



Deep Learning in Cloud-Edge AI Systems

Wei Wen

COMPUTATIONAL EVOLUTIONARY INTELLIGENCE LAB

Electrical and Computing Engineering Department

Duke University

www.pittnuts.com

AI Systems Landing

Into the Cloud:



Facebook Big Basin



AI Systems Landing

Upon the Edge:



iPhone X
Face ID



DJI Drone



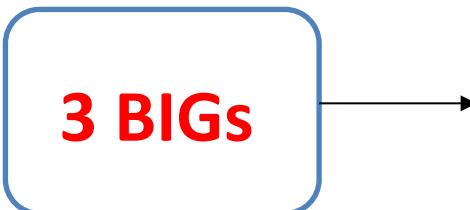
Autonomous Driving

Challenges in landing???

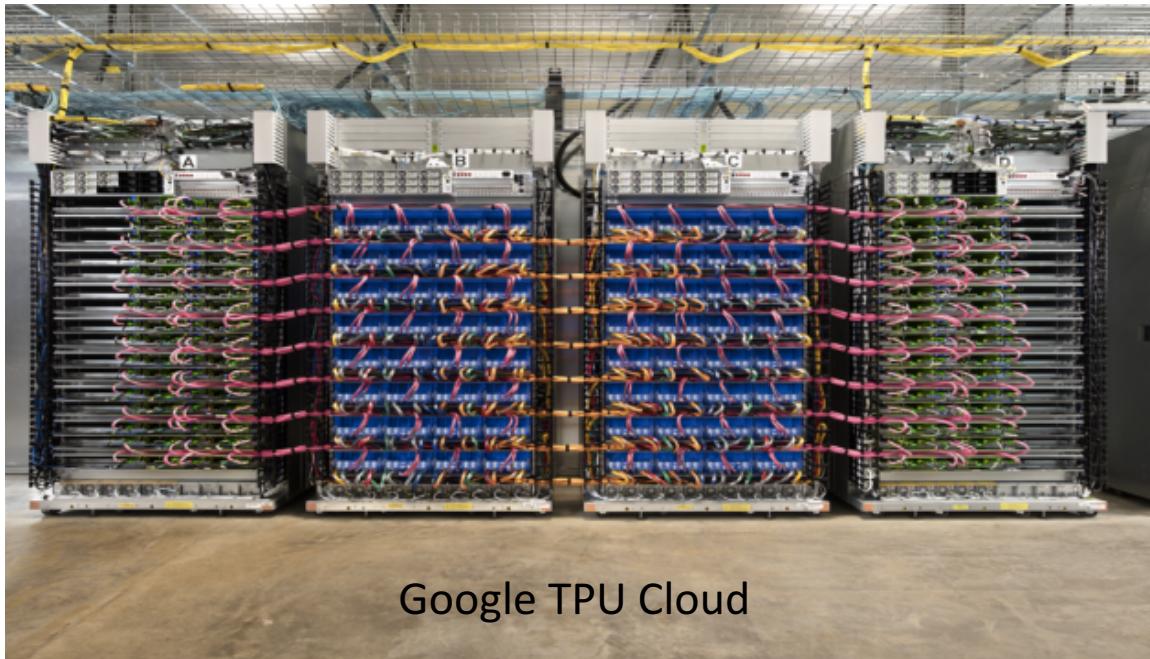
The Rocket Analogy of Deep Learning

An analogy by Andrew Ng

Rocket: **Big** Neural Networks
Engine: **Big** Computing
Fuel: **Big** Data



Deep Learning in the Cloud



Tons of Cables!!!

Communication
Bottleneck!!!

Parallelism

Deep Learning on the Edge



Model Size and Inference Speed Matter!!!

Research Highlights

- Distributed Training in the Cloud
 - TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning, NIPS 2017 (oral)
 - On-going work on large-batch training
- Efficient Inference on the Edge
 - Structurally Sparse DNNs (NIPS 2016 & ICLR 2018)
 - Lower-rank DNNs (ICCV 2017)
 - A Compact DNN: Approaching GoogLeNet-Level Accuracy of Classification and Domain Adaptation (CVPR 2017)
 - Direct Sparse Convolution (ICLR 2017)



TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning

NIPS 2017 (Oral)

Wei Wen¹, Cong Xu², Feng Yan³, Chunpeng Wu¹,
Yandan Wang⁴, Yiran Chen¹, Hai (Helen) Li¹

Duke University¹, Hewlett Packard Labs²,
University of Nevada - Reno³, University of Pittsburgh⁴
<https://github.com/wenwei202/terngrad>

Background – Stochastic Gradient Descent

Minimization target: $C(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n Q(\mathbf{z}_i, \mathbf{w})$

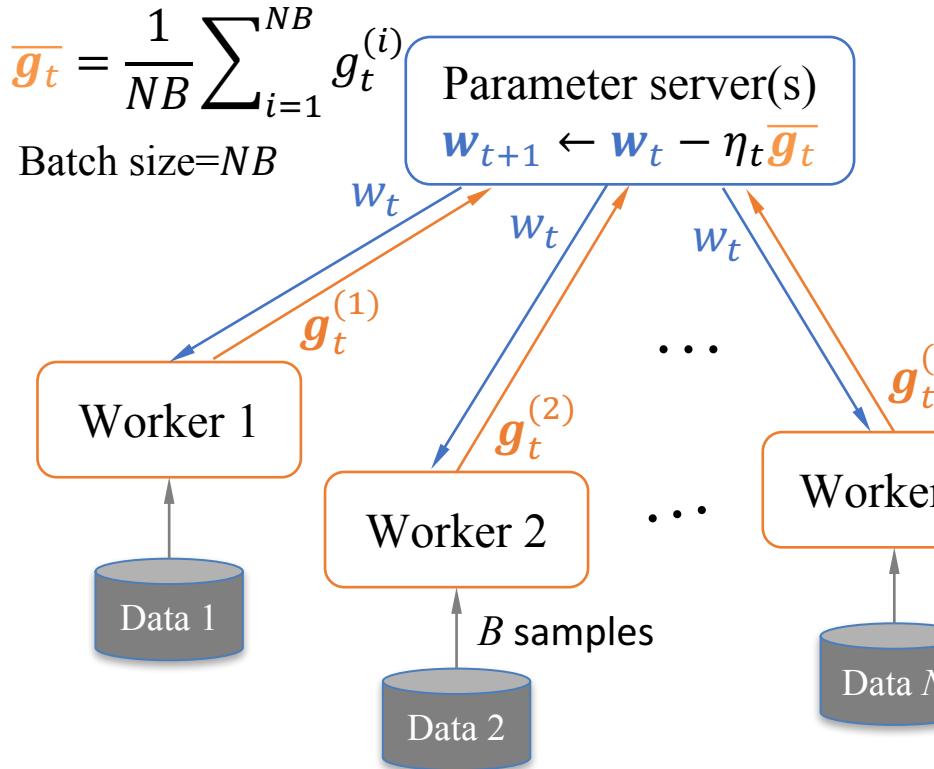
Batch gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta_t}{n} \sum_{i=1}^n g_t^{(i)}$ Computation expensive
 $g_t^{(i)} = \nabla Q(\mathbf{z}_i, \mathbf{w}_t)$

Mini-batch stochastic gradient descent (SGD):

$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta_t}{B} \sum_{b=1}^B g_t^{(b)}$ Computation cheap

$B << n$ samples are randomly drawn from training dataset

Background - Distributed Deep Learning



Synchronized Data Parallelism for Stochastic Gradient Descent (SGD):

1. Training data is split to N subsets
2. Each worker has a model replica (copy)
3. Each replica is trained on a data subset
4. Synchronization in parameter server(s)

Scalability:

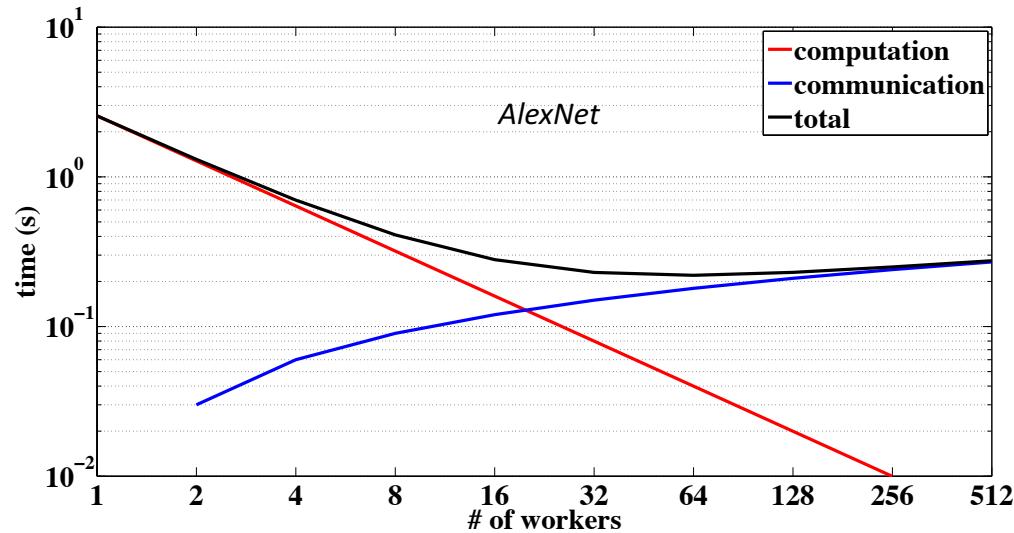
1. Computing time decreases with N
2. Communication can be the bottleneck
3. This work: quantizing gradient to three (i.e., *ternary*) levels $\{-1, 0, 1\}$ (<2bits)

