



LJMU Research Online

Joksas, D, Freitas, P, Chai, Z, Ng, WH, Buckwell, M, Li, C, Zhang, WD, Xia, QF, Kenyon, AJ and Mehonic, A

Committee Machines—A Universal Method to Deal with Non-Idealities in Memristor-Based Neural Networks

<http://researchonline.ljmu.ac.uk/id/eprint/13408/>

Article

Citation (please note it is advisable to refer to the publisher's version if you intend to cite from this work)

Joksas, D, Freitas, P, Chai, Z, Ng, WH, Buckwell, M, Li, C, Zhang, WD, Xia, QF, Kenyon, AJ and Mehonic, A Committee Machines—A Universal Method to Deal with Non-Idealities in Memristor-Based Neural Networks. Nature Communications. ISSN 2041-1723 (Accepted)

LJMU has developed **LJMU Research Online** for users to access the research output of the University more effectively. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in LJMU Research Online to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

The version presented here may differ from the published version or from the version of the record. Please see the repository URL above for details on accessing the published version and note that access may require a subscription.

For more information please contact researchonline@ljmu.ac.uk

<http://researchonline.ljmu.ac.uk/>

Committee Machines—A Universal Method to Deal with Non-Idealities in Memristor-Based Neural Networks

D. Joksas¹, P. Freitas², Z. Chai², W. H. Ng¹, M. Buckwell¹,
C. Li³, W. D. Zhang², Q. Xia³, A. J. Kenyon¹, and A. Mehonic¹

¹*Department of Electronic and Electrical Engineering,
University College London, London (United Kingdom)*

²*Department of Electronics and Electrical Engineering,
Liverpool John Moores University, Liverpool (United Kingdom)*

³*Department of Electrical and Computer Engineering,
University of Massachusetts Amherst (United States of America)*

Abstract

Artificial neural networks are notoriously power- and time-consuming when implemented on conventional von Neumann computing systems. Consequently, recent years have seen an emergence of research in machine learning hardware that strives to bring memory and computing closer together. A popular approach is to realise artificial neural networks in hardware by implementing their synaptic weights using memristive devices. However, various device- and system-level non-idealities usually prevent these physical implementations from achieving high inference accuracy. We suggest applying a well-known concept in computer science—committee machines—in the context of memristor-based neural networks. Using simulations and experimental data from three different types of memristive devices, we show that committee machines employing ensemble averaging can successfully increase inference accuracy in physically implemented neural networks that suffer from faulty devices, device-to-device variability, random telegraph noise and line resistance. Importantly, we demonstrate that the accuracy can be improved even without increasing the total number of memristors.

25 I. INTRODUCTION

26 Artificial neural networks (ANNs), with all of their variants, are now the main tools in
27 machine learning tasks, such as classification. The vast amounts of data being constantly
28 produced have enabled successful training and operation of ANNs. However, to achieve
29 high inference accuracy, it is usually necessary for neural networks to have a large number of
30 parameters. This results in both training [1] and inference [2] stages being time- and power-
31 consuming. This is largely caused by the need to transfer data from memory to computing
32 units—physical separation of memory and computing is the essence of any von Neumann
33 system.

34 One of the most promising solutions to these problems is the paradigm of non-von Neu-
35 mann computing and, specifically, analogue implementations of synapses (weights) in phys-
36 ical ANNs. Because there are many more synapses than there are neurons in ANNs, the
37 matrix-vector multiplications, in which the synaptic weight values are used, are the costli-
38 est operations in these networks, both in terms of power and time. Computing directly in
39 memory would minimise data transfers from off-chip memory, thus the most popular ap-
40 proach is using analogue memory devices as proxies for synaptic weights of ANNs (both
41 fully connected and their variants [3, 4]). A common technique is to arrange such devices
42 in a structure, called crossbar array, in which every device (or a pair of devices) is used to
43 represent a single synaptic weight or, more generally, an entry in a matrix [5]. Memristive
44 devices, such as phase-change memories (PCMs) [6, 7] or resistive random-access memories
45 (RRAMs) [8, 9], have been considered as candidates for such tasks. Although here we fo-
46 cus on ex-situ training, such systems have been successfully utilised for in-situ training too
47 [10, 11].

48 In memristive implementations of ANNs, the main concern is that various non-idealities
49 associated with these devices can prevent these systems from achieving high accuracy [12,
50 13]. Examples of non-idealities affecting inference accuracy include, but are not limited
51 to, devices not being able to electroform, devices stuck in one of the resistance states after
52 electroforming, device-to-device (D2D) variability and random telegraph noise (RTN). When
53 training analogue systems in-situ, limited endurance and non-linear resistance modulation
54 too have to be taken into account. To mitigate the effects of these device non-idealities, it is
55 often necessary to modify device structure [9], to use more advanced programming schemes

56 [14] or to use additional circuitry [15] or high-precision processing units [16] in conjunction
57 with memristive elements. On the system level, there is an issue of line resistance which
58 affects the distribution of currents and thus decreases the accuracy. These line resistance
59 effects can be partially compensated for algorithmically [17] or partially mitigated by using
60 multiple smaller crossbar arrays [18]. Examples of past efforts at dealing with these and
61 other non-idealities of memristive devices and systems are listed in Table I; most of these
62 non-idealities are still the main focus of the research in the neuromorphic community.

63 We propose a simple way to mitigate the effects of all types of non-idealities during
64 inference. We suggest combining several non-ideal memristor-based neural networks into
65 committees to achieve better accuracy. The committee machine (CM) method we propose
66 significantly increases the inference accuracy and does not increase the computation time
67 because memristive ANNs in such committees work in parallel.

68 In this work, we firstly explain the simulation setup—what networks were trained,
69 how they were simulated and how they were combined into CMs. After that, follows
70 the experimental part. We investigate three different types of memristor technology—
71 tantalum/hafnium oxide-based (Ta/HfO₂), tantalum oxide-based (Ta₂O₅), and amorphous
72 vacancy modulated conductive oxide-based (aVMCO) devices. By exploring their non-
73 idealities relevant to inference—faulty devices, D2D variability, RTN, and line resistance—
74 we use the experimental data to simulate memristive ANNs working individually and in
75 committees.

76 II. RESULTS

77 A. Simulation setup

78 Fully connected ANNs were trained in software to recognise handwritten digits (using
79 MNIST data base [19]). Architectures with one hidden layer were explored. Unless stated
80 otherwise, the simulations used networks with 25 hidden neurons. However, networks with
81 50, 100 and 200 hidden neurons were additionally employed to evaluate the effectiveness of
82 the proposed method while controlling for the total number of memristors required. Follow-
83 ing training, weights of ANNs were mapped onto pairs of conductances using proportional
84 mapping scheme (see [20]) to simulate memristor-based ANNs. Finally, these memristive

85 networks were disturbed using experimental data to reflect the effect of device- and system-
86 level non-idealities.

87 After simulating physical non-idealities, the networks were combined into CMs that em-
88 ployed ensemble averaging (EA) [21]. The principle of EA is shown in Figure 1A—several
89 networks are combined in parallel and then their outputs are averaged. After that, the
90 prediction is made using the averaged vector—the prediction is the label corresponding to
91 the largest entry in the vector.

92 CM methods are frequently used even with conventional ANNs. Methods, such as EA,
93 often produce better accuracy than that of the best individual network in a committee [22].
94 Although there are other types of CMs besides EA, they often rely on training additional
95 gating networks or boosting networks during the training stage. Using a gating network in
96 this scenario would produce additional problems—to avoid it acting as a performance bottle-
97 neck, it too would have to be implemented on crossbar arrays. Various non-idealities would
98 decrease the effectiveness of this gating network which is responsible for making the deci-
99 sions about the whole committee of ANNs. Likewise, we speculate that boosting of networks
100 would not be feasible in ex-situ training because it requires information about where indi-
101 vidual ANNs perform poorly—this cannot be known precisely until they are implemented
102 physically on crossbar arrays and the non-idealities manifest themselves. To authors' best
103 knowledge, the application of boosting in the context of memristive neural networks seems
104 to have been explored only once before [23]; as expected, it requires training each memristive
105 implementation differently because non-idealities manifest themselves differently in different
106 crossbar arrays.

107 There exist modifications of EA algorithm that could potentially perform better. One
108 example of this is generalized ensemble method (GEM) which, instead of using equal weight-
109 ings for each network during averaging (as in EA), uses a different one for each network [21].
110 These weightings are analytically determined by considering correlation of errors between
111 different networks. But because [21] only considered networks with mean square error loss
112 function (while our networks used cross-entropy loss function), this work does not explore
113 GEM. Instead, we investigated whether it is possible to achieve a better performance by
114 optimising the weightings numerically. This method, like GEM and others previously men-
115 tioned, might be impractical because, firstly, these weightings could be determined only after
116 the ANNs are physically implemented on crossbars, and, secondly, the devices could change

117 throughout their lifetimes thus affecting the optimal weightings.

118 Even with the assumption that the devices would have perfect retention, we found that
119 optimisation of weightings achieves effectively the same performance. Because of these rea-
120 sons, we focus only on EA in the main text, but present our results of optimising weightings
121 in Supplementary Figure S5. We stress that we are open to the idea that other CM methods
122 besides EA could be utilised successfully for ex-situ training in the context of memristive
123 ANNs. However, in this work we focus on demonstrating that CMs can be used to improve
124 the accuracy of memristor-based ANNs in general.

125 With EA, we find that even when the memristive ANNs, which go into a committee, all
126 use the same digital weights that are mapped onto crossbar arrays (see Figure 1B), committee
127 of memristor-based networks can still achieve higher accuracy than just a single non-ideal
128 network. Although all networks have the same *digital* weights before mapping, their physical
129 implementations (which we call "disturbances" in Figures 1B, C because they can usually
130 be represented by the modification of individual weights) will be different. For example, in
131 one crossbar array, a certain set of devices will be faulty, while in the other crossbar array, it
132 will be a different set. This will result in different physical implementations having slightly
133 different learned representations of the data set, or, to paraphrase, different networks will
134 be "damaged" differently by the non-idealities. This means that these committees will be
135 able to combine different representations, and thus achieve higher accuracy. However, by
136 definition, such approach would almost certainly not yield a committee accuracy that is
137 higher than the accuracy of a single digitally implemented network.

138 A better approach is to use different digital networks for different physical implementa-
139 tions that go into a committee (see Figure 1C). This approach much more resembles the
140 conventional application of EA in computer science. In the context of memristive crossbar
141 arrays, it would not only help to mitigate the effects of the non-idealities (as in the case
142 of Figure 1B), but would also allow to combine the representations of digital networks that
143 were different even before the mapping stage. Most importantly, this method allows for a
144 committee to achieve higher accuracy which is sometimes even higher than that of individual
145 networks with digitally implemented weights. We thus used this method in this analysis.
146 An example comparison of these two approaches is presented in Supplementary Figure S8.

147 In this work, any given committee used only one network architecture but each network
148 was initialised differently before training, thus trained networks had different sets of weights.

149 Although it was not explored in this work, combining different network architectures in a
150 committee of memristor-based networks might be advantageous. Furthermore, in this work
151 we focus on fully connected ANNs but CMs could be applied to other variants of neural
152 networks as well. Due to the simplicity of EA, it could, for example, be employed in con-
153 volutional neural networks (CNNs) [24], which are often used for image classification. This
154 might be of interest as CNNs have been successfully implemented using crossbar arrays re-
155 cently [25]. However, crossbar implementations are naturally more suited to fully connected
156 networks, therefore we limit ourselves to this architecture but are open to exploring the
157 effectiveness of EA with memristive CNNs in the future.

158 B. Ta/HfO₂ RRAM

159 With array-level data available, Ta/HfO₂ experiments provide the most complete pic-
160 ture of device- and system-level non-idealities. In this subsection, we present not only the
161 analysis of faulty devices and D2D variability, but also careful consideration of the line resis-
162 tance effects. Ta/HfO₂ memristors do not exhibit apparent RTN and overall have excellent
163 retention properties [26], and thus are perfect candidates for inference application.

164 1. *Faulty devices and device-to-device variability*

165 The most energy-efficient procedure to modulate the conductance of memristors is by
166 the application of voltage pulses. In an ideal scenario, one would apply identical pulses
167 and observe constant increases in conductance with each pulse. This is rarely the case
168 in practise, but, fortunately, this type of behaviour is more relevant for in-situ training
169 where it is necessary to ensure linear adjustment of ANN’s weights [27]. In ex-situ training,
170 conductance verification schemes can be used to program the devices precisely. Because the
171 devices would have to be programmed only once, one can spend additional resources to do so
172 accurately by applying SET (potentiation) and RESET (depression) pulses until a desirable
173 conductance state is achieved.

174 Even with this approach, there remain two obstacles—faulty devices and D2D variability.
175 It is observed in most memristor technologies that at least a small fraction of the devices
176 tends to get stuck in a particular conductance state. Additionally, even if not stuck, different

177 devices might behave differently; for example, they might have different conductance ranges.
178 Figure 2A shows conductance changes in Ta/HfO₂ RRAM devices (in a 128 × 64 crossbar
179 array) when they are applied with voltage pulses. We can see from the median values
180 that overall the devices' conductance tends to increase as more SET pulses are applied.
181 However, the wider bottom regions of the violin plots indicate that some devices are stuck
182 around high resistance state (HRS) and cannot set entirely no matter how many voltage
183 pulses are applied. There also exist devices that are stuck in low resistance state (LRS), or
184 simply do not span the full conductance range.

185 Figure 2A combines data from multiple SET cycles for each of the memristors, thus it
186 is important to understand how do these devices behave individually. Figures 2B-F show
187 conductance of 5 (out of 8,192) devices over 11 SET and RESET cycles. In the five dia-
188 grams, the radial component represents the conductance (in mS) and the angular component
189 represents the number of applied pulses. Figure 2B shows an example of preferable (and
190 typical) device behaviour—conductance changes in a continuous fashion and spans a wide
191 range of conductance values, from ~0.1 mS to ~1.0 mS. Although RESET cycles tend to
192 feature abrupt decreases in conductance, one can always repeat a cycle and exploit the more
193 predictable behaviour of SET cycles.

194 When encoding continuous numbers into crossbar devices' conductances, it is often prefer-
195 able to choose a large enough conductance range. Using data from Figure 2A, one could,
196 for example, choose the range between the first and the last median points (from ~0.1 mS
197 to ~1.0 mS). Device, whose behaviour is presented in Figure 2B, could be easily set to any
198 conductance within that range, as we have seen before. On the other hand, device, whose
199 behaviour is presented in Figure 2C, although operating in a predictable fashion, has smaller
200 conductance range. We can see that in all cycles, its conductance does not exceed 0.8 mS.
201 This is an example of D2D variability that can make it difficult to choose optimal operating
202 range and set the conductance of all devices precisely.

203 Device, whose behaviour is presented in Figure 2D, shows high cycle-to-cycle variability.
204 Although that could prove to be a problem in some applications, this specific device might
205 perfectly serve its purpose in ex-situ training of ANNs. We can observe that this device
206 spans the same conductance range as device from Figure 2B, even if in an unpredictable
207 manner. Because all states in the full range are, in theory, achievable, one can cycle the
208 device multiple times until it is set to the required conductance level.

209 Lastly, we have devices whose negative effect is most difficult to mitigate—faulty devices.
210 Figure 2E shows behaviour of a device stuck at high conductance values, while Figure 2F
211 shows behaviour of a device stuck at low conductance values. No matter how many pulses
212 the devices are applied with or how many times they are cycled, they exhibit almost no
213 conductance variation and thus, in most cases, cannot be used to encode information.

214 Knowing that some devices perform like the ones whose behaviour is shown in Fig-
215 ures 2C,E,F, it is important to minimise their negative effect. If the conductance that a
216 device has to be set to is outside that device’s range, it is sensible to set it to the closest
217 achievable conductance. Although there is little that can be done about fully stuck memris-
218 tors, it is possible to optimise the behaviour of devices like the one in Figure 2C that simply
219 have smaller conductance range. For example, if such a device has to be set to 0.9 mS, one
220 would set it to the highest achievable conductance (~ 0.8 mS). In the following simulations
221 involving faulty devices and D2D variability, operating range between the first and the last
222 median points was used, the devices were chosen randomly from the 128×64 crossbar and
223 set to the most desirable states, as described in this paragraph.

224 2. *Line resistance*

225 The effect of line resistance can be extremely detrimental in many crossbar-based im-
226 plementations of ANNs. That is especially the case if the crossbars used are large and the
227 resistance of the interconnects is high (compared to memristors’ resistance). Because in a
228 neural network many of the inputs are non-zero at any given time, a lot of current accumu-
229 lates in the bit lines which results in significant voltage drops across the interconnects, and
230 thus the current distribution across the crossbar is affected in a major way.

231 The Ta/HfO₂ crossbar has shape 128×64 and so this shape was chosen for all the simula-
232 tions involving line resistance. Even relatively small ANNs of architecture $784(+1):25(+1):10$
233 would need $2 \times (785 \times 25 + 26 \times 10) = 39,770$ memristors to be implemented. Even if not
234 all the inputs were used at any given time, it would not be possible to fit all the memristors
235 onto a single crossbar of shape 128×64 . To overcome this, we decided to simulate multiple
236 crossbars, each of which would implement a subset of the synaptic weights, but, for a given
237 synaptic layer, would all compute in parallel. Because $\lceil 785/128 \rceil = 7$, seven crossbars were
238 used to implement the first synaptic layer; the first crossbar utilized bottom 113 word lines,

239 while the other six crossbars used bottom 112 word lines because $113 + 6 \times 112 = 785$. The
240 second synaptic layer was implemented using eighth crossbar utilizing its bottom 26 word
241 lines.

242 Figure 3A shows an example of how the first synaptic layer of $784(+1):25(+1):10$ neural
243 network could be implemented. Specifically, it shows how the first subset of weights would
244 be implemented using one of the crossbars. Because we use proportional mapping scheme,
245 positive and negative weights would be implemented in different bit lines. In Figure 3A,
246 memristors designated to implement positive weights are coloured in blue, memristors desig-
247 nated to implement negative weights are coloured in orange and unelectroformed memristors
248 are coloured in black. Because simulations were constrained by experimental data, some of
249 the devices were left unused and assumed to be unelectroformed. In practise, the crossbars
250 could be manufactured to fit the geometry of the ANNs.

251 In each synaptic layer, the corresponding output currents from each of the crossbars would
252 be added together. Additionally, output currents at the bit lines implementing negative
253 weights would be subtracted from the output currents at the neighbouring bit lines (to their
254 left) implementing positive weights. For example, in the example configuration of Figure 3A,
255 output current at the 2nd bit line would be subtracted from the output current at the 1st bit
256 line, etc.

257 Unfortunately, even when using multiple smaller crossbars, the interconnects can signif-
258 icantly disturb current distribution in the crossbar. Average output current decreases due
259 to line resistance in all seven crossbars of Ta/HfO₂ devices (whose resistance ranges from
260 $\sim 1 \text{ k}\Omega$ to $\sim 11 \text{ k}\Omega$, and their interconnect resistance is 0.35Ω and 0.32Ω in the word and bit
261 lines, respectively), are shown in the heatmap in Figure 3B. We can see that the current
262 decreases can range from $\sim 12\%$ at the outputs nearest to the applied voltages to $\sim 16\%$ at
263 the outputs in the rightmost bit lines that are used. In the supplementary information, we
264 provide a possible strategy of mitigating line resistance effects in supervised learning. This
265 scheme was not employed in the simulations described in the main text because we wanted
266 to find out how well the CM method would deal with noticeable line resistance effects.

