

Designing Robust DNN Models That Exploit Energy-Reliability Tradeoffs

François Leduc-Primeau

École Polytechnique de Montréal

Accelerate Al Workshop 2023 - Thursday May 4, 2023

Energy consumption of AI in perspective

- Large language models
 - 4 Wh / query × 3 queries/min
 = 720 W

[A. S. Luccioni et al. Estimating the Carbon Footprint of BLOOM, a 176B Parameter Language Model, 2022]



Energy consumption of AI in perspective

- Large language models
 - 4 Wh / query × 3 queries/min
 = 720 W

[A. S. Luccioni et al. Estimating the Carbon Footprint of BLOOM, a 176B Parameter Language Model, 2022]

Self-driving cars

NVIDIA "Robotaxi" platform:
 2× Orin SOCs + 2× Ampere GPUs
 = 800 W (TDP)





Energy consumption of AI in perspective

- Large language models
 - 4 Wh / query × 3 queries/min
 = 720 W

[A. S. Luccioni et al. Estimating the Carbon Footprint of BLOOM, a 176B Parameter Language Model, 2022]

Self-driving cars

NVIDIA "Robotaxi" platform:
 2× Orin SOCs + 2× Ampere GPUs
 = 800 W (TDP)







• Human brain: \approx 20 W

Digital IC Context

Issues:

- Energy efficiency of CMOS is not improving
- Data movement consumes most of the energy

(Partial) solutions:

- Reduce data movement
- Increase on-chip memory
- Use specialized hardware
- In-memory computing

SRAM Energy vs Reliability





[Di Mauro et al., TCAS-I, 2020]

Memristor crossbar :



- $G_{i,j}$ is the conductance value
- ► X_i is the input voltage
- Z_j is the output voltage

Memristor crossbar :



- \blacktriangleright $G_{i,i}$ is the conductance value
- X_i is the input voltage
- Z_j is the output voltage

Memristor crossbar :



 $Z_j = r \sum_{i=1}^L G_{i,j} X_i$

- $G_{i,j}$ is the conductance value
- X_i is the input voltage
- Z_j is the output voltage

Memristor crossbar :



 $Z_j = r \sum_{i=1}^L G_{i,j} X_i$

- $G_{i,j}$ is the conductance value
- X_i is the input voltage
- \blacktriangleright Z_j is the output voltage

Memristor crossbar : X, No. G_{1,1} G_{1,3} G_{1,M} G, X₂_ $\mathbf{G}_{2,M}$ G, ` G23 X \sim G G_{12} G $\mathbf{G}_{\mathrm{L,M}}$

 $Z_j = r \sum_{i=1}^L G_{i,j} X_i$

- $G_{i,j}$ is the conductance value
- X_i is the input voltage
- \blacktriangleright *Z_j* is the output voltage

Advantages:

- High density
- Non volatile
- Operates at low voltage
- Computation within memory



 $Z_j = r \sum_{i=1}^L G_{i,j} X_i$

- $G_{i,j}$ is the conductance value
- X_i is the input voltage
- Z_j is the output voltage

Advantages:

- High density
- Non volatile
- Operates at low voltage
- Computation within memory

Challenges:

- Writing can be slow
- Cannot have a negative resistance
- Device-to-device and cycle-to-cycle variations



- $G_{i,j}$ is the conductance value
- X_i is the input voltage
- Z_j is the output voltage

The conductance values are noisy



[Source: Joshi et al. Accurate deep neural network inference using computational phase-change memory. Nature Communications, 2020.]

