

All-optical spiking neurosynaptic networks with self-learning capabilities

J. Feldmann¹, N. Youngblood², C.D. Wright³, H. Bhaskaran² and W.H.P. Pernice^{1}*

¹Institute of Physics, University of Muenster, Heisenbergstr. 11, 48149 Muenster, Germany

²Department of Materials, University of Oxford, Parks Road, OX1 3PH Oxford, UK

³Department of Engineering, University of Exeter, Exeter, EX4 QF, UK

*Correspondence to: wolfram.pernice@uni-muenster.de.

Software-implementations of brain-inspired computing underlie many important computational tasks, from image processing to speech recognition, artificial intelligence and deep learning applications. Yet, unlike real neural tissue, traditional computing architectures physically separate the core computing functions of memory and processing, making fast, efficient and low-energy computing difficult to achieve. To overcome such limitations, an attractive alternative is to design hardware that mimics neurons and synapses which, when connected in networks or neuromorphic systems, process information in a way more analogous to brains. Here we present an all-optical version of such a neurosynaptic system capable of supervised and unsupervised learning. We exploit wavelength division multiplexing techniques to implement a scalable circuit architecture for photonic neural networks, successfully demonstrating pattern recognition directly in the optical domain. Such photonic neurosynaptic networks promise access to the high speed and bandwidth inherent to optical systems, attractive for the direct processing of optical telecommunication and visual data..

25 **Introduction**

26 In our everyday life, artificial neural networks (ANNs) are already heavily active behind the
27 scenes, for example carrying out tasks such as face and speech recognition that are frequently
28 performed on our mobile phones¹. Thinking of more complex applications, such as medical
29 diagnostics² and autonomous driving^{1 3}, high-speed data-analysis will become even more
30 important in the future. However, fulfilling this demand for fast and efficient processing
31 using traditional computation techniques is problematic, due to speed and energy
32 inefficiencies⁴. Traditional computers are built following the von-Neumann architecture,
33 having two separate units for memory and processor and operating in a sequential way one
34 command at a time. Compared to the massively parallel signal processing of the brain, it
35 becomes clear why simulating a neural network in software on a machine based on the von-
36 Neumann architecture and limited by the transfer of data between memory and processor,
37 cannot be efficient⁵. A more radical approach – neuromorphic computing – seeks to overcome
38 the limitations of carrying out brain-like processing using conventional computers by
39 developing hardware mimics of the basic building blocks of biological brains, i.e. neurons and
40 synapses, and combining these into suitably-scaled networks and arrays. Such an approach
41 could, for example, enable the efficient processing and analysis of data in parallel directly on-
42 chip, so finding widespread utility in power-critical situations such as for mobile devices and
43 so-called “edge computing” applications⁶.

44 Recently, a number of different concepts for realizing hardware (i.e. neuromorphic)
45 implementations for artificial intelligence (AI) have been proposed in the electrical domain^{7 8}
46 but optical approaches are very much in their infancy⁹⁻¹³. A most promising candidate for
47 photonic neuromorphic computing is however that based around phase-change materials and
48 devices, since these have been shown to exhibit an intrinsic ability to provide in hardware the
49 basic integrate-and-fire functionality of neurons and the plastic weighting operation of

50 synapses^{14 15 16 17 18}. However, a fully optical, integrated and scalable neuromorphic
51 framework for implementing spiking neural networks using phase-change materials has—to
52 the best of our knowledge—not yet been demonstrated. In this work, therefore, we propose
53 and fabricate an all-optical spiking neuron circuit, with integrated all-optical synapses, and
54 demonstrate that such a system is capable of the prototypical AI task of pattern recognition.
55 Moreover, training/learning in the system is implemented in both a supervised and an
56 unsupervised way, both cases having a wide range of applications but relying on different
57 learning rules. In the first case, where training sets with pairs of known inputs and outputs are
58 present, supervised learning rules can be applied, such as the well-known backpropagation
59 algorithm^{19 20}. The second case requires unsupervised learning, meaning that the network
60 adapts on its own to specific repeating features and patterns that are unknown beforehand²¹.
61 We devise a scalable, layered architecture, based on wavelength division multiplexing
62 (WDM), for realizing such complex integrated photonic systems, presenting a photonic neural
63 network consisting of four neurons and sixty synapses (and 140 optical elements in total)
64 which is able to successfully recognize letters presented to it. By implementing an all-optical
65 spiking neural network on a nanophotonic chip, we provide a first step towards optical
66 neuromorphic systems, which benefit from the high bandwidth and fast signaling properties
67 that come with operating fully in the optical domain^{22 23}.

