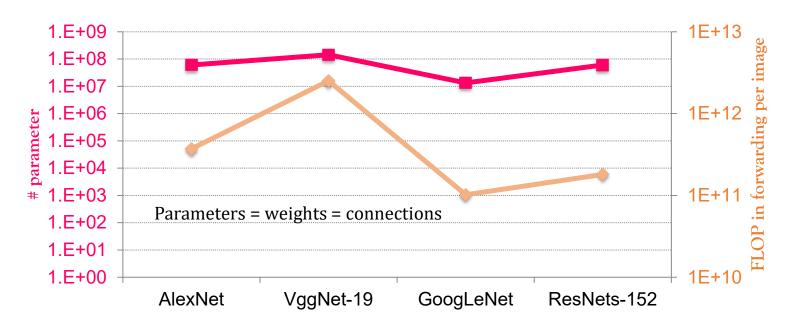
Circuit Name	FPR			TPI	R (%)	
Circuit Name	(%)	TR	GR	LR	<b>SVM</b>	RouteNet
des_perf	0.54	17	56	54	42	74
edit_dist	1.00	25	36	38	28	64
fft	0.30	21	45	54	31	71
matrix_mult_a	0.21	13	30	34	12	49
matrix_mult_b	0.24	13	37	41	20	53
Average	0.46	18	41	44	27	62

DNN-enabled Electronic Design Automation (EDA): Placement & route, High-level synthesis, Timing analysis, Data argumentation, ...

# Learning Structured Sparsity in Deep Neural Networks



## **Complexity of Deep Neural Networks**



Winners of ImageNet Challenge in recent years

Fewer parameters, fewer computation (FLOP: Floating Point Operation)

How to reduce the number of parameters in DNN so as to reduce FLOP, meanwhile maintain the classification accuracy?

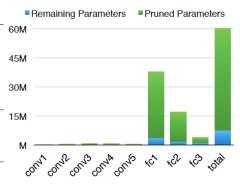
# Non-structurally Sparse DNNs

- State-of-the-art methods to reduce the number of parameters
  - Weight regularization (L1-norm)
  - Connection pruning

Layer	conv1	conv2	conv3	conv4	conv5
Sparsity%	0.927	0.95	0.951	0.942	0.938
Theoretical speedup	2.61	7.14	16.12	12.42	10.77

Sparsity: the ratio of zeros remained

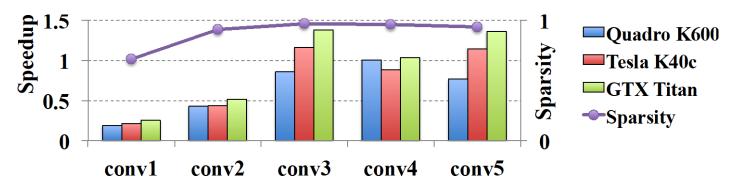
Layer	Weights	FLOP	Act%	Weights%	FLOP%
conv1	35K	211M	88%	84%	84%
conv2	307K	448M	52%	38%	33%
conv3	885K	299M	37%	35%	18%
conv4	663K	224M	40%	37%	14%
conv5	442K	150M	34%	37%	14%
fc1	38M	75M	36%	9%	3%
fc2	17M	34M	40%	9%	3%
fc3	4M	8M	100%	25%	10%
Total	61M	1.5B	54%	11%	30%
	'				



