



# Synchronous subnanosecond clock and data recovery for optically switched data centres using clock phase caching

Kari A. Clark<sup>1</sup>✉, Daniel Cletheroe<sup>2</sup>, Thomas Gerard<sup>1</sup>, Istvan Haller<sup>2</sup>, Krzysztof Jozwik<sup>2</sup>, Kai Shi<sup>2</sup>, Benn Thomsen<sup>2</sup>, Hugh Williams<sup>2</sup>, Georgios Zervas<sup>1</sup>, Hitesh Ballani<sup>2</sup>, Polina Bayvel<sup>1</sup>, Paolo Costa<sup>2</sup>✉ and Zhixin Liu<sup>1</sup>✉

**The rapid growth in the amount of data being transferred within data centres, combined with the slowdown in Moore's Law, creates challenges for the future scalability of electronically switched data-centre networks. Optical switches could offer a future-proof alternative, and photonic integration platforms have been demonstrated with nanosecond-scale optical switching times. End-to-end switching time is, however, currently limited by the clock and data recovery time, which typically takes microseconds, removing the benefits of nanosecond optical switching. Here we show that a clock phase caching technique can provide clock and data recovery times of under 625 ps (16 symbols at 25.6 Gb s<sup>-1</sup>). Our approach uses the measurement and storage of clock phase values in a synchronized network to simplify clock and data recovery versus conventional asynchronous approaches. We demonstrate the capabilities of our technique using a real-time prototype with commercial transceivers and validate its resilience against temperature variation and clock jitter.**

The rate of data transmitted between servers within data centres has increased rapidly over the last few years<sup>1</sup>, driven by cloud adoption and data-intensive cloud workloads such as data analytics and machine learning. Cloud providers have been able to accommodate this fast growth by relying on Moore's Law for networking, whereby, every two years, electronic switch integrated circuits double their bandwidth at the same cost and power. The long-term sustainability of this trend, however, is being questioned by two upcoming challenges. First, similar to the situation for processor integrated circuits, scaling the transistor density on electronic switch integrated circuits is fundamentally limited by power dissipation as few-nanometre transistor sizes are approached<sup>2</sup>. Second, electronic high-speed serial transceiver data rates are predicted to be hard to scale beyond 112 Gb s<sup>-1</sup> due to the steep increase in dielectric loss when operating at high frequencies<sup>3,4</sup>. Consequently, increasing the aggregate switch capacity will require a proportional increase in the number of serial transceivers surrounding the chip, resulting in greater power density and packaging complexity.

Although continued bandwidth scaling in the near future could be supported by architectural optimizations such as co-packaged optics<sup>5</sup>, preserving cost neutrality in the medium to long term appears very challenging. This uncertainty has motivated research in optical switches as a viable alternative to electronic switches<sup>6</sup>. Optical switches simply redirect the incoming signals onto output ports without any optical/electronic conversion or digital processing. Accordingly, they do not suffer from the limitations of transistor or transceiver technology. They could, therefore, provide a future-proof solution for bandwidth scaling within data centres<sup>7</sup>.

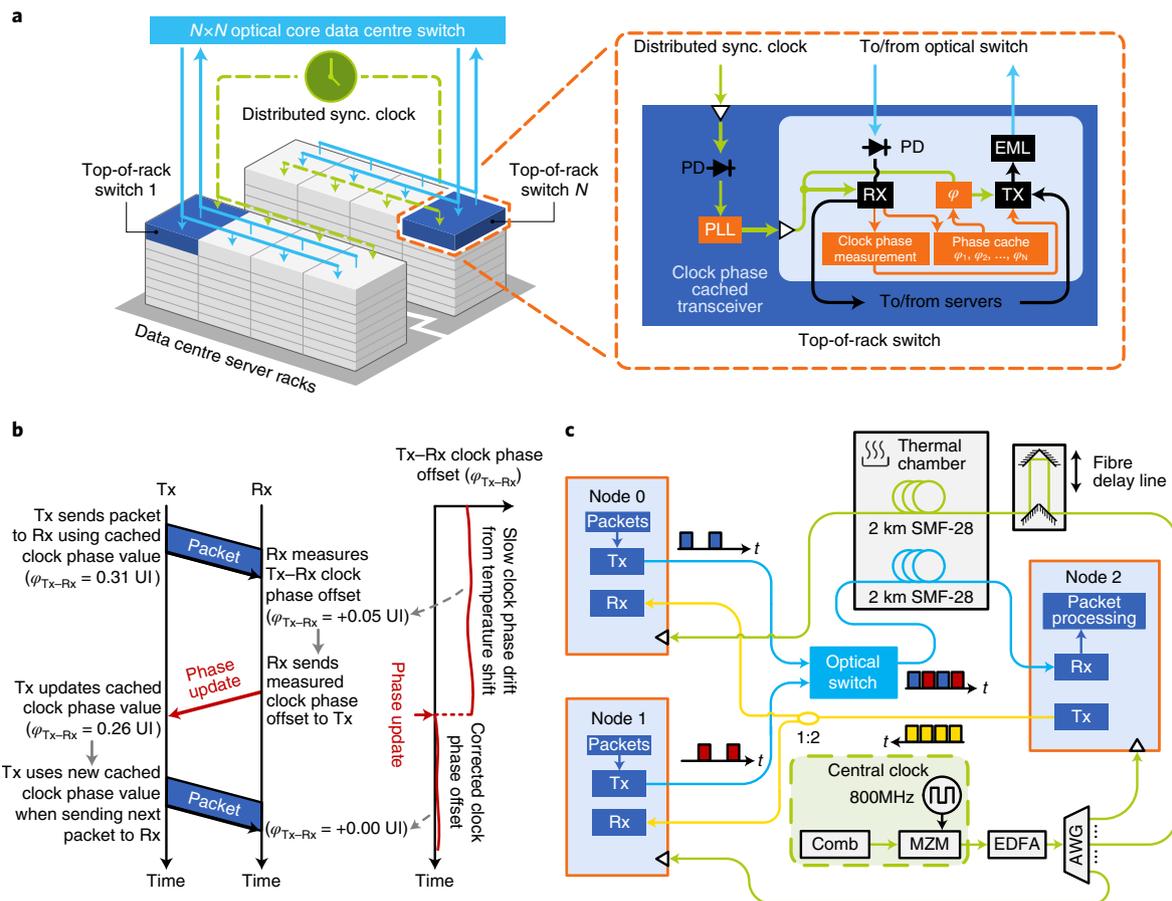
Optical switches would need to support nanosecond-granularity switching, because it is critical to supporting modern data-centre