68 **Photonic implementation of an artificial neuron**

69 A sketch of our optical spiking neuron circuit, also incorporating all-optical synapses, and
70 how it is integrated on a nanophotonic platform is shown in Figure 1a-b. The neuron
71 represents a system comprising N input (pre-synaptic) neurons, one output (post-synaptic)
72 neuron and N interconnecting synapses. Each connection between the pre-synaptic neurons
73 and the post-synaptic neuron has a certain weight w_i . In the configuration of Figure 1a, optical
74 pulses from the pre-synaptic neurons are fed from the left into the connecting synapses,

75 thence to the post-synaptic neuron itself. The synapses are built of optical waveguides (see
76 Figure 1b) and weighting is achieved via phase-change material cells (here of area $3.6 \mu\text{m}^2$
77 and shown as red squares) integrated on top of the waveguides, which can modify the
78 propagating optical mode in a controlled manner. Phase-change materials (PCMs) are
79 commonly used in re-writable optical disc technologies, such as Blu-ray RE, and exhibit a
80 large contrast in the absorption of light between their amorphous and crystalline states
81 (phases) of matter^{24 25 26}. With the PCM-cell in the amorphous state, the synaptic waveguide
82 is highly transmissive, representing a strong connection between two neurons. In the
83 crystalline state, however, most of the light is absorbed leading to a weak connection. After
84 the input pulses have been weighted, they are combined into a single waveguide using
85 wavelength division multiplexing (WDM)^{27,28} and guided to the (output) spiking neuron
86 circuit. This is composed of a ring resonator with its own integrated PCM-cell that can be
87 switched (between crystalline and amorphous states) by the incoming combined pulses.
88 Switching the neuronal PCM-cell in turn changes the optical resonance condition of the ring
89 and its propagation loss. When the neuronal PCM-cell is in the crystalline state, a suitable
90 probe pulse sent along the ‘output’ waveguide couples strongly into the ring resonator and so
91 no output pulse (spike) will be observed. However, if the instantaneous combined power of
92 the weighted input pulses from the pre-synaptic neurons is high enough to switch the neuronal
93 PCM-cell to its amorphous state, the probe pulse is no longer on resonance with the ring and
94 will be transmitted past the ring, so generating an output neural spike. As the switching of the
95 PCM-cell only occurs above a certain threshold power, the neuron only generates an output
96 pulse (spike) if the weighted sum of the input power exceeds this threshold. Thus, the system
97 naturally emulates the basic integrate-and-fire functionality of a biological neuron, with the
98 distinction that the artificial neuron integrates over the optical power at a fixed time, as
99 opposed to its biological counterpart that integrates incoming pulses over time. This artificial
100 neuron, shown in photonic circuit form in Figure 1c, serves as a building block in layered

101 photonic networks (described later) suited to the scalable implementation of neurosynaptic
102 systems.

103 The neurosynaptic system described so far in the above provides the basic structure needed
104 for supervised learning tasks, where the weights of the inputs are set by an external
105 supervisor. However, in order to also be capable of unsupervised learning, we add a feedback
106 waveguide channelling parts of the neuron's output spike back to the synaptic PCM-cells. In
107 this way (and as described in more detail later), the connections from all inputs that
108 contributed to a particular output spike will be enhanced, while those that did not contribute
109 will be weakened – or in Hebbian terms, “neurons that fire together, wire together.”

110 Figure 1d shows an optical micrograph of the actual implementation (x3) of a single-neuron
111 neurosynaptic system, fabricated via electron-beam lithography on a silicon-nitride on silicon-
112 oxide platform. Several input waveguides each with a synaptic PCM-cell on top (indicated by
113 red ovals) are fed to the upper waveguide using small ring resonators as a simple multiplexing
114 device (the ring resonators have four different radii increasing linearly from 40 μm to 55 μm ,
115 an optical Q of around 10000 and provide insertion loss of 1.5 dB - see supplementary section
116 S5 for more details. This upper waveguide then leads the light to the neuronal (large) ring
117 resonator (radius 60 μm) with its own integrated PCM-cell (area 9 μm^2). Probe pulses sent to
118 the waveguide lying below the neuronal ring resonator either couple to it or generate an
119 output spike depending, as described previously, on whether the neuronal PCM-cell is in the
120 crystalline or amorphous state. The light is coupled onto and off the chip using grating
121 couplers which provide access to multiple optical fibres in the measurement setup.

122 **Optical performance of a single neuron device**

123 Figure 2a shows in detail the photonic signal processing and optical operation of a single
124 neuron. Each input pulse is sent on a different wavelength λ_i and firstly partly absorbed

125 (weighted) by the relevant synaptic PCM-cell. After weighting, the individual input
126 waveguides are combined with a multiplexer to a single waveguide, summing up the input
127 powers. If this power is high enough to switch the neuronal PCM-cell of the large ring
128 resonator (see Figure 2b), an output spike is generated and in the case of unsupervised
129 learning the synaptic weights are adjusted using the feedback loop. In this specific
130 device/example an output pulse is generated if the summed power here exceeds 430 pJ (see
131 section 9 in the supplementary materials). In Figure 2b a more detailed scanning electron
132 micrograph of the ring resonator used to deliver the spiking neuron function can be seen. The
133 neuronal PCM-cell used to tune the resonance condition is deposited on top of a waveguide
134 crossing specially designed for low optical losses (0.23 dB²⁹). The second waveguide
135 crossing (without a PCM-cell) is only used for testing purposes and offers the ability of a bias
136 input.

137 In order to characterise the integrate-and-fire type response (or activation function) of our
138 neuron circuit, several transmission spectra for the resonator were obtained after sending
139 pulses of different energy to the neuronal PCM-cell via the crossing waveguide (Figure 2c),
140 see Methods section. By plotting the transmission of the ring resonator after different
141 excitation pulses vs. the pulse energy at a fixed wavelength, an activation function as shown
142 in Figure 2d is obtained. Depending on which wavelength is chosen, the contrast and
143 maximum transmission level of the output function can be adjusted. Figure 2d shows the
144 activation at 1553.4 nm (i.e. at the dashed line in Figure 2c), representing the operation with
145 the highest contrast of 9 dB between the output states. It can clearly be seen that only above a
146 threshold energy of 60% of the maximum pulse energy is a significant output generated. This
147 non-linear response resembling the rectified linear unit (ReLU) function is crucial for neural
148 activation functions, since it projects complex input data to higher dimensions enabling linear
149 separation by the output neurons³⁰.