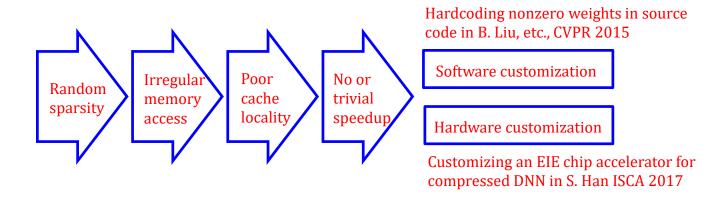
AlexNet, S. Han, et al., NIPS 2015

AlexNet, B. Liu, et al., CVPR 2015

## Theoretical Speedup ≠ Practical Speedup

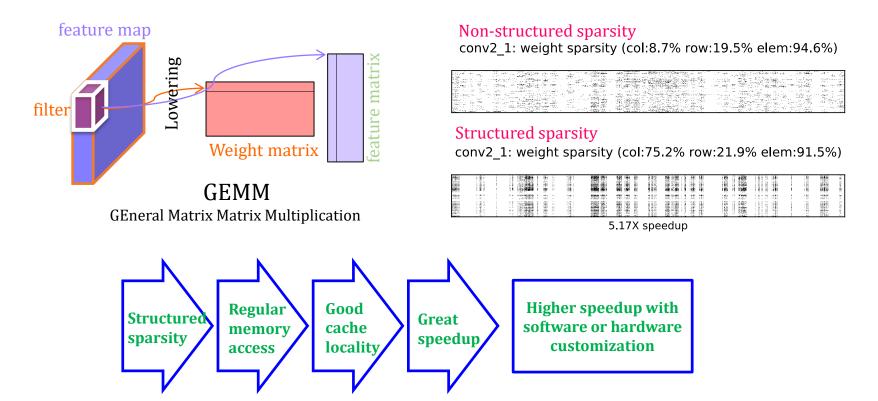


Forwarding speedups of AlexNet on GPU platforms and the sparsity. Baseline is GEMM of cuBLAS. The sparse matrixes are stored in the format of Compressed Sparse Row (CSR) and accelerated by cuSPARSE.



## Theoretical Speedup ≈ Practical Speedup

**Example**: Removing rows/columns in GEMM (row/column-wise sparsity)



# Structured Sparsity Regularization

Group Lasso regularization in ML model

$$\underset{\mathbf{w}}{\operatorname{argmin}} \left\{ E\left(\mathbf{w}\right) \right\} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ E_{D}\left(\mathbf{w}\right) + \lambda_{g} \cdot R_{g}\left(\mathbf{w}\right) \right\}$$

$$\underset{\mathbf{w}}{\operatorname{argmin}} \left\{ E\left(\mathbf{w}\right) \right\} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ E_{D}\left(\mathbf{w}\right) \right\}$$

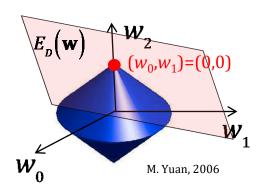
$$s.t. R_{g}\left(\mathbf{w}\right) \leq \eta_{g}$$

#### Example:

$$R_g([w_0, w_1], [w_2]) = \sqrt{w_0^2 + w_1^2} + \sqrt{w_2^2} \le \eta_g$$

#### Many groups will be zeros

$$R_g(\mathbf{w}) = \sum_{g=1}^{G} ||\mathbf{w}^{(g)}||_g,$$
$$||\mathbf{w}^{(g)}||_g = \sqrt{\sum_{i=1}^{|\mathbf{w}^{(g)}|} (w_i^{(g)})^2}$$



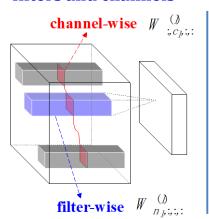
# **SSL: Structured Sparsity Learning**

Group Lasso regularization in DNNs:

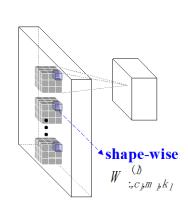
$$E(\mathbf{W}) = E_D(\mathbf{W}) + \lambda \cdot R(\mathbf{W}) + \lambda_g \cdot \sum_{l=1}^{L} R_g \left( \mathbf{W}^{(l)} \right)$$
$$R_g(\mathbf{w}) = \sum_{g=1}^{G} ||\mathbf{w}^{(g)}||_g$$

Learned structured sparsity is determined by the way of splitting groups

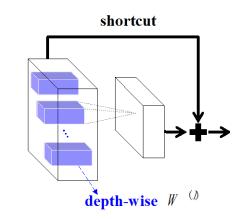
Penalize unimportant filters and channels



**Learn filter shapes** 



Learn the depth of layers



## Penalizing unimportant filters & channels

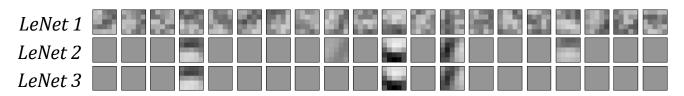
#### LeNet on MNIST

Table 1: Results after penalizing unimportant filters and channels in *LeNet* 

LeNet #	Error	Filter # §	Channel # §	FLOP §	Speedup §
1 (baseline)	0.9%	20—50	1—20	100%—100%	1.00×—1.00×
2	0.8%	5—19	1—4	25%—7.6%	1.64×—5.23×
3	1.0%	3—12	1—3	15%—3.6%	1.99×—7.44×

<sup>§</sup>In the order of conv1—conv2

#### conv1 filters (gray level 128 represents zero)



Fewer but more natural patterns

# Learning smaller filter shapes

Table 2: Results after learning filter shapes in *LeNet* 

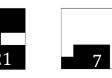
LeNet#	Error	Filter size §	Channel #	FLOP	Speedup
1 (baseline)	0.9%	25—500	1—20	100%—100%	1.00×—1.00×
4	0.8%	21—41	1—2	8.4%—8.2%	2.33×—6.93×
5	1.0%	7—14	1—1	1.4%—2.8%	5.19×—10.82×

<sup>§</sup> The sizes of filters after removing zero shape fibers, in the order of conv1—conv2

Learned shapes of conv1 filters:



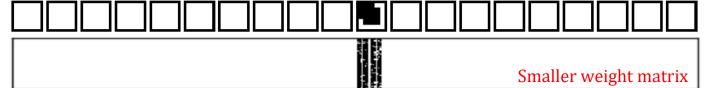
LeNet 1 LeNet 4



LeNet 5

Learned shape of conv2 filters @ LeNet 5

3D 20x5x5 filters is regularized to 2D filters!