Communication Bottleneck

$$t_{comp} = \frac{flop}{flops \cdot n}$$

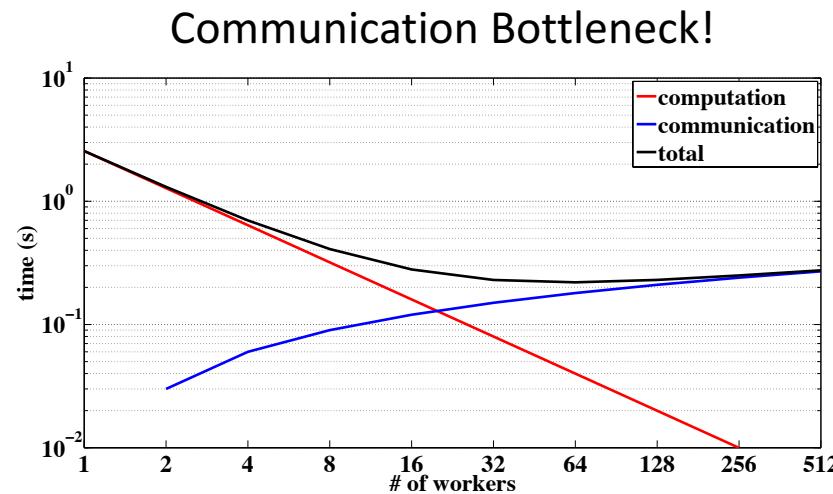
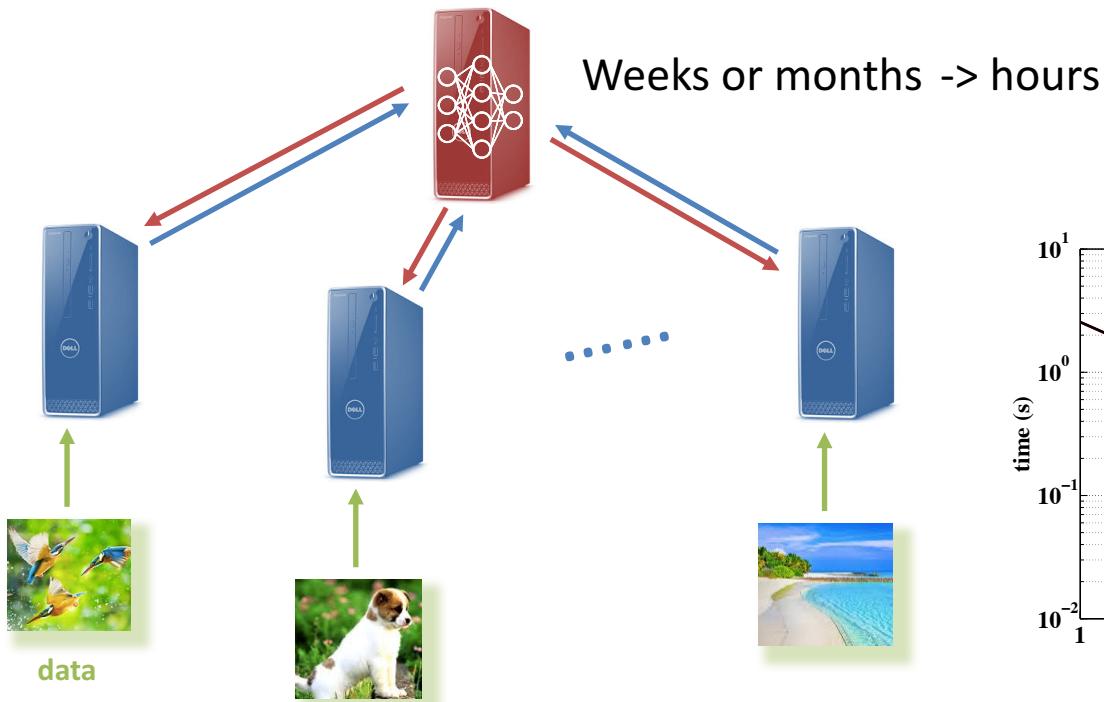
$$t_{comm} = \frac{2 \cdot |\mathbf{W}| \cdot g(n)}{b}$$

$$g(n) = \log_2(n)$$

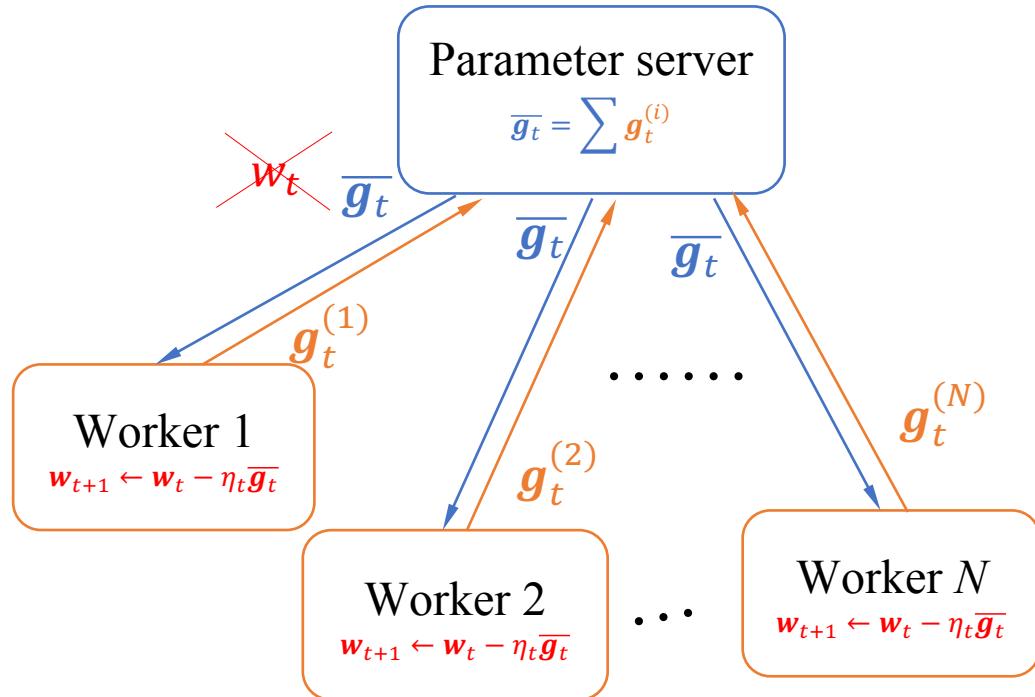


Credit: Alexander Ulanov

Background - Distributed Deep Learning

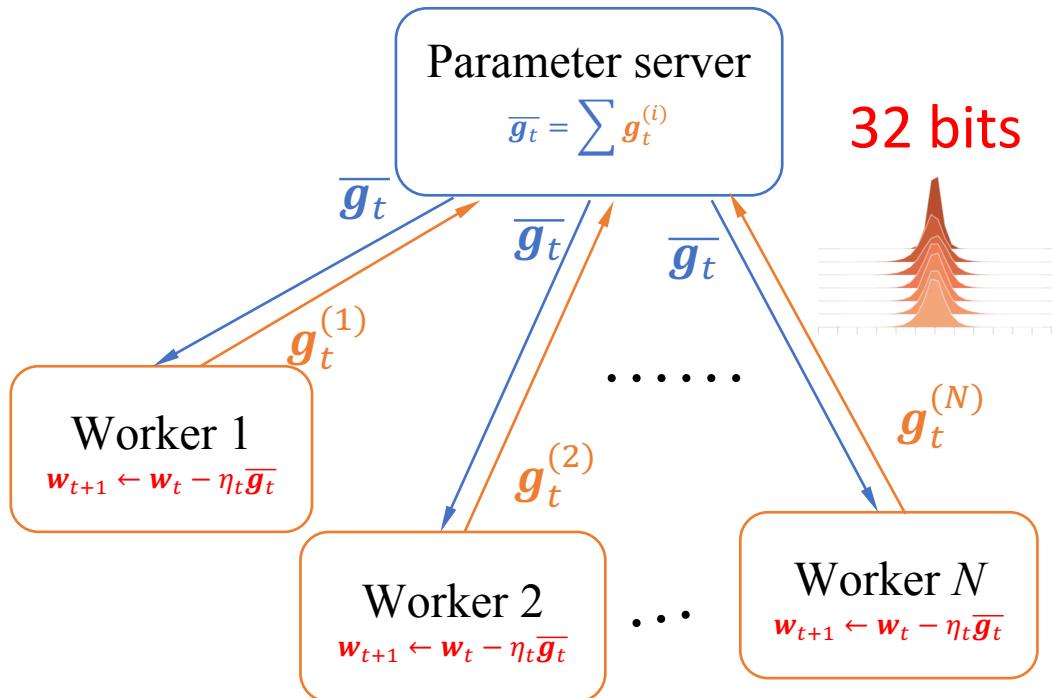


An Alternative Setting



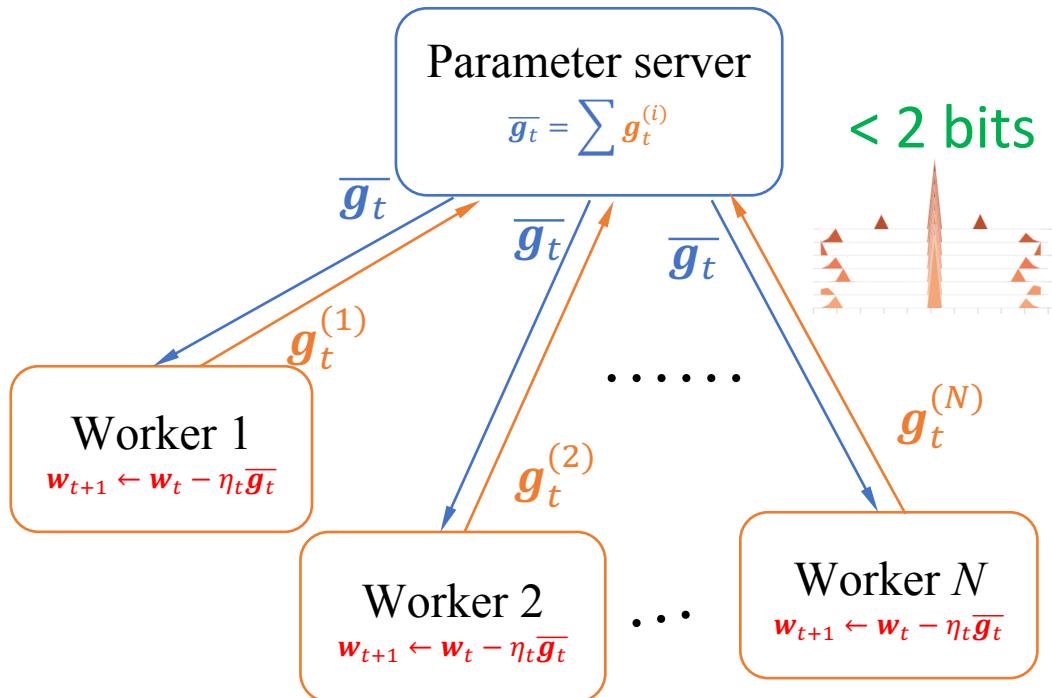
1. Only exchange gradients
2. Gradient quantization can reduce communication in both directions