267 3. *Inference accuracy*

268 Figure 4 shows the accuracy of individual networks, as well as of their committees; mem-
269 ristic ANNs were simulated by taking into account three non-idealities of Ta/HfO₂ crossbar
270 explored earlier—faulty devices, D2D variability and line resistance. As indicated by the
271 yellow box plot in Figure 4, individual networks implemented digitally achieve ~95.9% me-
272 dian accuracy. Networks disturbed to reflect the effect of non-idealities achieve ~91.0%
273 median accuracy, as indicated by the vermilion box plot. Although that is a substantial
274 drop in accuracy, we see that as more networks are added to the committee, the more the
275 accuracy increases. When 5 networks are used in a committee, median accuracy increases
276 up to ~95.7%, as indicated by the rightmost green box plot.

277 **C. Ta₂O₅ RRAM**

278 In order to explore the effectiveness of minimising adverse effects of RTN, we use another
279 memristor technology based on Ta₂O₅. To investigate RTN, measurements from a single
280 device were considered. To simulate line resistance effects, interconnect resistance from
281 Ta/HfO₂ was used and the same crossbar shape was assumed.

282 1. *Random telegraph noise*

283 Memristors often suffer from RTN resulting in a different accuracy at any given instant
284 in time. Ta₂O₅ device was characterised by measuring the current of 8 resistance states
285 multiple times. Figure 5 shows the cumulative probability plots for those resistance states,
286 together with lognormal fits modelling the nature of RTN. One of the things that the figure
287 reveals is that higher resistance states suffer from higher degree of RTN. Fits for every
288 resistance state, together with occurrence rates (see Supplementary Table SII), were used
289 to disturb the weights of ANNs in order to reproduce the effect of RTN.

290 2. *Inference accuracy*

291 The results combining RTN and line resistance effects for Ta₂O₅ device are shown in
292 Figure 6. From the difference in median accuracy between yellow and blue box plots, we can

293 notice that there is a significant drop in accuracy simply due to mapping of weights onto
 294 conductances. That is not surprising given that only 8 states were available for mapping.
 295 One can also observe that further drop in median accuracy due to non-idealities is not
 296 as severe—it drops to $\sim 94.1\%$. The RTN disturbance magnitude is limited to $<100\%$ in
 297 most cases, which possibly contributes to its smaller effect on accuracy. Additionally, Ta_2O_5
 298 device has much higher resistance (ranging from $25\text{ k}\Omega$ to $200\text{ k}\Omega$), thus line resistance is also
 299 less of a concern. When non-ideal networks are combined into committees of 5, the median
 300 accuracy jumps to $\sim 96.5\%$ —even higher than the software baseline of individual networks.
 301 This reveals additional trend seen in all the simulations performed—the higher the accuracy
 302 of the individual non-ideal memristive networks, the higher the accuracy of the committees
 303 that they are part of.

304 **D. aVMCO RRAM**

305 Further, we consider a third memristor technology—one based on aVCMO materials. We
 306 test the effects of RTN by considering measurements from a single device. Line resistance
 307 effects were simulated by using interconnect resistance and shape of Ta/HfO₂ crossbar array.

308 *1. Random telegraph noise*

309 Figure 7 shows the cumulative probability plots for 8 resistance states of an aVMCO
 310 device suffering from RTN. Like in Ta_2O_5 , we observe that higher resistance states experience
 311 RTN of higher magnitude. However, compared to Ta_2O_5 , the RTN magnitude is much more
 312 predictable. Fits for each of the 8 resistance states, together with occurrence rates (see
 313 Supplementary Table SIII), were used to simulate the effect of RTN in aVMCO-based neural
 314 networks.

315 *2. Inference accuracy*

316 The results combining RTN and line resistance are shown in Figure 8. As with Ta_2O_5 , we
 317 see a large drop due to mapping onto conductances—consequence of very few states available
 318 for mapping. More interestingly, the accuracy of individual memristor-based networks with

319 and without non-idealities is almost identical. That is because the occurrence rate of RTN
320 in aVMCO device is small and there is a much smaller probability of RTN having large
321 magnitude. Additionally, resistance of aVMCO device is even higher than that of Ta₂O₅
322 device—it ranges from 1 MΩ to 7.5 MΩ. Therefore, line resistance has even a smaller effect
323 in a hypothetical array of aVMCO devices. Due to median accuracy of individual non-ideal
324 memristor-based networks being higher ($\sim 94.6\%$), the median accuracy of committees is
325 higher too—in committees of size 5 it increases to $\sim 96.7\%$.

326 III. DISCUSSION

327 The results from the previous section suggest that the method of using committee ma-
328 chines to improve the accuracy of memristive neural networks is technology- and non-ideality-
329 agnostic. CMs can mitigate the effects of faulty devices, D2D variability, RTN and line
330 resistance in combination with each other. Although CM method is slightly less effective
331 with large line resistance (see discussion in the supplementary information), in all cases, we
332 observe that the accuracy of individual non-ideal networks largely determines the accuracy
333 of committees. That is consequential because it means that although committees always
334 increase the accuracy, there is still an incentive to optimise the devices and systems that
335 implement these networks—the higher the accuracy of individual networks, the higher the
336 accuracy of the committees.

337 It is also important to consider whether using larger networks, instead of committees of
338 smaller networks, would yield the same results if the same number of synapses (or mem-
339 ritors) was used in the large network as in the committee of smaller networks. In our
340 previous work we found that the accuracy of networks before disturbance (which we call
341 “starting accuracy”) has a huge effect on the robustness to non-idealities—the larger the
342 starting accuracy, the more robust the networks become [20]. One way to achieve higher
343 starting accuracy is to have larger networks, e.g. if we have a network with one hidden layer,
344 we might increase the number of neurons in that hidden layer, which would likely result in
345 higher accuracy after training and thus higher robustness.

346 Figure 9 shows a comparison of CMs of memristor-based networks disturbed using faulty
347 devices and D2D variability data from Ta/HfO₂ crossbar, when controlled for the total
348 number of memristors that is required to implement them (line resistance was not taken

349 into account due to long time required to simulate it in large networks). We can observe
350 that committees of two networks, each with 25 hidden neurons, (leftmost data point of
351 the orange curve) achieve $\sim 0.9\%$ higher median accuracy than individual networks with
352 50 hidden neurons (second data point from the left in the vermilion curve), despite both
353 requiring almost identical total number of memristors. Committees of two networks, each
354 with 100 hidden neurons, (third data point from the left in the orange curve) achieve $\sim 1.1\%$
355 higher median accuracy than individual networks with 200 hidden neurons (rightmost data
356 point in the vermilion curve), even though both require almost the same total number of
357 memristors. Even larger improvement is gained when committees of four networks, each with
358 50 hidden neurons, (second data point from the left in the blue curve) are used instead—
359 then the accuracy is improved by $\sim 1.5\%$, with almost the exact total number of memristors
360 used.

361 For different non-idealities and even different training schemes of the ANNs, the equiv-
362 alents of Figure 9 might be different, but there are a few common characteristics in all of
363 them. In all cases, for a given total number of memristors used, there is an optimal number
364 of networks that should be used in a committee. Additionally, we observe that the more
365 severe a non-ideality is, the more apparent the effectiveness of committees becomes. Finally,
366 sometimes the committees (for a fixed total number of memristors) might achieve lower
367 accuracy than individual networks but only if the networks that they replace are very small
368 and the non-ideality is not very detrimental. If the networks that are being replaced with
369 committees of smaller networks, are sufficiently large, the committees will achieve higher
370 accuracy. An example of that is shown in Supplementary Figure S7 where aVMCO device
371 is minimally affected by the non-idealities and so the advantage of committees becomes
372 apparent only when replacing larger networks.

373 The reason why committees work in the context of non-ideal implementations and why
374 they work best when they are used to replace large networks might, to some extent, lie in
375 their training. When it comes to training fully connected networks, their accuracy tends to
376 saturate as more parameters are added. Supplementary Figure S4 shows that networks with
377 50 hidden neurons can be trained to achieve significantly higher accuracy than networks with
378 25 hidden neurons. However, networks with 200 hidden neurons achieve only slightly higher
379 accuracy than networks with 100 hidden neurons. This also means that networks with 200
380 hidden neurons will be only slightly more robust to non-idealities than networks with 100

381 hidden neurons. When such networks are affected by non-idealities, their accuracy drops
382 to similar values but the smaller network can work in a committee with other networks,
383 totalling almost the same number of memristors as the large network, but achieving higher
384 accuracy overall. This is the most likely reason why the committees of smaller networks are
385 effective at dealing with non-idealities, especially when replacing large networks.

386 In addition to the accuracy improvements, committees can provide flexibility in mem-
387 ristic implementations of neural networks. Digital implementations of ANNs have very
388 predictable behaviour due to the precision of digital logic. Analogue implementations, on
389 the other hand, can vary greatly even if they use the same weights before the mapping
390 onto conductances—that is a result of the stochastic nature of memristors that implement
391 these ANNs. The parallel and modular nature of committee machines makes memristive
392 systems much more flexible. For example, if the verification accuracy of one of the ANNs in
393 a memristor-based CM deteriorates below acceptable levels, its outputs could be disabled
394 to ensure higher accuracy of the rest of the committee.

395 Importantly, this introduced parallelism comes at almost no extra cost. For a fixed total
396 number of memristors, a committee of smaller networks, compared to a large individual
397 network, would only require a few additional output and bias neurons, and an averaging
398 functionality, which could potentially be implemented in hardware. For example, an ANN
399 with 50 hidden neurons would require 846 neurons in total, while a committee of two ANNs,
400 each with 25 hidden neurons (and thus requiring almost the same total number of memris-
401 tors), would require 857 neurons in total.

402 In summary, our simulations employing experimental data from three different types of
403 memristive devices show that committee machines employing ensemble averaging can be used
404 to mitigate the effects of device- and system-level non-idealities in memristor-based neural
405 networks. EA allows to achieve higher inference accuracy in physically implemented neural
406 networks that suffer from faulty devices, device-to-device variability, random telegraph noise,
407 and even line resistance. This method is a universal way to deal with the most common
408 non-idealities and is straightforward to implement during the fabrication stage. Increased
409 modularity of these memristive neural network systems will increase not only their inference
410 accuracy, but also their robustness and flexibility, even without the need to sacrifice area.
411 Although some level of non-idealities in memristors is unavoidable, CM method allows us
412 to deal with these on the system level and is agnostic to a particular technology or, to some

413 degree, type of the non-ideality.

414 METHODS

415 Experiments

416 Ta/HfO₂ RRAM 1T1R array consists of NMOS transistors fabricated in a commercial
417 fab (feature size of 2 μm) and Pt/HfO₂/Ta devices. The bottom electrode was deposited by
418 evaporation of 20 nm Pt layer on top of a 2 nm tantalum (Ta) adhesive layer; the electrode
419 was patterned by photolithography and a lift-off process. A 5 nm HfO₂ switching layer was
420 deposited by atomic layer deposition using water and tetrakis(dimethylamido)hafnium as
421 precursors at 250 °C. Sputter-deposited Ta of 50 nm thickness followed by 10 nm Pd was
422 used in a lift-off process to serve as the top electrode. The filamentary based Ta₂O₅ device
423 consists of a TiN/4nm stoichiometric Ta₂O₅/20 nm nonstoichiometric TaO_x/10 nm TaN/TiN
424 stack with a cross-sectional area of 75 nm × 75 nm, while the non-filamentary-based aVMCO
425 has a cross-sectional area of 135 nm × 135 nm and is composed of a TiN/8 nm amorphous-
426 Si/8 nm anatase TiO₂/TiN stack. Ta₂O₅ and aVMCO devices were fabricated by imec. The
427 detailed fabrication process parameters can be found in references [11, 28, 29] for Ta/HfO₂,
428 Ta₂O₅ and aVMCO RRAMs respectively.

429 The conductance of Ta/HfO₂ devices was modulated by applying SET pulses (500 μs @
430 2.5 V and gate voltage increasing from 0.6 V to 1.6 V). After each of the 11 cycles, RESET
431 pulses were applied (5 μs @ 0.9 V increasing to 2.2 V and gate voltage of 5 V). The voltage
432 was being increased linearly throughout the 100 pulses. All electrical tests for Ta₂O₅ and
433 aVMCO devices were done with a Keysight B1500A. The RTN data is extracted by switching
434 the device into 8 uniformly distributed resistance levels between 25 kΩ and 200 kΩ, and 8
435 nearly uniformly distributed resistance levels between 1 MΩ and 7.5 MΩ with incremental
436 RESET DC sweeps [30] for Ta₂O₅ and aVMCO respectively. RTN measurement is then
437 carried out at each resistance level at a 0.1 V and 3 V read-out for Ta₂O₅ and aVMCO
438 respectively, with a sampling time of 2 ms/point and 10,000 sampling point per resistance
439 level for an RTN measurement period of 20 s.

440 Simulations

441 In this work, feed-forward ANNs with fully connected layers and continuous weights were
442 trained to recognise handwritten digits using the MNIST data base. All 60,000 MNIST
443 training images were used during the training stage; training set consisted of 50,000 images
444 and verification set consisted of 10,000 images. All 10,000 test images were used to evaluate
445 the inference accuracy of ANNs. Networks used 784 input neurons representing pixel inten-
446 sities of MNIST images of 28×28 pixel size, as well as one bias neuron. 10 output neurons
447 were used; they represented the ANNs' predictions of 10 handwritten digits. Hidden layers
448 used sigmoid activation function, while the output layer used softmax activation function.
449 Weights were optimised by minimising cross-entropy error function using stochastic gradi-
450 ent descent. Learning rate of 0.01 and patience of 25 epochs were used. 25 networks were
451 trained for each architecture explored by initialising them differently. When numerically op-
452 timising ANNs' weightings, optimisation was performed by employing verification set, while
453 the performance was evaluated using the test set. The code was implemented in Python.

454 Weights were mapped onto pairs of memristors' conductances using proportional map-
455 ping scheme—synaptic weights were made proportional to one of the conductances in the
456 pair, while the other was left unelectroformed. The zero weight was interpreted as given—
457 in practise, it would be implemented by not electroforming the device, thus resulting in its
458 negligible conductance. Although aVMCO devices do not have electroforming stage, for con-
459 sistency we assumed that additional insulating circuit elements could be used to implement
460 the zero weight. Negative weights would be implemented by placing certain memristors in
461 dedicated bit lines of the crossbars whose outputs would be subtracted from the outputs at
462 the corresponding bit lines implementing positive weights. Maximum weights after mapping
463 were optimised separately for each set of network architecture and conductance levels; in
464 each case this was done by excluding a certain proportion, p_L , of weights with largest abso-
465 lute values. What p_L values were used for each simulation is summarised in Supplementary
466 Table SI. More details on the mapping procedure can be found in our past work [20].

467 All non-idealities, except for line resistance, were simulated by disturbing the individual
468 conductances of memristor-based ANNs. To investigate line resistance, nodal analysis was
469 employed. By setting up simultaneous linear equations using Ohm's law and Kirchhoff's
470 current law, those were solved in sparse matrix representation using Python's library `scipy`.

471 After simulating memristor non-idealities, committees of different ANNs were composed.
472 Committees used EA, i.e. the outputs of individual networks in a committee were averaged
473 to produce a single output vector. In EA, the output vectors of individual networks can
474 simply be added together (if the weightings of different networks are the same, as we assume
475 in the main text); the label corresponding to the entry with the highest value would be
476 the prediction of the committee. This addition can be performed either in software, or, if
477 the activation function of the last neuronal layer can be implemented physically, it can be
478 performed by adding corresponding currents produced by the circuitry of this activation
479 function.

480 In the simulations, neural networks that go into a committee were chosen randomly.
481 This was done to reflect the most convenient strategy when manufacturing such systems—
482 because one does not need to selectively choose the networks, manufactured crossbars can be
483 easily programmed without the need to replace them if they perform poorly when working
484 individually (unless their effect is so detrimental that they have to be ignored which can
485 be made possible with this technique). Besides, devices might change over time, thus these
486 simulations, which show what happens when one does not selectively choose the networks,
487 are valuable to investigate conditions where it is not possible to replace the networks.

488 In the simulations, 25 base networks were used (each having different set of weights) for
489 each of the architectures. Then all of their weights were mapped onto pairs of conductances
490 using HRS/LRS values extracted from experiments. Finally, to reflect the effect of each of
491 the non-idealities, all networks were disturbed multiple times. In each disturbance iteration,
492 multiple combinations of networks were chosen and their performance as a committee of
493 certain size was evaluated. In total, for most simulations, 10,000 data points were recorded
494 for a committee of every size—these data captured the variations of base networks, their
495 combinations and different disturbance iterations. Only simulations involving line resistance
496 or numerical optimisation of weights had fewer data points for some committee sizes (due
497 to long simulation times).

498 **DATA AVAILABILITY**

499 The data that support the findings of this study are available from the corresponding
500 author upon reasonable request.

501 **AUTHOR CONTRIBUTIONS**

502 A.M. and D.J. conceived the idea and designed the study. A.M., P.F. and Z.C. per-
503 formed the experimental measurements. D.J. performed the simulations and analysed the
504 experimental and simulation results. C.L. and Q.X. provided the experimental data of the
505 programming of a Ta/HfO₂ 1T1R RRAM array. A.M., W.D.Z. and A.J.K. supervised the
506 research. D.J. wrote the initial manuscript. All authors contributed to the discussions of
507 the results and improved the text.

508 **COMPETING INTERESTS STATEMENT**

509 The authors declare that the research was conducted in the absence of any commercial
510 or financial relationships that could be construed as a potential conflict of interest.

511 **FUNDING**

512 A.M. acknowledges funding from the Royal Academy of Engineering under the Re-
513 search Fellowship scheme, A.J.K. acknowledges funding from the Engineering and Physi-
514 cal Sciences Research Council (EP/P013503/1) and the Leverhulme Trust (RPG-2016-135),
515 W.D.Z. acknowledges funding from the Engineering and Physical Sciences Research Council
516 (EP/S000259/1).

