# **CMOS-Compatible Electrochemical Synaptic Transistor**

# 2 Arrays for Deep Learning Accelerator

- 3 Jinsong Cui<sup>1</sup>, Fufei An<sup>1</sup>, Jiangchao Qian<sup>1</sup>, Yuxuan Wu<sup>1</sup>, Luke L. Sloan<sup>1</sup>, Saran Pidaparthy<sup>1</sup>,
- 4 Jian-Min Zuo<sup>1,2</sup>, and Oing Cao<sup>1-5\*</sup>
- <sup>1</sup>Department of Materials Science and Engineering, University of Illinois at Urbana-Champaign,
- 6 Urbana, IL, USA. <sup>2</sup>Materials Research Laboratory, University of Illinois at Urbana-Champaign,
- 7 Urbana, IL, USA. <sup>3</sup>Department of Electrical and Computer Engineering, University of Illinois at
- 8 Urbana-Champaign, Urbana, IL, USA. <sup>4</sup>Department of Chemistry, University of Illinois at
- 9 Urbana-Champaign, Urbana, IL, USA. <sup>5</sup>Holonyak Micro and Nanotechnology Laboratory,
- 10 University of Illinois at Urbana-Champaign, Urbana, IL, USA
- 11 \*e-mail: qingcao2@illinois.edu

12

1

13 Abstract: In-memory-computing architectures based on memristive crossbar arrays could offer 14 higher computing efficiency than traditional hardware in deep learning. However, the core memory devices must be capable of performing high-speed and symmetric analogue 15 16 programming with small variability. They should also be compatible with silicon technology 17 and scalable to nanometre-size footprints. Here we report an electrochemical synaptic transistor 18 that operates by shuffling protons between hydrogenated tungsten oxide channel and gate 19 through a zirconium dioxide protonic electrolyte. These devices are compatible with 20 complementary metal-oxide-semiconductor (CMOS) technology and can be scaled to lateral 21 dimensions of 150×150 nm<sup>2</sup>. They can be monolithically integrated with silicon transistors to 22 form pseudo-crossbar arrays as chip-area and energy efficient deep learning accelerators in 23 which parallel and symmetric programming of the channel conductance with small cycle-to-24 cycle variations are achieved via gate-voltage pulse. They can be programmed at frequencies approaching megahertz and offer endurance over 100 million read-write cycles. 25

#### Main

Advances in deep learning have allowed computers to outperform humans in a range of complicated tasks. These developments have though largely been achieved by increasing the depth and size of the neural network models used, which has exponentially increased the computational load and energy consumption needed to train and execute them<sup>1</sup>. In-memory-computing architectures based on crossbar arrays of non-volatile memory are of potential use as efficient deep learning accelerators<sup>2</sup>. They offer compatibility with the back-end-of-line (BEOL) integration with silicon logic circuits and excellent scalability, providing on-chip data storage with density and capacity high enough to eliminate off-chip memory access during data intensive computations, and capability to perform analogue matrix operations in parallel without incurring data movement. They could thus potentially offer higher computational efficiency and lower chip-area cost than microprocessors based on von Neumann architectures<sup>3</sup>.

To successfully train deep neural networks using memristive arrays, the weight movement in the training algorithms needs to be accurately mapped into the characteristics of the memory hardware. The device conductance change used to represent the weight modulation should be symmetric for positive and negative stimulus, and be precise with small cycle-to-cycle variability. If not, the training accuracy will be dramatically affected<sup>4-6</sup>. This requirement, however, remains a challenge as in most memristive devices — including those based on phase-change materials (PCMs)<sup>7</sup>, filament-forming metal oxides<sup>8-10</sup>, and ferroelectric oxides<sup>11</sup> — conductance increase (potentiation) and decrease (depression) characteristics are highly asymmetric with large variability due to the underlying material and physical device mechanisms<sup>3-5</sup>.

Metal-oxide memristors, for example, typically have an abrupt potentiation that is due to positive feedback between the filament growth and the electric field, and a large variability as the filament-formation process is intrinsically stochastic. Such non-ideal asymmetric programming can be compensated for by using advanced training algorithms, but it typically requires doubling

the chip area, and penalties in latency and energy consumption<sup>12</sup>. Recently, memristors with 1 2 improved symmetry have been developed by employing complex switching media<sup>13-16</sup>, but they 3 still suffer from large variability, and many of them are not compatible with direct BEOL integration<sup>15,16</sup>. 4 5 Electrochemical synaptic transistors were first used to construct artificial neural networks in 6 1960s<sup>17</sup>, and have recently resurfaced as a candidate to realize deep learning accelerators based 7 on in-memory-computing architectures<sup>18</sup>. As initially designed for neuromorphic computing 8 rather than digital information storage, their channel conductance can be precisely modulated by 9 the electrochemical intercalation reactions controlled by a bias applied to separate gate terminals. 10 They thus provide multistate analogue programming with high symmetry and low variability 11 compared to two-terminal memristive devices. However, current electrochemical random-access 12 memory (ECRAM) demonstrations are limited by poor compatibility with silicon 13 complementary metal-oxide-semiconductor (CMOS) technology, which must be monolithically 14 integrated with the memristive arrays as part of the accelerator to provide peripheral functions. The electrolytes used in most ECRAM devices are either liquid or organic polymers<sup>19-21</sup>. 15 16 Despite their good performance, it is difficult to incorporate these devices in circuits as scaled 17 memory cells with the long-term stability and reliability required for electronic applications, and 18 they are often only functional in a controlled environment. All-solid-state inorganic ECRAM

designs have been developed. However, some use PdHx, which suffer from their chemical instability towards air and moisture, as a protonic intercalant reservoir<sup>22,23</sup>. Many others use intercalation of lithium ions<sup>24,25</sup>, which is not compatible with silicon-CMOS as lithium solid electrolytes are generally air-sensitive and lithium ions can readily diffuse across various oxides into the silicon lattice. These lithium ions change the silicon charge carrier concentration as shallow interstitial dopants. The others adopt oxygen ions as the intercalant 26-28, whose sluggish motion limits the device programming speed typically to above several milliseconds<sup>29,30</sup>. In this Article, we report an all-solid-state inorganic ECRAM that operates through the

reversible insertion of protons into a HxWO3 channel from hydrogenated ZrO2 electrolyte and

19

20

21

22

23

24

25

26

H<sub>x</sub>WO<sub>3</sub> gate. These devices exhibit highly symmetric programming under gate-voltage pulses, low cycle-to-cycle and device-to-device spatiotemporal variability, reliable read-write operations above 100 million cycles, and energy consumption below femto-joule per transaction. Their channel conductance can be tuned over a wide range, from nano to micro-siemens, enabling the construction of large-size arrays with optimised power and performance<sup>3</sup>. The small radius of the proton intercalant leads to fast ionic dynamics for high-speed write-read programming with frequencies approaching one megahertz. WO<sub>3</sub> and ZrO<sub>2</sub> are compatible with silicon-CMOS technology and associated wafer-scale microfabrication techniques, which allows the ECRAM dimensions to be scaled down to 150×150 nm<sup>2</sup>. With amorphous HfO<sub>2</sub> as both the interlayer dielectric and the proton-diffusion barrier, the ECRAMs can be fabricated above silicon circuits without affecting the performance of the underlying logic transistors. They are functional in harsh conditions such as high vacuum and elevated temperature, and their non-volatile conductance exhibits low drift and long retention. To demonstrate parallel operation and monolithic integration, we fabricate ECRAM pseudo-crossbar arrays addressed with integrated silicon selector transistors. These ECRAM-CMOS hybrid in-memory-computing accelerators could achieve training accuracy comparable to that of digital accelerators based on static random-access memory (SRAM), but under drastically reduced chip-area cost and energy consumption.

