



Objective and Contributions

- ▶ A novel notion of **Boolean variation** and new mathematical framework of its calculus providing the chain rule similar to continuous gradient.
- ▶ A novel **Boolean logic backpropagation and optimization** method allowing for deep neural networks to **operate solely with Boolean logic** and to be **trained directly in Boolean domain**.
- ▶ **State-of-the-art results** compared to binarized neural networks, evaluated on challenging tasks with ConvNets, Transformers, ...
- ▶ **Significantly energy-efficient**: both training & inference.

Boolean Neural Networks

Boolean Neuron. Let L be a logic gate such as AND, OR, XOR, XNOR. The neuron's pre-activation output is given as follows:

$$s = w_0 + \sum_{i=1}^m L(w_i, x_i), \quad \text{where } w_i, x_i \in \mathbb{B} := \{\text{T}, \text{F}\}. \quad (1)$$

Mixed Boolean-Real Neuron. For flexibility, we can extend to Boolean weights with real-valued inputs, and real-valued weights with Boolean inputs.

Forward Activation. $y = \text{T}$ if $s \geq \tau$ and $y = \text{F}$ if $s < \tau$ where s is the preactivation, τ is a fixed or learnable scalar threshold.

Boolean Training

Algorithm 1: Illustration with a FC layer.

```

Input   : Learning rate  $\eta$ , number of iterations  $T$ ;
Initialize:  $m_{i,j}^{l,0} = 0$ ;  $\beta^0 = 1$ ;
1 for  $t = 0, \dots, T - 1$  do
2   Compute  $x^{l+1,t}$  following Eq. 2; ; /* 1. Forward pass */
3   Receive  $\frac{\delta \text{Loss}}{\delta x_{k,i}^{l+1,t}}$  from downstream layer; ; /* 2. Backward pass */
4   Compute and backpropagate  $g^{l,t}$  of Eq. 6; ; /* 2.1 Backpropagation */
5    $N_{\text{tot}} := 0, N_{\text{unchanged}} := 0$ ; ; /* 2.2 Weight update process */
6   foreach  $w_{i,j}^l$  do
7     Compute  $q_{i,j}^{l,t+1}$  following Eq. 5;
8     Update  $m_{i,j}^{l,t+1} = \beta^t m_{i,j}^{l,t} + \eta^t q_{i,j}^{l,t+1}$ ;
9      $N_{\text{tot}} \leftarrow N_{\text{tot}} + 1$ ;
10    if xnor( $m_{i,j}^{l,t+1}, w_{i,j}^l$ ) = T then
11       $w_{i,j}^{l,t+1} = \neg w_{i,j}^l, m_{i,j}^{l,t+1} = 0$ ; /* invert */
12    else
13       $w_{i,j}^{l,t+1} = w_{i,j}^l, N_{\text{unchanged}} \leftarrow N_{\text{unchanged}} + 1$ ; /* keep */
14    end
15  end
16  Update  $\eta^{t+1}, \beta^{t+1} = N_{\text{unchanged}}/N_{\text{tot}}$ ;
17 end

```

Boolean Variation

Definition 1. Order relations in \mathbb{B} are defined as: $\text{F} < \text{T}$, and $\text{T} > \text{F}$.

Definition 2. For $a, b \in \mathbb{B}$, the variation from a to b , denoted $\delta(a \rightarrow b)$, is defined as: $\delta(a \rightarrow b) \stackrel{\text{def}}{=} \text{T}$ if $b > a$, $\stackrel{\text{def}}{=} 0$ if $b = a$, and $\stackrel{\text{def}}{=} \text{F}$ if $b < a$.

Definition 3. For $f \in \mathcal{F}(\mathbb{B}, \mathbb{D})$, $\forall x \in \mathbb{B}$, write $\delta f(x \rightarrow \neg x) := \delta(f(x) \rightarrow f(\neg x))$. The variation of f w.r.t. x , i.e., $f'(x)$, is defined as: $f'(x) \stackrel{\text{def}}{=} \mathbf{xnor}(\delta(x \rightarrow \neg x), \delta f(x \rightarrow \neg x))$. Here, \mathbb{D} is either a logic set \mathbb{B} or a numeric set, e.g., \mathbb{R} or \mathbb{Z} .