151 **Supervised and un-supervised learning**

152 Having found the working point for our all-optical spiking neuron, supervised learning tests
153 are now carried out. In this case, the synaptic weights of the network are set by an external
154 supervisor (as for example done in software-based ANNs using the backpropagation
155 algorithm). Here, a training set of data consisting of pairs of input patterns and the expected
156 output is shown to the network. Depending on the deviation between the expected and actual
157 output, the synaptic weights are adjusted in an optimization process until the solution is best
158 approximated and the network is trained.

159 The experimental neural network used is composed of two single neurosynaptic systems (of
160 the type shown in Figure 1d) each consisting of four input (pre-synaptic) neurons connected
161 to one output (post synaptic) neuron by four PCM synapses. The weights of the first neuron
162 were set to the pattern “1010” meaning that the first and third PCM synapse were in a high
163 transmission state (contributing significantly to the activation energy) and the second and
164 fourth were in a low transmission state (contributing less to the activation energy). The second
165 neuron was trained in the same way to the pattern “1100”. In Figure 3a and b the post-
166 synaptic neuron output is plotted as a function of the input pattern. It can clearly be seen that
167 in both cases (i.e. for both input patterns) the neurons were trained successfully and, based on
168 the neuron’s output, it can be easily concluded which pattern was presented to the network.
169 Using only two output neurons on the same set of input neurons, our all-optical neuromorphic
170 system can already solve simple image recognition tasks. By increasing the number of inputs
171 per neuron and the number of neurons, more complex images can be processed and more
172 difficult tasks, such as letter (or digit) recognition or language identification can be solved
173 using this same basic approach, as we show later (in experiment for letter recognition, by
174 simulation for digit recognition and language identification).

175 Above we illustrated an example of supervised learning. This learning technique is feasible
176 for many tasks but has the limitation that a training set with tuples of input patterns and
177 expected outputs must be present. If the output is unknown, for example if an unknown but
178 repeating pattern must be found from a data stream, then supervised learning is not applicable
179 and unsupervised learning procedures are necessary. In an unsupervised approach, the
180 network updates its weights on its own and in this way adapts to a certain pattern over time,
181 without the need for an external supervisor.

182 In order to do this, an update rule needs to be defined. A common concept in unsupervised
183 learning is spike-timing-dependent plasticity (STDP) following Hebb's postulate³². Here the
184 change in the synaptic weight after an output spike depends on the relative timing between the
185 input and the output spike of a neuron (i.e. the timing difference between pre- and post-
186 synaptic neuron firings). If an input signal arrives right before an output spike was generated,
187 that input signal is likely to have contributed to reaching the firing threshold and the
188 corresponding weight will be increased. If the input pulse arrives after the output spike
189 occurred, the synaptic weight will be decreased. The amount of potentiation (weight increase)
190 or depression (weight decrease) is a function of the time difference between input and output
191 spike, as described by Bi and Poo³³.

192 A similar but simplified learning rule is applied in our all-optical neuron approach. As the
193 timing between incoming pulses and output pulses in our case is fixed (because we operate the
194 neuron in a clocked way, one complete input pattern per time step – see Methods section),
195 there is no varying time delay between input and output events. We therefore adopt a
196 simplified learning rule, increasing the synaptic weights of all inputs that contributed to a
197 spike generation, and decrease the weights of all that did not. Experimentally we obtain this
198 behaviour by overlapping (in time) the output pulses with the input pulses (see Methods
199 section).

200 Figures 3d and e show the development over time of the four synaptic weights during
201 unsupervised learning by a single neuron. Initially all the PCM-synapses are in the amorphous
202 (high transmittance) state. When the input pattern ‘0110’ is repeated, the neuron adapts to it
203 over time, until the neuron has finally learned this pattern without any intervention from an
204 external supervisor. The neuron is now specialized to recognize this particular pattern. From
205 Figure 3d it is clear that the weights w_2 and w_3 , corresponding to the inputs three and four,
206 stay almost constant over time, as the overlapping input and feedback pulses preserve their
207 amorphous state. In contrast, weights w_1 and w_4 are depressed stepwise with each epoch.

208 **A scalable architecture for photonic artificial neural networks**

209 Having successfully demonstrated a single-neuron neurosynaptic system as a fundamental
210 building block for photonic neural networks, a way of connecting these artificial neurons into
211 larger networks is now developed. An architecture exploiting individually addressable,
212 interlinked, photonic layers is thus implemented, as shown schematically in Figure 4a. The
213 whole network consists of an input and an output layer which are optically connected via N
214 hidden layers. Each hidden layer takes the output of the previous layer as an input and passes
215 its outputs to the next layer. The input layer is the optical interface to the real world, taking
216 the data to be processed and distributing it to the next level in the network.