### **Device architecture and operations**

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

A device schematic is shown in Fig. 1a, with its gate stack revealed in a cross-sectional scanning transmission electron microscopy (STEM) micrograph (Fig. 1b). The detailed fabrication process is described in Methods and Extended Data Fig. 1. We choose the stoichiometric amorphous WO<sub>3</sub> deposited by the reactive sputtering from WO<sub>3</sub> target as the ECRAM channel as it gives both low base conductance and long proton retention (See Supplementary Note 1). Compared to VO<sub>2</sub><sup>31,32</sup>, the intercalation of protons into amorphous WO<sub>3</sub> does not involve the phase transformation of the crystal structure<sup>33</sup>, ensuring faster ECRAM operations with better device endurance. A hydrogen-spillover process is then utilized to

1 introduce precisely controlled concentration of protons into the WO3 on wafer scale (See 2 Supplementary Note 2), followed by the deposition of an ultrathin film of ZrO<sub>2</sub> by atomic layer 3 deposition (ALD) as the solid-state protonic gate electrolyte. The secondary-ion mass 4 spectrometry (SIMS) depth profile indicates that substantial amount of protons are diffused into 5 the ZrO<sub>2</sub> from the H-doped WO<sub>3</sub> channel during ALD (Fig. 1c), which help to passivate the 6 surface and grain boundaries of ZrO<sub>2</sub> with hydroxyl groups, forming pathways for fast protonic 7 conduction<sup>34</sup>. The gate is a H<sub>x</sub>WO<sub>3</sub> and metal bilayer, where the H-rich H<sub>x</sub>WO<sub>3</sub> serves as the 8 proton reservoir and helps to minimize the device built-in potential (See Supplementary Note 9  $3)^{21}$ . Device channel conductance G can be programed through 64 (6-bit) discrete states with 10 high-level of symmetry, small cycle-to-cycle variability, and a large dynamic range above 20 11 under voltage pulses (amplitude V<sub>G</sub>=±4V, Fig. 1d), in contrary to those built with elemental 12 metal gate which can only operate symmetrically under current pulses (Extended Data Fig. 2). 13 The zoom-in view shows three continuous conductance states out of the total 64 and confirms 14 that the CMOS-compatible protonic ECRAM prototype has low read noise with conductance 15 down to nanosiemens regime (Fig. 1e). The gate current during write is merely around a few 16 pico-amperes (Fig. 1f), corresponding to a low switching energy of about 10 pico-joule per step, 17 which can be further reduced to below femto-joule with the adoption of up to 106 times faster 18 write pulses down to microseconds. By adjusting the initial hydrogen concentration in the 19 H<sub>x</sub>WO<sub>3</sub> channel through varying parameters in the hydrogen spillover (See Supplementary Note 20 2), the device base conductance  $G_0$  and the absolute conductance change  $\Delta G$  per weight-update 21 step can be effectively modulated. Fig. 1g illustrates the programing of a ECRAM built on a 22 more hydrogen-rich H<sub>x</sub>WO<sub>3</sub> channel, whose G<sub>0</sub> is increased from 1 to 200 nS with the average 23  $\Delta G$  per step simultaneously increased by  $\sim 10$  times, likely assisted by defect-defect interactions and the increase of proton concentration in the hydrogenated ZrO<sub>2</sub> electrolyte<sup>20</sup>. This tunability 24 25 is critical for large-scale array implementations, where the memristive resistance needs to be 26 optimized to minimize the thermal noise, the voltage drop in the interconnect metal lines, and the 27 overhead on the peripheral circuits<sup>3</sup>, and balance the trade-off between the absolute  $\Delta G$  per step

and the on/off ratio. The channel conductance modulation also depends exponentially on the amplitude of the gate-voltage pulses, as higher gate bias increases the gate-current density and thus the total amount of injected protons per pulse (Fig. 1h). These ECRAM cells demonstrate good retention. The drift coefficient v, obtained by fitting the change of channel conductance as a function of time t to the power law of  $G(t)=G_{t0}(t/t_0)^{-\nu}$  where  $G_{t0}$  is the initial conductance at time  $t_0$  and v is the fitted drift coefficient, is ultra-low (Fig. 1i), especially compared to PCMs<sup>35</sup>. It shows that the slow self-discharging of the ECRAM under zero gate bias ensures that these discrete conductance levels are well preserved. States deviating further from  $G_0$  exhibits a slight increase of the drift coefficient, likely caused by the larger gradient for proton diffusion. Such low drift coefficient over at least 10<sup>3</sup> seconds is more than enough to ensure that the accuracy will not be affected when ECRAM-based accelerators are used to train multilayer perceptrons and even convolutional neural networks with weight sharing in convolutional layers, where weights stored on ECRAM cells are continuously updated during back propagation<sup>36</sup>. The device retention is further improved if measured with zero gate current, which is consistent with the proven capability of electrochromic windows built on H<sub>x</sub>WO<sub>3</sub> to retain their coloration density for years without external power supply<sup>37</sup>.

### Device operating speed, scalability, and endurance

High-speed programing of protonic ECRAMs was then demonstrated with fast gate-voltage pulses. With  $G_0$  around 1  $\mu$ S, a large  $\Delta G$  of 3  $\mu$ S can be induced by 32 potentiation-depression pulses with width down to 300 microseconds to afford an on/off ratio of 3 (Fig. 2a). Fig. 2b shows the reproducible cycling with 10 microsecond and 5 microsecond write pulses. Although the device dynamic range becomes smaller, the programming characteristics are still highly symmetric, and the  $\Delta G$  after each pulse remains larger than 1 nS with low variability and read noise (Extended Data Fig. 3) to enable the storage of 8-bit of information within a single ECRAM cell. The effect of the off-state current can be eliminated with the addition of a dummy column in the crossbar architecture<sup>38</sup>. Compared to state-of-the-art ECRAMs employing oxygen ions as the intercalant<sup>28</sup>, when operated under the 10 microsecond-wide write pulses to achieve

the same conductance-modulation ratio, our protonic ECRAMs require about 7 times lower number of weight-update steps under a ~40% lower gate-voltage amplitude, resulting from proton's higher ionic diffusivity. There is some upshift of the baseline  $G_0$ , indicating that more weight-update operations are required in the depression branch to bring the conductance back to its initial level. However, such asymmetry associated with the slightly different number of the intermediate states between the minimum and maximum conductance states will not degrade the training accuracy as predicted in simulation<sup>5</sup>, and can be compensated by applying voltage pulses with slightly different amplitude or width for the potentiation versus depression operations. Since the conductance change correlates with the amount of charge, i.e., the number of protons, injected or extracted,  $\Delta G$ , as well as the corresponding energy consumption per weight-update step, scales linearly with the pulse width (Fig. 2c). In addition to the time required for the weight update, the ECRAM device speed is also limited by the read transients (See Extended Data Fig. 4 for the measurement setup). For read without write, the ECRAM channel behaves as a resistor capacitively coupled to the device capacitance. We monitored the current flowing across the source-drain electrodes, showing that the settling time  $t_{\text{read}}$  after the application of the read-voltage pulse was less than 500 nanoseconds as limited by the resistor-capacitor (RC) delay (Fig. 2d). For read post update, we observed a fast recovery current pulse followed by a slower decay eventually leading toward a stable sense current (Fig. 2e). This transient behavior is qualitatively similar to that of the ECRAMs employing lithium or oxygen ions as the intercalants, and it is limited by the dielectric relaxation processes and the charge transfer from the center of the channel to contacts on the edge. But the quantitative settling time  $t_{\text{read-after-write}}$  for protonic ECRAMs is >100 times faster with proton's high ionic diffusivity<sup>29,30</sup>. Both t<sub>read</sub> and t<sub>read-after-write</sub> are expected to diminish with scaling, which reduces both the gate capacitance and the required lateral diffusion distance of injected protons. In CMOS-compatible, all-inorganic protonic ECRAMs with their channel lateral dimensions scaled down to 150×150 nm<sup>2</sup> (Fig. 2f), which represent the smallest lateral ECRAMs ever fabricated, their symmetric and