-
- 517 [1] E. Strubell, A. Ganesh, and A. McCallum, “Energy and policy considerations for deep learning
518 in NLP,” *arXiv preprint arXiv:1906.02243*, 2019.
- 519 [2] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with
520 pruning, trained quantization and huffman coding,” in *International Conference on Learning
521 Representations*, 2016, San Juan (Puerto Rico), arXiv preprint arXiv:1510.00149.
- 522 [3] C. Li, Z. Wang, M. Rao, D. Belkin, W. Song, H. Jiang, P. Yan, Y. Li, P. Lin, M. Hu, N. Ge,
523 J. P. Strachan, M. Barnell, Q. Wu, R. S. Williams, J. J. Yang, and Q. Xia, “Long short-term
524 memory networks in memristor crossbar arrays,” *Nature Machine Intelligence*, vol. 1, no. 1,
525 pp. 49–57, 2019, doi: [10.1038/s42256-018-0001-4](https://doi.org/10.1038/s42256-018-0001-4).
- 526 [4] Z. Wang, C. Li, W. Song, M. Rao, D. Belkin, Y. Li, P. Yan, H. Jiang, P. Lin, M. Hu, J. P.
527 Strachan, N. Ge, M. Barnell, Q. Wu, A. G. Barto, Q. Qiu, R. S. Williams, Q. Xia, and J. J.
528 Yang, “Reinforcement learning with analogue memristor arrays,” *Nature Electronics*, vol. 2,
529 no. 3, p. 115, 2019, doi: [10.1038/s41928-019-0221-6](https://doi.org/10.1038/s41928-019-0221-6).
- 530 [5] Z. Sun, G. Pedretti, E. Ambrosi, A. Bricalli, W. Wang, and D. Ielmini, “Solving matrix
531 equations in one step with cross-point resistive arrays,” *Proceedings of the National Academy
532 of Sciences*, vol. 116, no. 10, pp. 4123–4128, 2019, doi: [10.1073/pnas.1815682116](https://doi.org/10.1073/pnas.1815682116).
- 533 [6] S. R. Nandakumar, M. Le Gallo, I. Boybat, B. Rajendran, A. Sebastian, and E. Eleftheriou,
534 “A phase-change memory model for neuromorphic computing,” *Journal of Applied Physics*,
535 vol. 124, no. 15, p. 152135, 2018, doi: [10.1063/1.5042408](https://doi.org/10.1063/1.5042408).
- 536 [7] S. Ambrogio, P. Narayanan, H. Tsai, R. M. Shelby, I. Boybat, C. D. Nolfo, S. Sidler, M. Gior-
537 dano, M. Bodini, N. C. P. Farinha, B. Killeen, C. Cheng, Y. Jaoudi, and G. W. Burr,
538 “Equivalent-accuracy accelerated neural-network training using analogue memory,” *Nature*,
539 vol. 558, no. 7708, pp. 60–67, 2018, doi: [10.1038/s41586-018-0180-5](https://doi.org/10.1038/s41586-018-0180-5).
- 540 [8] S. Yu, Z. Li, P. Y. Chen, H. Wu, B. Gao, D. Wang, W. Wu, and H. Qian, “Binary neu-
541 ral network with 16 Mb RRAM macro chip for classification and online training,” in *In-
542 ternational Electron Devices Meeting*. IEEE, 2016, San Francisco (United States), doi:
543 [10.1109/IEDM.2016.7838429](https://doi.org/10.1109/IEDM.2016.7838429).
- 544 [9] J. Woo, K. Moon, J. Song, S. Lee, M. Kwak, J. Park, and H. Hwang, “Improved synap-
545 tic behavior under identical pulses using $\text{AlO}_x/\text{HfO}_2$ bilayer RRAM array for neuromor-

- 546 phic systems,” *IEEE Electron Device Letters*, vol. 37, no. 8, pp. 994–997, 2016, doi:
547 [10.1109/LED.2016.2582859](https://doi.org/10.1109/LED.2016.2582859).
- 548 [10] M. Prezioso, F. Merrih-Bayat, B. D. Hoskins, G. C. Adam, K. K. Likharev, and D. B.
549 Strukov, “Training and operation of an integrated neuromorphic network based on metal-
550 oxide memristors,” *Nature*, vol. 521, no. 7550, pp. 61–64, 2015, doi: [10.1038/nature14441](https://doi.org/10.1038/nature14441).
- 551 [11] C. Li, D. Belkin, Y. Li, P. Yan, M. Hu, N. Ge, H. Jiang, E. Montgomery, P. Lin, Z. Wang,
552 W. Song, J. P. Strachan, M. Barnell, Q. Wu, R. S. Williams, J. J. Yang, and Q. Xia, “Ef-
553 ficient and self-adaptive in-situ learning in multilayer memristor neural networks,” *Nature*
554 *communications*, vol. 9, no. 1, p. 2385, 2018, doi: [10.1038/s41467-018-04484-2](https://doi.org/10.1038/s41467-018-04484-2).
- 555 [12] A. Chen and M. R. Lin, “Variability of resistive switching memories and its impact on cross-
556 bar array performance,” in *2011 International Reliability Physics Symposium*. IEEE, 2011,
557 Monterey (United States), doi: [10.1109/IRPS.2011.5784590](https://doi.org/10.1109/IRPS.2011.5784590).
- 558 [13] J. Kang, Z. Yu, L. Wu, Y. Fang, Z. Wang, Y. Cai, Z. Ji, J. Zhang, R. Wang, and Y. Yang,
559 “Time-dependent variability in RRAM-based analog neuromorphic system for pattern recogni-
560 tion,” in *International Electron Devices Meeting*. IEEE, 2017, San Francisco (United States),
561 doi: [10.1109/IEDM.2017.8268340](https://doi.org/10.1109/IEDM.2017.8268340).
- 562 [14] L. Xia, W. Huangfu, T. Tang, X. Yin, K. Chakrabarty, Y. Xie, Y. Wang, and H. Yang,
563 “Stuck-at fault tolerance in RRAM computing systems,” *IEEE Journal on Emerging and*
564 *Selected Topics in Circuits and Systems*, vol. 8, no. 1, pp. 102–115, 2017, doi: [10.1109/JET-](https://doi.org/10.1109/JETCAS.2017.2776980)
565 [CAS.2017.2776980](https://doi.org/10.1109/JETCAS.2017.2776980).
- 566 [15] C. Li, M. Hu, Y. Li, H. Jiang, N. Ge, E. Montgomery, J. Zhang, W. Song, N. Dávila, C. E.
567 Graves, Z. Li, J. P. Strachan, P. Lin, Z. Wang, M. Barnell, Q. Wu, S. Williams, J. Yang,
568 and Q. Xia, “Analogue signal and image processing with large memristor crossbars,” *Nature*
569 *Electronics*, vol. 1, no. 1, pp. 52–59, 2018, doi: [10.1038/s41928-017-0002-z](https://doi.org/10.1038/s41928-017-0002-z).
- 570 [16] M. Le Gallo, A. Sebastian, R. Mathis, M. Manica, H. Giefers, T. Tuma, C. Bekas, A. Curioni,
571 and E. Eleftheriou, “Mixed-precision in-memory computing,” *Nature Electronics*, vol. 1, no. 4,
572 p. 246, 2018, doi: [10.1038/s41928-018-0054-8](https://doi.org/10.1038/s41928-018-0054-8).
- 573 [17] M. Hu, J. P. Strachan, Z. Li, and S. R. William, “Dot-product engine as computing mem-
574 memory to accelerate machine learning algorithms,” in *17th International Symposium on Quality*
575 *Electronic Design*, 2016, Santa Clara (United States), doi: [10.1109/ISQED.2016.7479230](https://doi.org/10.1109/ISQED.2016.7479230).

- 576 [18] Q. Xia and J. J. Yang, “Memristive crossbar arrays for brain-inspired computing,” *Nature*
577 *materials*, vol. 18, no. 4, p. 309, 2019, doi: [10.1038/s41563-019-0291-x](https://doi.org/10.1038/s41563-019-0291-x).
- 578 [19] Y. LeCun, C. Cortes, and C. J. C. Burges, “The MNIST database of handwritten digits,”
579 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist>
- 580 [20] A. Mehonic, D. Joksas, W. H. Ng, M. Buckwell, and A. J. Kenyon, “Simulation of inference
581 accuracy using realistic RRAM devices,” *Frontiers in Neuroscience*, vol. 13, p. 593, 2019, doi:
582 [10.3389/fnins.2019.00593](https://doi.org/10.3389/fnins.2019.00593).
- 583 [21] M. P. Perrone and L. N. Cooper, “When networks disagree: Ensemble methods for hybrid
584 neural networks,” in *Artificial Neural Networks for Speech and Vision*. Chapman and Hall,
585 1993, pp. 126–142.
- 586 [22] S. Hashem and B. Schmeiser, “Improving model accuracy using optimal linear combinations of
587 trained neural networks,” *IEEE Transactions on Neural Networks*, vol. 6, no. 3, pp. 792–794,
588 1995, doi: [10.1109/72.377990](https://doi.org/10.1109/72.377990).
- 589 [23] B. Li, L. Xia, P. Gu, Y. Wang, and H. Yang, “Merging the interface: Power, area and accuracy
590 co-optimization for RRAM crossbar-based mixed-signal computing system,” in *Proceedings of*
591 *the 52nd Annual Design Automation Conference*, 2015, San Francisco (United States), doi:
592 [10.1145/2744769.2744870](https://doi.org/10.1145/2744769.2744870).
- 593 [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional
594 neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105,
595 Lake Tahoe (United States), doi: [10.1145/3065386](https://doi.org/10.1145/3065386).
- 596 [25] Z. Wang, C. Li, P. Lin, M. Rao, Y. Nie, W. Song, Q. Qiu, Y. Li, P. Yan, J. P. Strachan,
597 N. Ge, N. McDonald, Q. Wu, M. Hu, H. Wu, R. S. Williams, Q. Xia, and J. J. Yang, “In situ
598 training of feed-forward and recurrent convolutional memristor networks,” *Nature Machine*
599 *Intelligence*, vol. 1, no. 9, pp. 434–442, 2019, doi: [10.1038/s42256-019-0089-1](https://doi.org/10.1038/s42256-019-0089-1).
- 600 [26] H. Jiang, L. Han, P. Lin, Z. Wang, M. H. Jang, Q. Wu, M. Barnell, J. J. Yang, H. L. Xin, and
601 Q. Xia, “Sub-10 nm ta channel responsible for superior performance of a HfO₂ memristor,”
602 *Scientific reports*, vol. 6, p. 28525, 2016, doi: [10.1038/srep28525](https://doi.org/10.1038/srep28525).
- 603 [27] G. W. Burr, R. M. Shelby, S. Sidler, C. Di Nolfo, J. Jang, I. Boybat, R. S. Shenoy,
604 P. Narayanan, K. Virwani, E. U. Giacometti, B. N. Kurdi, and H. Hwang, “Experimen-
605 tal demonstration and tolerancing of a large-scale neural network (165 000 synapses) using
606 phase-change memory as the synaptic weight element,” *IEEE Transactions on Electron De-*

- 607 *vices*, vol. 62, no. 11, pp. 3498–3507, 2015, doi: [10.1109/TED.2015.2439635](https://doi.org/10.1109/TED.2015.2439635).
- 608 [28] Y. Fan, L. Zhang, D. Crotti, T. Witters, M. Jurczak, and B. Govoreanu, “Direct evidence
609 of the overshoot suppression in Ta₂O₅-based resistive switching memory with an integrated
610 access resistor,” *IEEE Electron Device Letters*, vol. 36, no. 10, pp. 1027–1029, 2015, doi:
611 [10.1109/LED.2015.2470081](https://doi.org/10.1109/LED.2015.2470081).
- 612 [29] B. Govoreanu, D. Crotti, S. Subhechha, L. Zhang, Y. Chen, S. Clima, V. Paraschiv, H. Hody,
613 C. Adelman, M. Popovici, O. Richard, and M. Jurczak, “A-VMCO: A novel forming-free, self-
614 rectifying, analog memory cell with low-current operation, nonfilamentary switching and excel-
615 lent variability,” in *Symposium on VLSI Technology*, 2015, Kyoto (Japan), doi: [10.1109/VLSIT.2015.7223717](https://doi.org/10.1109/VLSIT.2015.7223717).
- 617 [30] Z. Chai, W. Zhang, P. Freitas, F. Hatem, J. F. Zhang, J. Marsland, B. Govoreanu, L. Goux,
618 G. S. Kar, S. Hall, P. Chalker, and J. Robertson, “The over-reset phenomenon in Ta₂O₅
619 RRAM device investigated by the RTN-based defect probing technique,” *IEEE Electron Device
620 Letters*, vol. 39, no. 7, pp. 955–958, 2018, doi: [10.1109/LED.2018.2833149](https://doi.org/10.1109/LED.2018.2833149).
- 621 [31] C. Sung, S. Lim, H. Kim, T. Kim, K. Moon, J. Song, J.-J. Kim, and H. Hwang, “Effect
622 of conductance linearity and multi-level cell characteristics of TaO_x-based synapse device on
623 pattern recognition accuracy of neuromorphic system,” *Nanotechnology*, vol. 29, no. 11, p.
624 115203, 2018, doi: [10.1088/1361-6528/aaa733](https://doi.org/10.1088/1361-6528/aaa733).
- 625 [32] Y. Fang, Z. Yu, Z. Wang, T. Zhang, Y. Yang, Y. Cai, and R. Huang, “Improvement of HfO_x-
626 based RRAM device variation by inserting ALD TiN buffer layer,” *IEEE Electron Device
627 Letters*, vol. 39, no. 6, pp. 819–822, 2018, doi: [10.1109/LED.2018.2831698](https://doi.org/10.1109/LED.2018.2831698).
- 628 [33] B. Govoreanu, A. Redolfi, L. Zhang, C. Adelman, M. Popovici, S. Clima, H. Hody,
629 V. Paraschiv, I. Radu, A. Franquet, J. C. Liu, J. Swerts, O. Richard, H. Bender, L. Altimime,
630 and M. Jurczak, “Vacancy-modulated conductive oxide resistive RAM (VMCO-RRAM): An
631 area-scalable switching current, self-compliant, highly nonlinear and wide on/off-window re-
632 sistive switching cell,” in *International Electron Devices Meeting*. IEEE, 2013, Washington
633 (United States), doi: [10.1109/IEDM.2013.6724599](https://doi.org/10.1109/IEDM.2013.6724599).
- 634 [34] A. J. Kenyon, M. S. Munde, W. H. Ng, M. Buckwell, D. Joksas, and A. Mehonic, “The
635 interplay between structure and function in redox-based resistance switching,” *Faraday Dis-
636 cussions*, vol. 213, pp. 151–163, 2019, doi: [10.1039/C8FD00118A](https://doi.org/10.1039/C8FD00118A).

- 637 [35] W. Wu, H. Wu, B. Gao, P. Yao, X. Zhang, X. Peng, S. Yu, and H. Qian, “A methodology
638 to improve linearity of analog RRAM for neuromorphic computing,” in *Symposium on VLSI*
639 *Technology*. IEEE, 2018, Honolulu (United States), doi: [10.1109/VLSIT.2018.8510690](https://doi.org/10.1109/VLSIT.2018.8510690).
- 640 [36] Z. Chai, P. Freitas, W. Zhang, F. Hatem, J. F. Zhang, J. Marsland, B. Govoreanu, L. Goux,
641 and G. S. Kar, “Impact of RTN on pattern recognition accuracy of RRAM-based synaptic
642 neural network,” *IEEE Electron Device Letters*, vol. 39, no. 11, pp. 1652–1655, 2018, doi:
643 [10.1109/LED.2018.2869072](https://doi.org/10.1109/LED.2018.2869072).

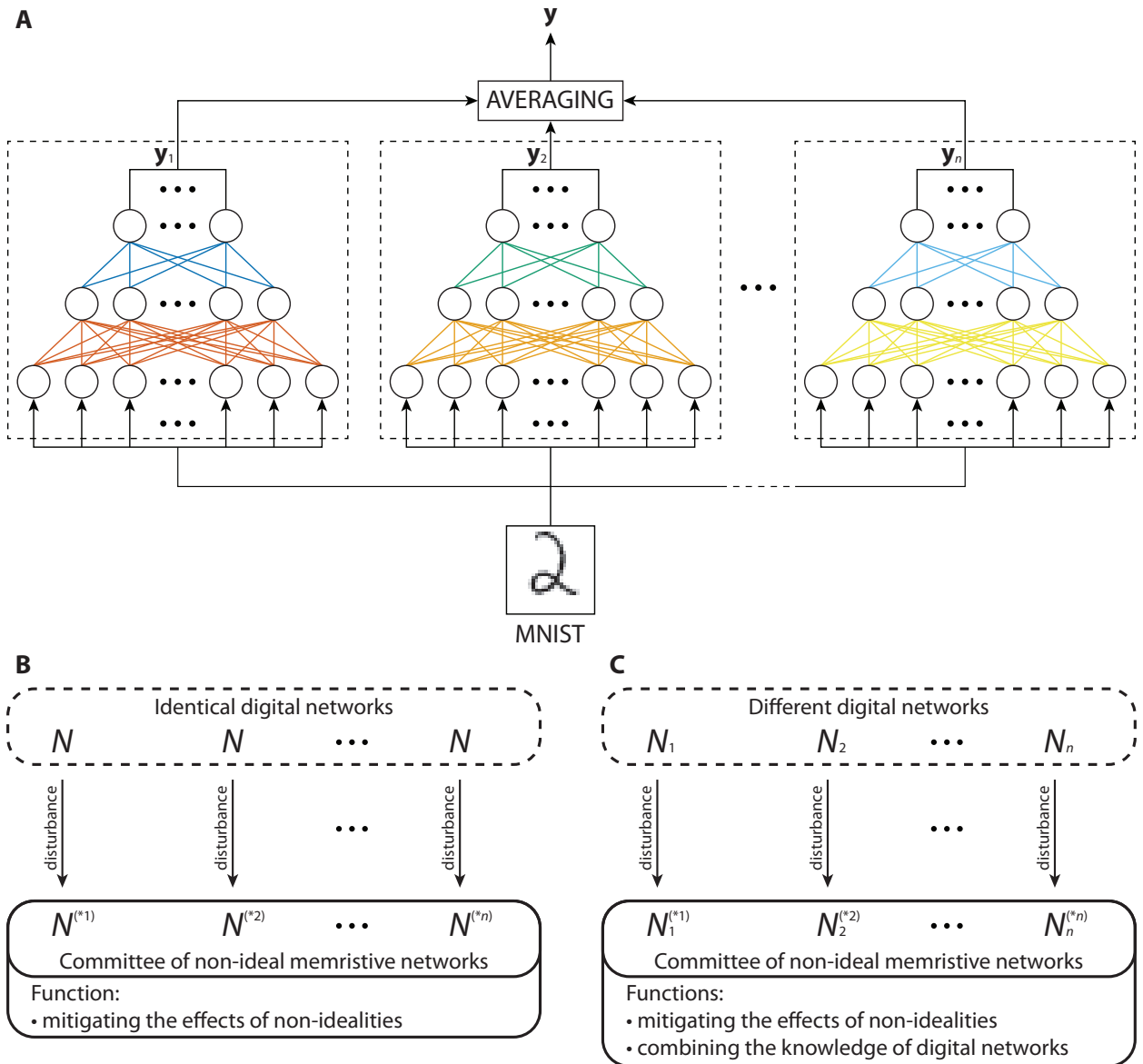


Figure 1. Using multiple neural networks to improve inference accuracy. **A)** The principle of EA. **B)** Using identical digital networks when implementing committees of memristive neural networks only helps to deal with the damage to the networks caused by the non-idealities. **C)** Using different digital networks when implementing committees of memristive neural networks both helps to deal with the damage to the networks caused by the non-idealities and allows to combine the knowledge about the data set acquired by individual digital networks.