Outline



- 2 Training robust models: Overview
- **3** Training robust models: Sharpness-aware training



4 Characterizing robustness: MemSE

Outline



2 Training robust models: Overview

3 Training robust models: Sharpness-aware training



Characterizing robustness: MemSE

F. Leduc-Primeau (Polytechnique Montréal)

Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)

"Bit-masking" faults:



- Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)
- "Bit-masking" faults:



fault is detected on the sign bit of a value

- Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)
- "Bit-masking" faults:



value replaced by zero

Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)

"Bit-masking" faults:



- Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)
- "Bit-masking" faults:



fault is detected on any other bit of a value

- Energy-reliability tradeoff: p(η) = e^{-aη}
 (η normalized energy, p bit flip probability)
- "Bit-masking" faults:



Affected bit value is replaced with the sign bit

DNN architecture comparison (CIFAR-10)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)



G. B. Hacene, et al. Training modern deep neural networks for memory-fault robustness. In IEEE ISCAS, 2019.

F. Leduc-Primeau (Polytechnique Montréal)

Outline



2 Training robust models: Overview

3 Training robust models: Sharpness-aware training



Characterizing robustness: MemSE

Sharpness-aware training

► Optimize for the loss value and the loss sharpness: "sharpness-aware minimization (SAM)" [Foret et al. ICLR 2021]

Sharpness-aware training

- ► Optimize for the loss value and the loss sharpness: "sharpness-aware minimization (SAM)" [Foret et al. ICLR 2021]
- Minimize the loss of a worst-case perturbation ε within a neighborhood region of radius ρ of a weight w:

$$L_{\mathsf{SAM}}(w) = \min_{w} \max_{\|\epsilon\|_2 \leq
ho} L(w + \epsilon)$$

Sharpness-aware training

- ► Optimize for the loss value and the loss sharpness: "sharpness-aware minimization (SAM)" [Foret et al. ICLR 2021]
- Minimize the loss of a worst-case perturbation ε within a neighborhood region of radius ρ of a weight w:

$$L_{\mathsf{SAM}}(w) = \min_{w} \max_{\|\epsilon\|_2 \le
ho} L(w + \epsilon)$$

- ► Adaptive sharpness-aware minimization (ASAM) builds on SAM: [Kwon et al. ICML 2021]
 - Sharpness is scale-invariant
 - Allows larger ρ

$$L_{\text{ASAM}}(w) = \min_{w} \max_{\|\epsilon/|w|\|_2 \le \rho} L(w+\epsilon)$$

SAMSON = Sharpness-aware minimization scaled by outlier normalization

Gonçalo Mordido, Sébastien Henwood, Sarath Chandar, François Leduc-Primeau. SAMSON: Sharpness-Aware Minimization Scaled by Outlier Normalization for Improving DNN Generalization and Robustness, arXiv 2211.11561, 2023.

F. Leduc-Primeau (Polytechnique Montréal)

SAMSON = Sharpness-aware minimization scaled by outlier normalization

Targeting <u>both</u> generalization and model robustness

Gonçalo Mordido, Sébastien Henwood, Sarath Chandar, François Leduc-Primeau. SAMSON: Sharpness-Aware Minimization Scaled by Outlier Normalization for Improving DNN Generalization and Robustness, arXiv 2211.11561, 2023.

F. Leduc-Primeau (Polytechnique Montréal)

- SAMSON = Sharpness-aware minimization scaled by outlier normalization
- Targeting <u>both</u> generalization and model robustness
- No noise simulation and no knowledge of the expected noise is required.

Gonçalo Mordido, Sébastien Henwood, Sarath Chandar, François Leduc-Primeau. SAMSON: Sharpness-Aware Minimization Scaled by Outlier Normalization for Improving DNN Generalization and Robustness, arXiv 2211.11561, 2023.

SAMSON = Sharpness-aware minimization scaled by outlier normalization

- Targeting <u>both</u> generalization and model robustness
- No noise simulation and no knowledge of the expected noise is required.

Loss function

$$L_{\text{SAMSON}}(w) = \min_{w} \max_{\|\epsilon\| \boldsymbol{w} \|_{\rho} / \|w\|_{2} \le \rho} L(w + \epsilon),$$

Gonçalo Mordido, Sébastien Henwood, Sarath Chandar, François Leduc-Primeau. SAMSON: Sharpness-Aware Minimization Scaled by Outlier Normalization for Improving DNN Generalization and Robustness, arXiv 2211.11561, 2023.