217 A single layer of the network consists of a collector, a distributor and its neurosynaptic
218 elements. The collector gathers all the outputs from the previous layer, which are then equally
219 distributed to the N neurons within the layer (fully connected network) by the distributor. The
220 photonic neurons themselves operate as described in detail before: a phase-change synapse
221 weights the inputs and a WDM multiplexer builds the sum, which is passed to the activation
222 unit that decides if a neuronal output pulse is transmitted. In this architecture, each layer is
223 addressed optically with its own waveguide for generating the probe signal. Therefore, the

224 optical power in the layer is not limited by the transmitted optical response from a previous
225 layer.

226 Figure 4b describes how the constituent parts of a layer translate into the actual photonic
227 circuit. The outputs from a previous layer are multiplexed onto a single waveguide using ring
228 resonators (thus building the collector). This signal is then equally distributed to the neurons
229 within this layer, again using ring resonators for demultiplexing (thus building the
230 distributor). By choosing the gap between the feeding waveguide and ring resonator, the
231 coupling efficiency can be tuned (see supplementary figure S8.) Following the formula for the
232 coupling efficiency $c_{\text{eff},i}=1/(N'+1-i)$ with N' neurons on the layer and the neuron position i
233 then, for example, in a layer with four neurons this means, that $1/(4+1-1)=1/4$ of the light is
234 coupled to the first neuron, $1/3$ of the remaining light is passed to the second neuron, $1/2$ to
235 the third and the residual light to the fourth neuron. The circuit-diagram of the actual photonic
236 neuron, the neuro-synaptic system, was shown in Figure 1c and is the same as used in the
237 experiments described above. The output pulses of a layer can then be connected to the
238 collector of the next layer.

239 We note that using the above approach no waveguide crossings are needed for distributing the
240 signal to the neurons, thus preventing crosstalk and losses. Because the output pulses are
241 generated for each layer individually, there is also no accumulation of errors and signal
242 contamination over subsequent layers. This fact also simplifies the timing of the network as
243 each layer can be processed step by step: First, the input pulses are sent, and the activation
244 units are switched where appropriate. Second, the output probe pulses are sent and transmitted
245 to the next layer (if the threshold for neuron switching was reached). In a final step, the PCM-
246 cells on the rings have to be returned to their initial state.

247

248 **Realization of a single-layer neurosynaptic system**

249 Figure 5 shows the experimental implementation of a full layer of the proposed neural
250 network design consisting of four neurons with 15 synapses each. The full device is composed
251 of more than 140 optical elements; optical micrographs of the photonic circuit are presented
252 in section 2 of the supplementary materials. This network is capable, by way of example, to
253 differentiate between four 15-pixel images, here representing the four letters A, B, C and D.
254 In this system, the neurons are optically fed via an integrated WDM distributor with 15 ring
255 resonators per neuron, while the collector is implemented off-chip using fiber-based WDM
256 components (see Methods section). As desired, all four output neurons are only activated
257 when the learned pattern is shown (Figure 5b): Neuron 1 only fires if pattern ‘A’ is shown,
258 neuron 2 only reacts to pattern ‘B’ and so forth. Thus, the network is able to successfully
259 classify the four 15-pixel images. (A more complex exemplar task of language recognition,
260 using a larger network with the same architecture, is discussed in the supplementary
261 materials). We note, that in the all-optical implementation of this architecture all artificial
262 neurons need to be recrystallized after each spiking event. Therefore the number of operation
263 cycles is eventually limited by the endurance of the PCM-cells. While individual PCM
264 devices in endurance experiments have already shown 10^{12} switching cycles¹⁶, further
265 improvements in material design and device engineering are needed for high-speed and long-
266 term switching operation.

267 Integrated phase-change photonic networks, designed and implemented as described above,
268 are capable of simple pattern recognition tasks and can adapt to specific patterns. When
269 operated with a waveguide feedback loop, they are capable of learning without an operator
270 needed, and can do this in a non-volatile fashion using phase-change materials. The large
271 contrast in absorption of light between their amorphous and crystalline state of matter makes
272 phase-change materials an attractive and simple solution to be integrated as synaptic

273 weighting mechanism. Compared to conventional computers that can only simulate the
274 parallelism of neural networks, our all-optical neurons are intrinsically suited for mimicking
275 biological neural networks. Compared to speeds of biological neural networks (~
276 milliseconds) our proposed neurons could operate several orders of magnitude faster, giving
277 rise to substantial potential in dealing with large amounts of data in a short amount of time.
278 Working exclusively in the optical domain, the spiking neurosynaptic network benefits from
279 high-bandwidth and fast data transfer rates intrinsic to light. Moreover, using a layered circuit
280 architecture, we present a pathway to scaling our network to more complex systems which
281 could be realized with foundry processing. This way also the off-chip components (used here
282 for experimental expediency only) such as laser sources, optical amplifiers and modulators,
283 can be integrated into a full system. Our integrated and novel design combines, via
284 wavelength division multiplexing techniques, the outputs of multiple phase-change synapses
285 to excite layered spiking phase-change neurons and holds promise for realizing all-optical
286 neural networks capable of addressing the upcoming challenges of big data and deep learning.