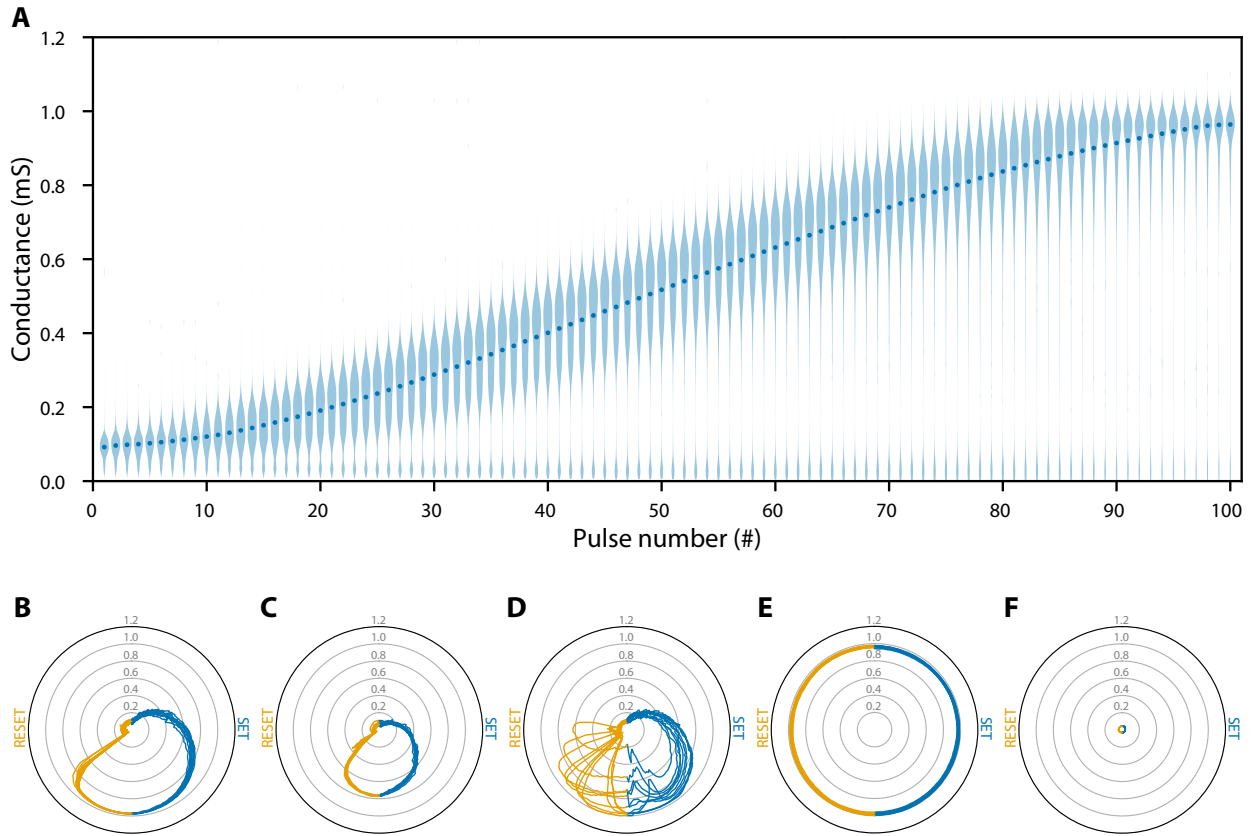


Figure 2. Experimental data of Ta/HfO₂ RRAM crossbar array of shape 128 × 64. **A)** Modulation of devices' conductance over 11 SET cycles, each consisting of a 100 potentiating pulses. Violin plots of gradual conductance changes are shown for all Ta/HfO₂ devices, with dots representing median conductance after a certain number of pulses. 100 points were used for Gaussian kernel density estimation. All violin plots have their maximum widths normalised. **B-F)** Examples of devices with their conductance (in mS) **B)** spanning the full range, **C)** spanning part of the full range, **D)** exhibiting cycle-to-cycle variability, **E)** stuck at high values, **F)** stuck at low values. These diagrams show conductance of five devices from Ta/HfO₂ crossbar array over 11 SET and RESET cycles. The radial component represents the conductance, while the angular component represents the number of applied pulses. The first SET cycle starts at the top of each of the diagrams. The conductance (in blue) over 100 SET pulses is displayed in a clockwise fashion across the right half of each of the diagrams. Following that, conductance (in orange) over 100 RESET pulses (starting at the bottom) is displayed across the left half of each of the diagrams, after which the next cycle is displayed. Cartesian version of these plots is shown in Supplementary Figure S9.

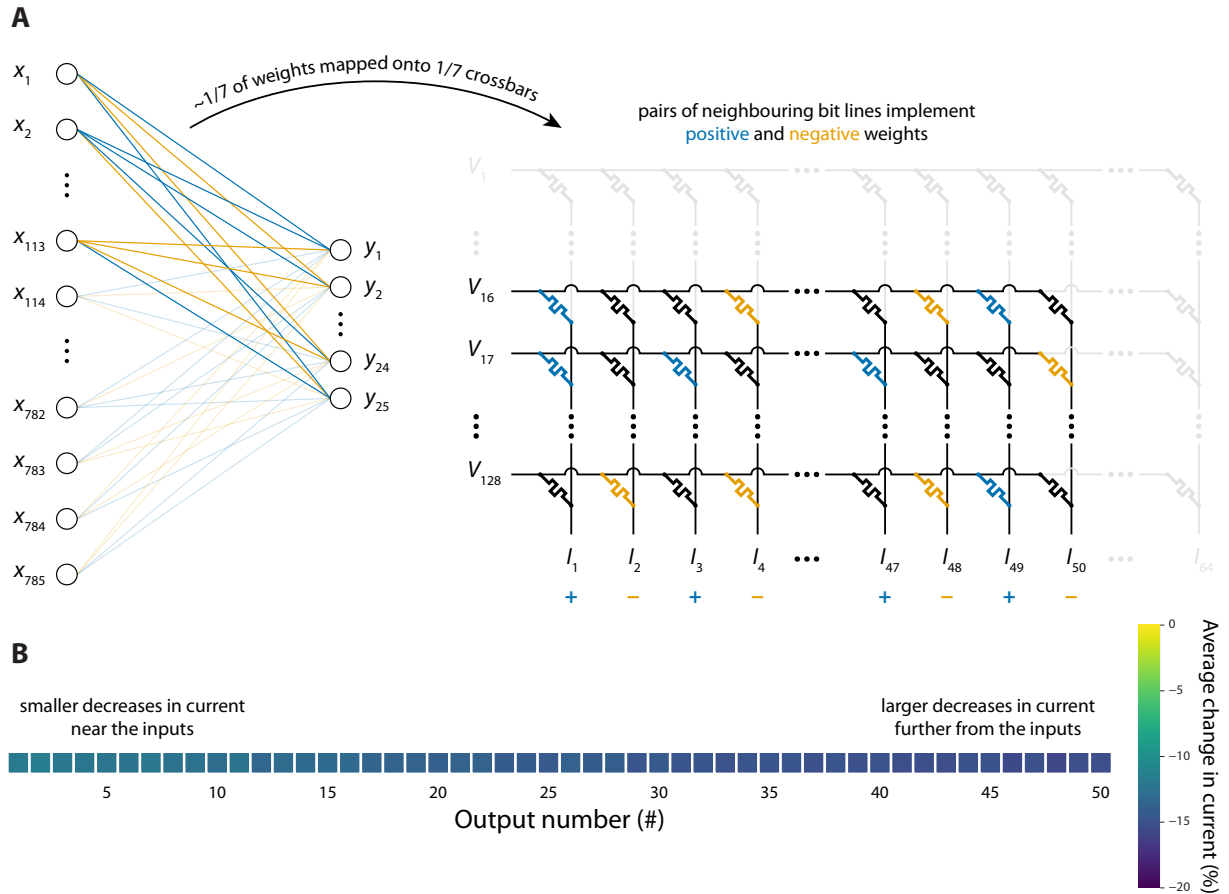


Figure 3. Theoretical implementation of a synaptic layer of shape 785×25 using crossbars of shape 128×64 . **A)** Mapping the first subset of weights onto one of the seven crossbars used to implement the whole synaptic layer. Positive weights and negative weights are mapped onto memristors in different bit lines. **B)** Heatmap of average changes in output currents due to line resistance (in all seven Ta/HfO₂ crossbars). For this particular simulation, it was assumed that Ta/HfO₂ devices can be programmed perfectly.

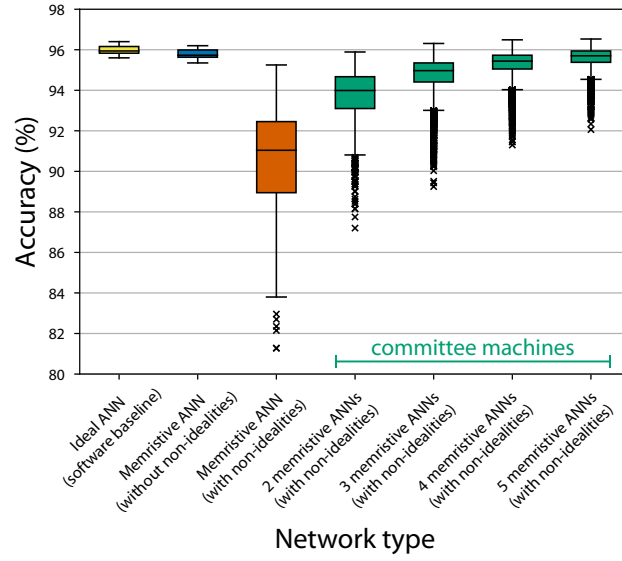


Figure 4. Accuracy achieved by individual networks and their committees when faulty devices, D2D variability data and line resistance of Ta/HfO₂ crossbar are taken into account. The maximum whisker length is set to $1.5 \times \text{IQR}$.

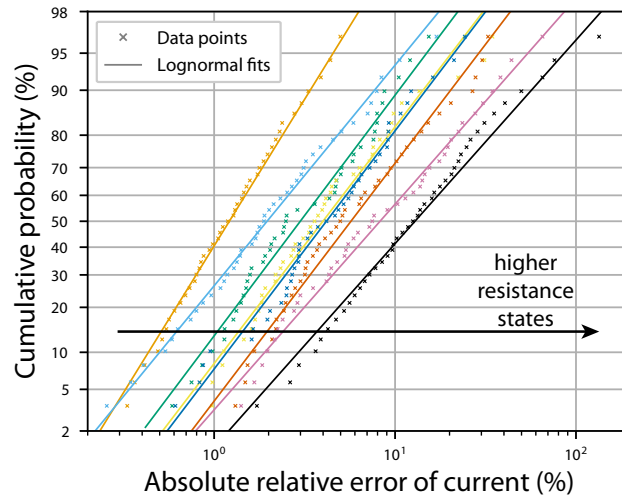


Figure 5. Cumulative probability plots of RTN-induced relative current deviations for all 8 resistance states of a Ta₂O₅ RRAM device. Lognormal fits are shown for each resistance state.

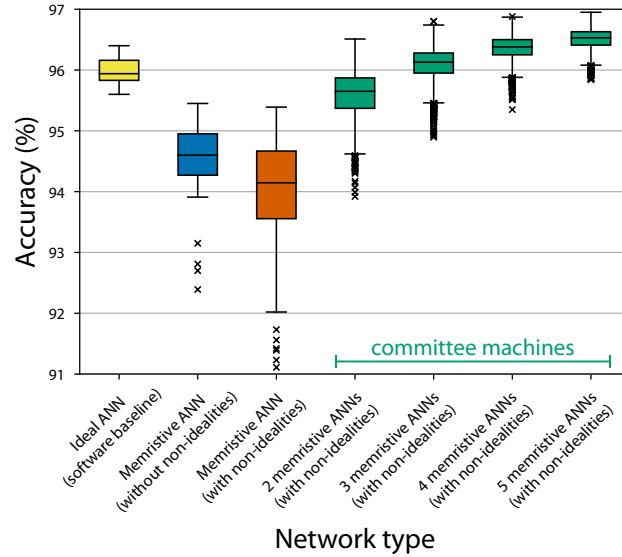


Figure 6. Accuracy achieved by individual networks and their committees when RTN data of a Ta_2O_5 device are taken into account. Additionally, interconnect resistance of $0.35\ \Omega$ and $0.32\ \Omega$ in the word and bit lines, respectively, (from Ta/HfO_2 array) was used to include line resistance effects. The maximum whisker length is set to $1.5 \times \text{IQR}$.

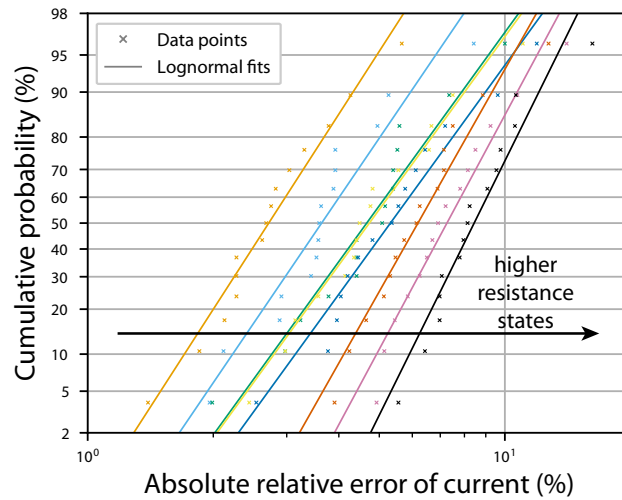


Figure 7. Cumulative probability plots of RTN-induced relative current deviations for all 8 resistance states of aVMCO RRAM device. Lognormal fits are shown for each resistance state.

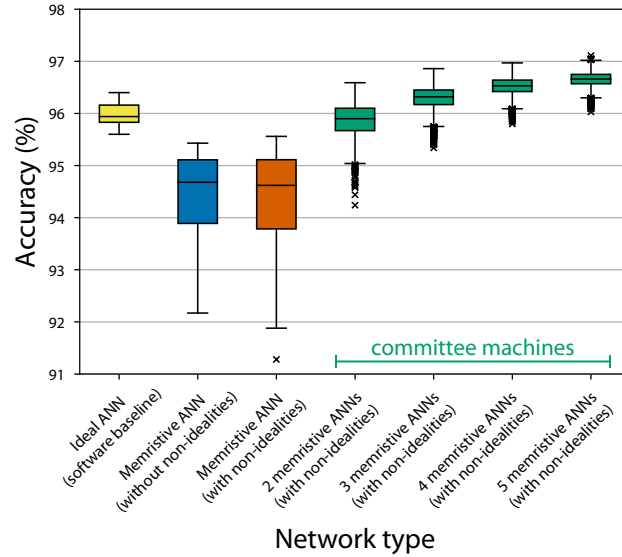


Figure 8. Accuracy achieved by individual networks and their committees when RTN data of an aVMCO device are taken into account. Additionally, interconnect resistance of $0.35\ \Omega$ and $0.32\ \Omega$ in the word and bit lines, respectively, (from Ta/HfO₂ array) was used to include line resistance effects. The maximum whisker length is set to $1.5 \times \text{IQR}$.

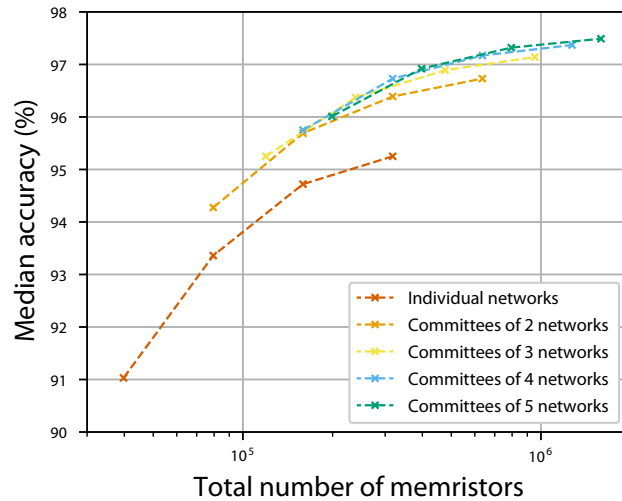


Figure 9. Median accuracy achieved by individual one-hidden-layer memristor-based networks and their committees, when controlled for total number of memristors required. The networks contained 25, 50, 100 or 200 hidden neurons and were disturbed using faulty devices and D2D variability data from Ta/HfO₂ crossbar.

First author (year)	Non-ideality	Device type	Proposed solution
C. Sung (2018) [31]	Current/voltage non-linearity	TaO _x RRAM	Hot-forming step is adopted
C. Li (2018) [15]	Current/voltage non-linearity	Ta/HfO ₂ RRAM	1T1R architecture is adopted
Y. Fang (2018) [32]	Device-to-device variability	HfO _x RRAM	Ultra-thin ALD-TiN buffer layer is introduced
B. Govoreanu (2013) [33]	Device-to-device variability	Al ₂ O ₃ /TiO ₂ (VMCO) RRAM	Non-filamentary RRAM is adopted
A. J. Kenyon (2019) [34]	Device-to-device variability	SiO _x RRAM	The roughness of bottom electrodes is increased
L. Xia (2017) [14]	Faulty devices	-	A modified mapping algorithm and redundancy schemes are used
S. Ambrogio (2018) [7]	Limited dynamic range	PCM	Two pairs of conductance of varying significance for every synaptic weight are used
M. Hu (2016) [17]	Line resistance	-	Advanced mapping algorithms are used to compensate for line resistance effects
W. Wu (2018) [35]	Programming non-linearity	HfO _x RRAM	Electro-thermal modulation layer is deposited on the switching layer
J. Woo (2016) [9]	Programming non-linearity	HfO ₂ RRAM	Bilayer structure is adopted
S. Ambrogio (2018) [7]	Programming non-linearity	PCM	PCM devices are used together with CMOS transistors
Z. Chai (2018) [36]	Random telegraph noise	TiO ₂ /a-Si (aVMCO) RRAM	Non-filamentary RRAM is adopted

Table I. Examples of past efforts at dealing with non-idealities of memristive devices and their systems.

Committee Machines—A Universal Method to Deal with Non-Idealities in Memristor-Based Neural Networks

D. Joksas¹, P. Freitas², Z. Chai², W. H. Ng¹, M. Buckwell¹,
C. Li³, W. D. Zhang², Q. Xia³, A. J. Kenyon¹, and A. Mehonic¹

¹*Department of Electronic and Electrical Engineering,
University College London, London (United Kingdom)*

²*Department of Electronics and Electrical Engineering,
Liverpool John Moores University, Liverpool (United Kingdom)*

³*Department of Electrical and Computer Engineering,
University of Massachusetts Amherst (United States of America)*

Abstract

Artificial neural networks are notoriously power- and time-consuming when implemented on conventional von Neumann computing systems. ~~Recent~~ Consequently, recent years have seen an emergence of research in machine learning hardware that strives to ~~break the bottleneck of von Neumann architecture and optimise the data flow, namely, to~~ bring memory and computing closer together. ~~One of the most often suggested solutions is the physical implementation of~~ A popular approach is to realise artificial neural networks in ~~which hardware by implementing~~ their synaptic weights ~~are realised with memristive devices, such as resistive random access memory~~ using memristive devices. However, various device- and system-level non-idealities usually prevent these physical implementations from achieving high inference accuracy. We suggest applying a well-known concept in computer science—committee ~~machine—in~~ machines—in the context of memristor-based neural networks. Using simulations and experimental data from three different types of memristive devices, we show that committee machines employing ensemble averaging can successfully increase inference accuracy in physically implemented neural networks that suffer from faulty devices, device-to-device variability, random telegraph noise and line resistance. Importantly, we ~~show~~ demonstrate that the accuracy can be improved even without increasing the total number of memristors.