Gradient Quantization for Communication Reduction



Only exchange
quantized
gradients

Gradient Quantization for Communication Reduction



Only exchange
quantized
gradients

Stochastic Gradients without Bias

Batch Gradient Descent

$$C(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n Q(\mathbf{z}_i, \mathbf{w})$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta_t}{n} \sum_{i=1}^n g_t^{(i)}$$

SGD

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \cdot g_t^{(I)}$$

I is randomly drawn from $[1, n]$

$$\mathbb{E}\{g_t^{(I)}\} = \nabla C(\mathbf{w})$$

No bias

TernGrad

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \cdot \text{ternarize}(g_t^{(I)})$$
$$\mathbb{E}\{\text{ternarize}(g_t^{(I)})\} = \nabla C(\mathbf{w})$$

No bias

TernGrad is Simple

$$\tilde{\mathbf{g}}_t = \text{ternarize}(\mathbf{g}_t) = s_t \cdot \text{sign}(\mathbf{g}_t) \circ \mathbf{b}_t$$

$$s_t \triangleq \|\mathbf{g}_t\|_\infty \triangleq \max(\text{abs}(\mathbf{g}_t))$$

$$\begin{cases} P(b_{tk} = 1 \mid \mathbf{g}_t) = |g_{tk}| / s_t \\ P(b_{tk} = 0 \mid \mathbf{g}_t) = 1 - |g_{tk}| / s_t \end{cases}$$

Example:

$$\mathbf{g}_t^{(i)}: [0.30, -1.20, \dots, 0.9]$$

$$s_t: 1.20$$

$$\text{Signs: } [1, -1, \dots, 1]$$

$$P(b_{tk} = 1 | g_t): [\frac{0.3}{1.2}, \frac{1.2}{1.2}, \dots, \frac{0.9}{1.2}]$$

$$\mathbf{b}_t: [0, 1, \dots, 1]$$

$$\widetilde{\mathbf{g}_t^{(i)}}: [0, -1, \dots, 1] * 1.20$$

$$\mathbf{E}_{\mathbf{z}, \mathbf{b}} \{\tilde{\mathbf{g}}_t\} = \mathbf{E}_{\mathbf{z}, \mathbf{b}} \{s_t \cdot \text{sign}(\mathbf{g}_t) \circ \mathbf{b}_t\}$$

$$= \mathbf{E}_{\mathbf{z}} \{s_t \cdot \text{sign}(\mathbf{g}_t) \circ \mathbf{E}_{\mathbf{b}} \{\mathbf{b}_t | \mathbf{z}_t\}\} = \mathbf{E}_{\mathbf{z}} \{\mathbf{g}_t\} = \nabla_{\mathbf{w}} C(\mathbf{w}_t)$$

No bias

TernGrad is Simple

```
146 def stochastical_binarize_gradients(grads_and_vars, scalers):
147     """Stochastically binarize gradients."""
148     gradients, variables = zip(*grads_and_vars)
149     binarized_gradients = []
150     for gradient, scaler in zip(gradients, scalers):
151         if gradient is None:
152             binarized_gradients.append(None)
153             continue
154         if isinstance(gradient, tf.IndexedSlices):
155             gradient_shape = gradient.dense_shape
156         else:
157             gradient_shape = gradient.get_shape()
158
159         zeros = tf.zeros(gradient_shape)
160         abs_gradient = tf.abs(gradient)
161         sign_gradient = tf.sign(gradient)
162         rnd_sample = tf.random_uniform(gradient_shape, 0, scaler)
163         where_cond = tf.less(rnd_sample, abs_gradient)
164         binarized_gradient = tf.cond(tf.size(gradient) < FLAGS.size_to_binarize,
165                                     lambda: gradient,
166                                     lambda: tf.where(where_cond, sign_gradient * scaler, zeros))
```

Convergence

Standard SGD almost truly converges under assumptions (Fisk 1965, Metivier 1981&1983, Bottou 1998)

Assumption 1:

$C(\mathbf{w})$ has a single minimum \mathbf{w}^* and $\forall \epsilon > 0, \inf_{\|\mathbf{w} - \mathbf{w}^*\|^2 > \epsilon} (\mathbf{w} - \mathbf{w}^*)^T \nabla_{\mathbf{w}} C(\mathbf{w}) > 0$

Assumption 2:

Learning rate γ_t decreases neither very fast nor very slow

$$\begin{cases} \sum_{t=0}^{+\infty} \gamma_t^2 < +\infty \\ \sum_{t=0}^{+\infty} \gamma_t = +\infty \end{cases}$$

Assumption 3 (gradient bound):

$$\mathbf{E} \{ \| \mathbf{g} \|^2 \} \leq A + B \| \mathbf{w} - \mathbf{w}^* \|^2$$

Standard SGD *almost-truly* converges

Assumption 3 (gradient bound):

$$\mathbf{E} \{ \| \mathbf{g} \|_\infty \cdot \| \mathbf{g} \|_1 \} \leq A + B \| \mathbf{w} - \mathbf{w}^* \|^2$$

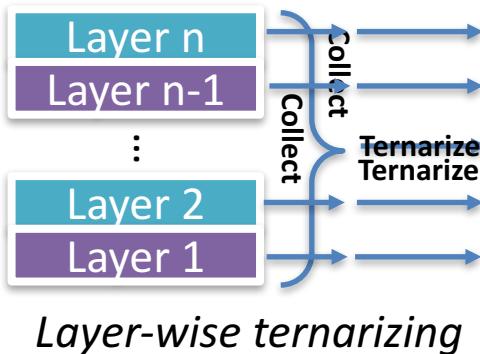
TernGrad *almost-truly* converges

$$\mathbf{E} \{ \| \mathbf{g} \|^2 \} \leq \mathbf{E} \{ \| \mathbf{g} \|_\infty \cdot \| \mathbf{g} \|_1 \} \leq A + B \| \mathbf{w} - \mathbf{w}^* \|^2$$

Stronger gradient bound in *TernGrad*

Closing Bound Gap

Two methods to push the gradient bound of *TernGrad* closer to the bound of standard SGD



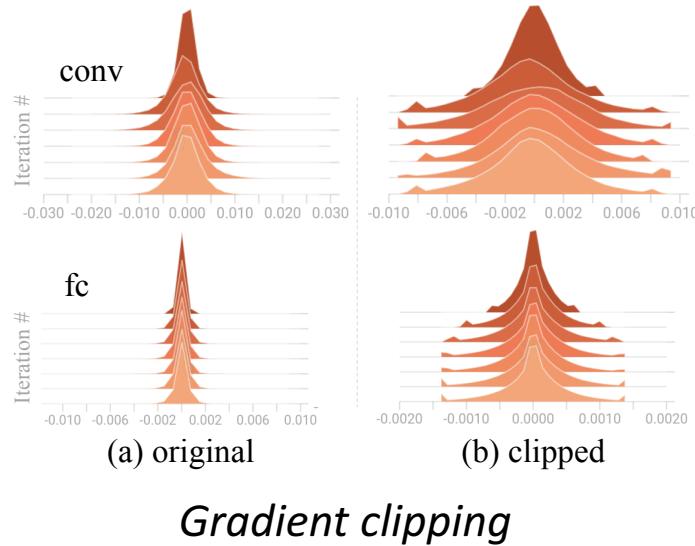
Approach:

1. Split gradients to buckets
2. Do TernGrad bucket by bucket

When bucket size == 1, TernGrad is floating SGD

Closing Bound Gap

Two methods to push the gradient bound of *TernGrad* closer to the bound of standard SGD



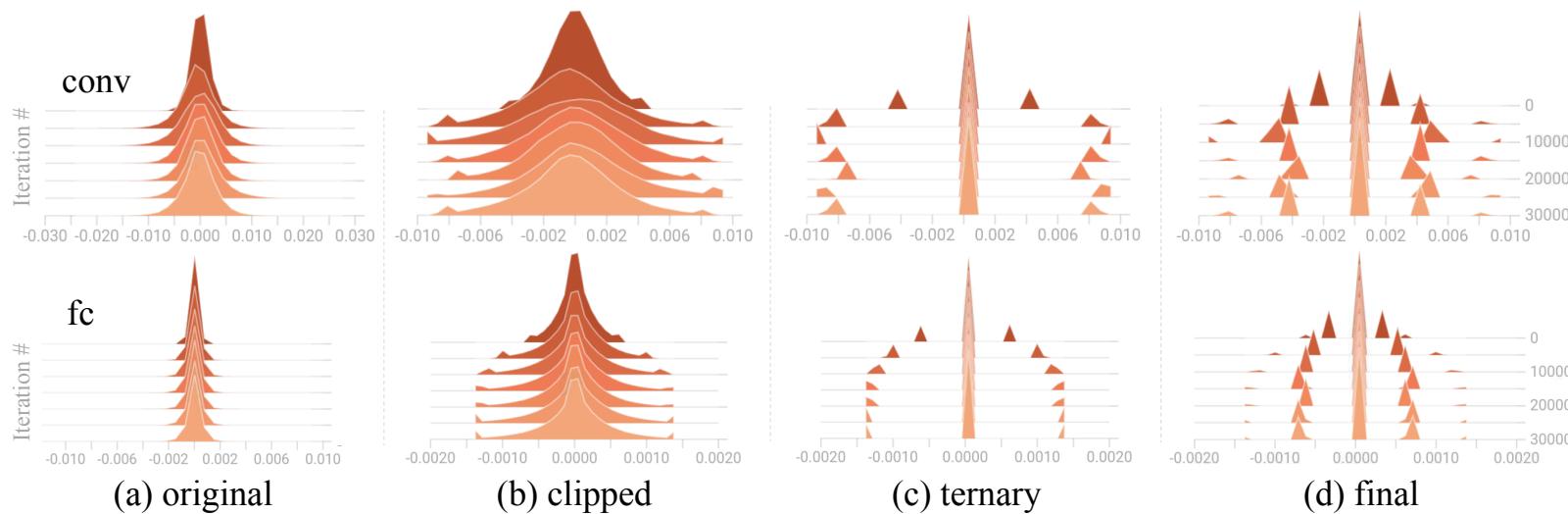
$$f(g_i) = \begin{cases} g_i & |g_i| \leq c\sigma \\ sign(g_i) \cdot c\sigma & |g_i| > c\sigma \end{cases}$$

$c=2.5$ works well for all tested datasets, DNNs and optimizers.