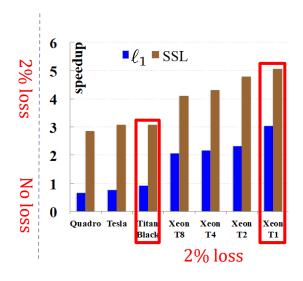
# AlexNet@ImageNet

Learning row-wise and column-wise sparsity:

14 4 24 2 24 2	2011	31	223				625 (W	1	200	ii ii	BLILLIAN S	14 6 6 14 6 6 14 6 6		Š.	100			Martin III				HAR HIS HIS	100 000
180.0	1 80 0		1444030	: :	A) II	100.1	HOME THE			18.1	4 5 2		. 1181	181	1867		T Shrift		4.5	:	***	171 M SF	10
100 F 100 F 100 F		100	1111	7		II.	5/622 184 5/622 184 1826 181	7	1000		1000000	12 × 32		1	121	î	150		100			HPE HEL	II divan

Table 4: Sparsity and speedup of AlexNet on ILSVRC 2012

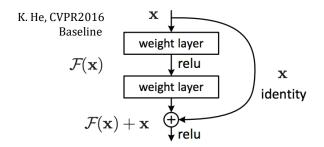
	These is opening and opening of the 2012													
#	Method	Top1 err.	Statistics	conv1	conv2	conv3	conv4	conv5						
1	$\ell_1$	44.67%	sparsity CPU × GPU ×	67.6% 0.80 0.25	92.4% 2.91 0.52	97.2% 4.84 1.38	96.6% 3.83 1.04	94.3% 2.76 1.36						
2	SSL	44.66%	column sparsity row sparsity CPU × GPU ×	0.0% 9.4% 1.05 1.00	63.2% 12.9% 3.37 2.37	76.9% 40.6% 6.27 4.94	84.7% 46.9% 9.73 4.03	80.7% 0.0% 4.93 3.05						
3	pruning[6]	42.80%	sparsity	16.0%	62.0%	65.0%	63.0%	63.0%						
4	$\ell_1$	42.51%	sparsity CPU × GPU ×	14.7% 0.34 0.08	76.2% 0.99 0.17	85.3% 1.30 0.42	81.5% 1.10 0.30	76.3% 0.93 0.32						
5	SSL	42.53%	column sparsity CPU × GPU ×	0.00% 1.00 1.00	20.9% 1.27 1.25	39.7% 1.64 1.63	39.7% 1.68 1.72	24.6% 1.32 1.36						



- 1. Non-structured sparsity method even slows down in some layers
- 2. layer-wise 5.1X /3.1X on CPU/GPU with 2% accuracy loss
- 3. layer-wise 1.4X on both CPU and GPU w/o accuracy loss
- 4. Higher speedups than non-structured speedups

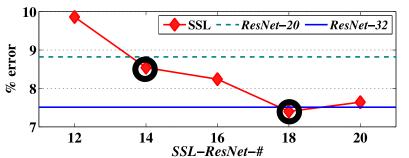
# Regularizing the depth of DNNs

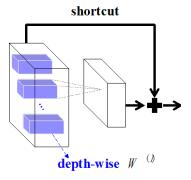
Experiments of ResNets on CIFAR-10



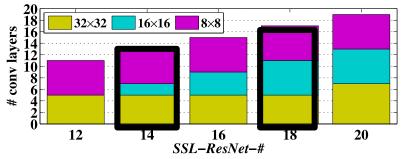
ResNet-20/32: baseline with 20/32 layers

SSL-ResNet-#: Ours with # layers after learning depth of ResNet-20





	# layers	error	# layers	error
ResNet	20	8.82%	32	7.51%
SSL-ResNet	14	8.54%	18	7.40%



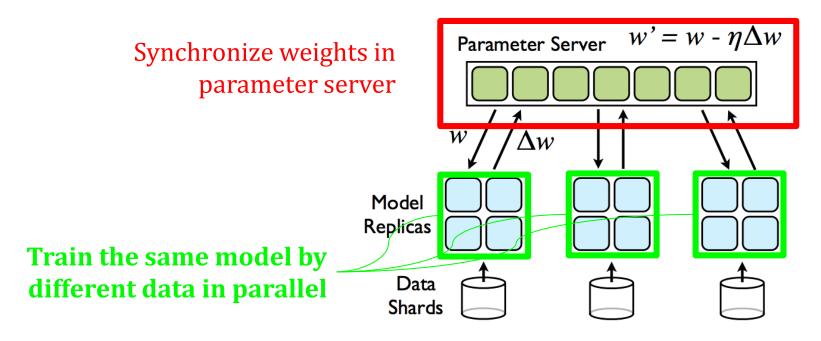
# TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning