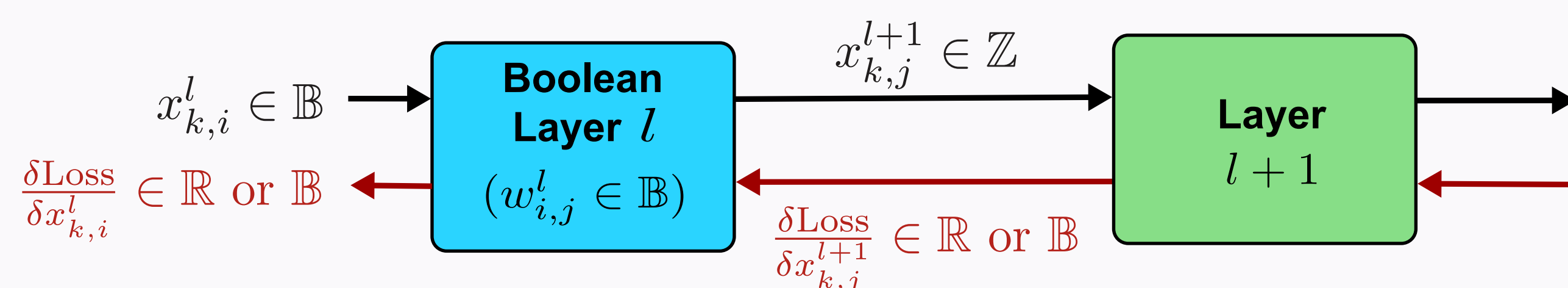
Remark. The variation of f w.r.t. x is T if f varies in same direction with x .

Definition 4. For $f \in \mathcal{F}(\mathbb{Z}, \mathbb{D})$, the variation of f w.r.t. $x \in \mathbb{Z}$ is defined as $f'(x) \stackrel{\text{def}}{=} \delta f(x \rightarrow x + 1)$, where δf is the variation defined in \mathbb{D} .

Theorem 5. The following properties hold:

- For $\mathbb{B} \xrightarrow{f} \mathbb{B} \xrightarrow{g} \mathbb{D}$: $(g \circ f)'(x) = \mathbf{xnor}(g'(f(x)), f'(x))$, $\forall x \in \mathbb{B}$.
- For $\mathbb{B} \xrightarrow{f} \mathbb{Z} \xrightarrow{g} \mathbb{D}$, $x \in \mathbb{B}$, if $|f'(x)| \leq 1$ and $g'(f(x)) = g'(f(x) - 1)$, then: $(g \circ f)'(x) = \mathbf{xnor}(g'(f(x)), f'(x))$.

Boolean Backpropagation



Consider the Boolean l -th layer, which is assumed a fully-connected layer

$$x_{k,j}^{l+1} = w_{0,j}^l + \sum_{i=1}^m L(x_{k,i}^l, w_{i,j}^l), \quad 1 \leq j \leq n, \quad (2)$$

where k is sample index, m and n are input and output sizes.

Remark. Layer l is connected to layer $l + 1$ that can be an activation layer, a batch normalization, an arithmetic layer, or any others.

Atomic Variation. Consider $L = \mathbf{xnor}$

$$q_{i,j,k}^l := \frac{\delta \text{Loss}}{\delta w_{i,j}^l} \Big|_k = \mathbf{xnor}\left(\frac{\delta \text{Loss}}{\delta x_{k,j}^{l+1}}, \frac{\delta x_{k,j}^{l+1}}{\delta w_{i,j}^l}\right) = \mathbf{xnor}\left(\frac{\delta \text{Loss}}{\delta x_{k,j}^{l+1}}, x_{k,i}^l\right), \quad (3)$$

$$g_{k,i,j}^l := \frac{\delta \text{Loss}}{\delta x_{k,i}^l} \Big|_j = \mathbf{xnor}\left(\frac{\delta \text{Loss}}{\delta x_{k,j}^{l+1}}, \frac{\delta x_{k,j}^{l+1}}{\delta x_{k,i}^l}\right) = \mathbf{xnor}\left(\frac{\delta \text{Loss}}{\delta x_{k,j}^{l+1}}, w_{i,j}^l\right). \quad (4)$$

Aggregation. Using the chain rules given by Theorem 5, we have

$$q_{i,j}^l := \frac{\delta \text{Loss}}{\delta w_{i,j}^l} = \sum_k \mathbf{1}(q_{i,j,k}^l = \text{T}) |q_{i,j,k}^l| - \sum_k \mathbf{1}(q_{i,j,k}^l = \text{F}) |q_{i,j,k}^l|, \quad (5)$$

$$g_{k,i}^l := \frac{\delta \text{Loss}}{\delta x_{k,i}^l} = \sum_j \mathbf{1}(g_{k,i,j}^l = \text{T}) |g_{k,i,j}^l| - \sum_j \mathbf{1}(g_{k,i,j}^l = \text{F}) |g_{k,i,j}^l|. \quad (6)$$