$$\blacktriangleright \quad \boldsymbol{g}_{i}^{\ell} = \boldsymbol{w}_{i}^{\ell} \times \frac{G_{\max}}{W_{\max}^{\ell}} \times \delta_{i}, \qquad \delta_{i} \sim \mathcal{N}(1, \sigma_{c}^{2})$$



$$\blacktriangleright \quad \boldsymbol{g}_{i}^{\ell} = \boldsymbol{w}_{i}^{\ell} \times \frac{G_{\max}}{W_{\max}^{\ell}} \times \delta_{i}, \qquad \delta_{i} \sim \mathcal{N}(1, \sigma_{c}^{2})$$



$$\blacktriangleright \quad \boldsymbol{g}_{i}^{\ell} = \boldsymbol{w}_{i}^{\ell} \times \frac{\boldsymbol{G}_{\max}}{\boldsymbol{W}_{\max}^{\ell}} \times \delta_{i}, \qquad \delta_{i} \sim \mathcal{N}(1, \sigma_{c}^{2})$$



F. Leduc-Primeau (Polytechnique Montréal)

$$\blacktriangleright \quad \boldsymbol{g}_{i}^{\ell} = \boldsymbol{w}_{i}^{\ell} \times \frac{\boldsymbol{G}_{\max}}{\boldsymbol{W}_{\max}^{\ell}} \times \delta_{i}, \qquad \delta_{i} \sim \mathcal{N}(1, \sigma_{c}^{2})$$



F. Leduc-Primeau (Polytechnique Montréal)

$$\blacktriangleright \quad \boldsymbol{g}_{i}^{\ell} = \boldsymbol{w}_{i}^{\ell} \times \frac{G_{\max}}{W_{\max}^{\ell}} \times \delta_{i}, \qquad \delta_{i} \sim \mathcal{N}(1, \sigma_{c}^{2})$$



F. Leduc-Primeau (Polytechnique Montréal)

Outline



- 2 Training robust models: Overview
- 3 Training robust models: Sharpness-aware training



Objective of MemSE

1. Develop analytical formulas of the mean squared error (MSE) of a neural network implemented using (noisy) memristors.

J. Kern, S. Henwood, G. Mordido, E. Dupraz, A. Aïssa-El-Bey, Y. Savaria, and F. Leduc-Primeau. MemSE: Fast MSE prediction for noisy memristor-based DNN accelerators. In IEEE AICAS, 2022.

Objective of MemSE

- 1. Develop analytical formulas of the mean squared error (MSE) of a neural network implemented using (noisy) memristors.
- 2. Hopefully they are faster to evaluate than Monte-Carlo simulations,

J. Kern, S. Henwood, G. Mordido, E. Dupraz, A. Aïssa-El-Bey, Y. Savaria, and F. Leduc-Primeau. MemSE: Fast MSE prediction for noisy memristor-based DNN accelerators. In IEEE AICAS, 2022.

Objective of MemSE

- 1. Develop analytical formulas of the mean squared error (MSE) of a neural network implemented using (noisy) memristors.
- 2. Hopefully they are faster to evaluate than Monte-Carlo simulations,

3. ... and differentiable.

J. Kern, S. Henwood, G. Mordido, E. Dupraz, A. Aïssa-El-Bey, Y. Savaria, and F. Leduc-Primeau. MemSE: Fast MSE prediction for noisy memristor-based DNN accelerators. In IEEE AICAS, 2022.

Noisy conductance values : $G_{i,j} = g_{i,j} + \delta_{i,j}^{q} + \epsilon_{i,j}^{v}$

Noisy conductance values : $G_{i,j} = g_{i,j} + \delta_{i,j}^q + \epsilon_{i,j}^v$

▶ g_{i,j} is the target conductance value 0 < g_{i,j} < G_{max}

Noisy conductance values : $G_{i,j} = g_{i,j} + \delta_{i,j}^{q} + \epsilon_{i,j}^{v}$

► g_{i,j} is the target conductance value 0 < g_{i,j} < G_{max}

• $\delta_{i,j}^{q}$ is the quantization error Uniform quantization with a resolution of $\Delta = \frac{G_{\text{max}}}{N}$

Noisy conductance values : $G_{i,j} = g_{i,j} + \delta_{i,j}^{q} + \epsilon_{i,j}^{v}$