# **Distributed Deep Learning**

When parallelism increase, communication is the bottleneck

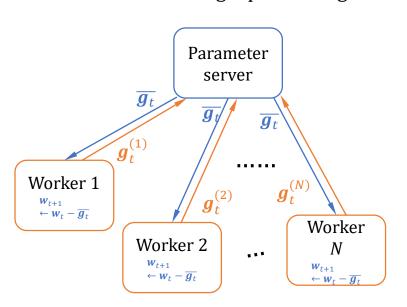


DistBelief by Google

# TernGrad – distributed training with ternary gradients

Key ideas to reduce communication:

- 1. Randomly quantize gradients to only three levels  $(0, \pm 1)$
- 2. The expectations of quantized gradients equal original values
- 3. Exchange quantized gradients instead of floating weights



**Algorithm 1** *TernGrad*: distributed SGD training using ternary gradients.

#### Worker: i = 1, ..., N

- Input  $z_t^{(i)}$ , a part of a mini-batch of training samples  $z_t$
- 2 Compute gradients  $\boldsymbol{g}_t^{(i)}$  under  $\boldsymbol{z}_t^{(i)}$
- Ternarize gradients to  $ilde{m{g}}_t^{(i)} = ternarize(m{g}_t^{(i)})$
- 4 Push ternary  $\tilde{\boldsymbol{g}}_{t}^{(i)}$  to the server
- Pull averaged gradients  $\overline{g_t}$  from the server
- Update parameters  $oldsymbol{w}_{t+1} \leftarrow oldsymbol{w}_t \eta \cdot \overline{oldsymbol{g}_t}$

#### Server:

7 Average ternary gradients  $\overline{m{g}_t} = \sum_i \tilde{m{g}}_t^{(i)}/N$ 

$$\tilde{g}_{t} = ternarize(g_{t}) = s_{t} \cdot sign(g_{t}) \circ b_{t}$$

$$s_{t} \triangleq max(abs(g_{t}))$$

$$\begin{cases} P(b_{tk} = 1 \mid g_{t}) = |g_{tk}|/s_{t} \\ P(b_{tk} = 0 \mid g_{t}) = 1 - |g_{tk}|/s_{t} \end{cases}$$

## **TernGrad - Convergence**

#### Mathematically guarantee the convergence of *TernGrad*

L. Bottou 1998

**Assumption 1.** C(w) has a single minimum  $w^*$  and gradient  $-\nabla_w C(w)$  always points to  $w^*$ , i.e.,

$$\forall \epsilon > 0, \inf_{\|\boldsymbol{w} - \boldsymbol{w}^*\|^2 > \epsilon} (\boldsymbol{w} - \boldsymbol{w}^*)^T \nabla_{\boldsymbol{w}} C(\boldsymbol{w}) > 0.$$
(8)

Convexity is a subset of Assumption 1, and we can easily find non-convex functions satisfying it.

**Assumption 2.** Learning rate  $\gamma_t$  is positive and constrained as  $\sum_{t=0}^{+\infty} \gamma_t^2 < +\infty$  and  $\sum_{t=0}^{+\infty} \gamma_t = +\infty$ , which ensures  $\gamma_t$  decreases neither very fast nor very slow respectively.

We assume the gradient is bounded as

**Assumption 3** (Gradient Bound). The gradient g is bounded as  $\mathbf{E} \{ max(abs(g)) \cdot ||g||_1 \} \le A + B ||w - w^*||^2$ , where A, B > 0 and  $||\cdot||_1$  is  $\ell_1$  norm.

**Theorem 1.** When online learning systems update as  $w_{t+1} = w_t - \gamma_t (s_t \cdot sign(g_t) \circ b_t)$  using stochastic ternary gradients, they converge almost surely toward minimum  $w^*$ , i.e.,  $P(\lim_{t\to+\infty} w_t = w^*) = 1$ .