Boolean Optimizer. With $q_{i,j}^l$ obtained in Eq. 5, the rule for optimizing $w_{i,j}^l$ subjected to making the loss decreased is simply given according to its definition as:

$$w_{i,j}^{l,t+1} = \neg w_{i,j}^l \text{ if } \mathbf{xnor}(q_{i,j}^l, w_{i,j}^l) = \text{T}. \quad (7)$$

Binarized Neural Networks — An Approximated Binary

They learn binary weights, \mathbf{w}_{bin} , through **full-precision latent weights**, \mathbf{w}_{fp} using common gradient-descent methods, i.e.

$$\mathbf{w}_{\text{bin}} = \text{sign}(\mathbf{w}_{\text{fp}} - \eta \cdot \mathbf{g}_{\mathbf{w}_{\text{fp}}})$$

- ✗ No gains during the training!
- ✗ Sub-optimal performance due the approximations of latent weight gradient!

Experimental Results

Our method, B \oplus LD, significantly outperforms the SOTA binarized neural networks (BNNs) both in terms of predictive performance and energy efficiency.

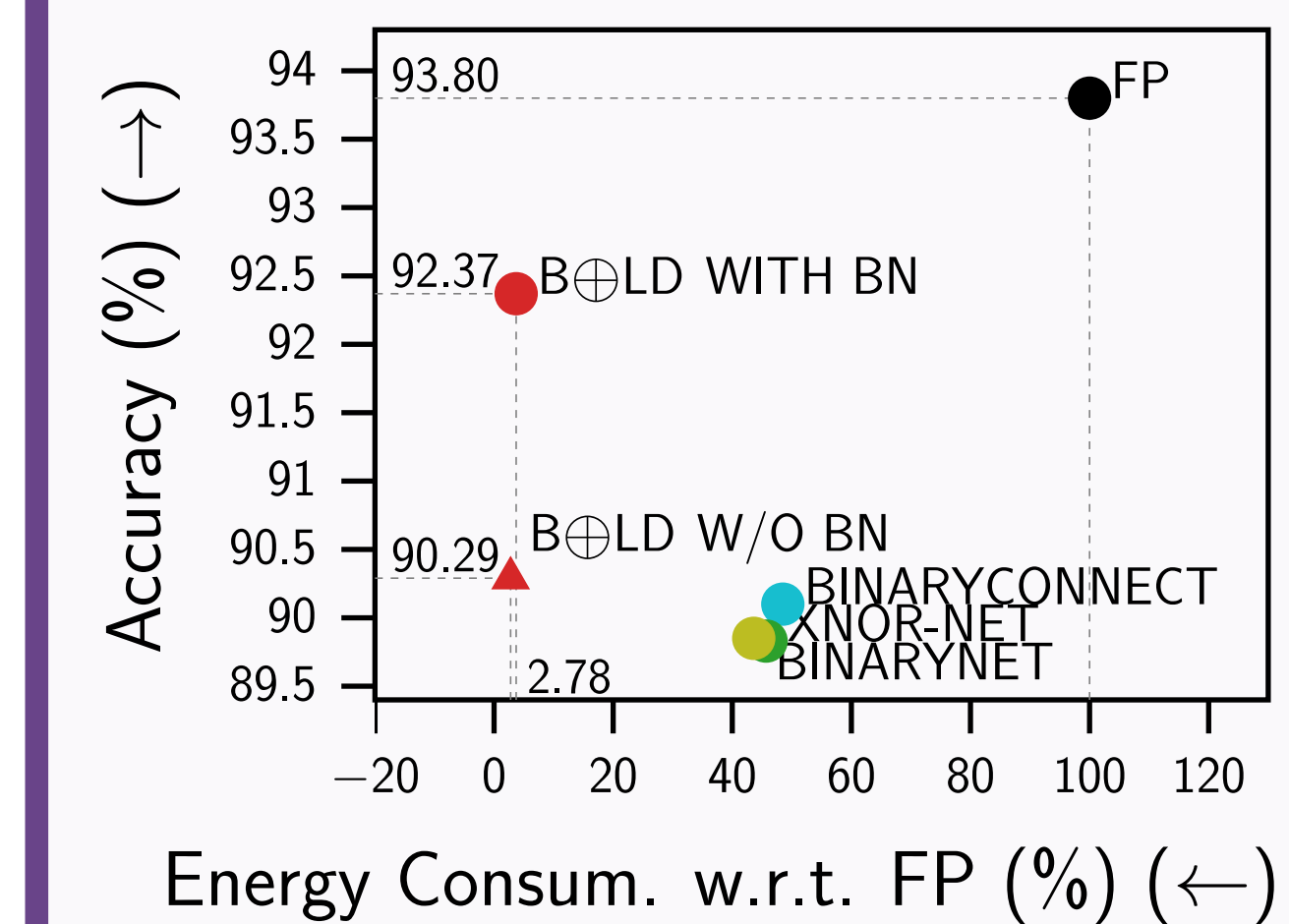


Table: Natural language understanding results with BERT models.