Suppose Gaussian distribution:

1. Change length 1.0%-1.5%
2. Change direction 2° - 3°
3. Small bias with variance reduced

Gradient Histograms



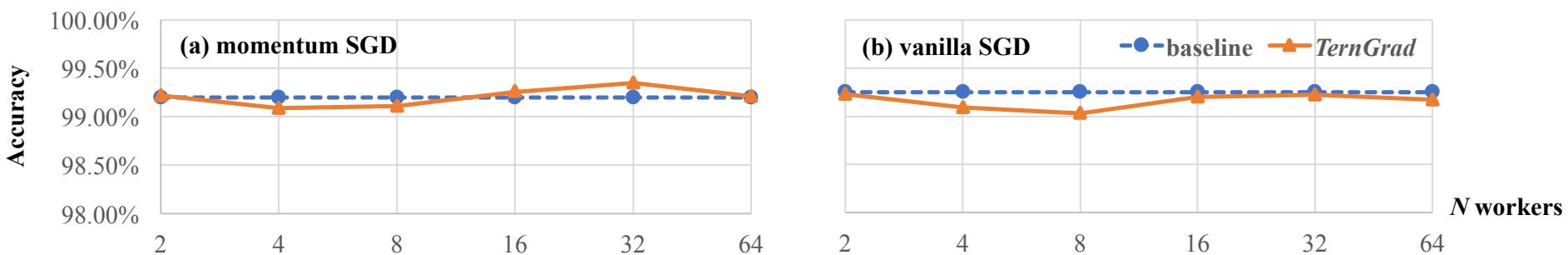
Average gradients are not ternary, but quantized with $2N+1$ levels (N : worker number).

Communication reduction: $\frac{32}{\log_2(2N+1)} > 1$, unless $N \geq 2^{30}$

Integration with Manifold Optimizers

(All experiments: All hyper-parameters are tuned for standard SGD and fixed in *TernGrad*)

LeNet (total mini-batch size 64): close accuracy & randomness in *TernGrad* results in small variance



Integration with Manifold Optimizers

CIFAR-10, mini-batch size 64 per worker

SGD	base LR	total mini-batch size	iterations	gradients	workers	accuracy
Adam	0.0002	128	300K	floating <i>TernGrad</i>	2 2	86.56% 85.64% (-0.92%)
Adam	0.0002	2048	18.75K	floating <i>TernGrad</i>	16 16	83.19% 82.80% (-0.39%)

Adam: D. P. Kingma, 2014

Tune hyper-parameters specifically for *TernGrad* may reduce accuracy gap

TernGrad works with manifold optimizers: Vanilla SGD, Momentum, Adam

Scaling to Large-scale Deep Learning

TernGrad: Randomness & regularization

- { (1) decrease randomness in dropout or
(2) use smaller weight decay
No new hyper-parameters added

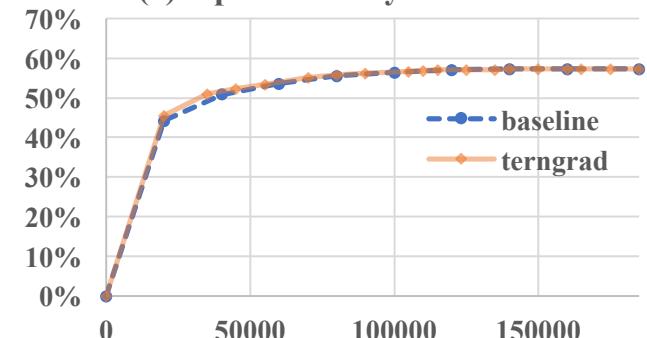
AlexNet

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

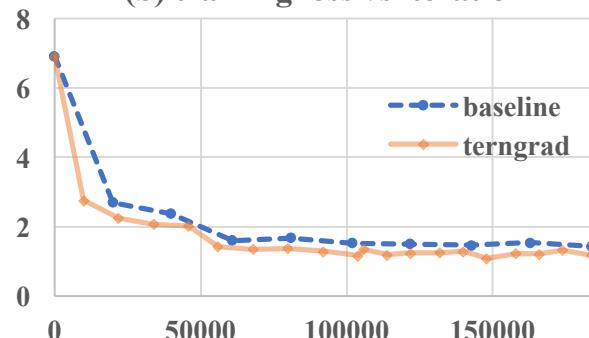
† DR: dropout ratio, the ratio of dropped neurons. ‡ *TernGrad* without gradient clipping.

Convergence Curves

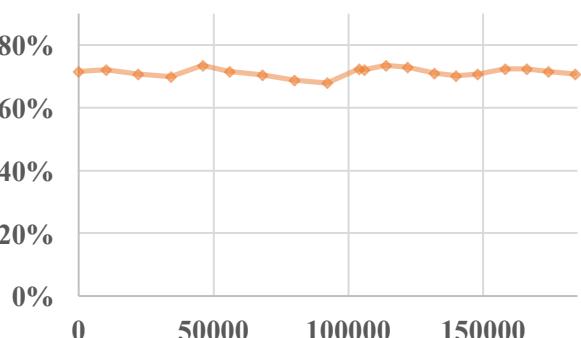
(a) top-1 accuracy vs iteration



(b) training loss vs iteration



(c) gradient sparsity of terngrad in fc6



AlexNet trained on 4 workers with mini-batch size 512

Coverages within the same epochs under the same base learning rate

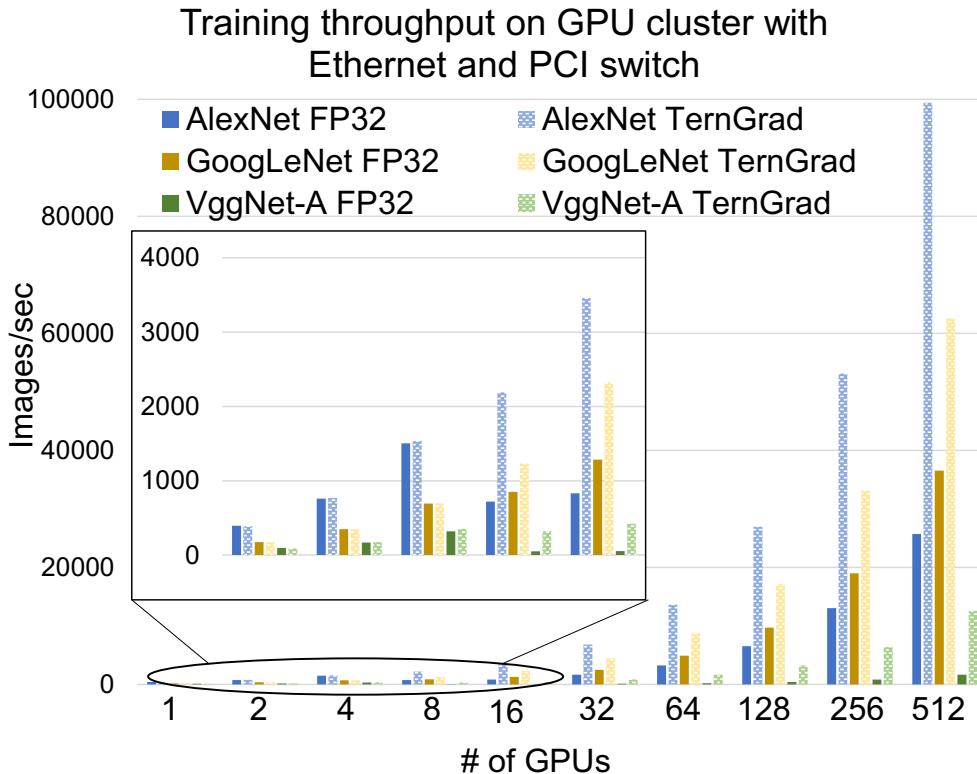
Scaling to Large-scale Deep Learning

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

GoogLeNet accuracy loss is <2% on average.

Tune hyper-parameters specifically for *TernGrad* may reduce accuracy gap

Performance Model



TernGrad gives higher speedup when

1. using more workers
2. using smaller communication bandwidth (Ethernet vs InfiniBand)
3. training DNNs with more fully-connected layers (VggNet vs GoogLeNet)

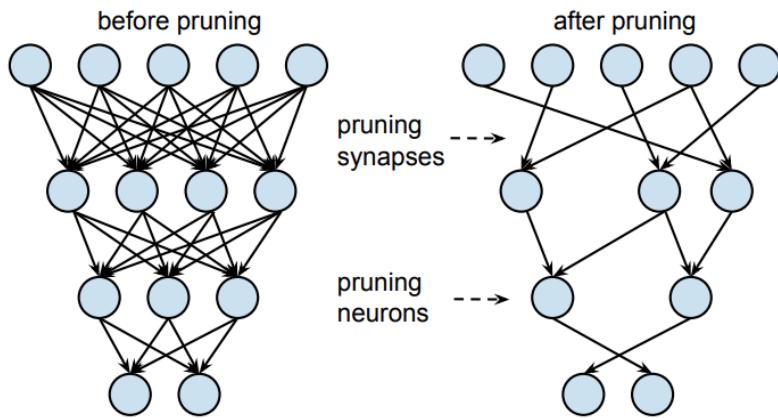


Structurally Sparse Deep Neural Networks

Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, Hai Li, “Learning Structured Sparsity in Deep Neural Networks”, NIPS, 2016. [[paper](#)][[code](#)][[poster](#)]