# **TernGrad - Gradients Histograms**

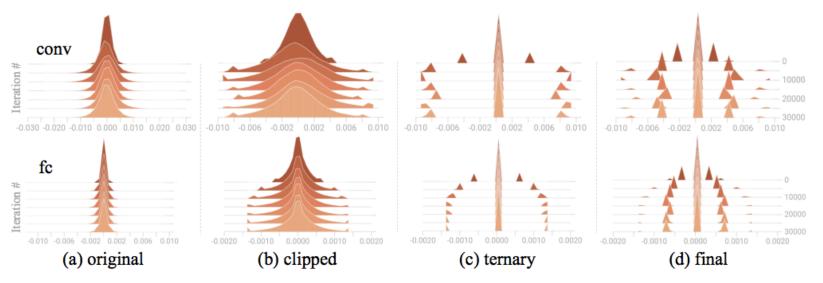


Figure 2: Histograms of (a) original floating gradients, (b) clipped gradients, (c) ternary gradients and (d) final averaged gradients. Visualization by TensorBoard. The DNN is *AlexNet* distributed on two workers, and vertical axis is the training iteration. Top row visualizes the third convolutional layer and bottom one visualizes the first fully-connected layer.

## TernGrad - AlexNet

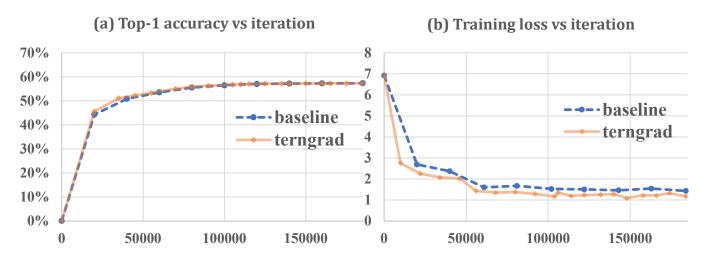


Table 2: Accuracy comparison for *AlexNet*.

base LR	mini-batch size	workers	iterations	gradients	weight decay	DR <sup>†</sup>	top-1	top-5
0.01	256	2	370K	floating TernGrad TernGrad-noclip ‡	0.0005 0.0005 0.0005	0.5 0.2 0.2	57.33% 57.61% 54.63%	80.56% 80.47% 78.16%
0.02	512	4	185K	floating <i>TernGrad</i>	0.0005 0.0005	0.5 0.2	57.32% 57.28%	80.73% 80.23%
0.04	1024	8	92.5K	floating TernGrad	0.0005 0.0005	0.5 0.2	56.62% 57.54%	80.28% 80.25%

<sup>&</sup>lt;sup>†</sup> DR: dropout ratio, the ratio of dropped neurons. <sup>‡</sup> TernGrad without gradient clipping.

# TernGrad - GoogLeNet

Table 3: Accuracy comparison for *GoogLeNet*.

base LR	mini-batch size	workers	iterations	gradients	weight decay	DR	top-5
0.04	128	2	600K	floating TernGrad	4e-5 1e-5	0.2 0.08	88.30% 86.77%
0.08	256	4	300K	floating TernGrad	4e-5 1e-5	0.2 0.08	87.82% 85.96%
0.10	512	8	300K	floating <i>TernGrad</i>	4e-5 2e-5	0.2 0.08	89.00% 86.47%

## TernGrad - Speedup

A performance model to evaluate the speed distributed training.

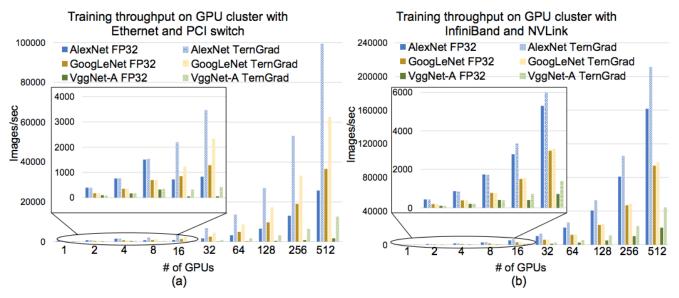


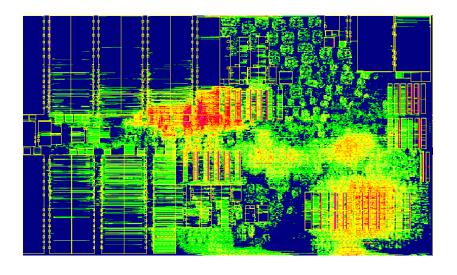
Figure 5: Training throughput on two different GPUs clusters: (a) 128-node GPU cluster with 1Gbps Ethernet, each node has 4 NVIDIA GTX 1080 GPUs and one PCI switch; (b) 128-node GPU cluster with 100 Gbps InfiniBand network connections, each node has 4 NVIDIA Tesla P100 GPUs connected via NVLink. Mini-batch size per GPU of *AlexNet*, *GoogLeNet* and *VggNet-A* is 128, 64 and 32, respectively