workloads such as key-value stores, which are dominated by small packets<sup>8</sup>. Fast optical switches with nanosecond optical switching time have been demonstrated using different techniques, including Mach-Zehnder modulators (MZMs)<sup>9</sup>, semiconductor optical amplifiers<sup>10</sup> and tunable lasers<sup>11</sup>. However, the overall end-to-end switching latency also includes the clock and data recovery (CDR) locking time in addition to the optical switching time. CDR locking time is incurred because, unlike the continuous point-to-point interconnections between electronic packet switches, a new physical-layer link is established every time an optical switch is reconfigured to transmit a new data packet. The receiver therefore needs to recover a new clock every time this occurs. In commercially available devices (for example, phase-interpolator-based CDR modules), the locking time is typically on the order of microseconds<sup>12</sup>, as it is necessary to compare the phase of a reference clock to the phase of the embedded data clock using an early-late vote counter. Therefore, the overall end-to-end switching latency would be limited by the CDR locking time rather than by the nanosecond switching time.

Long end-to-end switching time greatly reduces network utilization, defined as the percentage of time that the network is used to send data (as opposed to the time spent reconfiguring the network or for CDR locking). This occurs because the data transmitted in data-centre networks is dominated by small data packets, thus removing the benefits of optical switching. To quantify the impact of CDR locking time when switching data centre workloads, we analysed network traces from a cloud production service. We investigated the impact on network utilization as we varied the CDR locking time. As detailed below, we concluded that a subnanosecond CDR locking time (under 16 symbols at 25.6 Gb s<sup>-1</sup>) is needed for a data centre to achieve over 90% network utilization.

<sup>1</sup>Department of Electronic and Electrical Engineering, University College London, London, UK. <sup>2</sup>Microsoft Research, Cambridge, UK.

✉e-mail: kari.clark.14@ucl.ac.uk; paolo.costa@microsoft.com; zhixin.liu@ucl.ac.uk



**Fig. 1 | Example clock phase cached data-centre architecture, operational principle and demonstration.** **a**, An example clock-frequency-synchronised optically switched data-centre architecture, using our clock phase caching technique to enable subnanosecond CDR locking time. **b**, Operational principle of clock phase caching, showing a single clock phase update for a single transmitter-to-receiver pair. In clock phase caching, clock phase updates are performed at a slow rate across all transmitter-to-receiver pairs. This aligns the clock phase of all data packets arriving at each receiver, irrespective of origin, simplifying subnanosecond CDR. UI, unit interval, equal to one symbol period. **c**, Proof-of-concept experimental demonstration of clock phase caching operating on a 2-to-1 optical switch, in which 2 km clock fibre and 2 km data fibre are placed in a thermal chamber. The signals in these fibres counter-propagate to study the effect of worst-case rates of change of clock phase due to temperature on clock phase caching. PLL, phase-locked loop; Tx, transmitter; Rx, receiver;  $\varphi$ , clock phase interpolator; EML, externally modulated laser; MZM, Mach-Zehnder modulator; PD, photodiode; EDFA, erbium-doped fibre amplifier; AWG, arrayed-waveguide grating.

The CDR locking time in 1 to 2 symbols has been demonstrated using gated voltage-controlled oscillator-based CDRs<sup>13,14</sup>. However, there are two significant fundamental implementation issues associated with this approach: (1) the need for a dedicated gated oscillator per channel leads to high power consumption, large silicon area and possible crosstalk between oscillators; (2) there is an inability to operate at multiple data rates without large performance degradation, which is important for data-centre operation<sup>15</sup>. All-optical clock recovery in optical-injection locked lasers has also been investigated<sup>16</sup>, but these have limited stability for practical application in data centres. These implementation issues are not present for clock phase interpolator-based CDRs, but commercial implementations typically have CDR locking times on the order of microseconds<sup>12</sup>, with state-of-the-art research prototypes still limited to ~325 symbols or longer<sup>17</sup> due to clock phase position metastability when there is a 0.5-symbol phase offset between the initial CDR phase and incoming data<sup>15,17</sup>. This metastability problem arises from CDRs needing to be able to lock to any arbitrary incoming clock phase, which in turn arises from not establishing frequency and phase synchronization between transmitters and receivers connected to an optical switch.