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

linear programing characteristics are well preserved (Fig. 2g), showing the excellent device scalability. The absolute conductance values and the dynamic range of modulation are slightly degraded likely due to the impact from the contact resistance (Extended Data Fig. 5), and the larger read noise could be caused by the miniaturization of the device volume together with the increase of the channel surface-to-volume ratio, which reduce the averaging effect and increase the impact from traps at material interfaces. Even though the gate capacitance is significantly reduced with  $>10^3$  times smaller device channel area, the  $t_{read}$  of the scaled ECRAMs is only marginally improved to 370 nanoseconds (Fig. 2h), indicating that the device capacitance is mainly limited by parasitics from, for example, the overlap regions between the source/drain and gate electrodes with the hydrogenated ZrO<sub>2</sub> sandwiched in between. However, the rate of the current stabilization in the read post update shows much more significant improvement, with the apparent t<sub>read-after-write</sub> less than 1 microsecond (Fig. 2i), showing the benefit of ECRAM channellength scaling to mitigate the device read transients associated with charge/ion transport. To improve the device latency down to the desired 50-100 nanosecond regime, it is required to further suppress the device read transients and increase the  $\Delta G$  per step in weight update, which can be accomplished by further scaling down the lateral length of the H<sub>x</sub>WO<sub>3</sub> channel, optimizing the device structure to minimize the parasitic capacitance, modifying the composition and nanostructures of the ZrO<sub>2</sub>-based electrolyte to increase its ionic conductivity<sup>34</sup>, and reducing the gate stack (H<sub>x</sub>WO<sub>3</sub> channel and the electrolyte) thickness (Supplementary Note 4). In addition to their fast operations and excellent scalability, the CMOS-compatible protonic ECRAMs exhibited excellent endurance when modulated in air with fast voltage pulses. No degradation in terms of programming symmetry, cycle-to-cycle variability, baseline channel conductance, or conductance modulation was observed after 108 write-read operations, which correspond to over one million switches over the full on/off range (Fig. 2j). The moderate dynamic range of ECRAM corresponds to limited modulation of the proton concentration in the H<sub>x</sub>WO<sub>3</sub> channel, which ensures that the devices will not degrade due to repetitive lattice expansion or ion trapping<sup>39</sup>. There could be slightly oxygen loss accompanying the proton

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

insertion and removal in the amorphous H<sub>x</sub>WO<sub>3</sub> channel, but it is insignificant even after million times of operations driven by the gate pulses (Extended Data Fig. 6), as measured within the accuracy of TEM-energy dispersive X-ray spectroscopy (EDS). The variation of the maximum and minimum conductance after 32 potentiation and depression pulses is small with a relative standard deviation around 7%, which indicates a standard deviation normalized to the entire weight range around  $7\%/\sqrt{32} = 1\%$  for each weight-update cycle. With the 20 nm thick HfO<sub>2</sub> or Al<sub>2</sub>O<sub>3</sub> passivation deposited by ALD at 120 °C, the device operation was not affected even when measured in high vacuum (10<sup>-6</sup> Torr, Fig. 2k). It is a significant advantage over previous protonic ECRAMs employing polymer/nanoporous electrolytes and/or PdH<sub>x</sub> gate, which are only functional in a humidity-controlled, utmost medium vacuum, environment or forming gas<sup>21</sup>-<sup>23,40,41</sup>. The stable and symmetric programming characteristics were maintained even after annealing at 200 °C, and the device can be successfully operated at 80–100 °C without degrading the data retention or device endurance (Extended Data Fig. 7), which further verifies the reliability of our all-inorganic protonic ECRAM prototype in harsh environment. Similar as other memristive devices 40,42, the device conductance and conductance modulation systematically shift with the increase of temperature due to thermal excitation and activation, which needs to be compensated on the software level with the chip temperature continuously monitored by built-in temperature sensors. Finally, the performance of the all-inorganic protonic ECRAM is benchmarked with other BEOL-compatible analogue memory technologies, and shows clear advantages in symmetry, cycle-to-cycle variability, and energy consumption (See Supplementary Note 5).

#### ECRAM pseudo-crossbar arrays

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

All-inorganic protonic ECRAMs were monolithically integrated with silicon metal-oxide-semiconductor field-effect transistors (MOSFETs) (See Methods and Extended Data Fig. 8 for the process flow), in the form of 1-transistor-1-ECRAM cells (Fig. 3a), where the silicon transistor acted as the selector to solve the write-disturbance problem (Fig. 3b, there is negligible ECRAM conductance change with the selector transistor turned off). Note that ECRAM arrays

can potentially operate without selectors using the half-selection bias scheme<sup>27,28</sup>, and here we mainly use the fabrication and characterization of 1-transistor-1-ECRAM cells to verify the CMOS compatibility. With the ALD HfO<sub>2</sub> deposited on top of the completed silicon MOSFETs and the metal interconnects as both the interlayer dielectric and proton diffusion barrier, the performance of these underlying silicon devices was not affected by the addition of the ECRAM layer on top (Extended Data Fig. 9). Beyond individual cells, an integrated pseudo-crossbar array was constructed (Fig. 3c,d). In the pseudo-crossbar array fabric (Fig. 3e), parallel row-byrow weight updates were successfully implemented, as in the programming patterns shown in Fig. 3f. The average weight-update in the selected cells was 210±20%, while those in the unselected rows were not disturbed with an average  $\Delta G$  of -1±4%. After parallel programing, the pseudo-crossbar arrays can be used to parallelly execute the vector-matrix multiplication, which is the core and the most expensive computing operation in many deep-learning and imageprocessing algorithms. We used the color transformation as an example, where a transformation matrix is used to generate modified red (R), green (G), and blue (B) pixel data for recoloring (Fig. 3g). Each element of the transformation matrix was first mapped to the ECRAM array. The color of each pixel represented by an 8-bit number for each R/G/B channel was then converted to a voltage input vector delivered to the bit lines of the fabric. The current measured at the source lines for the weight sum was the result of the dot product between the input voltage vector and the analogue conductance matrix following the Kirchhoff's law, and it was finally encoded back to generate the modified RGB values (Fig. 3h). Despite the non-idealities associated with the device variability and conductance drift, the color transformation result was comparable to what performed by the software with a narrow error band (Fig. 3i). With the functionality of ECRAM pseudo-crossbar arrays verified in experiment, we then simulated the performance of accelerators built on ECRAM arrays integrated with silicon-CMOS peripheral circuits<sup>43</sup>. Based on the experimentally measured dynamic range, non-linearity, symmetry, read-write voltage and speed, cycle-to-cycle weight-update variation, and device-to-

device variations (Extended Data Fig. 10) of CMOS-compatible, all-inorganic protonic

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

ECRAMs (See Supplementary Note 6), the modeling, which considers all the device nonidealities, indicates that such in-memory-computing accelerators can achieve similar level of accuracy (86-93%, the small degradation in accuracy for devices operated with faster pulses is caused by their smaller on/off ratio) in learning the classification of the Modified National Institute of Standards and Technology (MNIST) dataset compared to the SRAM-based digital accelerators or the software (94%), as enabled by ECRAMs' good linearity and symmetry, low variability, and their large number of available states (Fig. 3j). Meanwhile, the overall energy consumption, including contributions from the ECRAM array, the selector transistors, and the periphery circuit, is up to two times lower than SRAM-based accelerators (Fig. 3k), and the chiparea cost, with the analogue and high-resistance ECRAMs monolithically integrated on top of the silicon circuits, is reduced by over 10 times. The latency is 50 times longer due to the much slower operation of our current ECRAM prototypes compared to SRAM (Fig. 31), but it is competitive against accelerators built on PCM and metal-oxide memristors (Supplementary Note 7)<sup>43</sup>. ECRAM accelerators can also be used to train the convolutional neural networks, achieving similar level of accuracy (~85%) as benchmarked against their SRAM-based counterparts or software (Fig. 3m), but superior to those employing other emerging non-volatile memories<sup>44</sup>.

### **Conclusions**

We have reported a CMOS-compatible, all-inorganic protonic ECRAM technology that offers highly symmetric switching characteristics, low cycle-to-cycle variability, high read-write pulse endurance, and good retention with low drift. Our approach uses a hydrogenated ZrO<sub>2</sub> as a protonic electrolyte to ensure CMOS-compatibility and provides high diffusivity of protons for fast device operation approaching megahertz. It adopts a symmetric gate stack composed of H<sub>x</sub>WO<sub>3</sub> channel and gate to enable symmetric voltage programming. The smallest lateral ECRAMs have an active device area of 150×150nm<sup>2</sup>, verifying excellent device scalability. We fabricated a pseudo-crossbar array composed of ECRAM synaptic memory cells monolithically integrated on top of silicon selector transistors, which can perform both parallel analogue programming and vector-matrix multiplication. The results illustrate the technological potential

- of our electrochemical synaptic transistor arrays in energy- and cost-efficient in-memory
- 2 computing accelerators for deep learning. The key remaining challenge is to further reduce the
- device speed down to 50–100 nanoseconds through both scaling and improving the solid-state
- 4 protonic electrolyte.