# RouteNet: Routability Prediction Using CNN



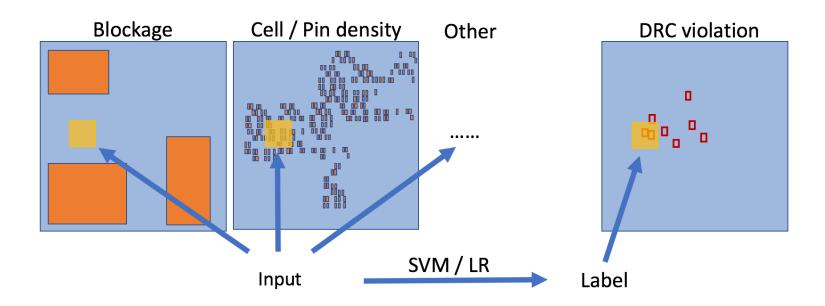
## **Early Routability Prediction**

- Routability: post-routing design rule violations
- Early prediction at placement stage
- Analytical techniques
  - Very fast
  - Not enough fidelity
- Trial routing
  - Acceptable fidelity
  - Not fast enough



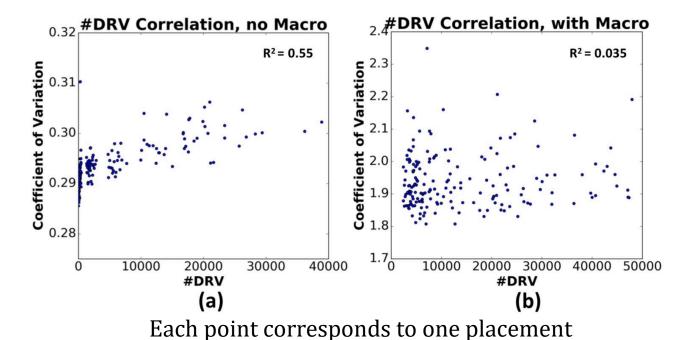
### **Previous ML Approaches**

• Learning on *small* cropped regions



### **Challenges of Macros**

- Layout is less homogeneous
- Correlation between pin density and #DRV becomes weak

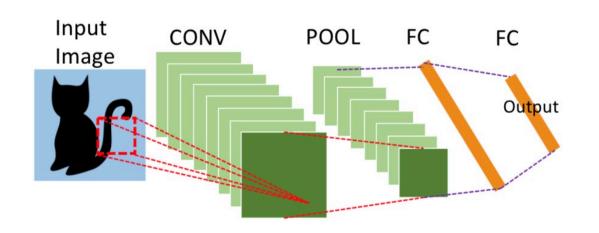


#### **Problem Formulations**

- Predicting overall number of design rule violations (#DRV)
  - Given two placement solutions, tell which is more routable with high fidelity

- DRV hotspot detection
  - Given a relatively routable placement solution, pinpoint DRV hotspots such that mitigation measures are well targeted

#### **CNN for #DRV Prediction**

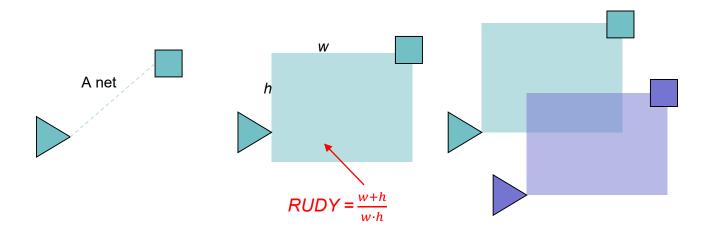


Convolutional (CONV) layers
Pooling (POOL) layers
Fully Connected (FC) layers
Widely used in image classification

- Given a cell placement, classify it among four routability levels, c0, c1, c2, c3
- c0 has the least #DRVs

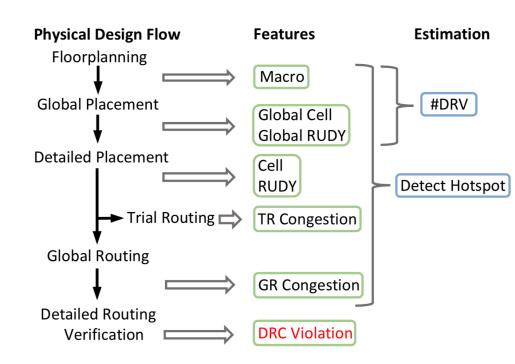
### **An Important Feature**

- RUDY (Rectangular Uniform wire DensitY) (P. Spinder et al. DATE07)
- RUDY at a point is superposition of RUDYs of multiple nets