In this Article, we report a technique that can achieve subnanosecond locking time through the measurement and storage of clock phase values in a synchronized network, to simplify clock and data recovery versus conventional asynchronous approaches, using commercial off-the-shelf transceivers. Optical switches lack packet buffers and thus require the network transmissions to be accurately synchronized at a nanosecond timescale to avoid conflicts<sup>18</sup>; this, in turn, requires that the network endpoints must have their clocks' frequencies synchronized. Network-wide frequency synchronization can be achieved in a controlled environment such as a data centre using well-known mechanisms for robust control-plane distribution of a reference clock<sup>19,20</sup>. This partly simplifies the CDR task, as only the clock phase (rather than both frequency and phase) needs to be recovered. Nevertheless, even with a known clock frequency, clock phase recovery can still take up to hundreds of nanoseconds<sup>8</sup>.

However, the clock phase shift across multiple transmissions between the same (transmitter, receiver) pair is relatively constant in a data-centre environment, changing only slowly with temperature due to the fibre time-of-flight change, which occurs due to fibre length expansion and fibre refractive index change<sup>21</sup>. We therefore propose 'caching' or storing the correct clock phase shift to be used

when transmitting between a given transmitter-to-receiver pair, so that the received clock phase is always constant, irrespective of which transmitter is communicating with the receiver.

To account for shifting of the correct clock phase values as the data-centre temperature changes, we designed a low-overhead clock phase update mechanism that all endpoints periodically execute. The clock phase synchronization achieved by clock phase caching is analogous to the clock synchronization of transistors across integrated circuits but on distance scales of up to 2 km. Clock phase updates are required at the data-centre scale because a temperature change across 2 km of optical fibre leads to a clock delay change on the order of  $80 \text{ ps } ^\circ\text{C}^{-1}$  (ref. 21). Across data-centre temperature ranges of up to  $40 \text{ }^\circ\text{C}$  (ref. 22) and data-centre distance scales of 2 km, this, in the worst case, leads to clock phase changes on the order of nanoseconds, requiring clock phase updates to maintain clock phase synchronization. In contrast, in electronic integrated circuits, temperature change causes negligible clock delay change<sup>23</sup>, and so periodic clock phase updates are not required.

### Clock phase caching technique

An example data-centre architecture using clock phase caching is illustrated in Fig. 1a. We consider a typical data-centre cluster comprising thousands of racks ( $\sim 48$  servers per rack) and propose to interconnect all the electronic top-of-rack switches through an all-optical switch fabric. Some of these switches can be used as Internet Protocol (IP) gateways to provide inter-cluster connectivity and enforce routing policies, such as firewalls. Each rack receives a synchronized clock (Fig. 1a), which is used to drive the transceiver logic. This clock could be distributed from multiple synchronized sources<sup>19</sup>, so it is tolerant to clock device failure, as represented by the dotted lines in Fig. 1a. Clock phase caching could also be used in other data-centre architectures, such as the direct interconnection of thousands of servers (instead of top-of-rack switches) from a single all-optical switch fabric.

A clock phase cache is located within each transceiver on each top-of-rack switch and contains a set of values (one value per receiver) corresponding to the phase shifts that need to be applied to the synchronized clock before each packet is sent from its transmitter. At start-up, every transmitter exchanges a packet with all receivers connected through the optical switch in the data centre. Every receiver then measures the phase offset of the received packet and feeds back this information to the transmitter to populate its phase cache. This process is then repeated periodically to account for clock phase drifts due to temperature variation. A single clock phase update between a transmitter and receiver is shown in Fig. 1b.

This technique allows subnanosecond CDR locking time to be consistently achieved across a wide range of temperature and jitter variation while only requiring a slow update frequency ( $\leq 10 \text{ Hz}$ ). The additional logic required to measure the phase offset, store the phase values and apply the phase shift (orange components, Fig. 1a) could be embedded in the switch die or in the transceiver itself. In our prototype implementation, clock phase caching is performed at the transmitter side. However, clock phase caching can also be performed at the receiver side (Supplementary Section 4). This would avoid the need to send the phase value updates back to the transmitters, but it would require the receiver to know the source of each incoming packet ahead of time.

In our approach, we distribute an optical clock via a control plane so that the clock frequency is synchronized. Although optical fibre-based clock synchronization, including clock phase and frequency transfer between nodes, has been used for metrology<sup>24</sup> and optical time domain division multiplexing<sup>25,26</sup>, it has not been investigated for burst-mode data communications. Furthermore, our method for measuring and controlling the clock phase is implemented on a digital CDR module and clock phase interpolator, respectively. The capability of phase tracking is therefore not limited

by optical or electronic delay lines, avoiding the cost, large size and limited clock phase tracking range associated with these devices.

### Impact of CDR locking time on network utilization

To assess the impact of CDR locking time for cloud data-centre networks, we analysed network traces from a cloud production service<sup>8</sup>. As shown in Fig. 2b, over 34% of the packets comprise fewer than 128 bytes, which means switching every 11 ns (with  $100 \text{ Gb s}^{-1}$  ports), while 98% of the packets comprise 576 bytes or fewer<sup>8</sup>. This is consistent with a similar study from Facebook, which found that 91% of the packets generated by their in-memory distributed cache service are 576 bytes or fewer<sup>27</sup>.