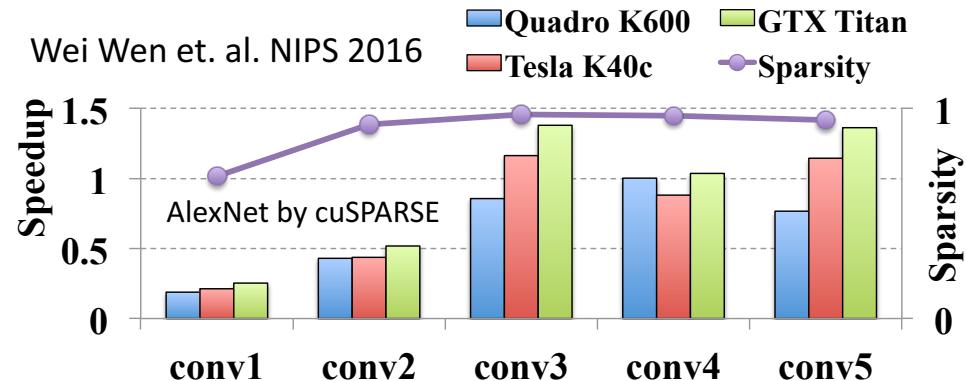
Wei Wen, Yuxiong He, Samyam Rajbhandari, Minjia Zhang, Wenhan Wang, Fang Liu, Bin Hu, Yiran Chen, Hai Li, “Learning Intrinsic Sparse Structures within Long Short-Term Memory”, ICLR, 2018. [[paper](#)][[code](#)]

Sparse Convolutional Neural Networks (SCNNs)

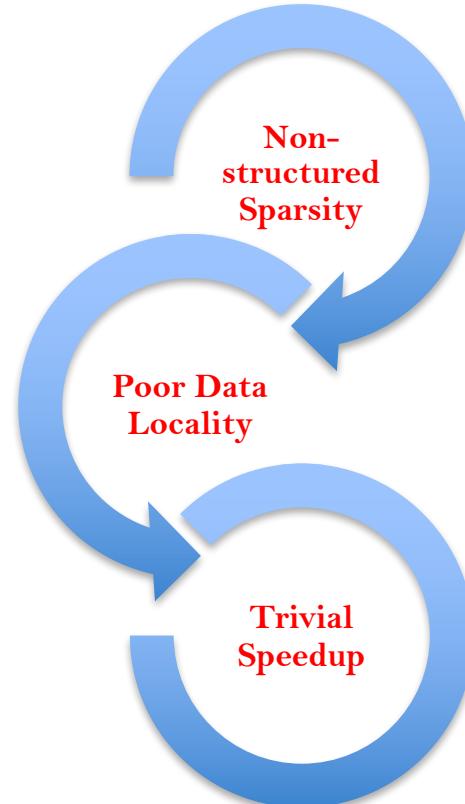
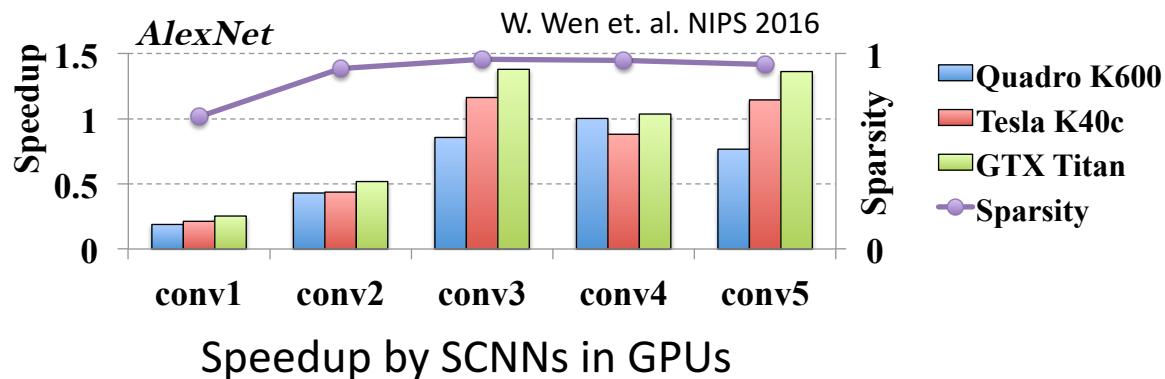


Significantly reduces the storage size of DNNs [1]
Good speedup with customized hardware [2]

- [1] S. Han et. al. NIPS 2015
[2] S. Han et. al. ISCA 2016.

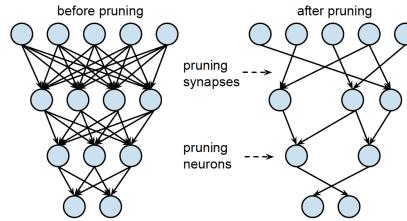
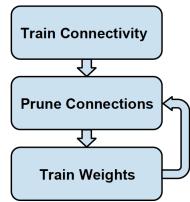


Inefficiency of SCNNs



- Format: Sparse matrixes are formatted as Compressed Sparse Row (CSR)
- Acceleration library: cuSPARSE

Structured Sparsity in DNNs

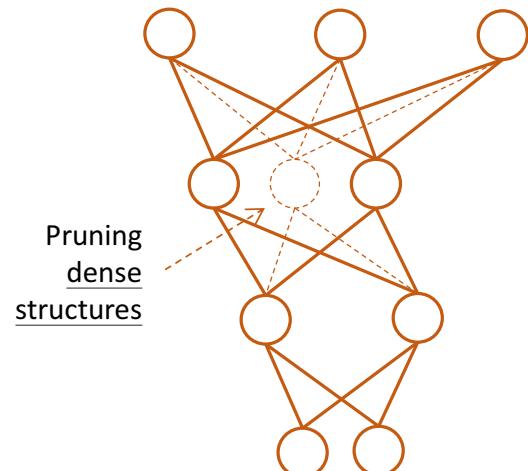
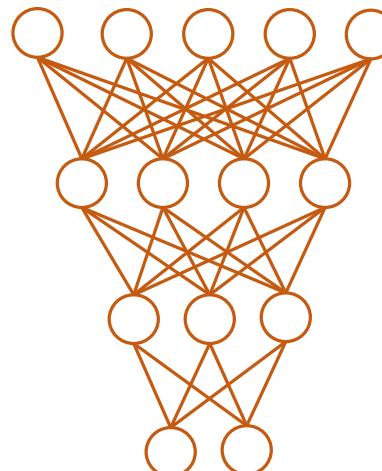


Faster Training speed
Higher practical computing speed

Wei Wen et. al. NIPS 2016
Wei Wen et. al. ICLR 2018

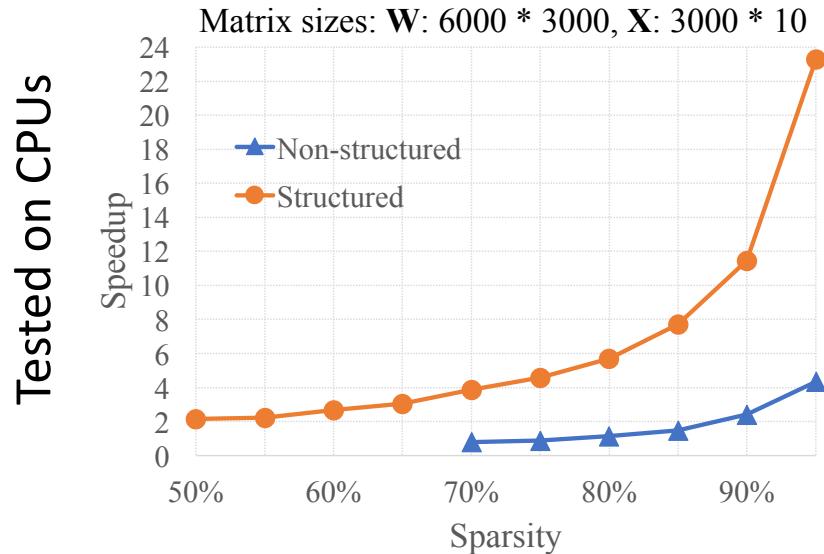
Train weights with
Group Lasso
Regularization

One-shot pruning from scratch
(with the same epochs)

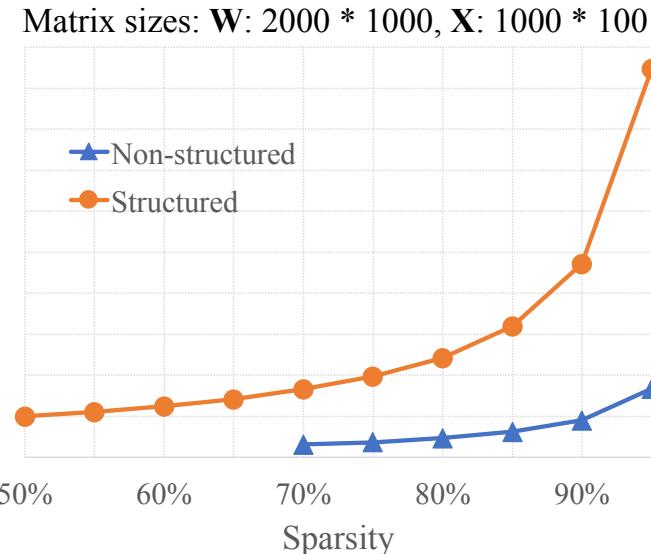


Efficiency of Structured Sparsity

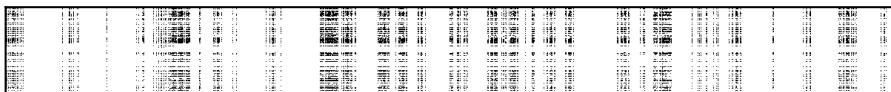
W. Wen et. al. ICLR 2018



Non-structured sparsity



Structured sparsity

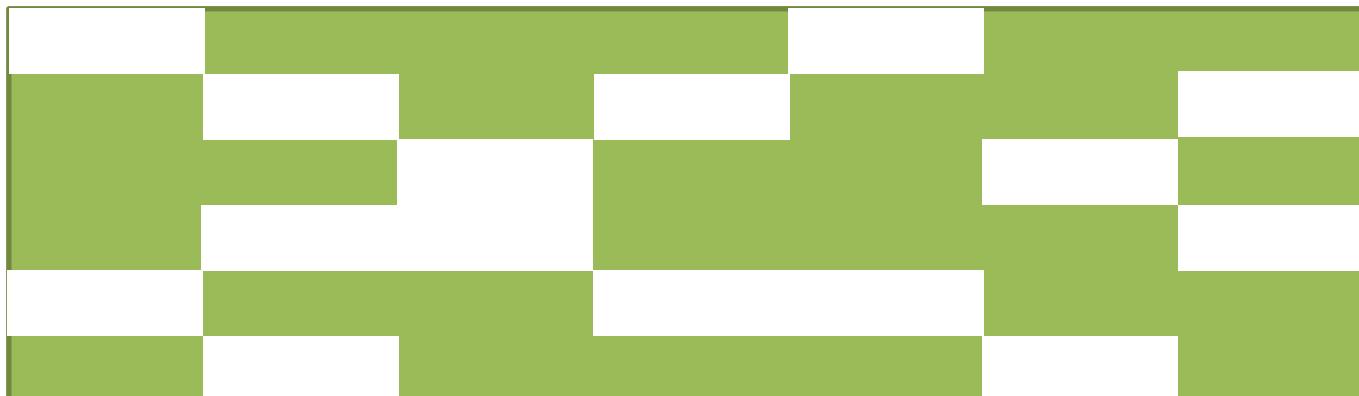


Structurally Sparse Deep Neural Networks

- What is Structurally Sparse Deep Neural Networks (SSDNNs)?