28 I. INTRODUCTION

29 Artificial neural networks (ANNs), with all of their variants, are now the main tools in
30 machine learning tasks, such as classification. The vast amounts of data being constantly
31 produced have enabled successful training and operation of ANNs. However, to achieve
32 high inference accuracy, it is usually necessary for neural networks to have a large number of
33 parameters. This results in both training [1] and inference [2] stages being time- and power-
34 consuming. This is largely caused by the need to transfer data from memory to computing
35 units—physical separation of memory and computing is the essence of any von Neumann
36 system.

37 One of the most promising solutions to these problems is the paradigm of non-von Neu-
38 mann computing and, specifically, analogue implementations of synapses (weights) in phys-
39 ical ANNs. Because there are many more synapses than there are neurons in ANNs, the
40 matrix-vector multiplications, in which the synaptic weight values are used, are the costli-
41 est operations in these networks, both in terms of power and time. Computing directly in
42 memory would minimise ~~costly~~ data transfers from off-chip memory, thus the most popular
43 approach is using analogue memory devices as proxies for synaptic weights of ANNs (both
44 fully connected and their variants [3, 4]). A common technique is to arrange such devices
45 in a structure, called crossbar array, in which every device (or a pair of devices) is used to
46 represent a single synaptic weight or, more generally, an entry in a matrix [5]. Memristive
47 devices, such as phase-change memories (PCMs) [6, 7] or resistive random-access memories
48 (RRAMs) [8, 9], have been considered as candidates for such tasks. Although here we fo-
49 cus on ex-situ training, such systems have been successfully utilised for in-situ training too
50 [10, 11].

51 In memristive implementations of ANNs, the main concern is that various non-idealities
52 associated with these devices can prevent these systems from achieving high accuracy [12,
53 13]. Examples of non-idealities affecting inference accuracy include, but are not limited
54 to, devices not being able to electroform, devices stuck in one of the resistance states after
55 electroforming, device-to-device (D2D) variability and random telegraph noise (RTN). When
56 training analogue systems in-situ, limited endurance and non-linear resistance modulation
57 too have to be taken into account. To mitigate the effects of these device non-idealities, it is
58 often necessary to modify device structure [9], to use more advanced programming schemes

59 [14] or to use additional circuitry [15] or high-precision processing units [16] in conjunction
60 with memristive elements. On the system level, there is an issue of line resistance which
61 affects the distribution of currents and thus decreases the accuracy. These line resistance
62 effects can be partially compensated for algorithmically [17] or partially mitigated by using
63 multiple smaller crossbar arrays [18]. Examples of past efforts at dealing with these and
64 other non-idealities of memristive devices and systems are listed in Table I; most of these
65 non-idealities are still the main focus of the research in the neuromorphic community.

66 We propose a simple way to mitigate the effects of all types of non-idealities during
67 inference. We suggest combining several non-ideal memristor-based neural networks into
68 committees to achieve better accuracy. The committee machine (CM) method we propose
69 significantly increases the inference accuracy and does not increase the computation time
70 because memristive ANNs in such committees work in parallel.

71 In this work, we firstly explain the simulation setup—what networks were trained,
72 how they were simulated and how they were combined into CMs. After that, follows
73 the experimental part. We investigate three different types of memristor technology—
74 tantalum/hafnium oxide-based (Ta/HfO₂), tantalum oxide-based (Ta₂O₅), and amorphous
75 vacancy modulated conductive oxide-based (aVMCO) devices. By exploring their non-
76 idealities relevant to inference—faulty devices, D2D variability, RTN, and line resistance—
77 we use the experimental data to simulate memristive ANNs working individually and in
78 committees.

79 II. RESULTS

80 A. Simulation setup

81 Fully connected ANNs were trained in software to recognise handwritten digits (using
82 MNIST data base [19]). Architectures with one hidden layer were explored. Unless stated
83 otherwise, the simulations used networks with 25 hidden neurons. However, networks with
84 50, 100 and 200 hidden neurons were additionally employed to evaluate the effectiveness of
85 the proposed method while controlling for the total number of memristors required. Follow-
86 ing training, weights of ANNs were mapped onto pairs of conductances using proportional
87 mapping scheme (see [20]) to simulate memristor-based ANNs. Finally, these memristive

88 networks were disturbed using experimental data to reflect the effect of device- and system-
89 level non-idealities.

90 After simulating physical non-idealities, the networks were combined into CMs that em-
91 ployed ensemble averaging (EA) [21]. The principle of EA is shown in Figure 1A—several
92 networks are combined in parallel and then their outputs are averaged. After that, the
93 prediction is made using the averaged vector—the prediction is the label corresponding to
94 the largest entry in the vector.

95 CM methods are frequently used even with conventional ANNs. Methods, such as EA,
96 often produce better accuracy than that of the best individual network in a committee [22].
97 Although there are other types of CMs besides EA, they often rely on training additional
98 gating networks or boosting networks during the training stage. Using a gating network in
99 this scenario would produce additional problems—to avoid it acting as a performance bottle-
100 neck, it too would have to be implemented on crossbar arrays. Various non-idealities would
101 decrease the effectiveness of this gating network which is responsible for making the deci-
102 sions about the whole committee of ANNs. Likewise, we speculate that boosting of networks
103 would not be feasible in ex-situ training because it requires information about where indi-
104 vidual ANNs perform poorly—this cannot be known precisely until they are implemented
105 physically on crossbar arrays and the non-idealities manifest themselves. To authors' best
106 knowledge, the application of boosting in the context of memristive neural networks seems
107 to have been explored only once before [23]; as expected, it requires training each memristive
108 implementation differently because non-idealities manifest themselves differently in different
109 crossbar arrays.

110 There exist modifications of EA algorithm that could potentially perform better. One
111 example of this is generalized ensemble method (GEM) which, instead of using equal weight-
112 ings for each network during averaging (as in EA), uses a different one for each network [21].
113 These weightings are analytically determined by considering correlation of errors between
114 different networks. But because [21] only considered networks with mean square error loss
115 function (while our networks used cross-entropy loss function), this work does not explore
116 GEM. Instead, we investigated whether it is possible to achieve a better performance by
117 optimising the weightings numerically. This method, like GEM and others previously men-
118 tioned, might be impractical because, firstly, these weightings could be determined only after
119 the ANNs are physically implemented on crossbars, and, secondly, the devices could change

120 throughout their lifetimes thus affecting the optimal weightings.

121 Even with the assumption that the devices would have perfect retention, we found that
122 optimisation of weightings achieves effectively the same performance. Because of these rea-
123 sons, we focus only on EA in the main text, but present our results of optimising weightings
124 in Supplementary Figure [S3S5](#). We stress that we are open to the idea that other CM meth-
125 ods besides EA could be utilised successfully for ex-situ training in the context of memristive
126 ANNs. However, in this work we focus on demonstrating that CMs can be used to improve
127 the accuracy of memristor-based ANNs in general.

128 With EA, we find that even when the memristive ANNs, which go into a committee, all
129 use the same ~~digitally implemented~~ digital weights that are mapped onto crossbar arrays
130 (see Figure 1B), committee of memristor-based networks can still achieve higher accuracy
131 than just a single non-ideal network. Although all networks have the same *digital* weights
132 before mapping, their physical implementations (which we call "disturbances" in Figures 1B,
133 C because they can usually be represented by the modification of individual weights) will
134 be different. For example, in one crossbar array, a certain set of devices will be faulty, while
135 in the other crossbar array, it will be a different set. This will result in different physical
136 implementations having slightly different learned representations of the data set, or, to
137 paraphrase, different networks will be "damaged" differently by the non-idealities. This
138 means that these committees will be able to combine different representations, and thus
139 achieve higher accuracy. However, by definition, such approach would almost certainly not
140 yield a committee accuracy that is higher than the accuracy of a single digitally implemented
141 network.

142 A better approach is to use different digital networks for different physical implementa-
143 tions that go into a committee (see Figure 1C). This approach much more resembles the
144 conventional application of EA in computer science. In the context of memristive crossbar
145 arrays, it would not only help to mitigate the effects of the non-idealities (as in the case
146 of Figure 1B), but would also allow to combine the representations of digital networks that
147 were different even before the mapping stage. Most importantly, this method allows for a
148 committee to achieve higher accuracy which is sometimes even higher than that of individual
149 networks with digitally implemented weights. We thus used this method in this analysis.
150 [An example comparison of these two approaches is presented in Supplementary Figure S8.](#)

151 In this work, any given committee used only one network architecture but each network

152 was initialised differently before training, thus trained networks had different sets of weights.
153 Although it was not explored in this work, combining different network architectures in a
154 committee of memristor-based networks might be advantageous. Furthermore, in this work
155 we focus on fully connected ANNs but CMs could be applied to other variants of neural
156 networks as well. Due to the simplicity of EA, it could, for example, be employed in con-
157 volutional neural networks (CNNs) [24], which are often used for image classification. This
158 might be of interest as CNNs have been successfully implemented using crossbar arrays re-
159 cently [25]. However, crossbar implementations are naturally more suited to fully connected
160 networks, therefore we limit ourselves to this architecture but are open to exploring the
161 effectiveness of EA with memristive CNNs in the future.

162 B. Ta/HfO₂ RRAM

163 With array-level data available, Ta/HfO₂ experiments provide the most complete pic-
164 ture of device- and system-level non-idealities. In this subsection, we present not only the
165 analysis of faulty devices and D2D variability, but also careful consideration of the line resis-
166 tance effects. Ta/HfO₂ memristors do not exhibit apparent RTN and overall have excellent
167 retention properties [26], and thus are perfect candidates for inference application.

168 1. *Faulty devices and device-to-device variability*

169 The most energy-efficient procedure to modulate the conductance of memristors is by
170 the application of voltage pulses. In an ideal scenario, one would apply identical pulses
171 and observe constant increases in conductance with each pulse. This is rarely the case
172 in practise, but, fortunately, this type of behaviour is more relevant for in-situ training
173 where it is necessary to ensure linear adjustment of ANN's weights [27]. In ex-situ training,
174 conductance verification schemes can be used to program the devices precisely. Because the
175 devices would have to be programmed only once, one can spend additional resources to do so
176 accurately by applying SET (potentiation) and RESET (depression) pulses until a desirable
177 conductance state is achieved.

178 Even with this approach, there remain two obstacles—faulty devices and D2D variability.
179 It is observed in most memristor technologies that at least a small fraction of the devices

180 tends to get stuck in a particular conductance state. Additionally, even if not stuck, different
181 devices might behave differently; for example, they might have different conductance ranges.
182 Figure 2A shows conductance changes in Ta/HfO₂ RRAM devices (in a 128 × 64 crossbar
183 array) when they are applied with voltage pulses. We can see from the median values
184 that overall the devices' conductance tends to increase as more SET pulses are applied.
185 However, the wider bottom regions of the violin plots indicate that some devices are stuck
186 around high resistance state (HRS) and cannot set entirely no matter how many voltage
187 pulses are applied. There also exist devices that are stuck in low resistance state (LRS), or
188 simply do not span the full conductance range.

189 Figure 2A combines data from multiple SET cycles for each of the memristors, thus it
190 is important to understand how do these devices behave individually. Figures 2B-F show
191 conductance of 5 (out of 8,192) devices over 11 SET and RESET cycles. In the five dia-
192 grams, the radial component represents the conductance (in mS) and the angular component
193 represents the number of applied pulses. Figure 2B shows an example of preferable (and
194 typical) device behaviour—conductance changes in a continuous fashion and spans a wide
195 range of conductance values, from ~0.1 mS to ~1.0 mS. Although RESET cycles tend to
196 feature abrupt decreases in conductance, one can always repeat a cycle and exploit the more
197 predictable behaviour of SET cycles.

198 When encoding continuous numbers into crossbar devices' conductances, it is often prefer-
199 able to choose a large enough conductance range. Using data from Figure 2A, one could,
200 for example, choose the range between the first and the last median points (from ~0.1 mS
201 to ~1.0 mS). Device, whose behaviour is presented in Figure 2B, could be easily set to any
202 conductance within that range, as we have seen before. On the other hand, device, whose
203 behaviour is presented in Figure 2C, although operating in a predictable fashion, has smaller
204 conductance range. We can see that in all cycles, its conductance does not exceed 0.8 mS.
205 This is an example of D2D variability that can make it difficult to choose optimal operating
206 range and set the conductance of all devices precisely.

207 Device, whose behaviour is presented in Figure 2D, shows high cycle-to-cycle variability.
208 Although that could prove to be a problem in some applications, this specific device might
209 perfectly serve its purpose in ex-situ training of ANNs. We can observe that this device
210 spans the same conductance range as device from Figure 2B, even if in an unpredictable
211 manner. Because all states in the full range are, in theory, achievable, one can cycle the

212 device multiple times until it is set to the required conductance level.

213 Lastly, we have devices whose negative effect is most difficult to mitigate—faulty devices.
214 Figure 2E shows behaviour of a device stuck at high conductance values, while Figure 2F
215 shows behaviour of a device stuck at low conductance values. No matter how many pulses
216 the devices are applied with or how many times they are cycled, they exhibit almost no
217 conductance variation and thus, in most cases, cannot be used to encode information.

218 Knowing that some devices perform like the ones whose behaviour is shown in Fig-
219 ures 2C,E,F, it is important to minimise their negative effect. If the conductance that a
220 device has to be set to is outside that device’s range, it is sensible to set it to the closest
221 achievable conductance. Although there is little that can be done about fully stuck memris-
222 tors, it is possible to optimise the behaviour of devices like the one in Figure 2C that simply
223 have smaller conductance range. For example, if such a device has to be set to 0.9 mS, one
224 would set it to the highest achievable conductance (~ 0.8 mS). In the following simulations
225 involving faulty devices and D2D variability, operating range between the first and the last
226 median points was used, the devices were chosen randomly from the 128×64 crossbar and
227 set to the most desirable states, as described in this paragraph.

228 2. Line resistance

229 The effect of line resistance can be extremely detrimental in many crossbar-based im-
230 plementations of ANNs. That is especially the case if the crossbars used are large and
231 the resistance of the interconnects are large is high (compared to memristors’ resistance).
232 Because in a neural network many of the inputs are non-zero at any given time, a lot of
233 current accumulates in the bit lines which results in significant voltage drops across the
234 interconnects, and thus the current distribution across the crossbar is affected in a major
235 way.

236 ~~Although there are many possible options for how to map synaptic weights onto crossbar~~
237 ~~arrays, the choice can determine the role of line resistance. It is often the case that synaptic~~
238 ~~layers of ANNs are large in size. However, that does not mean that the weights in those~~
239 ~~layers have to be mapped onto crossbars of equivalent shape; not only is that sometimes~~
240 ~~impossible, but it can also amplify the effect of line resistance. For example, if a synaptic~~
241 ~~layer with 785 input neurons (as is the case with the first layer of our ANNs) was mapped~~

242 ~~onto a crossbar with 785 word lines, massive amounts of current would accumulate in the~~
243 ~~bit lines.~~

244 The Ta/HfO₂ crossbar has shape 128 × 64 and so this shape was chosen for all the simula-
245 tions involving line resistance. Even relatively small ANNs of architecture 784(+1):25(+1):10
246 would need $2 \times (785 \times 25 + 26 \times 10) = 39,770$ memristors to be implemented. Even if not
247 all the inputs were used at any given time, it would not be possible to fit all the memristors
248 onto a single crossbar of shape 128 × 64. To overcome this, we decided to simulate multiple
249 crossbars, each of which would implement a subset of the synaptic weights, but, for a given
250 synaptic layer, would all compute in parallel. Because $\lceil 785/128 \rceil = 7$, seven crossbars were
251 used to implement the first synaptic layer; the first ~~six crossbars utilised all 128 crossbar~~
252 ~~utilized bottom 113~~ word lines, while the ~~last one used only the bottom 17~~ ~~other six crossbars~~
253 ~~used bottom 112~~ word lines because ~~$785 - 6 \times 128 = 17$~~ ~~$113 + 6 \times 112 = 785$~~ . The second
254 synaptic layer was implemented using eighth crossbar ~~utilising~~ ~~utilizing~~ its bottom 26 word
255 lines.

256 Figure 3A shows an example of how the first synaptic layer of 784(+1):25(+1):10 neural
257 network could be implemented. Specifically, it shows how the first subset of weights would
258 be implemented using one of the crossbars. Because we use proportional mapping scheme,
259 positive and negative weights would be implemented in different bit lines. In Figure 3A,
260 memristors designated to implement positive weights are coloured in blue, memristors des-
261 ignated to implement negative weights are coloured in orange and unelectroformed memris-
262 tors are coloured in black. Because simulations were constrained by experimental data, ~~the~~
263 ~~rightmost bit lines are some of the devices were left~~ unused and assumed to ~~contain only~~
264 ~~unelectroformed devices~~ ~~be unelectroformed~~. In practise, the crossbars could be manufactured
265 to fit the geometry of the ANNs.

266 In each synaptic layer, the corresponding output currents from each of the crossbars
267 would be added together. Additionally, output currents at the bit lines implementing neg-
268 ative weights would be subtracted from the output currents at the ~~corresponding bit lines~~
269 ~~neighbouring bit lines (to their left)~~ implementing positive weights. For example, in the ex-
270 ample configuration of Figure 3A, output current at the ~~26th~~ ~~2nd~~ bit line would be subtracted
271 from the output current at the 1st bit line, etc.

272 Unfortunately, even when using multiple smaller crossbars, the interconnects can signif-
273 icantly disturb current distribution in the crossbar. Average output current decreases due

274 to line resistance in all seven crossbars of Ta/HfO₂ devices (whose resistance ranges from
275 $\sim 1\text{ k}\Omega$ to $\sim 11\text{ k}\Omega$, and their interconnect resistance is $0.3\ \Omega$, $0.35\ \Omega$ and $0.32\ \Omega$ in the word
276 and bit lines, respectively), are shown in the ~~top heatmap of heatmap in~~ Figure 3B. We can
277 see that the current decreases can range from $\sim 1512\%$ at the outputs nearest to the applied
278 voltages to $\sim 1816\%$ at the outputs in the rightmost bit lines that are used. ~~Such large~~
279 ~~current decreases often result from large input voltages that are applied at the top part of~~
280 ~~the crossbar, far away from the outputs. Such inputs generate large amounts of current that~~
281 ~~flow through large portions of the bit lines and, with voltage drops across interconnects,~~
282 ~~disturb the overall current distribution in a major way.~~

283 ~~In some applications, such as supervised learning, it might be possible to strategically~~
284 ~~map certain inputs to certain word lines, so that the effect of line resistance is minimised.~~
285 ~~We propose intensity-aware reordering of ANN's inputs in which we record the average~~
286 ~~input intensities over training and verification sets, and then map inputs with highest~~
287 ~~average intensities to the word lines closest to the outputs of a crossbar. This makes~~
288 ~~it so that most of the current is generated near the outputs, while the currents in the~~
289 ~~top parts of the bit lines are disturbed minimally. Bottom heatmap in Figure 3B shows~~
290 ~~average current decreases when using such a scheme with an unseen test set—we observe~~
291 ~~significantly smaller decreases. Additionally, to make the influence of positive and negative~~
292 ~~weights (which are affected very differently in the naive mapping of Figure 3A) more equal~~
293 ~~and to increase the variability between different ANNs in a committee, we suggest random~~
294 ~~reordering of inputs and outputs. Both intensity-aware and random reordering were used~~
295 ~~in all the following simulations involving line resistance. The implementation of these~~
296 ~~methods individually and in combination with each other is explained in more detail in the~~
297 ~~supplementary information~~In the supplementary information, we provide a possible strategy
298 of mitigating line resistance effects in supervised learning. This scheme was not employed
299 in the simulations described in the main text because we wanted to find out how well the
300 CM method would deal with noticeable line resistance effects.