#### **Features for #DRV Prediction**

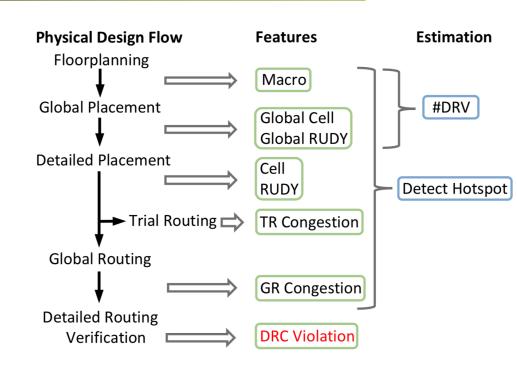
- Macro:
  - region occupied by macros
  - density of macro pins in each layer
- Cell:
  - density of cells
  - density of cell pins
- Global Cell:
  - cell features at global placement
- Global RUDY:
  - RUDY features calculated by global placement results



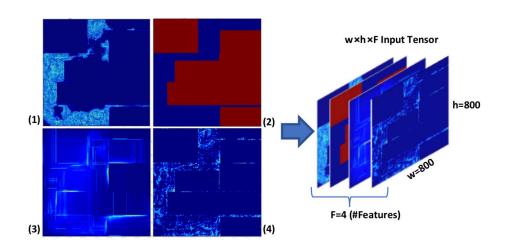
### **Additional Features for Hotspot Detection**

#### • RUDY

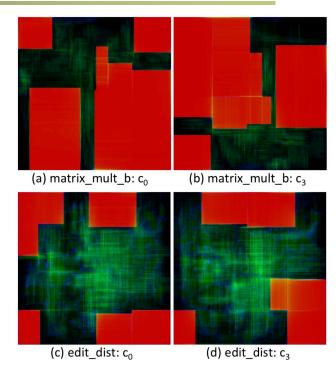
- long-range RUDY
  - RUDY from long-range nets
- short-range RUDY
  - DURY from short-range nets
- RUDY pins
  - pins with density value equal to the RUDY value of its net
- Congestion
  - trial global routing congestion
  - global routing congestion
- DRC violation
  - prediction target / label



#### **Feature Illustration**



Input tensor constructed by stacking 2D features: (1) Pin density, (2) macro (3) long-range RUDY, (4) RUDY pins



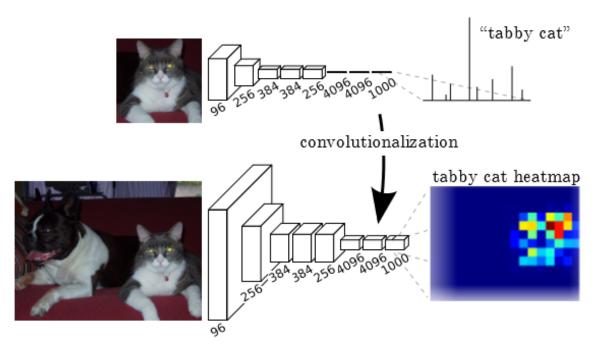
Input features for #DRV prediction.

Red: macro region

Green: global long-range RUDY

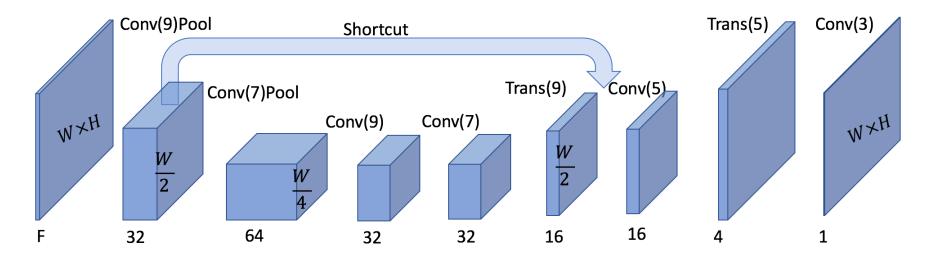
Blue: global RUDY pins

#### Fully Convolutional Network (FCN) for Hotspot Detection



Eliminate FC layers
May use transposed-convolutional to up-sample
Used in image segmentation, object detection

# **FCN** Architecture for Hotspot Detection



Filter size indicated in ()

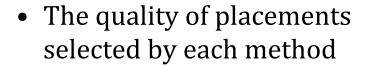
### **Experiment Setup**

- Five designs from ISPD 2015 placement contest
- ~300 different placements by placing macros in different ways
- Placement, routing and DRC are done by Cadence tool
- When a circuit is tested, the model trained with the other circuits
- SVM and Logistic Regression (LR) methods for comparison

Circuit Name	#Macros	#Cells	#Nets	Width (µm)	#Placements
des_perf	4	108288	110283	900	600
edit_dist	6	127413	131134	800	300
fft	6	30625	32088	800	300
matrix_mult_a	5	149650	154284	1500	300
matrix_mult_b	7	146435	151614	1500	300