Weights/connections are removed group by group.

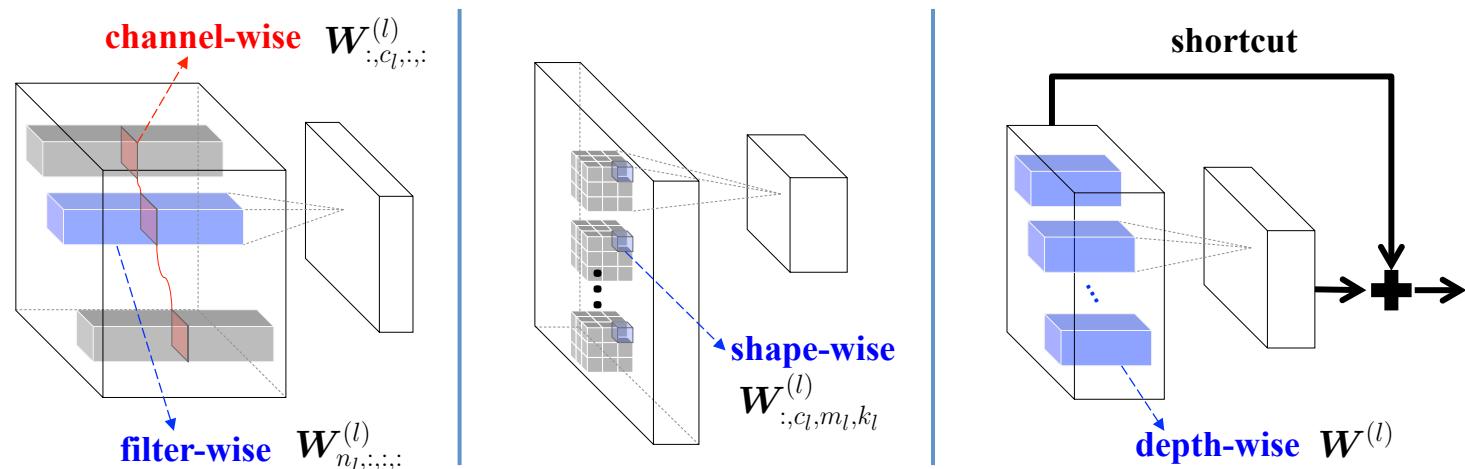
A group can be a rectangle block, a row, a column or even the whole matrix



Structurally Sparse Deep Neural Networks

In a nutshell, a group can be any form of a weight block, depending on what sparse structure you want to learn

In CNNs, a group of weights can be a channel, a 3D filter, a 2D filter, a filter shape fiber (i.e. a weight column), and even a layer (in ResNets).



Group Lasso regularization is all you need for SSDNNs

Step 1: Weights are split to G groups $\mathbf{w}^{(1..G)}$

e.g. $(w_1, w_2, w_3, w_4, w_5) \rightarrow$ group 1: (w_1, w_2, w_3) , group 2: (w_4, w_5)

Step 2: Add group Lasso on each group $\mathbf{w}^{(g)}$ $\|\mathbf{w}^{(g)}\|_g = \sqrt{\sum_{i=1}^{|w^{(g)}|} (w_i^{(g)})^2}$ i.e. vector length
 $\text{sqrt}(w_1^2 + w_2^2 + w_3^2), \text{sqrt}(w_4^2 + w_5^2)$

Step 3: Sum group Lasso over all groups as a regularization: $R_g(\mathbf{w}) = \sum_{g=1}^G \|\mathbf{w}^{(g)}\|_g$,
 $\text{sqrt}(w_1^2 + w_2^2 + w_3^2) + \text{sqrt}(w_4^2 + w_5^2)$

Step 4: SGD optimizing $\arg \min_{\mathbf{w}} \{E(\mathbf{w})\} = \arg \min_{\mathbf{w}} \{E_D(\mathbf{w}) + \lambda_g \cdot R_g(\mathbf{w})\}$

We refer to our method as ***Structured Sparsity Learning (SSL)***

Structured Sparsity Learning (SSL)

- Why can SSL learn to remove weights group by group?

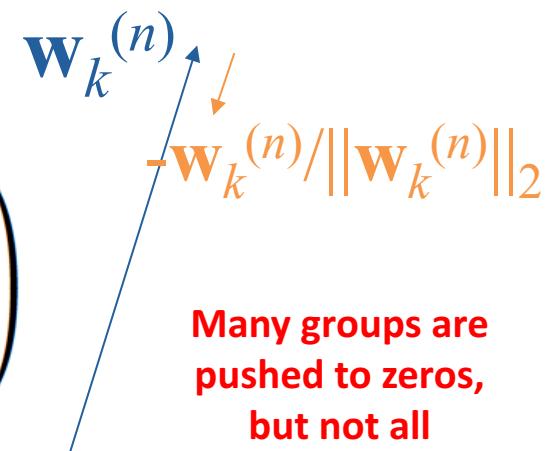
The gradient-based explanation:

$$\boxed{\mathbf{w}_k^{(n)}}$$

A group of
weights

$$\mathbf{w}_k^{(n)} \leftarrow \mathbf{w}_k^{(n)} - \eta \cdot \left(\boxed{\frac{\partial E(\mathbf{w})}{\partial \mathbf{w}_k^{(n)}}} + \lambda \cdot \boxed{\frac{\mathbf{w}_k^{(n)}}{\|\mathbf{w}_k^{(n)}\|_2}} \right)$$

Regular gradients to minimize error Additional gradients to learn sparsity



Many groups are
pushed to zeros,
but not all

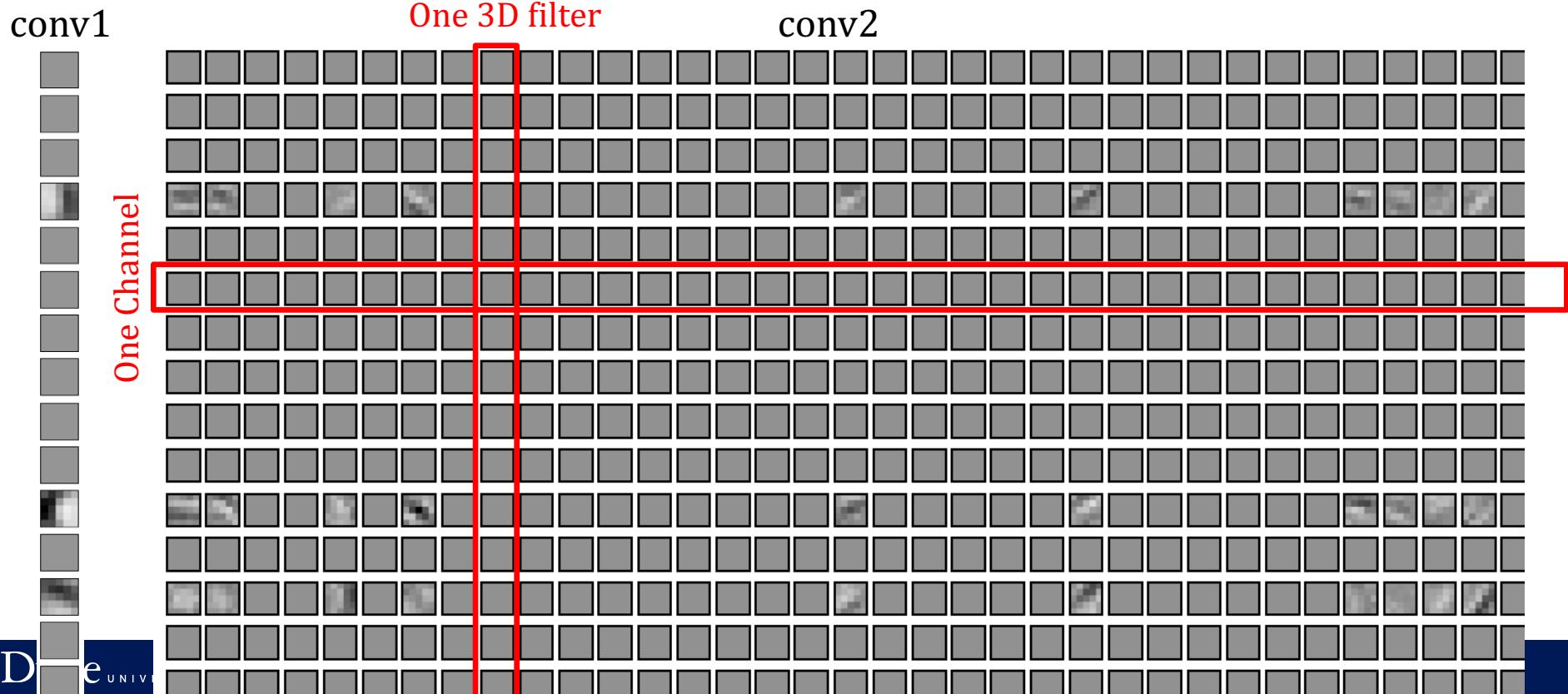
Easy to Implement

```
99 def add_dimen_grouplasso(var, axis=0):
100     with tf.name_scope("DimenGroupLasso"):
101         t = tf.square(var)
102         t = tf.reduce_sum(t, axis=axis) + tf.constant(1.0e-8)
103         t = tf.sqrt(t)
104         reg = tf.reduce_sum(t)
105         return reg
```

```
360     glasso_reg = add_dimen_grouplasso(train_var, axis=0)
361     self._regularization = self._regularization + glasso_reg * coef
```

```
408     grads, _ = tf.clip_by_global_norm(tf.gradients(cost + self._regularization, tvars),
409                                         config.max_grad_norm)
```