301 3. *Inference accuracy*

302 Figure 4 shows the accuracy of individual networks, as well as of their committees; mem-
303 ristive ANNs were simulated by taking into account three non-idealities of Ta/HfO₂ crossbar

304 explored earlier—faulty devices, D2D variability and line resistance. As indicated by the
 305 yellow box plot in Figure 4, individual networks implemented digitally achieve $\sim 95.9\%$ me-
 306 dian accuracy. Networks disturbed to reflect the effect of non-idealities achieve ~ 90.8 91.0%
 307 median accuracy, as indicated by the vermilion box plot. Although that is a substantial
 308 drop in accuracy, we see that as more networks are added to the committee, the more the
 309 accuracy increases. When 5 networks are used in a committee, median accuracy increases
 310 up to ~ 95.8 95.7% , as indicated by the rightmost green box plot.

311 C. Ta₂O₅ RRAM

312 In order to explore the effectiveness of minimising adverse effects of RTN, we use another
 313 memristor technology based on Ta₂O₅. To investigate RTN, measurements from a single
 314 device were considered. To simulate line resistance effects, interconnect resistance from
 315 Ta/HfO₂ was used and the same crossbar shape was assumed.

316 1. Random telegraph noise

317 Memristors often suffer from RTN resulting in a different accuracy at any given instant
 318 in time. Ta₂O₅ device was characterised by measuring the current of 8 resistance states
 319 multiple times. Figure 5 shows the cumulative probability plots for those resistance states,
 320 together with lognormal fits modelling the nature of RTN. One of the things that the figure
 321 reveals is that higher resistance states suffer from higher degree of RTN. Fits for every
 322 resistance state, together with occurrence rates (see Supplementary Table SII), were used
 323 to disturb the weights of ANNs in order to reproduce the effect of RTN.

324 2. Inference accuracy

325 The results combining RTN and line resistance effects for Ta₂O₅ device are shown in
 326 Figure 6. From the difference in median accuracy between yellow and blue box plots, we can
 327 notice that there is a significant drop in accuracy simply due to mapping of weights onto
 328 conductances. That is not surprising given that only 8 states were available for mapping.
 329 One can also observe that further drop in median accuracy due to non-idealities is not as

330 severe—it drops to $\sim 94.294.1\%$. The RTN disturbance magnitude is limited to $<100\%$ in
 331 most cases, which possibly contributes to its smaller effect on accuracy. Additionally, Ta_2O_5
 332 device has much higher resistance (ranging from $25\text{ k}\Omega$ to $200\text{ k}\Omega$), thus line resistance is also
 333 less of a concern. When non-ideal networks are combined into committees of 5, the median
 334 accuracy jumps to $\sim 96.5\%$ —even higher than the software baseline of individual networks.
 335 This reveals additional trend seen in all the simulations performed—the higher the accuracy
 336 of the individual non-ideal memristive networks, the higher the accuracy of the committees
 337 that they are part of.

338 D. aVMCO RRAM

339 Further, we consider a third memristor technology—one based on aVCMO materials. We
 340 test the effects of RTN by considering measurements from a single device. Line resistance
 341 effects were simulated by using interconnect resistance and shape of Ta/HfO₂ crossbar array.

342 1. Random telegraph noise

343 Figure 7 shows the cumulative probability plots for 8 resistance states of an aVMCO
 344 device suffering from RTN. Like in Ta_2O_5 , we observe that higher resistance states experience
 345 RTN of higher magnitude. However, compared to Ta_2O_5 , the RTN magnitude is much more
 346 predictable. Fits for each of the 8 resistance states, together with occurrence rates (see
 347 Supplementary Table SIII), were used to simulate [the](#) effect of RTN in aVMCO-based neural
 348 networks.

349 2. Inference accuracy

350 The results combining RTN and line resistance are shown in Figure 8. As with Ta_2O_5 , we
 351 see a large drop due to mapping onto conductances—consequence of very few states available
 352 for mapping. More interestingly, the accuracy of individual memristor-based networks with
 353 and without non-idealities is almost identical. That is because the occurrence rate of RTN
 354 in aVMCO device is small and there is a much smaller probability of RTN having large
 355 magnitude. Additionally, resistance of aVMCO device is even higher than that of Ta_2O_5

356 device—it ranges from 1 M Ω to 7.5 M Ω . Therefore, line resistance has even a smaller effect
357 in a hypothetical array of aVMCO devices. Due to median accuracy of individual non-ideal
358 memristor-based networks being higher ($\sim 94.794.6\%$), the median accuracy of committees
359 is higher too—in committees of size 5 it increases to $\sim 96.696.7\%$.

360 III. DISCUSSION

361 The results from the previous section suggest that the method of using committee ma-
362 chines to improve the accuracy of memristive neural networks is ~~technology-agnostic~~technology-
363 and non-ideality-agnostic. CMs can mitigate the effects of faulty devices, D2D variability,
364 RTN and line resistance in combination with each other. Although ~~line resistance is more~~
365 ~~difficult to deal with using committees due to the similar way in which all crossbars of~~
366 ~~different networks get affected, using random reordering can increase the effectiveness of~~
367 ~~ensembles of non-ideal memristive networks. In~~ CM method is slightly less effective with
368 large line resistance (see discussion in the supplementary information), in all cases, we
369 observe that the accuracy of individual non-ideal networks largely determines the accuracy
370 of committees. That is consequential because it means that although committees always
371 increase the accuracy, there is still an incentive to optimise the devices and systems that
372 implement these networks—the higher the accuracy of individual networks, the higher the
373 accuracy of the committees.

374 It is also important to consider whether using larger networks, instead of committees
375 of smaller networks, would yield the same results if the same number of synapses (or
376 memristors) was used in the large network as in the committee of smaller networks. In
377 our previous work we found that the accuracy of networks before disturbance (which we call
378 “starting accuracy”) has a huge effect on the robustness to non-idealities—the larger the
379 starting accuracy, the more robust the networks become [20]. One way to achieve higher
380 starting accuracy is to have larger networks, e.g. if we have a network with one hidden layer,
381 we might increase the number of neurons in that hidden layer, which would likely result in
382 higher accuracy after training and thus higher robustness.

383 Figure 9 shows a comparison of CMs of memristor-based networks disturbed using faulty
384 devices and D2D variability data from Ta/HfO₂ crossbar, when controlled for the total
385 number of memristors that is required to implement them (line resistance was not taken

386 into account due to long time required to simulate it in large networks). We can observe
387 that committees of two networks, each with 25 hidden neurons, (leftmost data point of
388 the orange curve) achieve $\sim 0.9\%$ higher median accuracy than individual networks with
389 50 hidden neurons (second data point from the left in the vermilion curve), despite both
390 requiring almost identical total number of memristors. Committees of two networks, each
391 with 100 hidden neurons, (third data point from the left in the orange curve) achieve $\sim 1.1\%$
392 higher median accuracy than individual networks with 200 hidden neurons (rightmost data
393 point in the vermilion curve), even though both require almost the same total number of
394 memristors. Even larger improvement is gained when committees of four networks, each with
395 50 hidden neurons, (second data point from the left in the blue curve) are used instead—
396 then the accuracy is improved by $\sim 1.5\%$, with almost the exact total number of memristors
397 used.

398 For different non-idealities and even different training schemes of the ANNs, the equiv-
399 alents of Figure 9 might be different, but there are a few common characteristics in all of
400 them. In all cases, for a given total number of memristors used, there is an optimal number
401 of networks that should be used in a committee. Additionally, we observe that the more
402 severe a non-ideality is, the more apparent the effectiveness of committees becomes. Finally,
403 sometimes the committees (for a fixed total number of memristors) might achieve lower
404 accuracy than individual networks but only if the networks that they replace are very small
405 and the non-ideality is not very detrimental. If the networks that are being replaced with
406 committees of smaller networks, are sufficiently large, the committees will achieve higher
407 accuracy. An example of that is shown in Supplementary Figure [S5-S7](#) where aVMCO de-
408 vice is minimally affected by the non-idealities and so the advantage of committees becomes
409 apparent only when replacing larger networks.

410 The reason why committees work in the context of non-ideal implementations and why
411 they work best when they are used to replace large networks might, to some extent, lie in
412 their training. When it comes to training fully connected networks, their accuracy tends to
413 saturate as more [weights-parameters](#) are added. Supplementary Figure [S2-S4](#) shows that
414 networks with 50 hidden neurons can be trained to achieve significantly higher accuracy
415 than networks with 25 hidden neurons. However, networks with 200 hidden neurons achieve
416 only slightly higher accuracy than networks with 100 hidden neurons. This also means that
417 networks with 200 hidden neurons will be only slightly more robust to non-idealities than

418 networks with 100 hidden neurons. When such networks are affected by non-idealities, their
419 accuracy drops to similar values but the smaller network can work in a committee with
420 ~~one more network~~other networks, totalling almost the same number of memristors as the
421 large network, but achieving higher accuracy overall. This is the most likely reason why the
422 committees of smaller networks are effective at dealing with non-idealities, especially when
423 replacing large networks.

424 In addition to the accuracy improvements, committees can provide flexibility in mem-
425 ristic implementations of neural networks. Digital implementations of ANNs have very
426 predictable behaviour due to the precision of digital logic. Analogue implementations, on
427 the other hand, can vary greatly even if they use the same weights before the mapping
428 onto conductances—that is a result of the stochastic nature of memristors that implement
429 these ANNs. The parallel and modular nature of committee machines makes memristive
430 systems much more flexible. For example, if the verification accuracy of one of the ANNs in
431 a memristor-based CM deteriorates below acceptable levels, its outputs could be disabled
432 to ensure higher accuracy of the rest of the committee.

433 Importantly, this introduced parallelism comes at almost no extra cost. For a fixed total
434 number of memristors, a committee of smaller networks, compared to a large individual
435 network, would only require a few additional output and bias neurons, and an averaging
436 functionality, which could potentially be implemented in hardware. For example, an ANN
437 with 50 hidden neurons would require 846 neurons in total, while a committee of two ANNs,
438 each with 25 hidden neurons (and thus requiring almost the same total number of memris-
439 tors), would require 857 neurons in total.

440 In summary, our simulations employing experimental data from three different types of
441 memristive devices show that committee machines employing ensemble averaging can be used
442 to mitigate the effects of device- and system-level non-idealities in memristor-based neural
443 networks. EA allows to achieve higher inference accuracy in physically implemented neural
444 networks that suffer from faulty devices, device-to-device variability, random telegraph noise,
445 and even line resistance. This method is a universal way to deal with the most common
446 non-idealities and is straightforward to implement during the fabrication stage. Increased
447 modularity of these memristive neural network systems will increase not only their inference
448 accuracy, but also their robustness and flexibility, even without the need to sacrifice area.
449 Although some level of non-idealities in memristors is unavoidable, CM method allows us

450 to deal with these on the system level and is agnostic to a particular technology or, to some
451 degree, type of the non-ideality.

452 METHODS

453 Experiments

454 Ta/HfO₂ RRAM 1T1R array consists of NMOS transistors fabricated in a commercial
455 fab (feature size of 2 μm) and Pt/HfO₂/Ta devices. The bottom electrode was deposited by
456 evaporation of 20 nm Pt layer on top of a 2 nm tantalum (Ta) adhesive layer; the electrode
457 was patterned by photolithography and a lift-off process. A 5 nm HfO₂ switching layer was
458 deposited by atomic layer deposition using water and tetrakis(dimethylamido)hafnium as
459 precursors at 250 °C. Sputter-deposited Ta of 50 nm thickness followed by 10 nm Pd was
460 used in a lift-off process to serve as the top electrode. The filamentary based Ta₂O₅ device
461 consists of a TiN/4nm stoichiometric Ta₂O₅/20 nm nonstoichiometric TaO_x/10 nm TaN/TiN
462 stack with a cross-sectional area of 75 nm × 75 nm, while the non-filamentary-based aVMCO
463 has a cross-sectional area of 135 nm × 135 nm and is composed of a TiN/8 nm amorphous-
464 Si/8 nm anatase TiO₂/TiN stack. Ta₂O₅ and aVMCO devices were fabricated by imec. The
465 detailed fabrication process parameters can be found in references [11, 28, 29] for Ta/HfO₂,
466 Ta₂O₅ and aVMCO RRAMs respectively.

467 The conductance of Ta/HfO₂ devices was modulated by applying SET pulses (500 μs @
468 2.5 V and gate voltage increasing from 0.6 V to 1.6 V). After each of the 11 cycles, RESET
469 pulses were applied (5 μs @ 0.9 V increasing to 2.2 V and gate voltage of 5 V). The voltage
470 was being increased linearly throughout the 100 pulses. All electrical tests for Ta₂O₅ and
471 aVMCO devices were done with a Keysight B1500A. The RTN data is extracted by switching
472 the device into 8 uniformly distributed resistance levels between 25 kΩ and 200 kΩ, and 8
473 nearly uniformly distributed resistance levels between 1 MΩ and 7.5 MΩ with incremental
474 RESET DC sweeps [30] for Ta₂O₅ and aVMCO respectively. RTN measurement is then
475 carried out at each resistance level at a 0.1 V and 3 V read-out for Ta₂O₅ and aVMCO
476 respectively, with a sampling time of 2 ms/point and 10,000 sampling point per resistance
477 level for an RTN measurement period of 20 s.

478 **Simulations**

479 In this work, feed-forward ANNs with fully connected layers and continuous weights were
480 trained to recognise handwritten digits using the MNIST data base. All 60,000 MNIST
481 training images were used during the training stage; training set consisted of 50,000 images
482 and verification set consisted of 10,000 images. All 10,000 test images were used to evaluate
483 the inference accuracy of ANNs. Networks used 784 input neurons representing pixel inten-
484 sities of MNIST images of 28×28 pixel size, as well as one bias neuron. 10 output neurons
485 were used; they represented the ANNs' predictions of 10 handwritten digits. Hidden [layer](#)
486 [layers](#) used sigmoid activation function, while the output layer used softmax activation func-
487 tion. Weights were optimised by minimising cross-entropy error function using stochastic
488 gradient descent. Learning rate of 0.01 and patience of 25 epochs were used. 25 networks
489 were trained for each architecture explored by initialising them differently. When numer-
490 ically optimising ANNs' weightings, optimisation was performed by employing verification
491 set, while the performance was evaluated using the test set. The code was implemented in
492 Python.

493 Weights were mapped onto pairs of memristors' conductances using proportional map-
494 ping scheme—synaptic weights were made proportional to one of the conductances in the
495 pair, while the other was left unelectroformed. The zero weight was interpreted as given—
496 in practise, it would be implemented by not electroforming the device, thus resulting in its
497 negligible conductance. Although aVMCO devices do not have electroforming stage, for con-
498 sistency we assumed that additional insulating circuit elements could be used to implement
499 the zero weight. Negative weights would be implemented by placing certain memristors in
500 dedicated bit lines of the crossbars whose outputs would be subtracted from the outputs at
501 the corresponding bit lines implementing positive weights. Maximum weights after mapping
502 were optimised separately for each set of network architecture and conductance levels; in
503 each case this was done by excluding a certain proportion, p_L , of weights with largest abso-
504 lute values. What p_L values were used for each simulation is summarised in Supplementary
505 Table SI. More details on the mapping procedure can be found in our past work [20].

506 All non-idealities, except for line resistance, were simulated by disturbing the individual
507 conductances of memristor-based ANNs. To investigate line resistance, [loop-nodal](#) analysis
508 was employed. By setting up simultaneous linear equations using [Ohm's law and Kirch-](#)

509 hoff's current ~~and voltage laws~~law, those were solved in sparse matrix representation using
510 Python's library `scipy`.

511 After simulating memristor non-idealities, committees of different ANNs were composed.
512 Committees used EA, i.e. the outputs of individual networks in a committee were averaged
513 to produce a single output vector. In EA, the output vectors of individual networks can
514 simply be added together (if the weightings of different networks are the same, as we assume
515 in the main text); the label corresponding to the entry with the highest value would be
516 the prediction of the committee. This addition can be performed either in software, or, if
517 the activation function of the last neuronal layer can be implemented physically, it can be
518 performed by adding corresponding currents produced by the circuitry of this activation
519 function.

520 In the simulations, neural networks that go into a committee were chosen randomly.
521 This was done to reflect the most convenient strategy when manufacturing such systems—
522 because one does not need to selectively choose the networks, manufactured crossbars can be
523 easily programmed without the need to replace them if they perform poorly when working
524 individually (unless their effect is so detrimental that they have to be ignored which can
525 be made possible with this technique). Besides, devices might change over time, thus these
526 simulations, which show what happens when one does not selectively choose the networks,
527 are valuable to investigate conditions where it is not possible to replace the networks.

528 In the simulations, 25 base networks were used (each having different set of weights) for
529 each of the architectures. Then all of their weights were mapped onto pairs of conductances
530 using HRS/LRS values extracted from experiments. Finally, to reflect the effect of each of
531 the non-idealities, all networks were disturbed multiple times. In each disturbance iteration,
532 multiple combinations of networks were chosen and their performance as a committee of
533 certain size was evaluated. In total, for ~~each simulation (except numerically optimised~~
534 ~~committees which used 1,000 points)~~most simulations, 10,000 data points were recorded
535 for a committee of every size—these data captured the variations of base networks, their
536 combinations and different disturbance iterations. Only simulations involving line resistance
537 or numerical optimisation of weights had fewer data points for some committee sizes (due
538 to long simulation times).

539 **DATA AVAILABILITY**

540 ~~All data generated or analysed during~~ The data that support the findings of this study
541 are ~~included in this published article (and its supplementary information file)~~ available from
542 the corresponding author upon reasonable request.

543 **AUTHOR CONTRIBUTIONS**

544 A.M. and D.J. conceived the idea and designed the study. A.M., P.F. and Z.C. per-
545 formed the experimental measurements. D.J. performed the simulations and analysed the
546 experimental and simulation results. C.L. and Q.X. provided the experimental data of the
547 programming of a Ta/HfO₂ 1T1R RRAM array. A.M., W.D.Z. and A.J.K. supervised the
548 research. D.J. wrote the initial manuscript. All authors contributed to the discussions of
549 the results and improved the text.

550 **COMPETING INTERESTS STATEMENT**

551 The authors declare that the research was conducted in the absence of any commercial
552 or financial relationships that could be construed as a potential conflict of interest.

553 **FUNDING**

554 A.M. acknowledges funding from the Royal Academy of Engineering under the Re-
555 search Fellowship scheme, A.J.K. acknowledges funding from the Engineering and Physi-
556 cal Sciences Research Council (EP/P013503/1) and the Leverhulme Trust (RPG-2016-135),
557 W.D.Z. acknowledges funding from the Engineering and Physical Sciences Research Council
558 (EP/S000259/1).