### 5 Methods

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

Fabrication of CMOS-compatible all-inorganic protonic ECRAMs. The process flow is schematically illustrated in Extended Data Figure 1. 40 nm of HfO<sub>2</sub> was first deposited as a proton-diffusion barrier and etching-stop layer by ALD (Veeco Nanotech ALD System) on a silicon substrate covered with 90 nm thermal oxide, using Tetrakis(dimethylamido)-hafnium and water as precursors at 200 °C. A photolithography step was then performed to define the sourcedrain electrode patterns into the photoresist (AZ 5214, exposure dose was 200mJ·cm<sup>-2</sup>), and an electron-beam evaporation (Temescal) was used to deposit 2 nm Cr/ 40 nm Pt, followed by liftoff in acetone to form the source-drain contacts. 80 nm of WO<sub>3</sub> was grown by reactive radiofrequency magnetron sputtering (Kurt J. Lesker PVD 75) of WO3 target (Kurt J. Lesker) by means of a plasma composed of argon as the carrier gas and oxygen as the reactive gas. After pumping down the sputtering-deposition chamber to high vacuum (5×10<sup>-7</sup> Torr), pure oxygen and argon were introduced with their relative flow rates adjusted to the desired ratio with the help of the mass-flow controllers to a stable chamber pressure of 1.5 millitorr. During deposition, the substrate was held at 300 °C and the radio-frequency forward input power was maintained at 140 W. After WO<sub>3</sub> deposition to the desired thickness of 50–100 nm, another photolithography step was performed to define the device channel footprint into the photoresist. A subsequent SF<sub>6</sub> reactive-ion etching (RIE, Oxford Mixed ICP-RIE system) removed the WO<sub>3</sub> in the unprotected area to stop on the HfO<sub>2</sub> or Pt metal surfaces. The inductively coupled plasma (ICP) power was 1,000 W, the RIE power was 100 W, the flow rate of SF<sub>6</sub> was 20 s.c.c.m., and the chamber pressure was maintained at 5 millitorr. A thin, i.e., 10 nm thick, Aluminum film was then blanketly deposited covering the whole substrate by sputtering (AJA Orion-8 Magnetron

Sputtering System). The deposited Al film was etched by soaking the substrate in 4 mol·L<sup>-1</sup> HCl 1 2 aqueous solution, and part of the hydrogen generated was incorporated into the WO<sub>3</sub> channel as 3 proton intercalants through hydrogen spillover. After further adjusting the proton concentration 4 in H<sub>x</sub>WO<sub>3</sub> by annealing at 150 °C in air, the substrate was immediately transferred to an ALD 5 chamber, where 15 nm of ZrO<sub>2</sub> was deposited using Tetrakis-(dimethylamino)-zirconium (Strem 6 Chemicals) and water vapor as precursors at 120 °C. The low ALD deposition temperature and 7 the adoption of water as precursor are critical to maintain the proton content within the ZrO<sub>2</sub> to enhance its ionic conductivity<sup>45,46</sup>. Photolithography was used again to pattern the gate electrode, 8 9 which has overlap with the source-drain contacts, into the photoresist. 20 nm of WO<sub>3</sub> was 10 deposited again by the radio-frequency reactive sputtering with the substrate held at room 11 temperature, and then converted to H<sub>x</sub>WO<sub>3</sub> by means of the same hydrogen spillover process as 12 described above. Electron-beam evaporation was used to deposit 2 nm Cr/40 nm Au followed by 13 lift-off of the photoresist in acetone to complete the gate electrode stack composed of 14 H<sub>x</sub>WO<sub>3</sub>/Cr/Au trilayer. Another 20 nm of HfO<sub>2</sub> or Al<sub>2</sub>O<sub>3</sub> was deposited by ALD as a passivation 15 layer to improve the device stability and endurance. A final lithography step was performed to 16 expose the probing pads for the source-drain electrodes, where the ZrO<sub>2</sub> and HfO<sub>2</sub>/Al<sub>2</sub>O<sub>3</sub> 17 dielectrics on top were etched by HF to complete the device fabrication flow. 18 Fabrication of nanoscale ECRAMs. To fabricate ECRAMs with their device channel length 19 and channel width scaled down to sub-micron regime, electron-beam lithography (EBL), instead 20 of photolithography, was used to define the device source-drain contacts with 150 nm channel 21 length, the device isolation pattern for the metal hard mask with width down to 150 nm, and the 22 composite H<sub>x</sub>WO<sub>3</sub>/metal gate, using 6% ethyl lactate (EL-6) and polymethyl methacrylate 23 (PMMA 950k A3, Kayaku Advanced Materials) bilayer as the photoresist. The exposure dose was 750 μC·cm<sup>-2</sup> and the developer was methyl isobutyl ketone (MIBK) 1:2 diluted with 24 25 isopropyl alcohol (IPA). In the device isolation step, 100 nm Au hardmask defined by liftoff was 26 used instead of the photoresist, which was stripped by soaking the substrate in commercial gold 27 etchant (Transene TFA) after the SF<sub>6</sub> RIE.

1 Monolithic integration of all-inorganic protonic ECRAMs with silicon MOSFETs. The 2 process, which is schematically illustrated in Extended Data Figure 8, started from fabricating 3 the silicon MOSFETs on a silicon-on-insulator (SOI) substrate (Soitec, 70 nm lightly p-doped 4 device-layer silicon on 2 µm buried oxide). A layer of thermal oxide was first grown by dry 5 oxidation at 1050 °C to a thickness of about 90 nm. Photolithography was performed to define 6 patterns of the heavily doped source-drain regions of the transistors into the photoresist (AZ-7 5214), followed by removing the SiO<sub>2</sub> in the exposed area with buffered oxide etchant. After 8 stripping the photoresist in acetone, a film of phosphorus-containing spin-on-dopant (Filmtronics 9 P509) was blanketly deposited by spin casting at 3,000 rpm for 30 seconds followed by a soft 10 bake at 110 °C for 3 min. Annealing at 850 °C for 20 minutes in a three-zone tube furnace 11 (Lindberg) with N<sub>2</sub> (2 L·min<sup>-1</sup>) and O<sub>2</sub> (1 L·min<sup>-1</sup>) flow caused the phosphorus to diffuse from 12 the spin-on-dopant into the underlying silicon to form the heavily n-doped source/drain contact 13 regions. After cooling down to room temperature, the wafer was immersed in HF to remove 14 both the spin-on-dopant and the thermal oxide mask, followed by piranha cleaning to remove the 15 residual phosphorus oxide. The device islands were then patterned by photolithography, and the 16 silicon in the exposed area was removed by timed ICP-RIE (20 mtorr, 20 s.c.c.m. CF4 flow, 300 17 W ICP power, and 100 W RIE power for 90 seconds) stopping on the buried oxide of the SOI 18 wafer. After removing the photoresist by acetone and cleaning the surface by piranha solution, 19 the source-drain contact electrodes were patterned by the third photolithography step, followed 20 by the electron-beam evaporation of 2 nm Cr/40 nm Au and lift-off in acetone. 40 nm of HfO<sub>2</sub> 21 was then deposited as the high- $\kappa$  gate dielectric by ALD, using Tetrakis(dimethylamido)hafnium and water as precursors at 200 °C. Another photolithography and the lift-off scheme 22 23 were subsequently performed again to define the metal gate electrodes composed of 2 nm Cr and 24 40 nm Au to complete the silicon MOSFET fabrication. Afterwards, another 40 nm HfO<sub>2</sub> was 25 deposited by ALD at 200 °C, serving as both the interlayer dielectric and the proton-diffusion 26 barrier. The ECRAM layer was then fabricated on top using the process described above. The 27 vias for interlayer interconnects were exposed by photolithography, with ICP-RIE utilized to