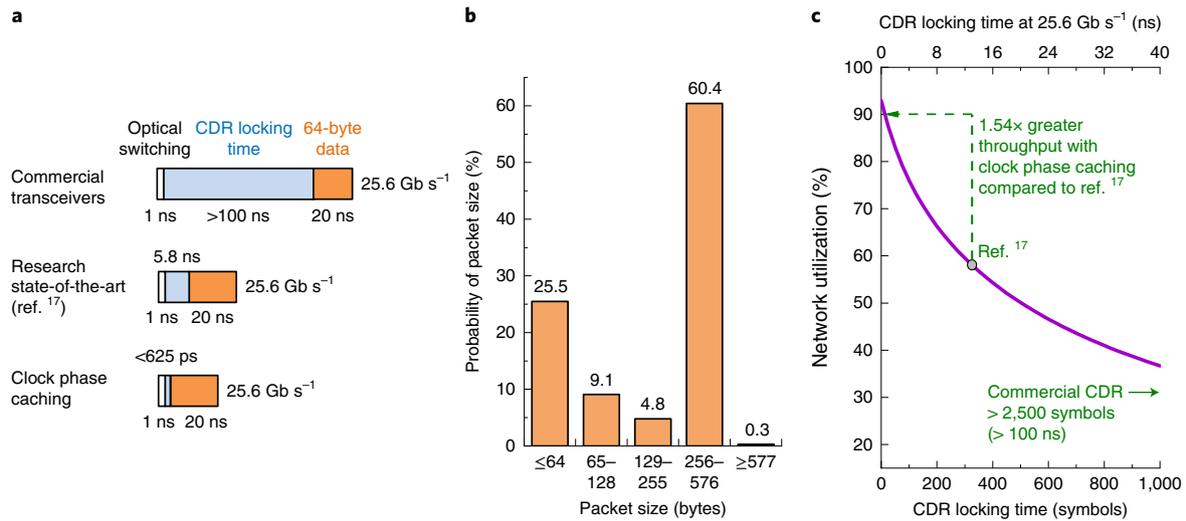
As we illustrate in Fig. 2a, with small packets, even if reconfiguring the optical switch takes only 1 ns, long locking time leads to low network utilization. The network utilization decreases as the CDR locking time increases because an increasing proportion of time is used for CDR locking rather than data packet reception. For example, assuming a  $25.6 \text{ Gb s}^{-1}$  transceiver, sending a 64-byte packet only takes 20 ns. Using a commercially available CDR with a locking time of 100 ns, the overall network utilization would be less than 17%. In Fig. 2c, we plot the results of a study to evaluate the impact of CDR locking time using the packet distribution in Fig. 2b, assuming  $4 \times 25 \text{ Gb s}^{-1}$  transceivers. As shown in the graph, reducing the CDR locking time from the state-of-the-art CDR locking time for phase interpolator CDRs of 325 symbols ( $12.7 \text{ ns}$  at  $25.6 \text{ Gb s}^{-1}$ )<sup>17</sup> to subnanosecond would result in a  $1.54\times$  improvement in network utilization.

### Experimental demonstration of clock phase caching

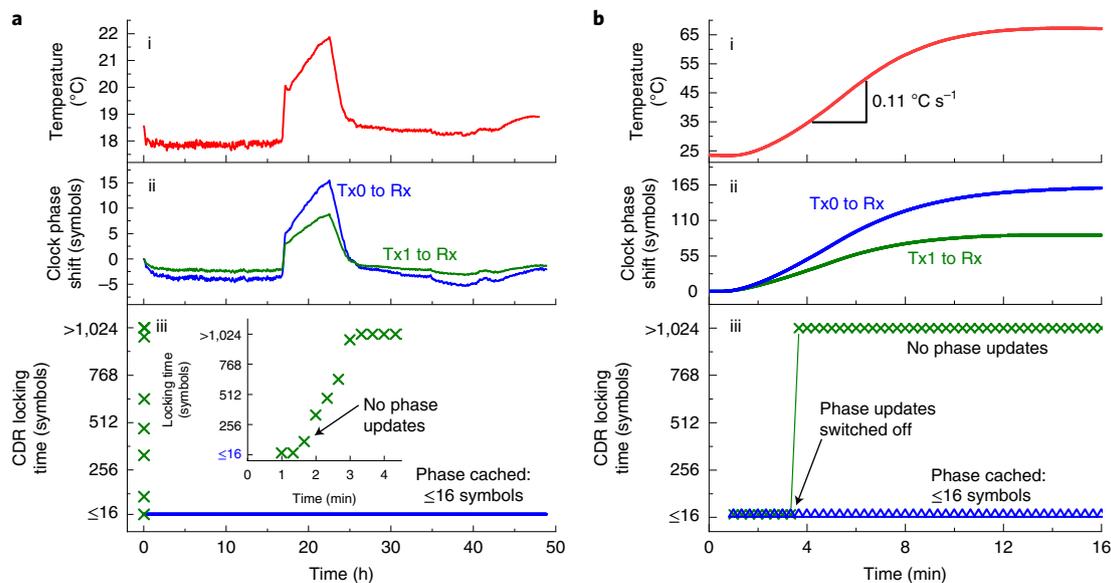
To demonstrate the feasibility of our phase caching technique and evaluate its performance, we built a prototype system using three field programmable gate arrays (FPGAs) and investigated the performance of our CDR technique in a  $2 \times 1$  optically switched network. As shown in Fig. 1c, the first two nodes (nodes 0 and 1) transmit 128-byte on-off keying (OOK) modulated packet payloads embedded in 60-ns packets at  $25.6 \text{ Gb s}^{-1}$ , via externally modulated lasers (EMLs), to node 2 through a  $2 \times 2$  LiNbO<sub>3</sub> optical switch with a switching time of 200 ps. After the switch, alternate packets from node 0 and node 1 propagate through 2 km of single-mode optical fibre (SMF-28) and are attenuated to  $-10.5 \text{ dBm}$  before reaching node 2 to emulate power budgets of intra-data-centre transmission standards<sup>28</sup>.