-
- 559 [1] E. Strubell, A. Ganesh, and A. McCallum, “Energy and policy considerations for deep learning
560 in NLP,” *arXiv preprint arXiv:1906.02243*, 2019.
- 561 [2] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with
562 pruning, trained quantization and huffman coding,” in *International Conference on Learning
563 Representations*, 2016, San Juan (Puerto Rico), arXiv preprint arXiv:1510.00149.
- 564 [3] C. Li, Z. Wang, M. Rao, D. Belkin, W. Song, H. Jiang, P. Yan, Y. Li, P. Lin, M. Hu, N. Ge,
565 J. P. Strachan, M. Barnell, Q. Wu, R. S. Williams, J. J. Yang, and Q. Xia, “Long short-term
566 memory networks in memristor crossbar arrays,” *Nature Machine Intelligence*, vol. 1, no. 1,
567 pp. 49–57, 2019, doi: [10.1038/s42256-018-0001-4](https://doi.org/10.1038/s42256-018-0001-4).
- 568 [4] Z. Wang, C. Li, W. Song, M. Rao, D. Belkin, Y. Li, P. Yan, H. Jiang, P. Lin, M. Hu, J. P.
569 Strachan, N. Ge, M. Barnell, Q. Wu, A. G. Barto, Q. Qiu, R. S. Williams, Q. Xia, and J. J.
570 Yang, “Reinforcement learning with analogue memristor arrays,” *Nature Electronics*, vol. 2,
571 no. 3, p. 115, 2019, doi: [10.1038/s41928-019-0221-6](https://doi.org/10.1038/s41928-019-0221-6).
- 572 [5] Z. Sun, G. Pedretti, E. Ambrosi, A. Bricalli, W. Wang, and D. Ielmini, “Solving matrix
573 equations in one step with cross-point resistive arrays,” *Proceedings of the National Academy
574 of Sciences*, vol. 116, no. 10, pp. 4123–4128, 2019, doi: [10.1073/pnas.1815682116](https://doi.org/10.1073/pnas.1815682116).
- 575 [6] S. R. Nandakumar, M. Le Gallo, I. Boybat, B. Rajendran, A. Sebastian, and E. Eleftheriou,
576 “A phase-change memory model for neuromorphic computing,” *Journal of Applied Physics*,
577 vol. 124, no. 15, p. 152135, 2018, doi: [10.1063/1.5042408](https://doi.org/10.1063/1.5042408).
- 578 [7] S. Ambrogio, P. Narayanan, H. Tsai, R. M. Shelby, I. Boybat, C. D. Nolfo, S. Sidler, M. Gior-
579 dano, M. Bodini, N. C. P. Farinha, B. Killeen, C. Cheng, Y. Jaoudi, and G. W. Burr,
580 “Equivalent-accuracy accelerated neural-network training using analogue memory,” *Nature*,
581 vol. 558, no. 7708, pp. 60–67, 2018, doi: [10.1038/s41586-018-0180-5](https://doi.org/10.1038/s41586-018-0180-5).
- 582 [8] S. Yu, Z. Li, P. Y. Chen, H. Wu, B. Gao, D. Wang, W. Wu, and H. Qian, “Binary neu-
583 ral network with 16 Mb RRAM macro chip for classification and online training,” in *In-
584 ternational Electron Devices Meeting*. IEEE, 2016, San Francisco (United States), doi:
585 [10.1109/IEDM.2016.7838429](https://doi.org/10.1109/IEDM.2016.7838429).
- 586 [9] J. Woo, K. Moon, J. Song, S. Lee, M. Kwak, J. Park, and H. Hwang, “Improved synap-
587 tic behavior under identical pulses using $\text{AlO}_x/\text{HfO}_2$ bilayer RRAM array for neuromor-

- 588 phic systems,” *IEEE Electron Device Letters*, vol. 37, no. 8, pp. 994–997, 2016, doi:
589 [10.1109/LED.2016.2582859](https://doi.org/10.1109/LED.2016.2582859).
- [10] M. Prezioso, F. Merrih-Bayat, B. D. Hoskins, G. C. Adam, K. K. Likharev, and D. B.
591 Strukov, “Training and operation of an integrated neuromorphic network based on metal-
592 oxide memristors,” *Nature*, vol. 521, no. 7550, pp. 61–64, 2015, doi: [10.1038/nature14441](https://doi.org/10.1038/nature14441).
- [11] C. Li, D. Belkin, Y. Li, P. Yan, M. Hu, N. Ge, H. Jiang, E. Montgomery, P. Lin, Z. Wang,
594 W. Song, J. P. Strachan, M. Barnell, Q. Wu, R. S. Williams, J. J. Yang, and Q. Xia, “Ef-
595 ficient and self-adaptive in-situ learning in multilayer memristor neural networks,” *Nature*
596 *communications*, vol. 9, no. 1, p. 2385, 2018, doi: [10.1038/s41467-018-04484-2](https://doi.org/10.1038/s41467-018-04484-2).
- [12] A. Chen and M. R. Lin, “Variability of resistive switching memories and its impact on cross-
598 bar array performance,” in *2011 International Reliability Physics Symposium*. IEEE, 2011,
599 Monterey (United States), doi: [10.1109/IRPS.2011.5784590](https://doi.org/10.1109/IRPS.2011.5784590).
- [13] J. Kang, Z. Yu, L. Wu, Y. Fang, Z. Wang, Y. Cai, Z. Ji, J. Zhang, R. Wang, and Y. Yang,
601 “Time-dependent variability in RRAM-based analog neuromorphic system for pattern recogni-
602 tion,” in *International Electron Devices Meeting*. IEEE, 2017, San Francisco (United States),
603 doi: [10.1109/IEDM.2017.8268340](https://doi.org/10.1109/IEDM.2017.8268340).
- [14] L. Xia, W. Huangfu, T. Tang, X. Yin, K. Chakrabarty, Y. Xie, Y. Wang, and H. Yang,
604 “Stuck-at fault tolerance in RRAM computing systems,” *IEEE Journal on Emerging and*
605 *Selected Topics in Circuits and Systems*, vol. 8, no. 1, pp. 102–115, 2017, doi: [10.1109/JET-](https://doi.org/10.1109/JETCAS.2017.2776980)
607 [CAS.2017.2776980](https://doi.org/10.1109/JETCAS.2017.2776980).
- [15] C. Li, M. Hu, Y. Li, H. Jiang, N. Ge, E. Montgomery, J. Zhang, W. Song, N. Dávila, C. E.
608 Graves, Z. Li, J. P. Strachan, P. Lin, Z. Wang, M. Barnell, Q. Wu, S. Williams, J. Yang,
609 and Q. Xia, “Analogue signal and image processing with large memristor crossbars,” *Nature*
610 *Electronics*, vol. 1, no. 1, pp. 52–59, 2018, doi: [10.1038/s41928-017-0002-z](https://doi.org/10.1038/s41928-017-0002-z).
- [16] M. Le Gallo, A. Sebastian, R. Mathis, M. Manica, H. Giefers, T. Tuma, C. Bekas, A. Curioni,
612 and E. Eleftheriou, “Mixed-precision in-memory computing,” *Nature Electronics*, vol. 1, no. 4,
613 p. 246, 2018, doi: [10.1038/s41928-018-0054-8](https://doi.org/10.1038/s41928-018-0054-8).
- [17] M. Hu, J. P. Strachan, Z. Li, and S. R. William, “Dot-product engine as computing mem-
614 ory to accelerate machine learning algorithms,” in *17th International Symposium on Quality*
615 *Electronic Design*, 2016, Santa Clara (United States), doi: [10.1109/ISQED.2016.7479230](https://doi.org/10.1109/ISQED.2016.7479230).

- 618 [18] Q. Xia and J. J. Yang, “Memristive crossbar arrays for brain-inspired computing,” *Nature*
619 *materials*, vol. 18, no. 4, p. 309, 2019, doi: [10.1038/s41563-019-0291-x](https://doi.org/10.1038/s41563-019-0291-x).
- 620 [19] Y. LeCun, C. Cortes, and C. J. C. Burges, “The MNIST database of handwritten digits,”
621 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist>
- 622 [20] A. Mehonic, D. Joksas, W. H. Ng, M. Buckwell, and A. J. Kenyon, “Simulation of inference
623 accuracy using realistic RRAM devices,” *Frontiers in Neuroscience*, vol. 13, p. 593, 2019, doi:
624 [10.3389/fnins.2019.00593](https://doi.org/10.3389/fnins.2019.00593).
- 625 [21] M. P. Perrone and L. N. Cooper, “When networks disagree: Ensemble methods for hybrid
626 neural networks,” in *Artificial Neural Networks for Speech and Vision*. Chapman and Hall,
627 1993, pp. 126–142.
- 628 [22] S. Hashem and B. Schmeiser, “Improving model accuracy using optimal linear combinations of
629 trained neural networks,” *IEEE Transactions on Neural Networks*, vol. 6, no. 3, pp. 792–794,
630 1995, doi: [10.1109/72.377990](https://doi.org/10.1109/72.377990).
- 631 [23] B. Li, L. Xia, P. Gu, Y. Wang, and H. Yang, “Merging the interface: Power, area and accuracy
632 co-optimization for RRAM crossbar-based mixed-signal computing system,” in *Proceedings of*
633 *the 52nd Annual Design Automation Conference*, 2015, San Francisco (United States), doi:
634 [10.1145/2744769.2744870](https://doi.org/10.1145/2744769.2744870).
- 635 [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional
636 neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105,
637 Lake Tahoe (United States), doi: [10.1145/3065386](https://doi.org/10.1145/3065386).
- 638 [25] Z. Wang, C. Li, P. Lin, M. Rao, Y. Nie, W. Song, Q. Qiu, Y. Li, P. Yan, J. P. Strachan,
639 N. Ge, N. McDonald, Q. Wu, M. Hu, H. Wu, R. S. Williams, Q. Xia, and J. J. Yang, “In situ
640 training of feed-forward and recurrent convolutional memristor networks,” *Nature Machine*
641 *Intelligence*, vol. 1, no. 9, pp. 434–442, 2019, doi: [10.1038/s42256-019-0089-1](https://doi.org/10.1038/s42256-019-0089-1).
- 642 [26] H. Jiang, L. Han, P. Lin, Z. Wang, M. H. Jang, Q. Wu, M. Barnell, J. J. Yang, H. L. Xin, and
643 Q. Xia, “Sub-10 nm ta channel responsible for superior performance of a HfO₂ memristor,”
644 *Scientific reports*, vol. 6, p. 28525, 2016, doi: [10.1038/srep28525](https://doi.org/10.1038/srep28525).
- 645 [27] G. W. Burr, R. M. Shelby, S. Sidler, C. Di Nolfo, J. Jang, I. Boybat, R. S. Shenoy,
646 P. Narayanan, K. Virwani, E. U. Giacometti, B. N. Kurdi, and H. Hwang, “Experimen-
647 tal demonstration and tolerancing of a large-scale neural network (165 000 synapses) using
648 phase-change memory as the synaptic weight element,” *IEEE Transactions on Electron De-*

- 649 *vices*, vol. 62, no. 11, pp. 3498–3507, 2015, doi: [10.1109/TED.2015.2439635](https://doi.org/10.1109/TED.2015.2439635).
- 650 [28] Y. Fan, L. Zhang, D. Crotti, T. Witters, M. Jurczak, and B. Govoreanu, “Direct evidence
651 of the overshoot suppression in Ta₂O₅-based resistive switching memory with an integrated
652 access resistor,” *IEEE Electron Device Letters*, vol. 36, no. 10, pp. 1027–1029, 2015, doi:
653 [10.1109/LED.2015.2470081](https://doi.org/10.1109/LED.2015.2470081).
- 654 [29] B. Govoreanu, D. Crotti, S. Subhechha, L. Zhang, Y. Chen, S. Clima, V. Paraschiv, H. Hody,
655 C. Adelman, M. Popovici, O. Richard, and M. Jurczak, “A-VMCO: A novel forming-free, self-
656 rectifying, analog memory cell with low-current operation, nonfilamentary switching and excel-
657 lent variability,” in *Symposium on VLSI Technology*, 2015, Kyoto (Japan), doi: [10.1109/VLSIT.2015.7223717](https://doi.org/10.1109/VLSIT.2015.7223717).
- 658 [30] Z. Chai, W. Zhang, P. Freitas, F. Hatem, J. F. Zhang, J. Marsland, B. Govoreanu, L. Goux,
659 G. S. Kar, S. Hall, P. Chalker, and J. Robertson, “The over-reset phenomenon in Ta₂O₅
660 RRAM device investigated by the RTN-based defect probing technique,” *IEEE Electron Device*
661 *Letters*, vol. 39, no. 7, pp. 955–958, 2018, doi: [10.1109/LED.2018.2833149](https://doi.org/10.1109/LED.2018.2833149).
- 662 [31] C. Sung, S. Lim, H. Kim, T. Kim, K. Moon, J. Song, J.-J. Kim, and H. Hwang, “Effect
663 of conductance linearity and multi-level cell characteristics of TaO_x-based synapse device on
664 pattern recognition accuracy of neuromorphic system,” *Nanotechnology*, vol. 29, no. 11, p.
665 115203, 2018, doi: [10.1088/1361-6528/aaa733](https://doi.org/10.1088/1361-6528/aaa733).
- 666 [32] Y. Fang, Z. Yu, Z. Wang, T. Zhang, Y. Yang, Y. Cai, and R. Huang, “Improvement of HfO_x-
667 based RRAM device variation by inserting ALD TiN buffer layer,” *IEEE Electron Device*
668 *Letters*, vol. 39, no. 6, pp. 819–822, 2018, doi: [10.1109/LED.2018.2831698](https://doi.org/10.1109/LED.2018.2831698).
- 669 [33] B. Govoreanu, A. Redolfi, L. Zhang, C. Adelman, M. Popovici, S. Clima, H. Hody,
670 V. Paraschiv, I. Radu, A. Franquet, J. C. Liu, J. Swerts, O. Richard, H. Bender, L. Altimime,
671 and M. Jurczak, “Vacancy-modulated conductive oxide resistive RAM (VMCO-RRAM): An
672 area-scalable switching current, self-compliant, highly nonlinear and wide on/off-window re-
673 sistive switching cell,” in *International Electron Devices Meeting*. IEEE, 2013, Washington
674 (United States), doi: [10.1109/IEDM.2013.6724599](https://doi.org/10.1109/IEDM.2013.6724599).
- 675 [34] A. J. Kenyon, M. S. Munde, W. H. Ng, M. Buckwell, D. Joksas, and A. Mehonic, “The
676 interplay between structure and function in redox-based resistance switching,” *Faraday Dis-*
677 *cussions*, vol. 213, pp. 151–163, 2019, doi: [10.1039/C8FD00118A](https://doi.org/10.1039/C8FD00118A).
- 678

- 679 [35] W. Wu, H. Wu, B. Gao, P. Yao, X. Zhang, X. Peng, S. Yu, and H. Qian, “A methodology
680 to improve linearity of analog RRAM for neuromorphic computing,” in *Symposium on VLSI
681 Technology*. IEEE, 2018, Honolulu (United States), doi: [10.1109/VLSIT.2018.8510690](https://doi.org/10.1109/VLSIT.2018.8510690).
- 682 [36] Z. Chai, P. Freitas, W. Zhang, F. Hatem, J. F. Zhang, J. Marsland, B. Govoreanu, L. Goux,
683 and G. S. Kar, “Impact of RTN on pattern recognition accuracy of RRAM-based synaptic
684 neural network,” *IEEE Electron Device Letters*, vol. 39, no. 11, pp. 1652–1655, 2018, doi:
685 [10.1109/LED.2018.2869072](https://doi.org/10.1109/LED.2018.2869072).

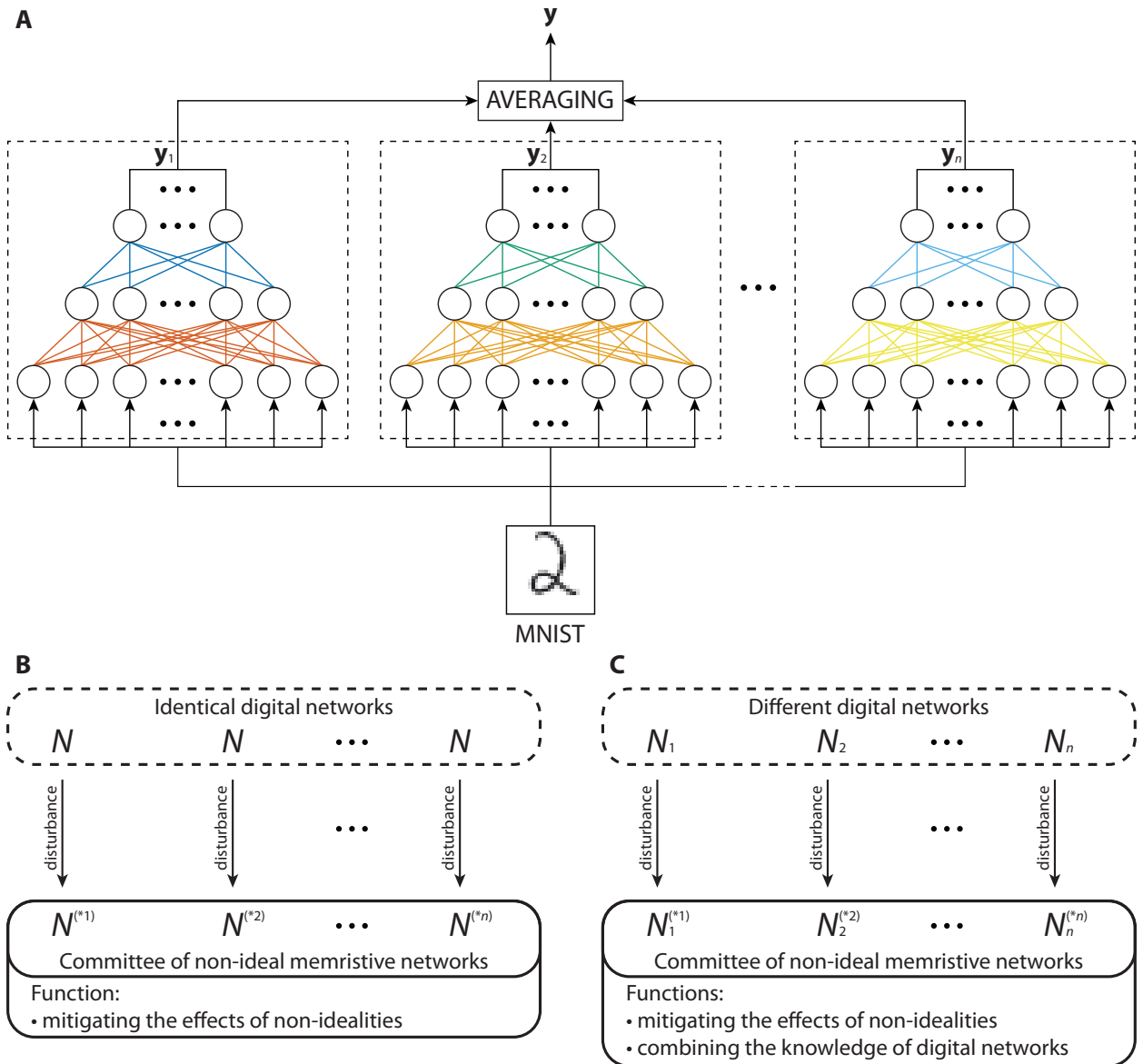


Figure 1. Using multiple neural networks to improve inference accuracy. **A)** The principle of EA. **B)** Using identical digital networks when implementing committees of memristive neural networks only helps to deal with the damage to the networks caused by the non-idealities. **C)** Using different digital networks when implementing committees of memristive neural networks both helps to deal with the damage to the networks caused by the non-idealities and allows to combine the knowledge of individual digital networks about the data set acquired by individual digital networks.