287

288 **References**

- 289 1. Lane, N. D. *et al.* Squeezing Deep Learning into Mobile and Embedded Devices. *IEEE*
290 *Pervasive Comput.* **16**, 82–88 (2017).
- 291 2. Amato, F. *et al.* Artificial neural networks in medical diagnosis. *J. Appl. Biomed.* **11**,
292 47–58 (2013).
- 293 3. Nawrocki, R. A., Voyles, R. M. & Shaheen, S. E. A Mini Review of Neuromorphic
294 Architectures and Implementations. *IEEE Trans. Electron Devices* **63**, 3819–3829
295 (2016).
- 296 4. Preissl, R. *et al.* Compass: A scalable simulator for an architecture for cognitive
297 computing. *Int. Conf. High Perform. Comput. Networking, Storage Anal. SC* (2012).
- 298 5. Neumann, J. von. *The Computer and The Brain.* (Yale University Press, 1958).
- 299 6. Wu, H., Yao, P., Gao, B. & Qian, H. Multiplication on the edge. *Nat. Electron.* **1**, 8–9
300 (2018).
- 301 7. Furber, S. Bio-Inspired Massively-Parallel Computation. *Parallel Comput. Road to*
302 *Exascale* **27**, 3–10 (2016).
- 303 8. Schmitt, S. *et al.* Neuromorphic hardware in the loop: Training a deep spiking network
304 on the BrainScaleS wafer-scale system. in *Proceedings of the International Joint*
305 *Conference on Neural Networks* 2227–2234 (2017).
- 306 9. Shen, Y. *et al.* Deep learning with coherent nanophotonic circuits. *Nat. Photonics* **11**,
307 441–446 (2017).
- 308 10. Vinckier, Q. *et al.* High-performance photonic reservoir computer based on a
309 coherently driven passive cavity. *Optica* **2**, 438 (2015).

- 310 11. Brunner, D., Soriano, M. C., Mirasso, C. R. & Fischer, I. Parallel photonic information
311 processing at gigabyte per second data rates using transient states. *Nat. Commun.* **4**,
312 1364–1367 (2013).
- 313 12. Ferreira De Lima, T., Shastri, B. J., Tait, A. N., Nahmias, M. A. & Prucnal, P. R.
314 Progress in neuromorphic photonics. *Nanophotonics* **6**, 577–599 (2017).
- 315 13. Vandoorne, K. *et al.* Experimental demonstration of reservoir computing on a silicon
316 photonics chip. *Nat. Commun.* **5**, 1–6 (2014).
- 317 14. Cheng, Z., Ríos, C., Pernice, W. H. P., Wright, C. D. & Bhaskaran, H. On-chip
318 photonic synapse. *Sci. Adv.* **2**, 1–7 (2017).
- 319 15. Kim, S. *et al.* NVM neuromorphic core with 64k-cell (256-by-256) phase change
320 memory synaptic array with on-chip neuron circuits for continuous in-situ learning. in
321 *International Electron Devices Meeting (IEDM), 2015 IEEE International* 17.1.1-
322 17.1.4 (2015).
- 323 16. Kuzum, D., Jeyasingh, R. G. D., Lee, B. & Wong, H. P. Nanoelectronic Programmable
324 Synapses Based on Phase Change Materials for Brain-Inspired Computing. *Nano Lett.*
325 **12**, 2179–2186 (2012).
- 326 17. Wright, C. D., Liu, Y., Kohary, K. I., Aziz, M. M. & Hicken, R. J. Arithmetic and
327 biologically-inspired computing using phase-change materials. *Adv. Mater.* **23**, 3408–
328 3413 (2011).
- 329 18. Pantazi, A., Woźniak, S., Tuma, T. & Eleftheriou, E. All-memristive neuromorphic
330 computing with level-tuned neurons. *Nanotechnology* **27**, 355205 (2016).
- 331 19. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-Based Learning Applied to
332 Document Recognition. *Proc. IEEE* **86**, 2278–2323 (1998).

- 333 20. Burr, G. W. *et al.* Experimental Demonstration and Tolerancing of a Large-Scale
334 Neural Network (165 000 Synapses) Using Phase-Change Memory as the Synaptic
335 Weight Element. *IEEE Trans. Electron Devices* **62**, 3498–3507 (2015).
- 336 21. Le, Q. V. *et al.* Building high-level features using large scale unsupervised learning. in
337 *29th International Conference on Machine Learning* (2012).
- 338 22. Alduino, A. & Paniccia, M. Interconnects: Wiring electronics with light. *Nat.*
339 *Photonics* **1**, 153–155 (2007).
- 340 23. Sun, C. *et al.* Single-chip microprocessor that communicates directly using light.
341 *Nature* **528**, 534–538 (2015).
- 342 24. Wuttig, M. & Yamada, N. Phase-change materials for rewriteable data storage. *Nat.*
343 *Mater.* **6**, 824–832 (2007).
- 344 25. Raoux, S., Xiong, F., Wuttig, M. & Pop, E. Phase change materials and phase change
345 memory. *MRS Bull.* **39**, 703–710 (2014).
- 346 26. Burr, G. W. *et al.* Recent Progress in Phase-Change Memory Technology. *IEEE J.*
347 *Emerg. Sel. Top. Circuits Syst.* **6**, 146–162 (2016).
- 348 27. Tait, A. N., Nahmias, M. A., Shastri, B. J. & Prucnal, P. R. Broadcast and Weight: An
349 Integrated Network For Scalable Phtonic Spike Processing. **32**, 3427–3439 (2014).
- 350 28. Tait, A. N. *et al.* Neuromorphic Silicon Photonic Networks. *Sci. Rep.* 1–10 (2016).
351 doi:10.1038/s41598-017-07754-z
- 352 29. Feldmann, J. *et al.* Calculating with light using a chip-scale all-optical abacus. *Nat.*
353 *Commun.* **8**, (2017).
- 354 30. Hahnloser, R. H. R., Sarpeshkar, R., Mahowald, M. A., Douglas, R. J. & Seung, H. S.