# **#DRV Prediction Fidelity**

• How methods recognize placements with the lowest #DRV level  $(c_0)$ 



 The best rank of top ten placements predicted to have least #DRV

	$c_0/c_1+c_2+c_3 \text{ accuracy (\%)}$				Best rank in top 10					
Circuit Name	SVM	LR	TR	GR	Route Net	SVM	LR	TR	GR	Route Net
des_perf	63	74	80	77	80	87 <sup>th</sup>	15 <sup>th</sup>	2 <sup>nd</sup>	1 <sup>st</sup>	2 <sup>nd</sup>
edit_dist	69	68	78	77	76	17 <sup>th</sup>	$17^{th}$	$3^{\rm rd}$	$3^{\rm rd}$	$2^{\mathrm{nd}}$
fft	66	62	73	70	75	6 <sup>th</sup>	$6^{th}$	$2^{nd}$	$33^{\rm rd}$	1 <sup>st</sup>
matrix_mult_a	66	65	78	74	72	30 <sup>th</sup>	5 <sup>th</sup>	1 <sup>st</sup>	1 <sup>st</sup>	5 <sup>th</sup>
matrix_mult_b	63	62	76	73	76	22 <sup>nd</sup>	93 <sup>rd</sup>	$4^{th}$	1 <sup>st</sup>	$4^{th}$
Average	65	66	77	74	76	32 <sup>nd</sup>	27 <sup>th</sup>	2 <sup>nd</sup>	8 <sup>th</sup>	3 <sup>rd</sup>



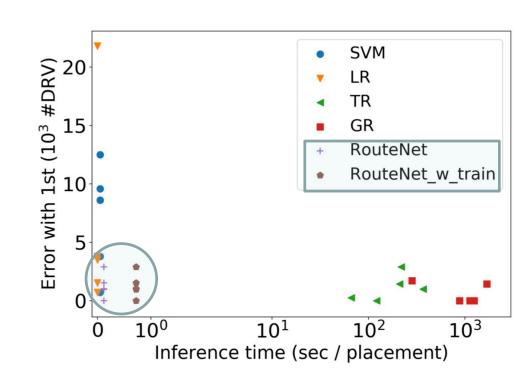
TR: Trial Routing GR: Global Routing

#### **#DRV Prediction Error and Runtime**

• Y: gap between the 'best in 10' and the actually 1st-ranked placement with least #DRV

• X: inference time taken for each method

 RouteNet achieves low inference time and high accuracy at the same time



# **DRV Hotspot Detection Evaluation**

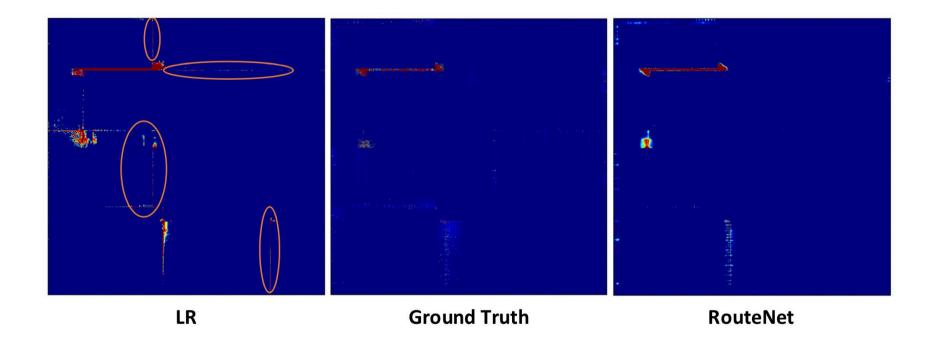
- Same decision threshold is used for all designs
- Slight different FPR, but all under 1%
- RouteNet is superior to all methods and improves global routing accuracy by 50%

Circuit Name	FPR	TPR (%)					
Circuit Name	(%)	TR	GR	LR	SVM	RouteNet	
des_perf	0.54	17	56	54	42	74	
edit_dist	1.00	25	36	38	28	64	
fft	0.30	21	45	54	31	71	
matrix_mult_a	0.21	13	30	34	12	49	
matrix_mult_b	0.24	13	37	41	20	53	
Average	0.46	18	41	44	27	62	

		Predicti	on Result	
		Positive	Negative	Evaluation
Label	Positive	TP	FN	$TPR = \frac{TP}{TP + FN}$
	Negative	FP	TN	$FPR = \frac{FP}{FP + TN}$

TPR (True Positive Rate) FPR (False Positive Rate)

### **DRC Hotspot Detection Demonstration**



### **Our Perspective**