1 etch through the ZrO<sub>2</sub> gate electrolyte of the ECRAM, the HfO<sub>2</sub> interlayer dielectric, and the 2 HfO<sub>2</sub> gate oxide of the silicon MOSFET to stop on the metal contact (5 mtorr, 10 s.c.c.m. CHF<sub>3</sub> 3 flow and 5 s.c.c.m. Ar flow, 100 W ICP power, and 40 W RIE power). A final lithography and 4 lift-off were performed to pattern the interlayer interconnects connecting the source electrodes of 5 the silicon selector transistors to the gate electrodes of the ECRAMs. 6 **Instrumentation.** The cross-sectional high-angle annular dark field (HAADF) STEM images 7 and the associated elemental configurations were obtained using the FEI Talos F200X G2 STEM 8 equipped with four-crystal EDS system (FEI Super-X). All STEM-EDS data were collected for 9 more than 30 minutes with a 10 microseconds dwell time and ~200 pA probe current, at an 10 accelerating voltage of 200 kV. The sample was prepared using the Thermo Scios2 dual-beam 11 focused-ion beam to a thickness around 50 nm under 5 kV. The SIMS elemental depth profile 12 was obtained using the Phi TRIFT III time-of-flight SIMS system with the material removal 13 performed using low-energy Cs<sup>+</sup> ion source. SEM micrographs were acquired using a Hitachi 14 S4800 microscope. The device electrical characterizations were performed either in ambient or 15 under high vacuum at the desired temperature using a manual probe station (LakeShore 16 Cryotronics CRX-6.5K) connected with a semiconductor parameter analyzer (Keysight B1500A) 17 equipped with integrated high-resolution source-measurement units (Keysight B1517A) and 18 waveform generator/fast-measurement unit (Keysight B1530A). The AZ 5214E photoresist was 19 patterned with a Heidelberg MLA150 aligner, and the EL6/PMMA photoresist was patterned 20 with a Elionix ELS-G150 150 keV EBL system. The Neurosim simulations were performed on the Illinois Campus Cluster HAL<sup>47</sup>. The X-ray diffraction (XRD) spectrum was recorded using 21 22 the Bruker D8 Advanced XRD system. The composition of deposited tungsten oxides was 23 determined by Rutherford backscattering spectrometry (RBS) using the NEC Pelletron 24 accelerator equipped with RBS chamber. Optical transmittances of WO<sub>3</sub> and H<sub>x</sub>WO<sub>3</sub> were 25 recorded using an Agilent Cary 5000 ultraviolet-visible absorbance spectrophotometer. XPS

spectra of WO<sub>3</sub> and H<sub>x</sub>WO<sub>3</sub> were acquired on a PHI Versa Probe III instrument with a

- 1 monochromatic Al Kα (1486.6 eV) source. The pass energy was 55 eV and spectra were
- 2 referenced to C1s peak (adventitious carbon) at 284.8 eV.

- 4 Data availability: All data are available in the main text or the supplementary materials. All
- 5 information and materials can be requested from the corresponding author.

6

- 7 Supplementary Information is linked to the online version of the paper at www.nature.com/natelectron/.
- 8 Acknowledgements This work was supported by the US National Science Foundation grant 1950182 (Q.C.) and
- 9 2139185 (Q.C. and J.-M.Z.). Significant aspects of the material characterizations and device fabrications were
- performed using the shared user facilities of the University of Illinois Materials Research Laboratory and Holonyak
- 11 Micro and Nanotechnology Laboratory. This work also made use of the Illinois Campus Cluster HAL, utilizes
- 12 resources supported by the National Science Foundation's Major Research Instrumentation program grant 1725729
- 13 and the University of Illinois at Urbana-Champaign. We thank Dr. Timothy Spila for assistance with RBS and
- 14 SIMS measurements, and Dr. John Baltrus for performing the XPS measurements.
- 15 Author Contributions Q.C. conceived and designed the experiments. J.-S.C., F.-F.A., Y.-X.W., and L.L.S.
- 16 performed the experiments. J.-C. Q., S.P. and J.-M.Z. performed the STEM/EDS analysis. Q.C. and J.-S.C. wrote
- 17 the manuscript. All authors discussed the results (all of which are reported in the main text and supplement) and
- 18 commented on the manuscript.
- 19 Competing Interests: J.-S.C. and Q.C. are inventors on a provisional patent application entitled "Solid-state
- electrochemical random access memory (ECRAM) and methods of making and operating a solid-state ECRAM"
- 21 (63/434,627) submitted by the Board of Trustees of the University of Illinois.
- 22 Author Information Reprints and permissions information is available at <a href="www.nature.com/reprints">www.nature.com/reprints</a>.
- 23 Correspondence and requests for materials should be addressed to Q.C. (email: qingcao2@illinois.edu).

2

#### References

- Anthony, L. F. W., Kanding, B. & Selvan, R. Carbontracker: Tracking and predicting the carbon footprint of training deep learning models. *arXiv*, 2007.03051 [cs.CY] (2020).
- Xia, Q. & Yang, J. J. Memristive crossbar arrays for brain-inspired computing.
   Nat. Mater. 18, 309-323 (2019).
- Haensch, W., Gokmen, T. & Puri, R. The next generation of deep learning hardware: Analog computing. *Proc. IEEE* **107**, 108-122 (2019).
- Yu, S. Neuro-inspired computing with emerging nonvolatile memorys. *Proc. IEEE* **106**, 260-285 (2018).
- 12 5 Xi, Y. *et al.* In-memory learning with analog resistive switching memory: A review and perspective. *Proc. IEEE* **109**, 14-42 (2021).
- Gokmen, T. & Vlasov, Y. Acceleration of deep neural network training with resistive cross-point devices: Design considerations. *Front. Neurosci.* **10**, 333 (2016).
- 7 Ambrogio, S. *et al.* Equivalent-accuracy accelerated neural-network training using analogue memory. *Nature* **558**, 60-67 (2018).
- Prezioso, M. *et al.* Training and operation of an integrated neuromorphic network based on metal-oxide memristors. *Nature* **521**, 61-64 (2015).
- Yao, P. *et al.* Fully hardware-implemented memristor convolutional neural network. *Nature* **577**, 641-646 (2020).
- Wan, W. *et al.* A compute-in-memory chip based on resistive random-access memory. *Nature* **608**, 504-512 (2022).
- Dutta, S. *et al.* Monolithic 3D integration of high endurance multi-bit ferroelectric FET for accelerating compute-in-memory. *IEDM Tech. Digest*, 36.34.31-36.34.34 (2020).
- 28 12 Gokmen, T. & Haensch, W. Algorithm for training neural networks on resistive device arrays. *Front. Neurosci.* **14**, 103 (2020).
- Woo, J. *et al.* Improved synaptic behavior under identical pulses using AlO<sub>x</sub>/HfO<sub>2</sub> bilayer RRAM array for neuromorphic systems. *IEEE Electron Device Lett.* **37**, 994-997 (2016).
- Wu, W. *et al.* A methodology to improve linearity of analog RRAM for neuromorphic computing. *VLSI Symp. Tech. Digest*, 103-104 (2018).
- Choi, S. *et al.* SiGe epitaxial memory for neuromorphic computing with reproducible high performance based on engineered dislocations. *Nat. Mater.* **17**, 335-340 (2018).
- Mou, X. *et al.* Analog memristive synapse based on topotactic phase transition for high-performance neuromorphic computing and neural network pruning. *Sci. Adv.* **7**. eabh0648 (2021).
- Widrow, B. Generalization and information storage in networks of adaline neurons. 435-461 (Spartan Books, 1962).
- 43 18 Fuller, E. J. *et al.* Redox transistors for neuromorphic computing. *IBM J. Res. Dev.* **63**, 9:1-9:9 (2019).