- ► g_{i,j} is the target conductance value 0 < g_{i,j} < G_{max}
- $\delta_{i,j}^{q}$ is the quantization error Uniform quantization with a resolution of $\Delta = \frac{G_{\text{max}}}{N}$
- $\epsilon_{i,j}^{v}$ is the noise due to variability in conductance programming $\epsilon_{i,j}^{v} \sim \mathcal{N}(0, \sigma^{2})$

Noisy conductance values : $G_{i,j} = g_{i,j} + \delta_{i,j}^{q} + \epsilon_{i,j}^{v}$

- ► g_{i,j} is the target conductance value 0 < g_{i,j} < G_{max}
- $\delta_{i,j}^{q}$ is the quantization error Uniform quantization with a resolution of $\Delta = \frac{G_{\text{max}}}{N}$
- $\epsilon_{i,j}^{v}$ is the noise due to variability in conductance programming $\epsilon_{i,j}^{v} \sim \mathcal{N}(0, \sigma^{2})$

For a neural network weight $w_{i,j}$, the conductance value corresponds to :

$$g_{i,j} = c w_{i,j}$$

Scaling factor
$$c = \frac{G_{\text{max}}}{W_{\text{max}}}$$

Noisy memristor computations

Conductance can only be positive \rightarrow we split *G* as $G = G^+ - G^-$

For a standard DNN :With the memristor-DNN : $z_j = \sum_{i=1}^{L} w_{i,j} x_i$ $Z_j = \frac{1}{c} (\sum_{i=1}^{L} rG_{i,j}^+ X_i - \sum_{i=1}^{L} rG_{i,j}^- X_i)$

We estimate the performance of the noisy network compared with its standard counterpart through the MSE:

$$\begin{split} \mathsf{MSE}[Z_j] &= \mathbb{E}[(Z_j - z_j)^2] \\ &= \mathbb{V}[Z_j] + (\mathbb{E}[Z_j] - z_j)^2 \,. \end{split}$$

Noisy memristor computations

Conductance can only be positive \rightarrow we split *G* as $G = G^+ - G^-$

For a standard DNN :With the memristor-DNN : $z_j = \sum_{i=1}^{L} w_{i,j} x_i$ $Z_j = \frac{1}{c} (\sum_{i=1}^{L} rG_{i,j}^+ X_i - \sum_{i=1}^{L} rG_{i,j}^- X_i)$

We estimate the performance of the noisy network compared with its standard counterpart through the MSE:

$$\begin{aligned} \mathsf{MSE}[Z_j] &= \mathbb{E}[(Z_j - z_j)^2] \\ &= \mathbb{V}[Z_j] + (\mathbb{E}[Z_j] - z_j)^2 \,. \end{aligned}$$

Noisy memristor computations

Conductance can only be positive \rightarrow we split *G* as $G = G^+ - G^-$

For a standard DNN :With the memristor-DNN : $z_j = \sum_{i=1}^{L} w_{i,j} x_i$ $Z_j = \frac{1}{c} (\sum_{i=1}^{L} rG_{i,j}^+ X_i - \sum_{i=1}^{L} rG_{i,j}^- X_i)$

We estimate the performance of the noisy network compared with its standard counterpart through the MSE:

$$\begin{split} \mathsf{MSE}[Z_j] &= \mathbb{E}[(Z_j - z_j)^2] \\ &= \mathbb{V}[Z_j] + (\mathbb{E}[Z_j] - z_j)^2 \,. \end{split}$$

Experiments on CIFAR-10



Experiments on CIFAR-10



MemSE Runtime



 Crucial to develop robust DNN models to target energy efficient hardware and be robust to attacks.

- Crucial to develop robust DNN models to target energy efficient hardware and be robust to attacks.
- Robustness is affected by DNN architecture and is not the same for every layer.

- Crucial to develop robust DNN models to target energy efficient hardware and be robust to attacks.
- Robustness is affected by DNN architecture and is not the same for every layer.
- Training for robustness yields significant improvements.

- Crucial to develop robust DNN models to target energy efficient hardware and be robust to attacks.
- Robustness is affected by DNN architecture and is not the same for every layer.
- Training for robustness yields significant improvements.
- Fast analytical characterization is possible.