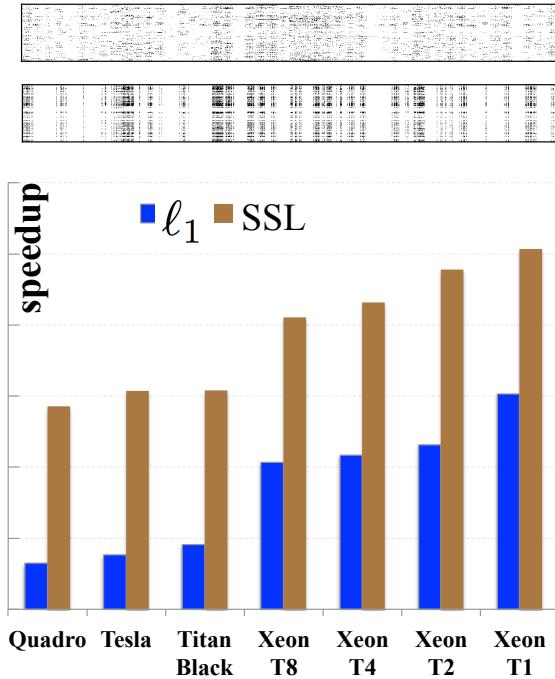
Structurally Sparse LeNet – removing channels and filters



Structurally Sparse *AlexNet* – removing columns and rows

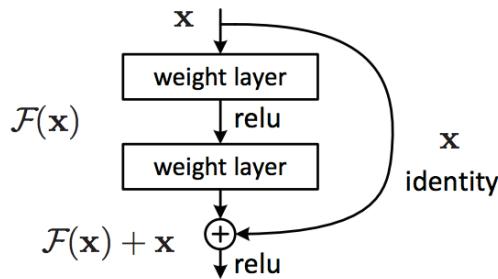
Table 4: Sparsity and speedup of *AlexNet* on ILSVRC 2012

#	Method	Top1 err.	Statistics	conv1	conv2	conv3	conv4	conv5
1 2% loss	ℓ_1	44.67%	sparsity	67.6%	92.4%	97.2%	96.6%	94.3%
			CPU ×	0.80	2.91	4.84	3.83	2.76
			GPU ×	0.25	0.52	1.38	1.04	1.36
2	SSL	44.66%	column sparsity	0.0%	63.2%	76.9%	84.7%	80.7%
			row sparsity	9.4%	12.9%	40.6%	46.9%	0.0%
			CPU ×	1.05	3.37	6.27	9.73	4.93
3	pruning [6]	42.80%	sparsity	16.0%	62.0%	65.0%	63.0%	63.0%
			sparsity	14.7%	76.2%	85.3%	81.5%	76.3%
			CPU ×	0.34	0.99	1.30	1.10	0.93
4 No loss	ℓ_1	42.51%	sparsity	0.08	0.17	0.42	0.30	0.32
			column sparsity	0.00%	20.9%	39.7%	39.7%	24.6%
			CPU ×	1.00	1.27	1.64	1.68	1.32
5	SSL	42.53%	sparsity	1.00	1.25	1.63	1.72	1.36



[6] S. Han, et. al. NIPS 2015

Structurally Sparse *ResNets* – removing layers



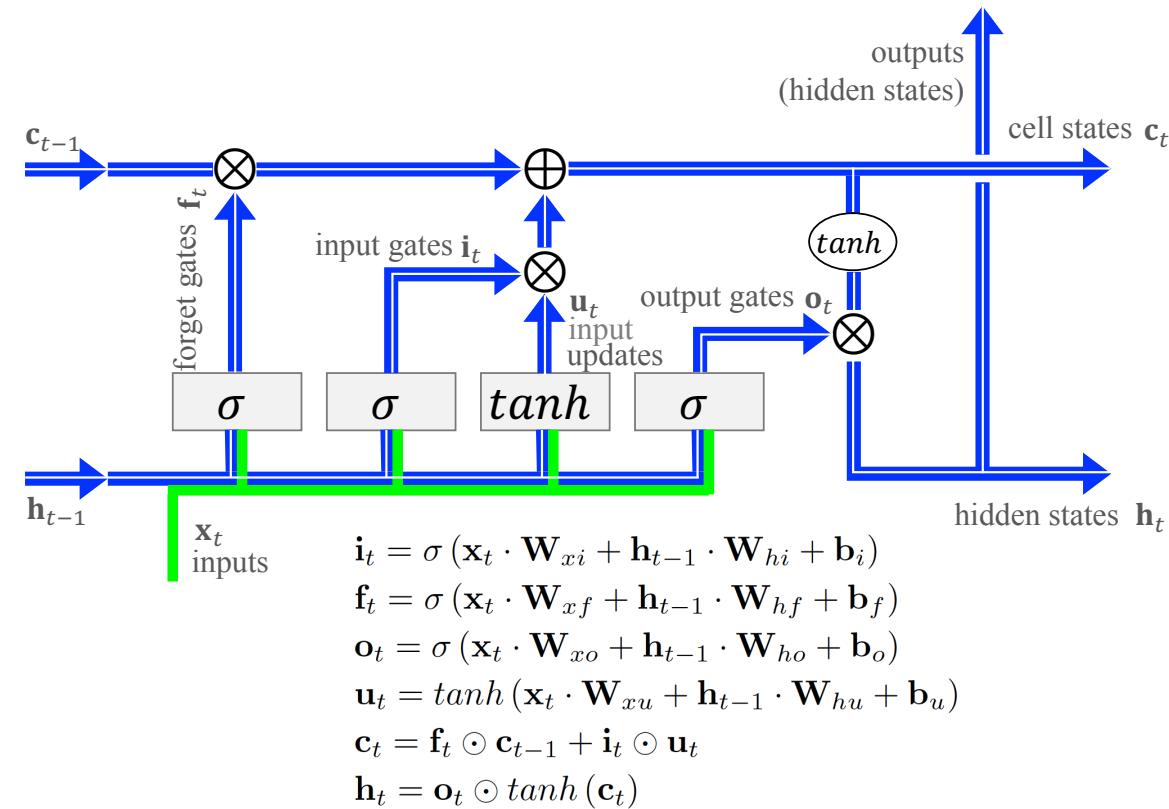
He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In *CVPR*. 2016.

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

SSL-ResNet-#: Ours with # layers after removing layers in ResNet-20

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

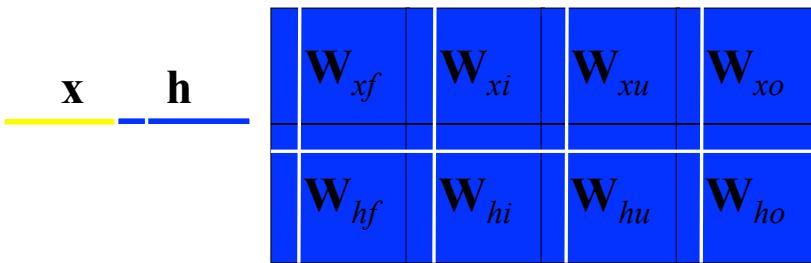
Hidden Structures in LSTMs



“Dimension Consistency”:

1. **Hidden Structures (in blue):** hidden states, input updates, cells, gates, outputs
2. All hidden structures must have the same dimension
3. Hidden structures cannot be independently removed

Learning Sparse Hidden Structures in LSTMs



Weights matrices in LSTM



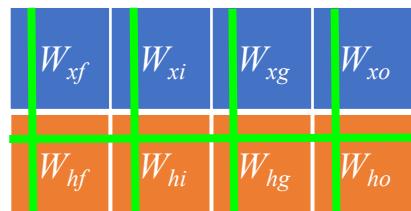
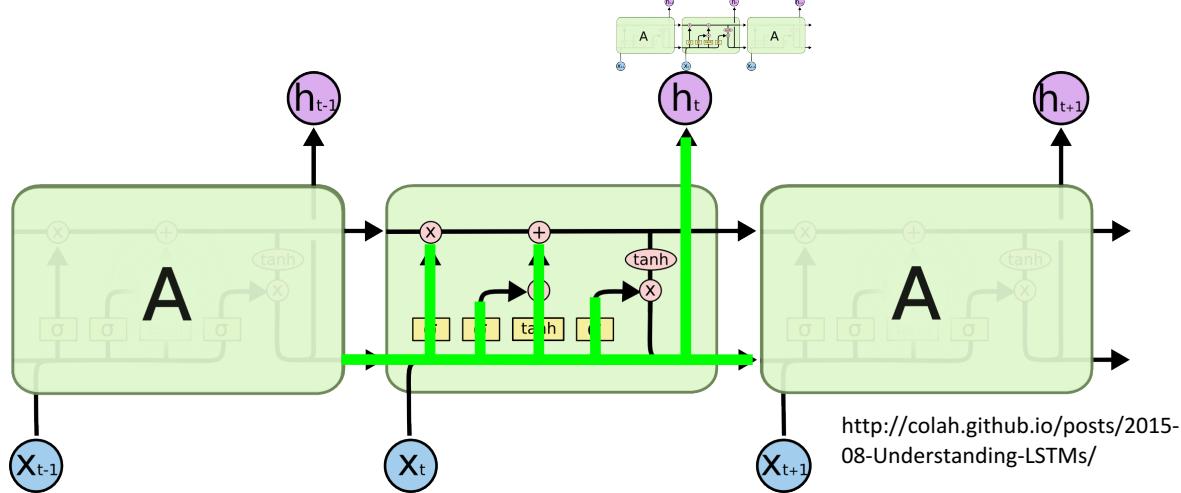
Weights in next layer(s)

$$\begin{aligned}\mathbf{i}_t &= \sigma(\mathbf{x}_t \cdot \mathbf{W}_{xi} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{hi} + \mathbf{b}_i) \\ \mathbf{f}_t &= \sigma(\mathbf{x}_t \cdot \mathbf{W}_{xf} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{hf} + \mathbf{b}_f) \\ \mathbf{o}_t &= \sigma(\mathbf{x}_t \cdot \mathbf{W}_{xo} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{ho} + \mathbf{b}_o) \\ \mathbf{u}_t &= \tanh(\mathbf{x}_t \cdot \mathbf{W}_{xu} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{hu} + \mathbf{b}_u) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{u}_t \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t)\end{aligned}$$