The clock-phase-cached transceiver architecture shown in Fig. 1a is implemented in the three FPGAs. As shown in Fig. 1b, node 2 periodically measures the clock phase offset using the FPGA CDR module and subsequently sends the clock phase offset values back to nodes 0 and 1 via the link shown in yellow, which emulates duplex interconnection. We then shift the transmitter clock phase using the clock phase interpolator in the transmitter. An 800-MHz reference clock was modulated onto a 24-channel frequency comb and distributed to all three nodes for frequency synchronization, emulating distributed frequency synchronization techniques such as White Rabbit<sup>29</sup> and Sync-E<sup>19</sup>.

To study the impact of varying environmental conditions observed in production data centres, we placed 2 km of SMF-28 for the data path and 2 km of SMF-28 for the clock path in a thermally controlled chamber. When switched on, the temperature started at  $25 \text{ }^\circ\text{C}$  and increased approximately linearly between 30 and  $50 \text{ }^\circ\text{C}$  at a rate of  $0.11 \text{ }^\circ\text{C s}^{-1}$ . We configured the system such that the signals propagated in opposite directions in the data and clock fibres. Variation in temperature therefore caused the phase shift in the two fibres to be additive, resulting in an overall phase shift at the receiver equal to 4 km SMF-28 ( $\sim 160 \text{ ps } ^\circ\text{C}^{-1}$ ), which allowed us to investigate the worst-case rate of clock phase shift in synchronous intra-data-centre interconnection. This contrasts with the case where clock and data signals propagate in the same direction, where



**Fig. 2 | Increase of data centre network utilization enabled by clock phase caching.** **a**, Receiver overhead caused by different CDR techniques when receiving minimum size (64 byte) data packets. **b**, Small-packet-dominated distribution of packet size in a measured production cloud data-centre traffic pattern. **c**, Gain in data centre network utilization from using clock phase caching versus long CDR locking time when handling the traffic pattern shown in **b** (assuming 1 ns optical switching time).

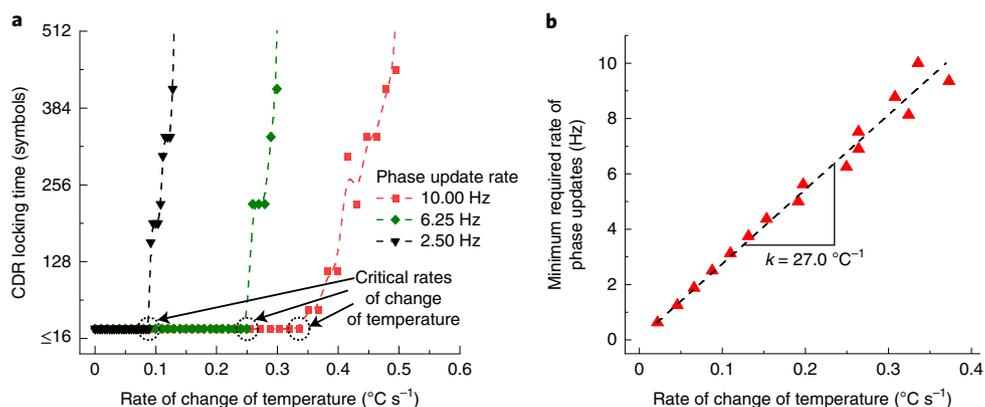


**Fig. 3 | Demonstration of the stability of clock phase caching.** **a**, Stability over 48 h of measurement: (i) recorded ambient temperature, with the air conditioner switched off between the 17th and 23rd hours; (ii) recorded receiver clock phase shift for packets originating from each transmitter; (iii) receiver CDR locking time. Blue line, with clock phase caching; green crosses, no clock phase caching. **b**, Stability of clock phase caching under a rapid rate of temperature change of 0.11 °C s<sup>-1</sup>, which is over three times greater than the worst-case rate of change of temperature we observed in a production cloud data centre: (i) temperature within a thermally controlled chamber; (ii) recorded clock phase shift for packets originating from each transmitter; (iii) receiver CDR locking time. Blue triangles, with clock phase caching; green crosses, clock phase caching turned off during the rapid temperature shift. A straight-line fit is shown as a guide to the eye.

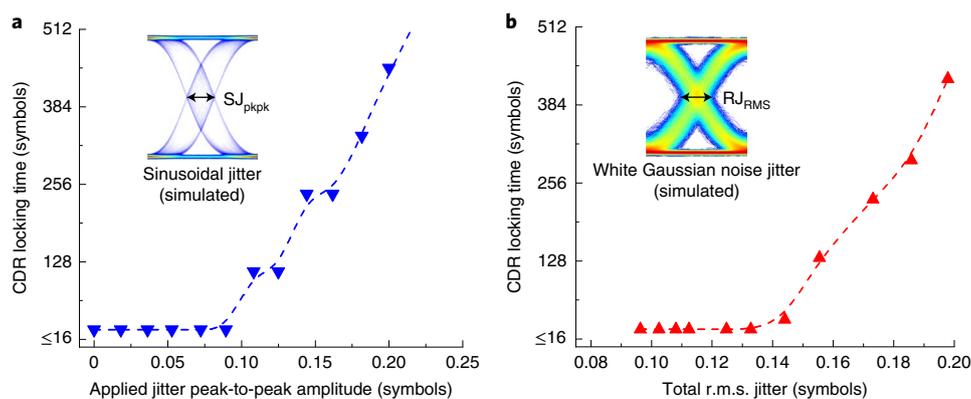
balanced changes in both fibres due to temperature variation would cause the phase shifts to cancel at the receiver.