- 355 Digital selection and analogue amplification coexist in a cortex- inspired silicon circuit.
356 *Nature* **405**, 947–951 (2000).
- 357 31. Rios, C. *et al.* Integrated all-photonics non-volatile multi-level memory. *Nat. Photonics*
358 **9**, 725–732 (2015).
- 359 32. Hebb, D. *The Organization of Behaviour*. (Wiley & Sons, 1949).
- 360 33. Bi, G. & Poo, M. Synaptic Modification by correlated activity: Hebb ’ s Postulate
361 Revisited. *Annu. Rev. Neurosci.* **24**, 139–166 (2001).
- 362 34 D. Goldhahn, T. Eckart & U. Quasthoff: Building Large Monolingual Dictionaries at
363 the Leipzig Corpora Collection: From 100 to 200 Languages. In: Proceedings of the 8th
364 International Language Resources and Evaluation (LREC'12), 2012.
- 365
- 366

367 **Acknowledgements**

368 This research was supported by EPSRC via grants EP/J018694/1, EP/M015173/1, and
369 EP/M015130/1 in the United Kingdom and the Deutsche Forschungsgemeinschaft (DFG)
370 grant PE 1832/5-1 in Germany. WHPP gratefully acknowledges support by the European
371 Research Council through grant 724707. We acknowledge funding from the European
372 Union's Horizon 2020 research and innovation programme under grant agreement No.
373 780848 (Fun-COMP).

374

375 **Author contribution**

376 WHPP, HB and CDW conceived the experiment. JF fabricated the devices with assistance
377 from NY. NY performed the deposition of the GST material. JF implemented the
378 measurement setup and carried out the measurements with help from NY. All authors
379 discussed the data and wrote the manuscript together.

380

381

382 **Competing interests**

383 The authors declare no competing interests.

384

385 **Figure 1. All-optical spiking neuronal circuits. a-b)** Schematic of the network realized in
386 this work consisting of several pre-synaptic input neurons and one post-synaptic output
387 neuron connected via PCM-synapses. The input spikes are weighted using PCM-cells and
388 summed up using a WDM multiplexer. If the integrated power of the postsynaptic spikes
389 surpasses a certain threshold, the PCM-cell on the ring resonator switches and an output pulse
390 (neuronal spike) is generated. **c)** Photonic circuit diagram of an integrated optical neuron with
391 symbol block shown in the inset (top right). Several of these blocks can be connected to larger
392 networks using the wavelengths inputs and outputs as described in more detail in Figure 5 **d)**
393 Optical micrograph of three fabricated neurons (B5, D1 and D2) showing four input ports.
394 The four small ring resonators on the left are used to couple light of different wavelengths
395 from the inputs to a single waveguide, which then leads to the phase-change material cell at
396 the crossing point with the large ring. The triangular structures on the bottom are grating
397 couplers used to couple light onto and off the chip.

398 **Figure 2. Spike generation and operation of the artificial neuron. a)** Schematic of the
399 photonic implementation of a phase-change neuron circuit. Light of different wavelength is
400 weighted by phase-change elements w_1 - w_4 and summed up by a multiplexer to a single
401 waveguide. If this activation energy surpasses a threshold, an output pulse is generated, and
402 the weights are updated. **b)** Scanning electron micrograph of a ring resonator used to
403 implement the activation function. By switching the PCM-cell on top of the waveguide
404 crossing, the resonance condition of the resonator can be tuned. The waveguide on the bottom
405 of the ring is used to probe the resonance and generate an output pulse. **c)** Transmission
406 measurement of the device in **b)** and its dependence of pulse energy. The resonance shifts
407 towards shorter wavelength with increasing pulse energy send to the PCM-cell on the ring. At
408 the same time the transmission increases because of reduced absorption in the PCM-cell and
409 thus changes the coupling between ring and waveguide. **d)** Normalized transmission to the

410 output at a fixed wavelength (dashed line in c)) showing the activation function used to define
411 the firing threshold of the neuron.

412 **Figure 3. Supervised and unsupervised learning with phase-change all-optical neurons.**

413 a) and b) show the neuron output of two individual neurons when presented with different
414 input patterns. Neuron one learned to recognize pattern ‘1010’, while neuron two generates an
415 output signal when ‘1100’ is shown. In this example the eight weights of the neural network
416 were set by an external supervisor. c) Schematic illustrating the unsupervised learning
417 mechanism in an all-optical neuron. If an output spike is generated, the synaptic-weights
418 where input and feedback pulses overlap in time are potentiated, while the weights that are
419 only hit by the single feedback pulse are depressed. d) Change of the four synaptic weights
420 over time when the pattern ‘0110’ is repeatedly shown starting from fully amorphous (high
421 transmitting) weights. The weights where input- and feedback pulse overlap stay almost
422 constant over several epochs. The other weights where only the feedback pulse is shown
423 decrease continuously. e) Development of the weights over time, clarifying that the
424 information of the pattern is encoded in the weights.

