

Fully parallel stochastic gradient descent on 1k metal-insulator-graphene (MIG) crossbar array

Xudong Zhuang¹, and Tania Roy^{1#}

¹Duke University, 101 Science Dr., Durham, NC 27705, US

In-memory stochastic gradient descent (SGD) through fully parallel weight update on memristive crossbar array shows great promise in accelerating deep learning training on analog compute-in-memory architecture, as it reduces the time complexity for weight update into time complexity of $O(1)$ as compared to the conventional row-/column-wise update scheme of $O(N)$ [1, 2]. However, it requires ultrahigh linearity and symmetry in conductance tuning curve on memristive devices. To relax this high requirement on device, various algorithms have been proposed but failed to fulfill the potential of fully-parallel in-memory SGD [3, 4]. Existed device solution shows large cell area and limited experimentally-demonstrated size [5]. In this work, we propose a metal-insulator-graphene (MIG) interfacial memristor stack with ultrahigh linear and symmetric conductance tuning curve that enables fully parallel in-memory stochastic gradient descent. It also shows on/off ratio over 10^3 , which benefits signal-to-noise ratio. We experimentally fabricate and characterize 32×32 MIG crossbar array to demonstrate the scalability of this device.

References

- [1] Q. Xia, J. J. Yang, Nat. Mater. 18 (2019) 209-323.
- [2] T. Gokmen, Y. Vlasov, Front. Neurosci. 10 (2016) 333.
- [3] T. Gokmen, W. Haensch, Front. Neurosci. 14 (2020) 103.
- [4] M. J. Rasch, F. Carta, O. Fagbohunge, T. Gokmen, Nat. Comm. 15 (2024) 7133.
- [5] E. J. Fuller, S. T. Keene, A. Melianas, Z. Wang, S. Agarwal, Y. Li, Y. Tuchman, C. D. James, M. J. Marinella, J. J. Yang, A. Salleo, A. A. Talin, Science 364 (2019) 570-574.

* Corresponding author: email: tania.roy@duke.edu

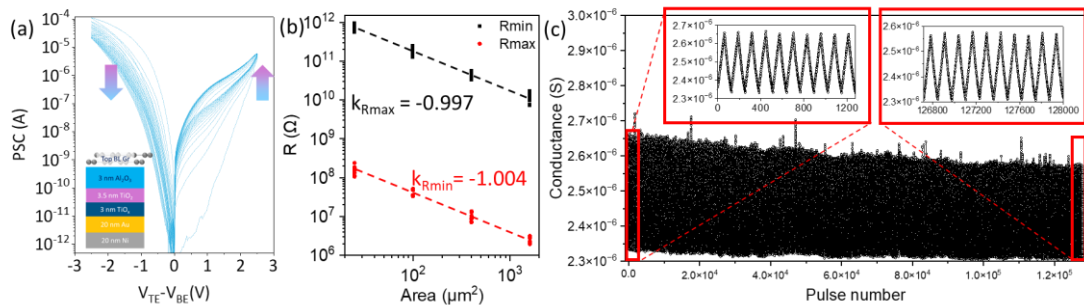


Fig. 1 (a) Consecutive DC double sweeps of MIG showing gradual programming properties. The MIG is potentiated when $V_{TE} > V_{BE}$ and depressed when $V_{TE} < V_{BE}$. (b) Maximum and minimum resistance of different device sizes change with area with slope closed to -1, proving non-filamentary properties. (c) Endurance of 128,000 read-write operations shows the linear and symmetric features can be maintained through multiple cycles. For each cycle, there are 64 potentiation and 64 depression pulses. The two insets are the zoom-in images of the 1st 10 and the last 10 cycles, showing high linearity and symmetry of conductance tuning curves.