To demonstrate the reliability of our technique, we first carried out a measurement over a 48 h period in a laboratory environment, where the temperature fluctuated within 5 °C. During this long-term stability test, we emulated cooling failure by turning off the laboratory air conditioner (Fig. 3a(i)), which led to a temperature rise of 2 °C at 17 h followed by a slow room-temperature

increase during the 17th and 23rd hours. After that, we turned on the air conditioner and cooled the room back to 18.5 °C. Our transceiver continuously monitored the phase shift of each transmitter, as shown in Fig. 3a(ii). The CDR locking time was determined by finding the first error-free 16-bit bin by averaging  $9.2 \times 10^8$  packets over 60 s, as detailed in the Methods. When clock phase caching was enabled, the system was error-free over the whole 48 h, resulting in a CDR locking time of under 16 symbols (under 625 ps), as shown in



**Fig. 4 | Impact of rate of change of temperature on clock phase caching.** **a**, CDR locking time for different rates of temperature and clock phase cache update rates across 2 km of SMF-28 clock fibre and 2 km of SMF-28 data fibre. In all cases, the rate of clock phase updates is small ( $<10$  Hz). B-spline fits (dashed lines) are shown as a guide to the eye. **b**, Minimum required rate of clock phase updates to achieve 625 ps (16 symbol) CDR locking time for different rates of temperature change for the fibre length described in **a**. Clock phase updates are required at a rate of 27.0 Hz for each  $1^\circ\text{C s}^{-1}$  of rate of change of temperature (note that the worst-case rate of change of temperature that we observed in a production cloud data centre was  $0.03^\circ\text{C s}^{-1}$ ). A linear regression fit (black dashed line) was used to calculate this proportionality constant.



**Fig. 5 | Impact of clock jitter on CDR locking time.** **a**, Impact of 1 MHz sinusoidal jitter on CDR locking time. The tolerance to applied peak-to-peak sinusoidal jitter ( $S_{J_{pkpk}}$ ) is 0.089 symbols, corresponding to 3.5 ps at the data rate of  $25.6\text{ Gb s}^{-1}$  used in our experiment. A B-spline fit (dashed blue line) is shown as a guide to the eye. **b**, Impact of white Gaussian noise jitter ( $R_{J_{RMS}}$ ) on CDR locking time. The r.m.s. jitter tolerance was measured to be 0.13 symbols, corresponding to 5.1 ps r.m.s. jitter at  $25.6\text{ Gb s}^{-1}$ . A B-spline fit (dashed red line) is shown as a guide to the eye.

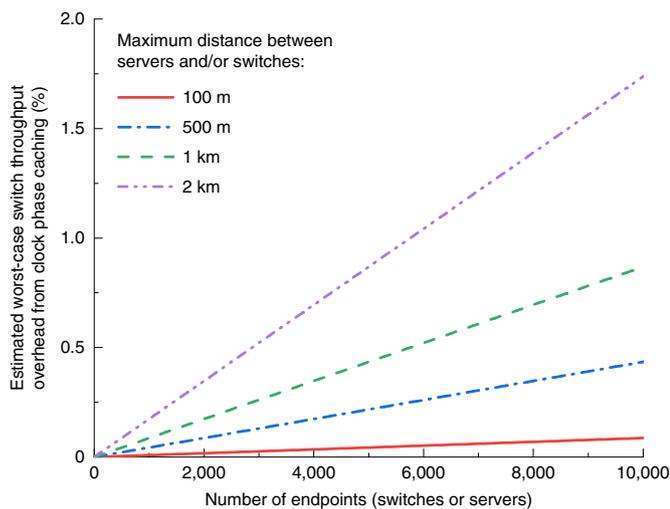
Fig. 3a(iii). When phase caching was disabled, we observed a quick degradation in bit error rate (BER) due to temperature drift, resulting in an increase in CDR locking time to over 40 ns within 4 min. Figure 3b shows the performance of our clock phase caching technique under an emulated worst-case rate of change of data-centre temperature. Figure 3b(i) and Fig. 3b(ii) show the measured temperature in the temperature chamber and the resulting total phase shifts. The  $0.11^\circ\text{C s}^{-1}$  increase in temperature resulted in a change of  $16\text{ ps s}^{-1}$  in fibre time of flight across  $2 \times 2\text{ km}$  of SMF-28. Figure 3b(iii) shows the measured CDR locking time with phase caching enabled at a rate of 10 Hz. Even with the rapid change of temperature, no errors were observed, and the CDR locking time was under 16 symbols (under 625 ps). When clock phase caching was initially run, and then switched off during the rate of temperature increase of  $0.11^\circ\text{C s}^{-1}$ , we observed a loss of clock phase alignment in less than one second due to the temperature-induced change in fibre time of flight.

The rate of clock phase drift is proportional to the rate of temperature change. This impacts the stability of our CDR technique, which requires a sufficient rate of recalibration of clock phases for

all transceiver pairs to maintain a CDR locking time of under 16 symbols (under 625 ps). To investigate the required phase update rate for different rates of temperature change, we employed an electronically controlled tunable optical delay line to emulate the temperature-induced change of time of flight across  $2 \times 2\text{ km}$  of SMF-28. The optical delay line is driven with triangular waveforms of different slopes to emulate different rates of temperature change of up to  $0.50^\circ\text{C s}^{-1}$ .

As we show in Fig. 4, the critical rate of temperature change at which the CDR locking time begins to increase to greater than under 16 symbols (under 625 ps) is proportional to the rate of clock phase updates. The proportionality constant,  $k$ , represents the clock phase update rate required to cope with the worst-case rate of change of temperature within the data centre and is  $27.0^\circ\text{C}^{-1}$  in our system.