425 **Figure 4. Scaling architecture for all-optical neural networks.**

426 a) The general neural network is composed of an input layer, an output layer and several hidden layers. Each of
427 these layers consist of a collector gathering the information from the previous layer, a
428 distributor that equally splits the signal to individual neurons and the neuronal and synaptic
429 elements of the layer itself. Each neuron has a weighting unit and a multiplexer to calculate
430 the weighted sum of the inputs. The sum is then fed to an activation unit which decides if an
431 output pulse is generated. b) Photonic implementation of a single layer from the network. The
432 collector unites the optical pulses from the previous layer using a WDM multiplexer. A
433 distributor made from the same rings as the collector but with adjusted coupling efficiency
434 equally distributes the input signal to the PCM synapses of each neuron. The letters “P”, “W”

435 and “R” denote the input ports used to probe the output, set the weights and return the
436 neuronal PCM to its initial state.

437 **Figure 5 Experimental realisation of a single layer spiking neural network.** a) The device
438 consists of four photonic neurons, each with 15 synapses. Each synapse corresponds to a pixel
439 in a 3x5 image (see b)) and is encoded in the wavelengths corresponding to the ring
440 multiplexers (see numbering in b)). The full device comprises an integrated photonic circuit
441 built up from 140 optical components. b) The change in output spike intensity is shown for
442 the four trained patterns illustrated on the right-hand side. The neural network successfully
443 recognizes the four patterns as each neuron only responds (spikes) to one of the patterns. The
444 error bars denote the standard deviation for n=5.

445

446

447 **Methods**

448 **Device fabrication**

449 The nanophotonic circuits used in this work are realized using electron-beam lithography
450 (EBL) with a 100-kV system (Raith EBPG 5150). In a first step, opening windows for lift-off
451 processing of alignment markers made from gold are exposed in the positive tone resist
452 Polymethylmethacrylat (PMMA) on a silicon wafer (Rogue Valley Microdevices) with a 3300
453 nm silicon oxide and 344 nm silicon nitride layer on top. After development in 1:3
454 MIBK:Isopropanol for 2 min, a stack of 5 nm chromium, 120 nm gold and 5 nm chromium
455 again are evaporated via electron-beam physical vapour deposition (PVD). The lift-off step to
456 remove the PMMA is performed in acetone, leaving the gold markers for the alignment in the
457 next EBL-steps.

458 In the second lithography step the photonic structures are defined. TI Prime is used as an
459 adhesion agent for the negative-tone ebeam resist maN 2403. The photonic circuitry is
460 developed in MF-319 for 60 s and afterwards placed on a hotplate at 105°C for two minutes
461 of reflow processing to reduce surface roughness. By reactive ion etching in a CHF₃/O₂
462 plasma, the resist mask is transferred into the sample till the silicon nitride is fully etched. The
463 remaining resist on the structures is removed in an oxygen plasma for 10 minutes.

464 The last EBL step consists of writing windows for the deposition of the phase-change material
465 and is executed in the same way as defining the marker windows. After development, 10 nm
466 of the phase-change material GST are sputter-deposited and covered by a 10 nm film of
467 indium tin oxide (ITO) to prevent oxidation of the GST. The GST and ITO capping layers
468 were deposited using RF sputtering with an argon plasma (5 mtorr working pressure, 15 sccm
469 Ar, 30 W RF power, and base pressure of 2×10^{-6} Torr). Finally, the GST is crystallized on a

470 hot plate for about 10 minutes at 210°C. The photonic circuits are composed of single mode
471 waveguides at 1550 nm with a width of 1.2 μm .

472 **Measurement setup**

473 The experimental setup used to operate the all-optical neurons comprises pattern generation
474 and read-out of the individual weights, as sketched in the Supplementary materials. The
475 optical read-out is achieved via transmission measurement using a continuous wave laser
476 (Santec, TSL 510) and four low-noise photodetectors D1-D4 (New Focus, Model 2011) that
477 are monitored on a computer. In order to couple light efficiently onto the chip, an optical fiber
478 array is aligned with respect to the on-chip grating couplers to provide multi-port input. For
479 optimal coupling efficiency to the chip the polarization is optimized with a set of polarization
480 controllers.

481 Pattern generation as input for the on-chip neuron is accomplished with four cw-lasers set to
482 different wavelength matching the on-chip multiplexer. Off-chip the four light paths are
483 combined using a fiber multiplexer and desired optical pulses are created using an electro-
484 optical modulator (EOM) and a computer controlled electrical pulse generator (Agilent, HP
485 8131A). After amplifying the pulses with an erbium-doped fiber amplifier (Pritel, EDFA) the
486 pulses are de-multiplexed again and guided to the on-chip device. Using circulators pump and
487 probe light are counter-propagating through the device enabling efficient separation of the
488 beam paths. Arbitrary input patterns are then selected by switching on and off the shutters of
489 the pump lasers.

490 Similar to the input pulses, the output pulses are also created from a cw-laser in combination
491 with an EOM and amplified by an EDFA. A small portion of the output is measured with
492 detector D0, while the remaining light is amplified and send to the ports F1-F4 as a feedback
493 pulse for weight adjustments. Turning off the feedback amplifier puts the device in a

494 supervised learning mode. The setup used for the four neuron-network is shown in
495 supplementary figure S2.