- 1 19 Gkoupidenis, P., Schaefer, N., Garlan, B. & Malliaras, G. G. Neuromorphic functions in PEDOT:PSS organic electrochemical transistors. *Adv. Mater.* **27**, 7176-7180 (2015).
- 4 20 Yao, X. *et al.* Protonic solid-state electrochemical synapse for physical neural networks. *Nat. Commun.* **11**, 3134 (2020).
- Fuller, E. J. *et al.* Parallel programming of an ionic floating-gate memory array for scalable neuromorphic computing. *Science* **364**, 570-574 (2019).
- Onen, M., Emond, N., Li, J., Yildiz, B. & del Alamo, J. A. CMOS-compatible protonic programmable resistor based on phosphosilicate glass electrolyte for analog deep learning. *Nano Lett.* **21**, 6111-6116 (2021).
- Onen, M. *et al.* Nanosecond protonic programmable resistors for analog deep learning. *Science* **377**, 539-543 (2022).
- Fuller, E. J. *et al.* Li-ion synaptic transistor for low power analog computing. *Adv. Mater.* **29**, 1604310 (2017).
- Tang, J. *et al.* ECRAM as scalable synaptic cell for high-speed, low-power neuromorphic computing. *IEDM Tech. Digest*, 13.11.11-13.11.14 (2018).
- Shi, J., Ha, S. D., Zhou, Y., Schoofs, F. & Ramanathan, S. A correlated nickelate synaptic transistor. *Nat. Commun.* **4**, 2676 (2013).
- Kim, S. *et al.* Metal-oxide based, CMOS-compatible ECRAM for deep learning accelerator. *IEDM Tech. Digest*, 35.37.31-35.37.34 (2019).
- 21 28 Lee, C., Choi, W., Kwak, M., Kim, S. & Hwang, H. Excellent synapse 22 characteristics of 50 nm vertical transistor with WO<sub>x</sub> channel for high density 23 neuromorphic system. *VLSI Symp. Tech. Digest*, 1-2 (2021).
- 24 29 Solomon, P. M. *et al.* Transient investigation of metal-oxide based, CMOS-25 compatible ECRAM. *IEEE Int. Reliab. Phys. Symp. Proc.*, 1-7 (2021).
- 26 30 Bishop, D. *et al.* Time-resolved conductance in electrochemical systems for neuromorphic computing. *2018 International Conference on Solid State Devices and Materials (SSDM)*, 23-24 (2018).
- Jo, M. *et al.* Gate-induced massive and reversible phase transition of VO<sub>2</sub> channels ssing solid-state proton electrolytes. *Adv. Funct. Mater.* **28**, 1802003 (2018).
- 32 Oh, C. *et al.* Deep proton insertion assisted by oxygen vacancies for long-term memory in VO<sub>2</sub> synaptic transistor. *Adv. Electron. Mater.* **7**, 2000802 (2021).
- 34 33 Leng, X. *et al.* Insulator to metal transition in WO<sub>3</sub> induced by electrolyte gating. *npj Quantum Mater.* **2**, 35 (2017).
- Meng, Y. *et al.* Review: Recent progress in low-temperature proton-conducting ceramics. *J. Mater. Sci.* **54**, 9291-9312 (2019).
- Ding, K. *et al.* Phase-change heterostructure enables ultralow noise and drift for memory operation. *Science* **366**, 210-215 (2019).
- 40 36 Li, Y. *et al.* Capacitor-based cross-point array for analog neural network with record symmetry and linearity. *VLSI Symp. Tech. Digest*, 25-26 (2018).
- 42 37 Zhang, J.-G., Benson, D., Tracy, C. E., Webb, J. & Deb, S. *Self-bleaching mechanism of electrochromic WO*<sub>3</sub> *films*. Vol. 2017 104-112 (SPIE, 1993).
- 44 38 Chen, P., Peng, X. & Yu, S. NeuroSim: A circuit-level macro model for benchmarking neuro-inspired architectures in online learning. *IEEE Trans.*
- 46 Comput. Aided Des. Integr. Circuits Syst. **37**, 3067-3080 (2018).

- Wen, R.-T., Granqvist, C. G. & Niklasson, G. A. Eliminating degradation and uncovering ion-trapping dynamics in electrochromic WO<sub>3</sub> thin films. *Nat. Mater.* **14**, 996-1001 (2015).
- 4 40 Melianas, A. *et al.* Temperature-resilient solid-state organic artificial synapses for neuromorphic computing. *Sci. Adv.* **6**, eabb2958 (2020).
- Katase, T., Onozato, T., Hirono, M., Mizuno, T. & Ohta, H. A transparent electrochromic metal-insulator switching device with three-terminal transistor geometry. *Sci. Rep.* **6**, 25819 (2016).
- Walczyk, C. *et al.* Impact of temperature on the resistive switching behavior of embedded HfO<sub>2</sub>-based RRAM devices. *IEEE Trans. Electron Devices* **58**, 3124-3131 (2011).
- 12 43 Chen, P., Peng, X. & Yu, S. NeuroSim+: An integrated device-to-algorithm 13 framework for benchmarking synaptic devices and array architectures. *IEDM* 14 *Tech. Digest*, 6.1.1-6.1.4 (2017).
- Peng, X., Huang, S., Jiang, H., Lu, A. & Yu, S. DNN+NeuroSim V2.0: An end-toend benchmarking framework for compute-in-memory accelerators for on-chip training. *IEEE Trans. Comput. Aided Des. Integr. Circuits Syst.* **40**, 2306-2319 (2021).
- 20 Zhou, Y. *et al.* Enhanced transmittance modulation of ITO/NiO<sub>x</sub>/ZrO<sub>2</sub>:H/WO<sub>3</sub>/ITO electrochromic devices. *Ionics* **22**, 25-32 (2016).
- Park, J. S. *et al.* Evidence of proton transport in atomic layer deposited Yttriastabilized zirconia films. *Chem. Mater.* **22**, 5366-5370 (2010).
- 23 47 Kindratenko, V. *et al.* in *Practice and Experience in Advanced Research* 24 *Computing (PEARC' 20).* 41–48 (Association for Computing Machinery).

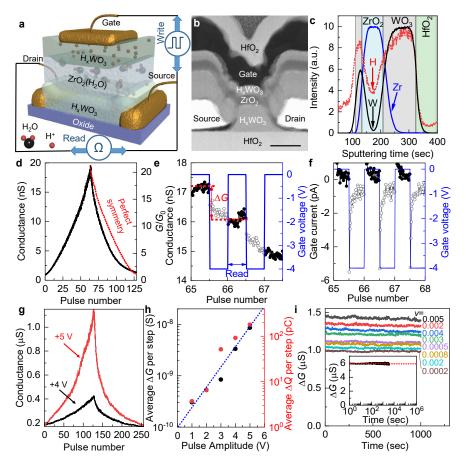


Figure 1. Structure and electrical characteristics of CMOS-compatible, all-inorganic protonic ECRAMs. a, Schematic illustrating the structure and read/write operations of the CMOS-compatible, all-inorganic protonic ECRAM prototype. b, Cross-sectional STEM micrograph showing the device gate stack. Scale bar: 50 nm. c, SIMS depth profiles for W (black), Zr (blue), and H (red) across the ECRAM device gate stack. Background color serves as visual guide to mark different layers including HfO<sub>2</sub> (green),  $H_xWO_3$  (grey), and  $ZrO_2$  (blue). The material interfaces are defined at positions where the major element (W and Zr for  $H_xWO_3$  and  $ZrO_2$  layers, respectively) composition has reached to 50% of the plateau. d, Programing of the CMOS-compatible, all-inorganic protonic ECRAM (channel length  $L_{ch}$ =10  $\mu$ m and width W=3  $\mu$ m) with gate-voltage pulses (64 potentiation then 64 depression pulses with amplitude of ±4 V and width of 3 sec). Red dotted line is the mirror image of the potentiation curve. The absolute device

channel conductance and the relative modulation with regard to the baseline conductance  $(G/G_0)$  are used for the left and right y-axis, respectively. **e-f**, Zoom-in views of the programing characteristics showing the readout of discrete conductance states with sufficient noise margin (part e, source-drain bias for read was 0.1 V) and the gate current during the weight-update (part f), driven by gate voltage pulses (blue, right axis). Red dotted lines serve as visual guide to mark the average channel conductance Current values measured during read and write operations are in each state. represented with solid and hollow symbols, respectively. **g,** Programing of the ECRAM fabricated on a more hydrogen-rich H<sub>x</sub>WO<sub>3</sub> channel with voltage pulses of identical width of 3 sec but different amplitude of ±4 V (black) and ±5 V (red), respectively. h, Average  $\Delta G$  per gate programing pulse (black, left axis) and the total charge injected by the gate current ( $\Delta Q$ , red, right axis) as a function of the pulse amplitude. Blue dotted line represents the linear fitting to the data. i, Retention of selected analog states (0.1 V read) under zero gate bias in ambient and the corresponding drift coefficient (v). Inset: Retention with the gate floating. Red dotted line represents the fitting to the power decay function of  $G(t)=G_{t0}(t/t_0)^{-\nu}$ .