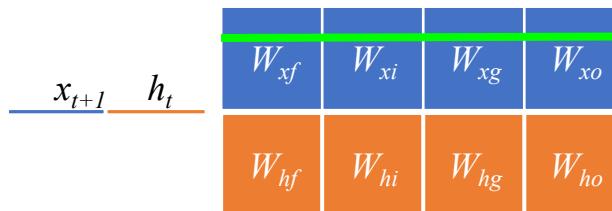
White strips (ISS: Intrinsic Sparse Structures) =
one group of weights in SSL

Removing one group = reduce hidden size by one

Structurally Sparse LSTMs



LSTM 1



LSTM 2

A group in SSL

Removing one group ==
reducing hidden size by one

Structurally Sparse LSTMs

Method	Dropout keep ratio	Perplexity (validate, test)	ISS # in (1st , 2nd) LSTM	Weight #	Total time*	Speedup	Mult-add reduction†
baseline	0.35	(82.57, 78.57)	(1500, 1500)	66.0M	157.0ms	1.00×	1.00×
ISS	0.60	(82.59, 78.65)	(373, 315)	21.8M	14.82ms	10.59×	7.48×
		(80.24, 76.03)	(381, 535)	25.2M	22.11ms	7.10×	5.01×
direct design	0.55	(90.31, 85.66)	(373, 315)	21.8M	14.82ms	10.59×	7.48×

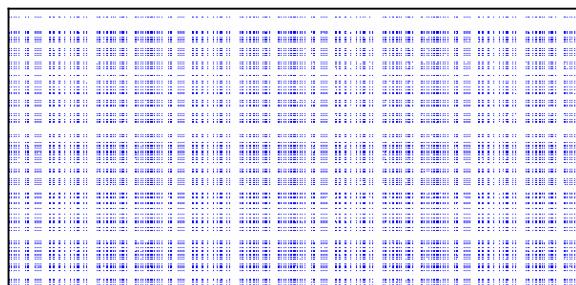
* Measured with 10 batch size and 30 unrolled steps.

† The reduction of multiplication-add operations in matrix multiplication. Defined as (original Mult-add)/(left Mult-add)

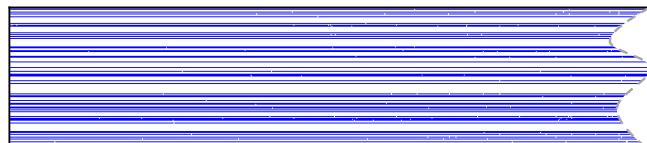
LSTM 1



LSTM 2



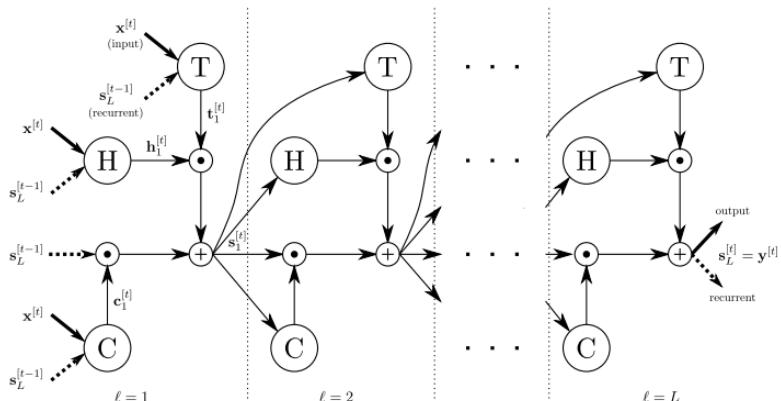
Output



Learned structurally
sparse LSTMs

Structurally Sparse Recurrent Highway Networks

What weights does a group have?



Zilly, Julian Georg, Rupesh Kumar Srivastava, Jan Koutník, and Jürgen Schmidhuber. "Recurrent highway networks." *arXiv:1607.03474* (2016).

Table 2: Learning ISS sparsity from scratch in RHNs.

Method	λ	Perplexity (validate, test)	RHN width	Parameter #
baseline	0.0	(67.9, 65.4)	830	23.5M
ISS*	0.004	(67.5, 65.0)	726	18.9M
Ours	0.005	(68.1, 65.4)	517	11.1M
ISS*	0.006	(70.3, 67.7)	403	7.6M
ISS*	0.007	(74.5, 71.2)	328	5.7M

* All dropout ratios are multiplied by 0.6×.

Question Answering with SSDNNs

Question Answering
in search engine:

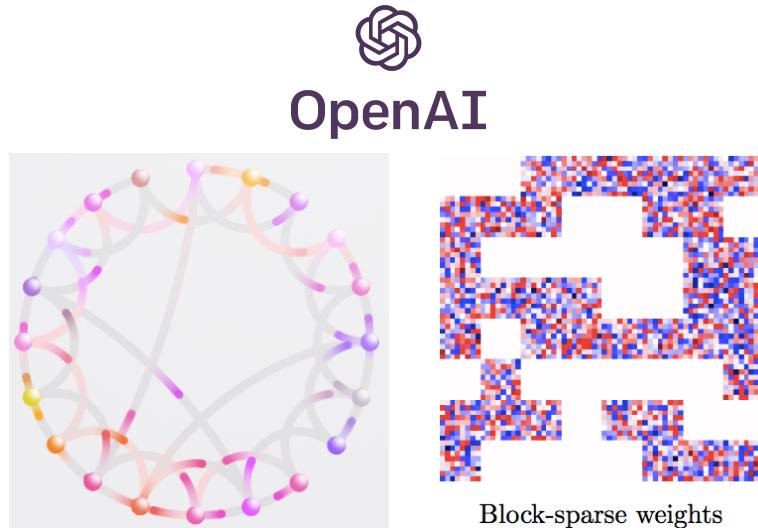


Table 4: Remaining ISS components in BiDAF by training from scratch.

	EM	F1	ModFwd1	ModBwd1	ModFwd2	ModBwd2	OutFwd	OutBwd	weight #	Total time*
Baseline	67.98	77.85	100	100	100	100	100	100	2.69M	6.20ms
Ours	67.36	77.16	87	81	87	92	74	96	2.29M	5.83ms
	66.32	76.22	51	33	42	58	37	26	1.17M	4.46ms
	65.36	75.78	20	33	40	38	31	16	0.95M	3.59ms
	64.60	74.99	23	22	35	35	25	14	0.88M	2.74ms

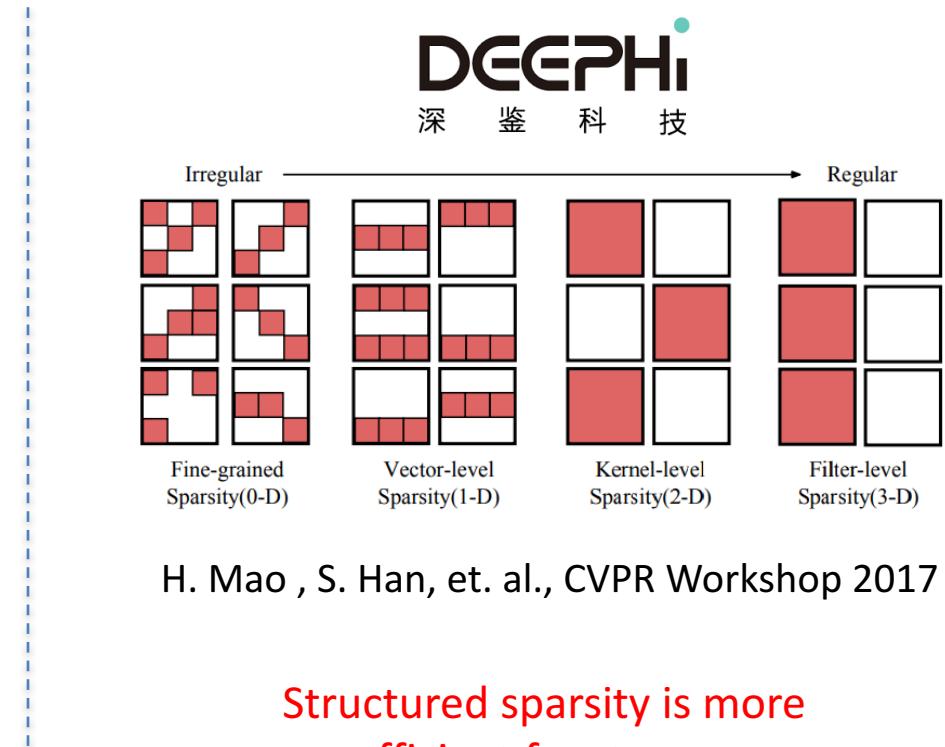
* Measured with batch size 1.

Efficiency of Structurally Sparse DNNs (after us)



S. Gray, A. Radford and D. P. Kingma, GPU Kernels for Block-Sparse Weights, OpenAI 2017

Structured sparsity is more efficient for computation



H. Mao , S. Han, et. al., CVPR Workshop 2017

Structured sparsity is more efficient for storage

Conclusion

- AI systems are landing into the cloud and upon the edge
- Communication bottleneck limits the scalability of distributed deep learning in the cloud
 - We propose *TernGrad* to quantize gradients to reduce communication volume.
- Lightweight computing engine restricts the deployment of cumbersome Deep Neural Networks on the edge
 - We propose Structurally Sparse Deep Neural Networks to simplify the “rocket” to lighten the load on the “engine”
- More done/on-going work on large-batch training and model acceleration
 - <http://www.pittnuts.com/#Publications>



THANKS! Q&A?

Wei Wen, Duke University
www.pittnuts.com