496 **Activation unit**

497 In order to obtain the ReLU function shown in Figure 2d pulses with different energies up to
498 approximately 700 pJ have been sent to the ring resonator resembling the activation unit.
499 Because of changes introduced in the structure of the PCM on the waveguide crossing, the
500 transmission spectrum undergoes a significant change between the initial crystalline state (0%
501 pulse energy) and the final amorphous state (maximum pulse energy). This is partially
502 explained by the resonance wavelength (at which the light is coupled into the ring and
503 therefore the transmission is minimal) slightly shifting to shorter wavelength with increasing
504 amount of amorphization (which causes a change in the real part of the refractive index of the
505 phase-change material, leading to a slightly shifted resonance condition (see Supplementary
506 Materials)). However, we also observe that the minimum transmission after amorphizing the
507 PCM-cell is much higher. This is a combined effect of the change in the imaginary part of the
508 refractive index (absorption) and a change in the extinction ratio of the resonator. The lower
509 absorption in the amorphous phase obviously leads to higher transmission but, equally
510 important, changing the loss per round trip in the ring affects the coupling between ring and
511 waveguide and therefore alters the extinction ratio.

512 **Estimation of the energy balance**

513 The maximum pulse energy used to switch the PCM-cell in this experiment was 710 pJ
514 employing optical pulses of 200 ns width. For setting the weights similar energies are used. In
515 the experiments we operate our neurons with relatively long optical pulses in order to
516 implement two-pulse switching²⁹ which relies on overlapping pulses in the time-domain. As
517 the output pulse was amplified off-chip and the optical path was therefore relatively long
518 (compared to on-chip waveguides), longer pulses had to be applied to ensure the overlap. We

519 note, however, that lower switching energies can be employed by moving towards ps
520 pulses^{29,31}. It is also important to note that this operation scheme does not require continuous
521 energy input for maintaining the state of the PCM weights due to their non-volatile response.
522 Therefore, the energy budget per neuron to perform one operation is given by the switching
523 energy for the ring resonator plus the energy required to return the PCM element to its
524 original fully-crystalline state.

525 **Update of weights in unsupervised learning**

526 Via a feedback waveguide, the neuron's output spike is guided back to the synaptic PCM-
527 cells (see Figure 3c). If the input pulses are now long enough in time such that they overlap
528 with the feedback pulse at the waveguide crossing where the synaptic PCM-cell is located,
529 then the overlapping pulses have enough energy to switch the synapse into its low-absorbing
530 amorphous state (weights w_2 and w_3 in Figure 3c). A feedback pulse that encounters a synapse
531 without an input pulse will partly crystallize the corresponding PCM-cell because of the lower
532 pulse energy, and therefore decrease the weight (weights w_1 and w_4). Due to the properties of
533 the phase-change material used, the PCM-cells can only be amorphized in a single step³¹,
534 meaning that if the neuron fires, all contributing inputs will always be potentiated completely.
535 Opposite to that, full crystallisation can be achieved in several steps²⁹ and the weights can
536 correspondingly be decreased stepwise. Successful unsupervised learning can be
537 accomplished using such weight update rules, as we show experimentally in Figure 3 for
538 small-scale systems, and in the supplementary information (section 10) by simulation for
539 larger-scale systems.

540 **Image encoding and implementation of the WDM distributor and collector**

541 To feed a certain pattern to the neural network, it has to be encoded in optical pulse patterns
542 which are presented to the on-chip network. The images shown in Figure 5b (corresponding to

543 the letters A-D) are encoded in the following way: each pixel corresponds to the resonance
544 wavelength of one of the ring resonators within a neuron, as indicated by the numbers
545 superimposed on the pixels in fig 5b. These wavelengths are aligned to WDM channels 27-41
546 in the telecommunication C-band. In the experiment we present the “white” pixels to the
547 network such that the pulse pattern corresponding to an ‘A’ is, for example, represented by an
548 optical pulse consisting of wavelengths 1, 3, 5, 8 and 14. These wavelengths are multiplexed
549 onto the input waveguide as described in more detail in the supplementary materials and
550 equally split to the synapses of the four neurons by the distributor. After adjusting the
551 synapses (PCM-cells) corresponding to the patterns ‘A’, ‘B’, ‘C’ and ‘D’ using optical pulses,
552 as described previously for the single neuron (Figure 3), the four different pulse patterns are
553 sent to the input of the device. Subsequently, the change in the output spike intensity is
554 observed for all four neurons as shown in Figure5b.

555 **Simulation of language identification with a two-layer network**

556 Using the scalable architecture, we further simulate the performance of a scaled-up,
557 multilayer version of the network of fig. 5a for carrying out a much more complex task of
558 language identification. The network for this task consists, as shown in Figure S16a),b) in the
559 supplementary materials of four input neurons, three hidden layer neurons and two output
560 neurons. The network is assembled using the scaling architecture described in Figure 4 and
561 built up using model representations of photonic neurons according to the measured
562 experimental data (activation function shown in Figure 2d)). This particular network is then
563 used to detect if the language of a given input text is either English or German (the sample
564 texts are taken from³⁴). In a first step the ratio of each vowel (“a”, “e”, “i”, “o”, “u”) and the
565 total number of characters in the input text is calculated (preprocessing). In a second step the
566 five obtained ratios are fed to the inputs of the neural network and the outputs are computed.
567 Each neuron in the simulated network uses the measured optical response from the on-chip

568 neurons (see supplementary materials section 10). Already with a count of about 35 words in
569 the input text, an accuracy above 90% for the language detection is attained. With 150 words
570 the accuracy reaches 99.6%.

571

572 **Data Availability Statement**

573 All data used in this study are available from the corresponding author upon reasonable
574 request.