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

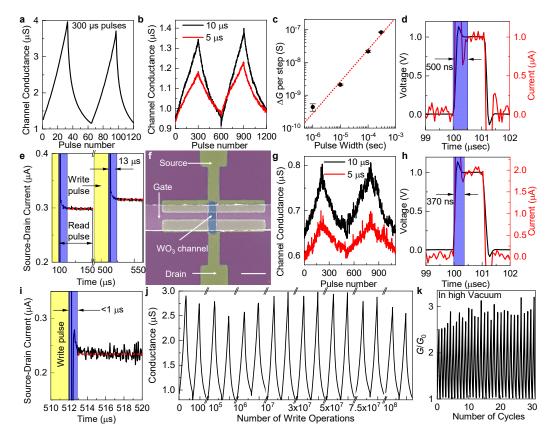


Figure 2. Speed, scaling, and endurance of CMOS-compatible, all-inorganic protonic ECRAMs. a, Programing of a micron-scale ( $L_{ch}$ =10  $\mu$ m, W=3  $\mu$ m) CMOS-compatible, all-inorganic protonic ECRAM with 300  $\mu$ sec gate-voltage pulses ( $\pm$ 4 V). b, Reproducible and highly symmetric programming with 10  $\mu$ sec (black) and 5  $\mu$ sec (red) write pulses (amplitude= $\pm$ 4 V). c, Double logarithmic scale plot showing the average  $\Delta G$  per weight-update step as a function of the pulse width with the same pulse amplitude of 4 V. Red dotted line represents the linear fitting to the data. Error bars represent the standard deviation. d, Measured transient variation of the sense current (red, right axis) during the ECRAM read operation performed with a voltage pulse (black, left axis) applied on the drain electrode and the source/gate grounded. The blue shading highlights the settling time  $t_{read}$ . e, Recovery waveform of the sense current of read post write pulses. The yellow shadings indicate the write pulses of 300  $\mu$ s, and the blue shadings highlight the settling time  $t_{read-after-write}$  required to reach steady states

whose averages are marked with red dashed lines. **f**, False-colored scanning-electron microscopy (SEM) micrograph of a nanometer-scale ( $L_{\rm ch}$ =150 nm, W=150 nm) CMOS-compatible, all-inorganic protonic ECRAM. Scale bar: 500 nm. **g**, Reproducible and symmetric switching characteristics of the scaled ( $L_{\rm ch}$ =W=150 nm) ECRAM modulated with 10 µsec (black) and 5 µsec (red) write pulses (amplitude=±4 V). **h**, Waveforms of the read voltage (black, left axis) and the sense current (red, right axis) of the scaled ECRAM to extract its  $t_{\rm read}$  as highlighted by the blue shading. **i**, Time-resolved source-drain current of the scaled ECRAM in a read performed right after a write pulse (yellow shadings), showing a faster settling time (blue shading) to steady state (red dashed line). **j**, Endurance test for 10<sup>8</sup> write-read pulses, showing no device degradation with the intermediate switching cycles plotted after 10<sup>5</sup>, 10<sup>6</sup>, 10<sup>7</sup>, 3×10<sup>7</sup>, 5×10<sup>7</sup>, 7.5×10<sup>7</sup>, and 10<sup>8</sup> pulses. **k**, Cycling the ECRAM passivated with HfO<sub>2</sub> between its high and low-conductance states in high vacuum of 10<sup>-6</sup> Torr.

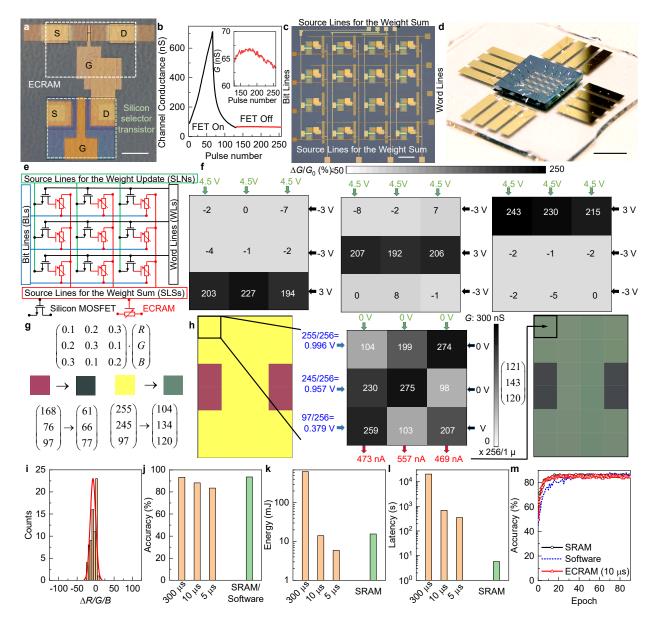


Figure 3. Monolithic integration of all-inorganic protonic ECRAMs with silicon transistors, and the operations of their pseudo-crossbar arrays as in-memory-computing accelerators. a, Optical image of a memory cell consisting of a ECRAM and a silicon MOSFET, which are fabricated on different layers. The white and green dashed lines serve as visual guide to mark the boundary of the ECRAM and the silicon selector transistor, respectively. Scale bar: 200 μm. b, Programing of the ECRAM cell by voltage-pulse trains composed of 64 potentiation followed by 64 depression steps, with the selector transistor turned either on (black) or off (red, inset). c-d, Optical

images of the fabricated ECRAM-cell array (part c, scale bar: 500 µm) and the bonded chip (part d, scale bar: 5 mm). e, Circuit diagram of the 3 by 3 ECRAM pseudocrossbar array. f, Parallel row-by-row programing of the ECRAM array, with the normalized conductance modulation ( $\Delta G/G$ ) indicated by grayscale. The voltages applied on the word lines (WLs) and the source lines for the weight update (SLNs) are shown in black and green, respectively. g, Color transformation performed by the dot product between the input RGB vector and the transformation matrix. transformation of a 4-pixel by 6-pixel image of letter I performed in parallel using the ECRAM array. Taken the top-left pixel as example, the encoded voltage inputs to the bit lines (BLs) are shown in blue, and the current outputs at the source lines for the weight sum (SLSs) measured in experiment are shown in red. The channel conductance G of each ECRAM after programing is illustrated in grayscale. Histogram showing the deviation of the modified RGB values calculated by the ECRAM array from those determined by software. j-l, Achievable accuracy (part j), energy consumption (part k), and latency (part I) for the ECRAM-CMOS hybrid in-memorycomputing accelerator to learn the classification of the MNIST dataset using different write-pulse widths (orange), as benchmarked against the SRAM-based digital accelerator or software (green). m, Simulated accuracy of ECRAM (red, operated with 10 usec gate-voltage pulses) and SRAM-based (black) accelerators to learn the classification of the CIFAR (Canadian Institute for Advanced Research)-10 dataset with the VGG (Visual Geometry Group)-8 network as a function of the training epochs performed. Results from learning in software (blue dotted line) are also displayed for comparison.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

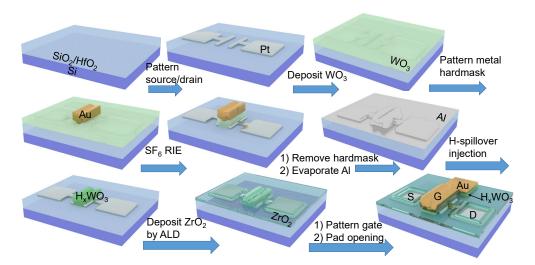
19

20

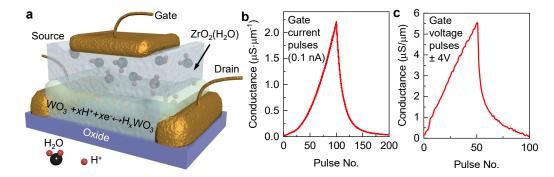
21

22

23

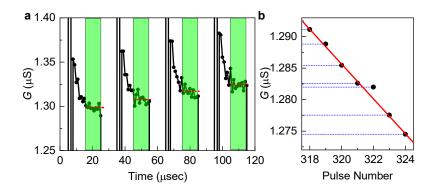


Extended Data Figure 1| Schematics illustrating the process flow to fabricate the CMOS-compatible, all-inorganic protonic ECRAM.

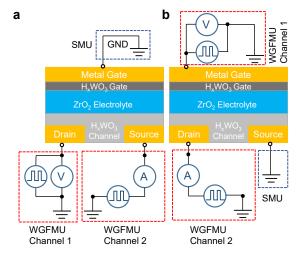


Extended Data Figure 2 All-inorganic protonic ECRAM with elemental metal gate.

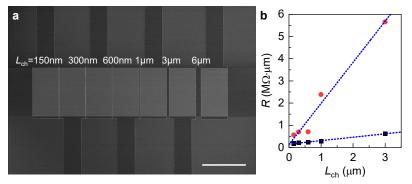
**a,** Schematic of the protonic ECRAM employing only Cr/Au as the gate electrode. **b,** Symmetric programing of the device channel conductance under gate-current pulses (100 potentiation then 100 depression with the gate-current pulse amplitude of 100 pA and width of 3 sec). **c,** Asymmetric programing of the same device but under gate-voltage pulses (50 potentiation then 50 depression with the gate-voltage pulse amplitude of ±4 V and width of 3 s), caused by the non-zero built-in open-circuit potential.