To evaluate the worst-case rate of change of temperature within data centres, we monitored the inlet and exhaust temperature for a rack in a production data centre with evaporative cooling over 228 days, and the largest temperature variation observed was  $9^\circ\text{C}$ . Even assuming this occurred across adjacent 5-min measurement intervals, it translates to a rate of change of temperature of  $0.03^\circ\text{C s}^{-1}$ .



**Fig. 6 | Estimated worst-case optical switch overhead from the clock phase caching technique.** The overhead resulting from clock phase caching is directly proportional to the number of clock phase cached endpoints connected to an optical switch. Clock phase caching supports all intra-data-centre distance scales with small network utilization overhead from performing clock phase cache updates. At the largest intra-data-centre distance scale, with 10,000 switches or servers interconnected by an optical switch with a maximum of 2 km of optical fibre between endpoints, the estimated overhead from clock phase caching is only 1.7%.

This is also consistent with recommendations for industrial data-centre design<sup>22</sup>. Based on Fig. 4b, a clock phase update rate of 1 Hz is sufficient to cope with this worst-case rate of change of temperature.

To demonstrate the jitter tolerance of our CDR technique, we separately applied different amplitudes of 1-MHz sinusoidal jitter and white noise jitter to the 800-MHz reference clock source. For both types of jitter, increasing the clock jitter amplitude eventually degraded the CDR locking time. The 2×2 km fibres were kept in an insulated chamber to minimize the impact of ambient temperature variation. In Fig. 5, we show the CDR locking time against root-mean-square (r.m.s.) jitter amplitude for 1-MHz sinusoidal jitter (Fig. 5a) and white noise jitter (Fig. 5b). The CDR locking time was under 16 symbols (under 625 ps) when the r.m.s. jitter was equal to or less than 5.2 ps. Standard reference oscillators have jitter of a few hundred femtoseconds, ensuring sufficient performance for our CDR technique for intra-data-centre interconnection. We also show the simulated effect of sinusoidal and white noise jitter on a non-return-to-zero on-off keyed (NRZ-OOK) signal.

### Estimated scalability of clock phase caching

Our clock phase caching technique must be scalable to support hyperscale data-centre interconnection. Based on the experimental results, we estimate the worst-case network utilization overhead resulting from clock phase caching at different data-centre distance scales. An estimate of the worst-case optical switch network utilization overhead resulting from clock phase caching is given by  $D(\frac{dT}{dt})kN(t_{\text{meas}} + t_{\text{update}})/2$  km, where  $D$  is the maximum distance between the servers and/or switches,  $\frac{dT}{dt}$  is the worst-case rate of change of data-centre temperature,  $k$  is the proportionality constant from Fig. 4b of  $27.0^\circ\text{C}^{-1}$ ,  $N$  is the number of nodes,  $t_{\text{meas}}$  is the time taken to measure the clock phase offset between each transmitter-to-receiver pair, as detailed in the Methods and  $t_{\text{update}}$  is the length of time each transmitter spends transmitting a packet to carry each phase update value.

In Fig. 6, we estimate the worst-case optical switch network utilization overhead of clock phase caching, based on our measured worst-case rate of change of data-centre temperature of  $0.03^\circ\text{C s}^{-1}$ , a time taken for each clock phase update per transmitter-to-receiver pair of  $2.08\ \mu\text{s}$  (as detailed in the Methods) and a time taken to transmit each update packet of  $0.065\ \mu\text{s}$ .

The estimated worst-case network utilization overhead resulting from clock phase caching is only 1.7% for 10,000 data-centre endpoints (switches or servers) separated by a maximum of 2 km. If 10,000 electronic top-of-rack switches, each networking 48 servers (480,000 servers in total), were connected by an optically switched fabric, our clock phase calibration technique would allow all possible transmitter-to-receiver pairs to communicate with each other with under 16 symbols (under 625 ps) CDR locking time with only a 1.7% optical switch worst-case network utilization overhead. For 10,000 endpoints, a traditional asynchronous CDR would need a clock recovery time of 25 symbols to have an equivalent network utilization to clock phase caching.

Further scalability improvements can be made by using optical fibre that experiences a smaller change in time of flight due to temperature change. For example, the change of the time of flight in hollow-core fibre is 20 times less than that in SMF-28<sup>21</sup>, which can potentially lead to an overhead reduction by a factor of 20, or scaling up of the number of supported nodes at the same overhead by a factor of 20. Specially coated fibre with three to four times lower thermal sensitivity is also commercially available for overhead reduction or scaling up the number of clock-synchronized nodes<sup>30</sup>.

### Conclusions

We have demonstrated a clock phase caching technique that enables a CDR locking time of under 625 ps for optically switched data-centre networks. Through a real-time prototype network, we show that the clock phase caching technique can tolerate a rate of change of temperature of  $0.11^\circ\text{C s}^{-1}$ , which is three times greater than the worst-case conditions in production cloud data centres, as well as a r.m.s. jitter tolerance of 0.135 symbols. Our technique allows for the utilization of an optically switched data-centre network to exceed 90%. We estimate that it can be scaled to support 10,000 data-centre endpoints (servers or switches) for realistic worst-case data-centre environmental temperature variation and fibre lengths with a worst-case small 1.7% overhead on network utilization.