Method	GLUE Benchmark (Accuracy, \uparrow)							Avg.	
	MNLI	QQP	QNLI	SST-2	COLA	SST-B	MRPC		RTE
FP BERT	84.9	91.4	92.1	93.2	59.7	90.1	86.3	72.2	83.9
BINARYBERT	35.6	66.2	51.5	53.2	0.0	6.1	68.3	52.7	41.0
BIBERT	66.1	84.8	72.6	88.7	25.4	33.6	72.5	57.4	63.2
BIT	77.1	82.9	85.7	87.7	25.1	71.1	79.7	58.8	71.0
B \oplus LD	75.6	85.9	84.1	88.7	27.1	68.7	78.4	58.8	70.9

Figure: Comparisons of our method against notable BNNs on CIFAR10 with VGG-SMALL.

Table: Image classification results with RESNET18 on IMAGENET. 'Base' is the mapping dimension of 1st layer.

Training Modality	Method	Acc. Cons. (%)	
		(%)	Tesla V100
FP BASELINE	RESNET18	69.7	100.00
FP SHORTCUT	BINARYNET	42.2	---
	XNOR-NET	51.2	---
	B \oplus LD + BN	51.8	3.87
	BI-REALNET:18	56.4	32.99
LARGE MODELS	BI-REALNET:34	62.2	58.24
	BI-REALNET:152	64.5	---
	MELIUS-NET:29	65.8	---
	B \oplus LD (Base 256)	66.9	24.45
KD: RESNET34	REAL2BINARY	65.4	---
	REACTNET-RESNET18	65.5	77.89
	BNEXT:18	68.4	37.51
	B \oplus LD + BN (Base 192)	65.9	16.91
KD: RESNET50	B \oplus LD (Base 256)	70.0	24.45
	POKEBNN-RESNET18	65.2	---

Table: Super-resolution results measured in PSNR (dB) (\uparrow), using the EDSR baseline.

Task	Method	SET5	SET14	BSD100	URBAN100	DIV2K
$\times 2$ FULL EDSR (FP)	FP BASELINE	38.11	33.92	32.32	32.93	35.03
	BINARYNET	38.01	33.63	32.19	31.60	34.67
	B \oplus LD	37.42	33.00	31.75	30.26	33.82
$\times 3$ FULL EDSR (FP)	FP BASELINE	34.65	30.52	29.25	---	31.26
	BINARYNET	34.37	30.24	29.10	---	30.93
	B \oplus LD	33.56	29.70	28.72	---	30.22
$\times 4$ FULL EDSR (FP)	FP BASELINE	32.46	28.80	27.71	26.64	29.25
	BINARYNET	32.17	28.53	27.62	26.14	29.04
	B \oplus LD	31.23	27.97	27.24	25.12	28.36

Table: Image segmentation results.

Dataset	Model	mIoU (%) (\uparrow)
CITYSCAPES	FP BASELINE	70.7
	BINARY DAD-NET	58.1
	B \oplus LD	67.4
PASCAL VOC 2012	FP BASELINE	72.1
	B \oplus LD	67.3

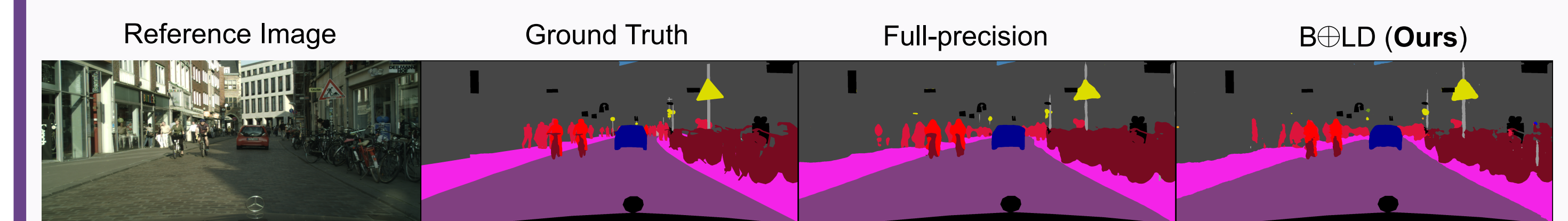


Figure: An example of CITYSCAPES.

References

- [1] Nguyen. "Boolean Variation and Boolean Logic BackPropagation". *Arxiv 2023*
- [2] Hubara et al. "Binarized Neural Networks". *NIPS 2016*