Extended Data Figure 3| Read noise and cycle-to-cycle variability of ECRAM operated with 10 μsec programming gate-voltage pulses. a, Zoom-in view of the programing characteristics showing the readout of four discrete conductance states with low read noise during the 10 μsec sensing period shaded in green, after 10 μsec settling time, out of the 300 conductance states shown in Fig. 2b. Red dotted lines serve as visual guide to mark the average channel conductance in each state. b, Zoom-in view of the read-out conductance (black circle) after the corresponding weight-update pulses. The red solid line is the expected conductance value obtained through fitting the depression curve. The small deviation of the experimental data to the expectation indicates small cycle-to-cycle variation. Blue dotted lines serve as visual guide to mark the measured conductance values on y axis.



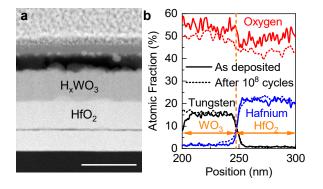
**Extended Data Figure 4| ECRAM pulse-measurement setups to extract the read transients. a,** Measurement configuration to extract *t*<sub>read</sub>. WGFMUs are connected to the ECRAM source and drain terminals. The gate is grounded. During measurement, the voltage pulse is applied on the drain side and the other channel of WGFMU reads the current signal. **b,** Measurement configuration to perform the microsecond pulse write operations and extract *t*<sub>read-after-write</sub>. WGFMUs are connected to the ECRAM gate and drain terminals. The source is grounded. During measurement, the voltage pulse is applied on the gate terminal for the weight-update and the other channel of WGFMU applies a small voltage pulse to reads the current flowing across the channel afterwards. SMU: source-measurement unit (Keysight B1517A). WGFMU: waveform generator/fast-measurement unit (Keysight B1530A). GND: ground.



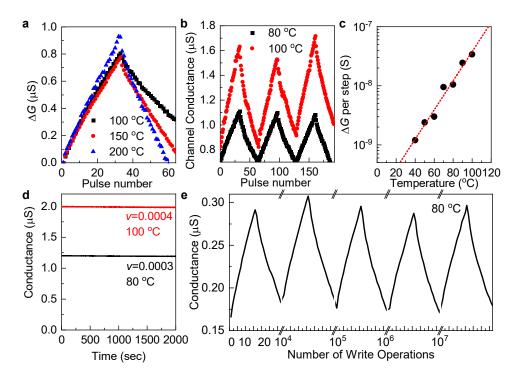
6 Extended Data Figure 5| Contact resistance of ECRAMs based on H<sub>x</sub>WO<sub>3</sub>. a, SEM

micrograph showing the transmission-line structure fabricated with  $L_{ch}$  varied from 6 µm down to 150 nm but identical W of 60 µm. Scale bar: 50 µm. **b**, Width normalized device resistance R as a function of  $L_{ch}$  for  $H_xWO_3$  with Pt contact switched between high (180  $M\Omega \cdot sq^{-1}$ ) and low (150  $k\Omega \cdot sq^{-1}$ ) sheet resistance corresponding to the operating dynamic range of ECRAM, showing a degradation of the device on/off ratio

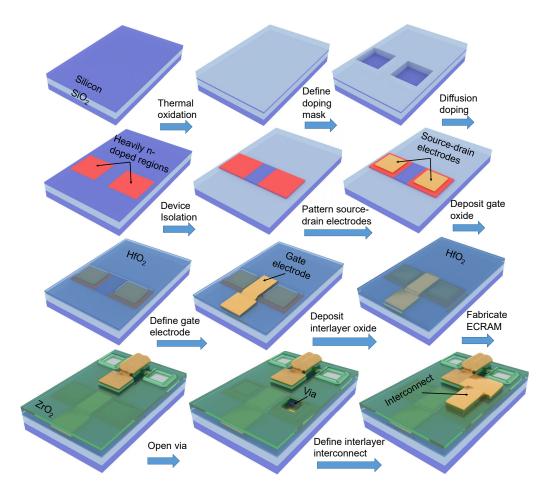
12 with the scaling of the  $L_{ch}$ .



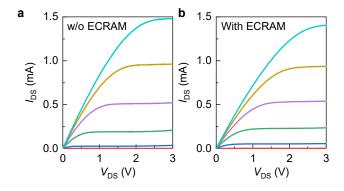
Extended Data Figure 6|  $H_xWO_3$  channel after the endurance test. a, STEM micrograph showing the gate stack of the ECRAM after operation with 100 million readwrite cycles. Scale bar: 100 nm. b, The depth profiles showing the atomic fractions of W (black), Hf (blue), and O (red) before (solid lines) and after (dotted lines) the  $10^8$  cycle endurance test, as measured by energy dispersive X-ray spectroscopy. The organ dashed line serves as a visual guide to mark the interface between  $WO_3$  and  $HfO_2$ .



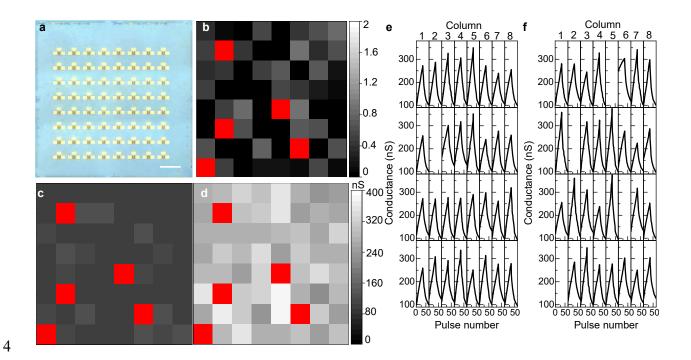
Extended Data Figure 7| Operation and reliability of all-inorganic protonic ECRAMs under elevated temperature. a, Programing of an all-inorganic protonic ECRAM with Al<sub>2</sub>O<sub>3</sub> passivation measured at the same temperature after being annealed at 100 °C (black), 150 °C (red), and 200 °C (blue), respectively. b, Programing an all-inorganic protonic ECRAM at 80 °C (black) and 100 °C (red).  $G_0$  and  $\Delta G$  increased due to the thermal excitation of additional carriers in the H<sub>x</sub>WO<sub>3</sub> channel and the increase of proton diffusivity, respectively. c, Logarithmic scale plot showing the average  $\Delta G$  per weight-update step as a function of the device operating temperature. Red dotted line represents the linear fitting to the data. d, Retention of the fully on-state conductance of the ECRAM cell at 80 °C (black) and 100 °C (red), and the corresponding drift coefficient ( $\nu$ ). e, ECRAM endurance test at 80 °C for 10<sup>7</sup> write-read pulses, showing no device degradation with the intermediate switching cycles plotted after 10<sup>4</sup>, 10<sup>5</sup>, 10<sup>6</sup>, and 10<sup>7</sup> pulses.



Extended Data Figure 8| Schematics illustrating the process flow to fabricate a 1-transistor-1-ECRAM memory cell, consisted of an all-inorganic protonic ECRAM and a silicon MOSFET in different layers monolithically integrated together on a SOI wafer substrate.



Extended Data Figure 9| Current-voltage characteristics of the silicon MOSFET selector measured before (part a) and after (part b) fabricating the protonic ECRAM layer on top with 40 nm HfO<sub>2</sub> as the interlayer dielectric.  $V_{\rm DS}$ : source-drain bias;  $I_{\rm DS}$ : source-drain current.



**Extended Data Figure 10| Device-to-device variation extracted from 8×8 ECRAM array. a**, Optical micrograph showing the completed 8 by 8 ECRAM array. Scale bar: 500 μm. **b-d**, Spatial mapping of the non-linearity (part **b**, standard deviation around 10%), minimum conductance (part **c**, standard deviation about 12%), and maximum conductance (part **d**, standard deviation about 13%) of the fabricated ECRAMs. Red indicated failed devices resulting from lithography or material defects. **e-f**, Collections of the programing characteristics of ECRAMs in row 1-4 (from top to bottom, part **e**) and row 5-8 (part **f**) of the array, respectively.

## References

2	45	Zhou, Y. et al. Enhanced transmittance modulation of ITO/NiOx/ZrO2:H/WO3/ITO
3		electrochromic devices. Ionics 22, 25-32 (2016).

- 4 46 Park, J. S. *et al.* Evidence of proton transport in atomic layer deposited Yttria-stabilized zirconia films. *Chem. Mater.* **22**, 5366-5370 (2010).
- stabilized zirconia films. *Chem. Mater.* **22**, 5366-5370 (2010).

  Kindratenko, V. *et al.* in *Practice and Experience in Advanced Research Computing (PEARC' 20).* 41–48 (Association for Computing Machinery).