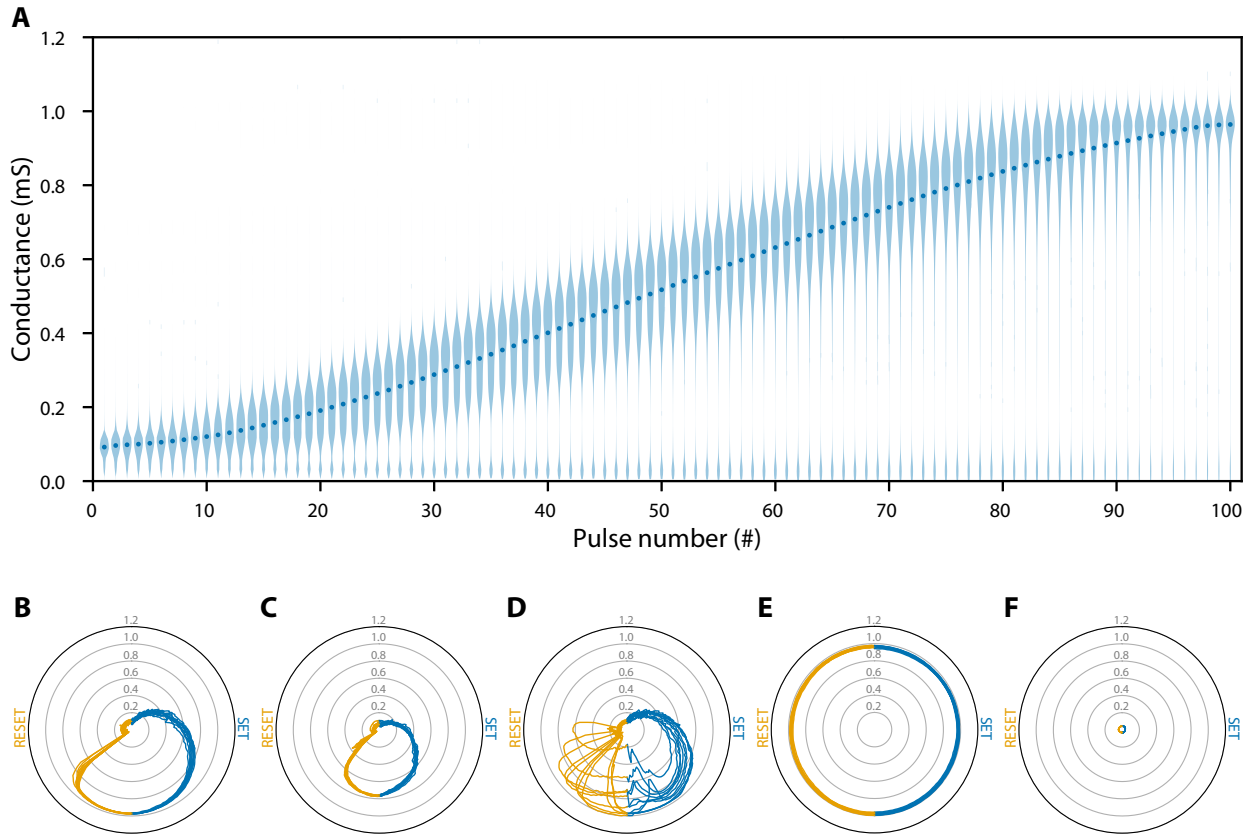


Figure 2. Experimental data of Ta/HfO₂ RRAM crossbar array of shape 128 × 64. **A)** Modulation of devices' conductance over 11 SET cycles, each consisting of a 100 potentiating pulses. Violin plots of gradual conductance changes are shown for all Ta/HfO₂ devices, with dots representing median conductance after a certain number of pulses. 100 points were used for Gaussian kernel density estimation. All violin plots have their maximum widths normalised. **B-F)** Examples of devices with their conductance (in mS) **B)** spanning the full range, **C)** spanning part of the full range, **D)** exhibiting cycle-to-cycle variability, **E)** stuck at high values, **F)** stuck at low values. These diagrams show conductance of five devices from Ta/HfO₂ crossbar array over 11 SET and RESET cycles. The radial component represents the conductance, while the angular component represents the number of applied pulses. The first SET cycle starts at the top of each of the diagrams. The conductance (in blue) over 100 SET pulses is displayed in a clockwise fashion across the right half of each of the diagrams. Following that, conductance (in orange) over 100 RESET pulses (starting at the bottom) is displayed across the left half of each of the diagrams, after which the next cycle is displayed. [Cartesian version of these plots is shown in Supplementary Figure S9.](#)

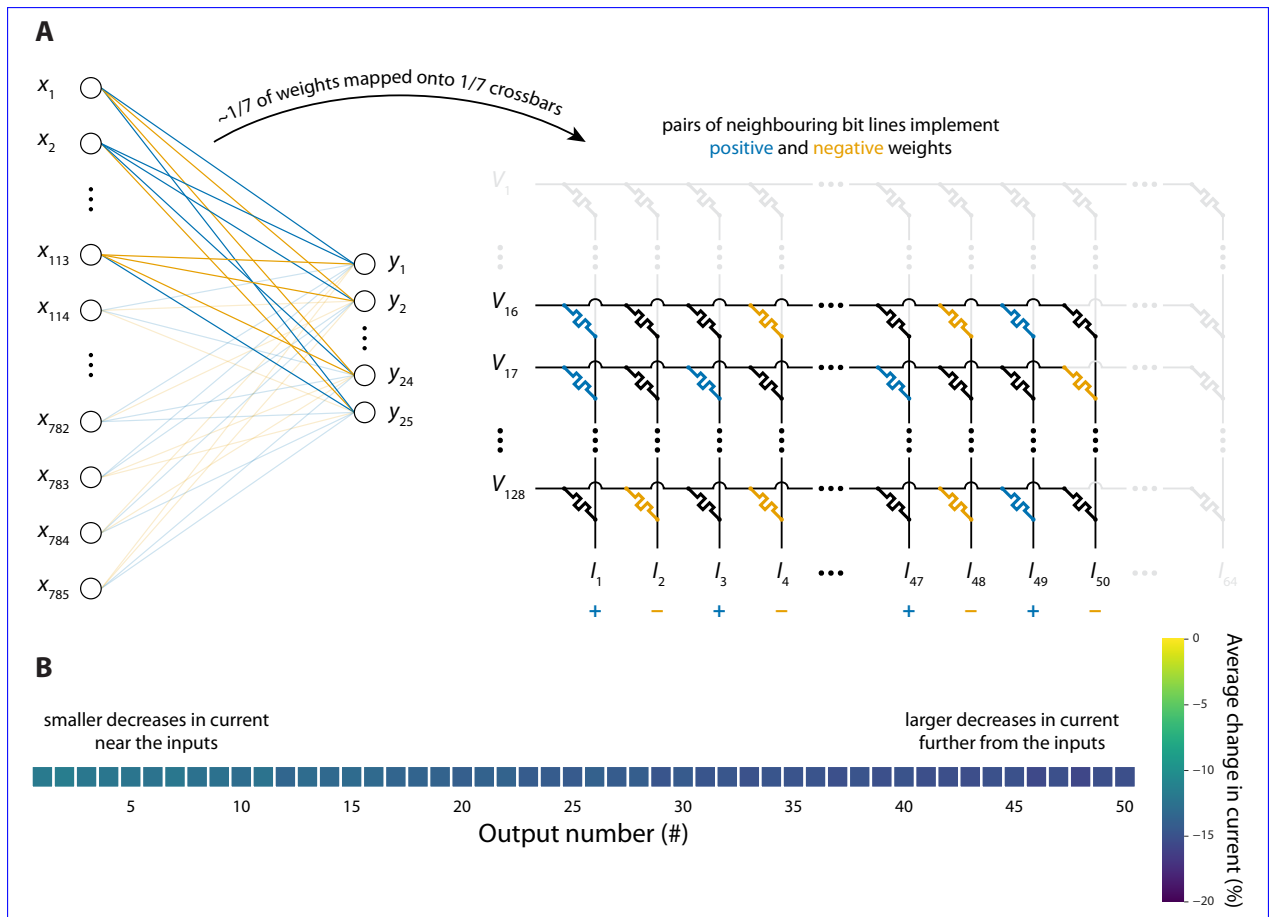


Figure 3. Theoretical implementation of a synaptic layer of shape 785×25 using crossbars of shape 128×64 . **A)** Mapping the first subset of weights onto one of the seven crossbars used to implement the whole synaptic layer. Positive weights and negative weights are mapped onto memristors in different bit lines. **B)** Heatmap of average changes in output currents due to line resistance (in all seven Ta/HfO₂ crossbars) ~~without and with a scheme that maps certain inputs onto certain word lines depending on expected average intensities of those inputs~~. For this particular simulation, it was assumed that Ta/HfO₂ devices can be programmed perfectly.

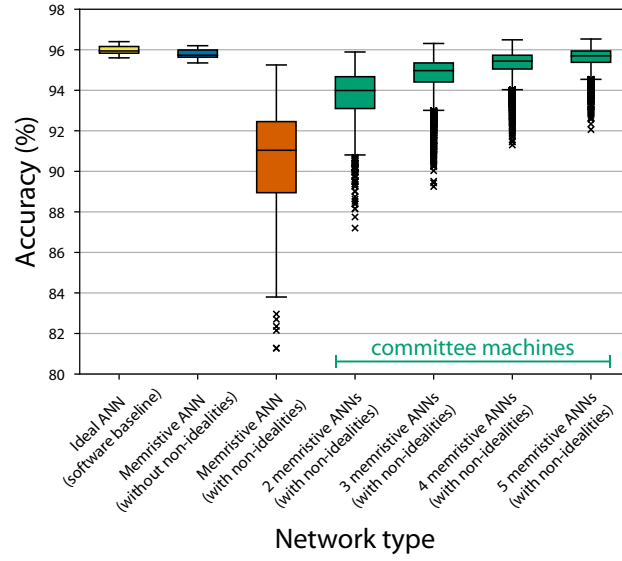


Figure 4. Accuracy achieved by individual networks and their committees when faulty devices, D2D variability data and line resistance of Ta/HfO₂ crossbar are taken into account. The maximum whisker length is set to $1.5 \times \text{IQR}$.

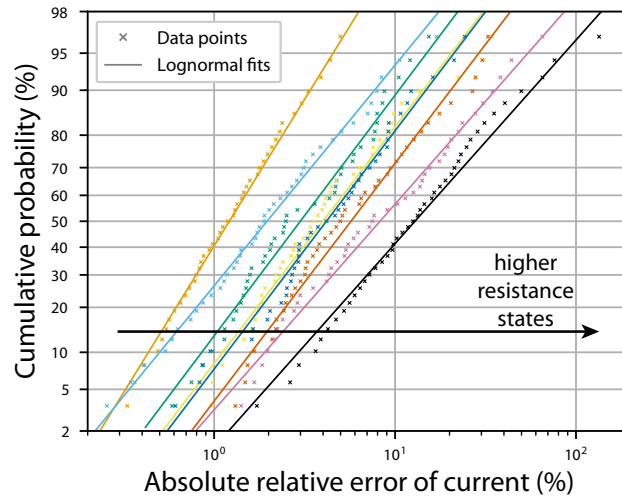


Figure 5. Cumulative probability plots of RTN-induced relative current deviations for all 8 resistance states of a Ta₂O₅ RRAM device. Lognormal fits are shown for each resistance state.

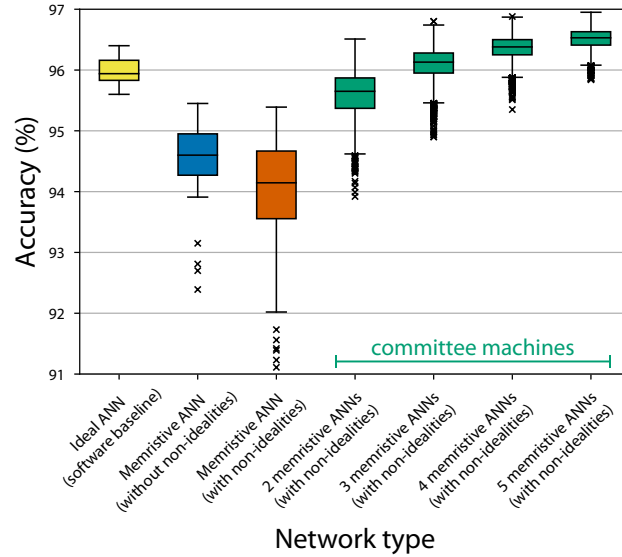


Figure 6. Accuracy achieved by individual networks and their committees when RTN data of a Ta_2O_5 device are taken into account. Additionally, interconnect resistance of $0.3\ \Omega$ – $0.35\ \Omega$ and $0.32\ \Omega$ in the word and bit lines, respectively, (from Ta/HfO₂ array) was used to include line resistance effects. The maximum whisker length is set to $1.5 \times \text{IQR}$.

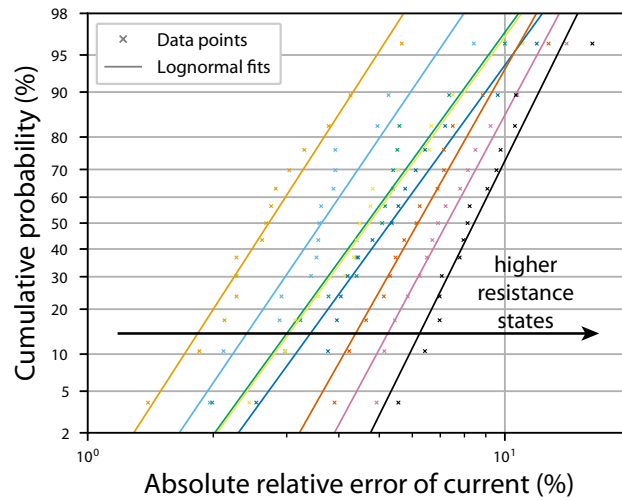


Figure 7. Cumulative probability plots of RTN-induced relative current deviations for all 8 resistance states of aVMCO RRAM device. Lognormal fits are shown for each resistance state.

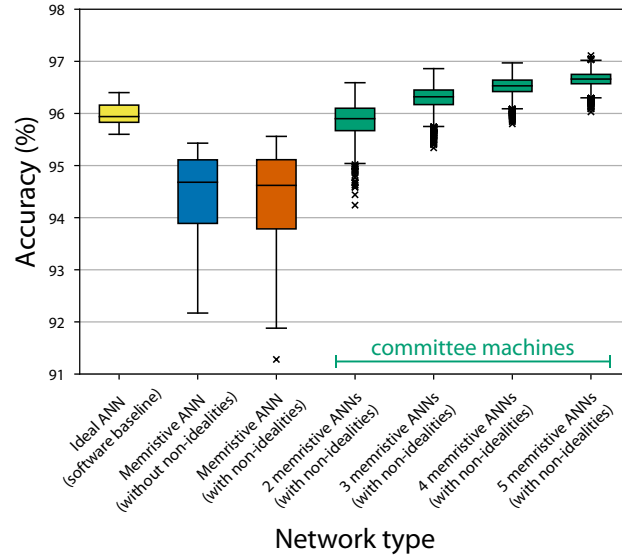


Figure 8. Accuracy achieved by individual networks and their committees when RTN data of an aVMCO device are taken into account. Additionally, interconnect resistance of $0.3\ \Omega$ $0.35\ \Omega$ and $0.32\ \Omega$ in the word and bit lines, respectively, (from Ta/HfO₂ array) was used to include line resistance effects. The maximum whisker length is set to $1.5 \times \text{IQR}$.

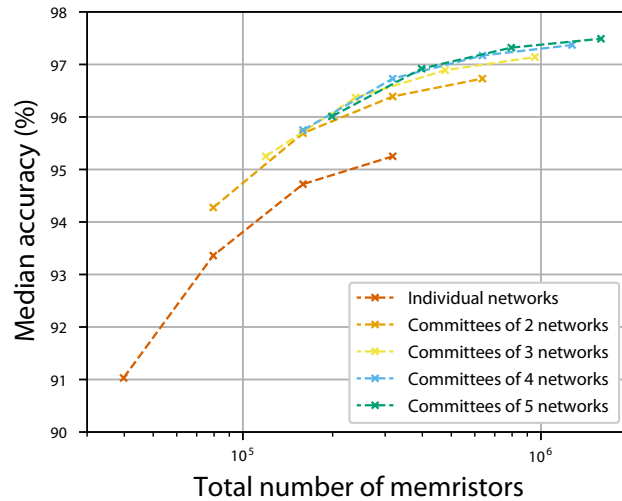


Figure 9. Median accuracy achieved by individual one-hidden-layer memristor-based networks and their committees, when controlled for total number of memristors required. The networks contained 25, 50, 100 or 200 hidden neurons and were disturbed using faulty devices and D2D variability data from Ta/HfO₂ crossbar.

First author (year)	Non-ideality	Device type	Proposed solution
C. Sung (2018) [31]	Current/voltage non-linearity	TaO _x RRAM	Hot-forming step is adopted
C. Li (2018) [15]	Current/voltage non-linearity	Ta/HfO ₂ RRAM	1T1R architecture is adopted
Y. Fang (2018) [32]	Device-to-device variability	HfO _x RRAM	Ultra-thin ALD-TiN buffer layer is introduced
B. Govoreanu (2013) [33]	Device-to-device variability	Al ₂ O ₃ /TiO ₂ (VMCO) RRAM	Non-filamentary RRAM is adopted
A. J. Kenyon (2019) [34]	Device-to-device variability	SiO _x RRAM	The roughness of bottom electrodes is increased
L. Xia (2017) [14]	Faulty devices	-	A modified mapping algorithm and redundancy schemes are used
S. Ambrogio (2018) [7]	Limited dynamic range	PCM	Two pairs of conductance of varying significance for every synaptic weight are used
M. Hu (2016) [17]	Line resistance	-	Advanced mapping algorithms are used to compensate for line resistance effects
W. Wu (2018) [35]	Programming non-linearity	HfO _x RRAM	Electro-thermal modulation layer is deposited on the switching layer
J. Woo (2016) [9]	Programming non-linearity	HfO ₂ RRAM	Bilayer structure is adopted
S. Ambrogio (2018) [7]	Programming non-linearity	PCM	PCM devices are used together with CMOS transistors
Z. Chai (2018) [36]	Random telegraph noise	TiO ₂ /a-Si (aVMCO) RRAM	Non-filamentary RRAM is adopted

Table I. Examples of past efforts at dealing with non-idealities of memristive devices and their systems.