Our clock phase caching technique could be of value in the various scientific and engineering applications that require high-throughput burst-mode transmission, and also in communities that require synchronization. For example, time-division multiplexing-based passive optical networks (PONs) could potentially benefit from the synchronization technology that we have demonstrated<sup>31</sup>. The optical network units in PONs could be synchronized in a similar manner to minimize the required inter-packet gap and data-detection latency, which is crucial for latency-sensitive applications such as the Internet of Skills, virtual reality and gaming<sup>32</sup>. Another example lies in the clock synchronization for quantum key distribution (QKD) systems, which require the time gate to be synchronized to the photon arrival time to identify the quantum signals correctly<sup>33</sup>. Our clock phase caching technique can track the drift of the time of flight in transmission links without the need for paired fibre, significantly reducing the complexity of QKD clock networks.

### Methods

**Optical clock distribution.** We have achieved frequency synchronization of the three FPGA nodes by modulating an 800-MHz reference clock onto a 24-tone optical frequency comb (spacing of 25 GHz) with a MZM. The frequency comb used was an optoelectronic comb source consisting of a 27-dBm continuous-wave light source emitting at 1,555 nm, a phase modulator and a MZM<sup>34</sup>. Both the phase modulator and the MZM were driven by a 25-GHz radiofrequency source, generating a 25-GHz spacing comb with 2-dB power flatness. The generated

comb was modulated and subsequently amplified to 18 dBm by an EDFA, yielding an average power of  $\sim 3$  dBm per tone. These optical clock signals were filtered by an arrayed waveguide grating (AWG) and detected by 18-GHz photodiodes with transimpedance amplifiers. In our experiment, the lowest receiver power required for each modulated tone was  $-11$  dBm. We thus emulated the use of a 1:8 optical splitter by attenuating the optical power from  $-0.5$  dBm (power per tone after AWG) to  $-11$  dBm, indicating that our current system can optically synchronize 192 nodes from a single clock source. In practice, relatively low-speed photodiodes (for example, 5-GHz bandwidth, which were unavailable at the time of the experiment) may provide higher sensitivities because of their slightly higher responsivity (for example, about  $1 \text{ A W}^{-1}$ ). This number of nodes can be easily increased to 3,072 if we use a high-power EDFA to amplify the modulated clock to 30 dBm. The clock signal between the AWG and node 1 additionally travels through 2 km of SMF-28. Noting that the clock distribution does not necessitate the use of a frequency comb, conventional laser arrays with channel spacing of 50 GHz or 100 GHz could also be used. Nevertheless, a frequency comb provides a compact source with well-defined wavelength spacing, which may ease the thermal and wavelength management in a clock-synchronized network with potentially low operating expenses. In addition, wideband and high optical signal-to-noise ratio (OSNR) frequency combs such as parametric combs<sup>35</sup> and thin-film LiNbO<sub>3</sub> (ref. <sup>36</sup>) combs can potentially allow the system to scale to more than 10,000 nodes.

**Transmitter data modulation.** The externally modulated lasers shown in Fig. 1c consist of 1,550 nm carriers modulated with data using 35-GHz-bandwidth electro-absorption modulators.

**Optical switch control.** The focus of this work is on demonstrating that clock phase caching can reduce the CDR locking time to subnanosecond in optical switches. To minimize experimental complexity, we used a  $2 \times 2$  optical switch, both transmitters were configured to output packets continuously with no central or edge scheduling, and frame alignment of packets from the two transmitters was performed manually. The  $2 \times 2$  optical switch (an 18-GHz LiNbO<sub>3</sub> MZM) was driven by a square wave clock signal (130-ns period: two 60-ns packets plus 5-ns inter-packet gap) from a Xilinx GTY transceiver (18 GHz bandwidth) from FPGA node 1.

**Experimental packet structure.** Our experimental packet structure used to measure the CDR locking time (shown in Supplementary Fig. 1) contains three De Bruijn sequences of  $2^9$  length. A media access control (MAC) layer protocol (64 bits) is embedded in the third sequence for frame alignment, clock phase offset feedback and packet identification. In a full system demonstration of clock phase caching, clock phase measurement packets (such as those used in our experimental demonstration) would be scheduled and transmitted periodically between standard Ethernet packets. In our implementation, for the purposes of simplifying implementation complexity as the aim of our experimentation is to demonstrate the proof of principle of clock phase caching, we allowed all packets transmitted to be potentially used for both clock phase measurement and data transmission, avoiding the need to implement a packet scheduler. The second and third sequences in the packets are used to measure clock phase offset when a clock phase update is required.

**CDR locking time measurement.** To assess the performance of phase caching, the first and second sequences act as the 128-byte payload and are divided into  $64 \times 16$ -bit bins. The number of bit errors falling in each bin across all packets arriving at the receiver is recorded in real time over intervals of 1 s. As shown in Supplementary Fig. 2, the CDR locking time was calculated as the first bin in the packet with a BER of under  $10^{-10}$ , if all following bins also had a BER of under  $10^{-10}$ . The BER was calculated by summing errors collected across all 64 bins.

**Transmitter clock phase shift.** The clock phase interpolator consists of a standard clock phase rotator that takes a half-rate clock from the output of the FPGA transceiver LC-Tank quad PLL, and outputs one clock phase selected from 128 evenly spaced phase steps split across two symbols (resolution of 64 steps per symbol). The selected clock phase is then used to drive the serial clock of the parallel-in serial-out (PISO) converter in the FPGA transmitter. The phase interpolator selected phase output is controlled by a PIPPM control circuit built into the Xilinx GTY transceiver. This control circuit allows the transmitter clock phase to be shifted by up to  $15/64$ th of a symbol every transmitter parallel clock cycle, which is 400 MHz in our proof-of-concept experiment. Within the 5-ns inter-packet gap, immediately before packet transmission to a receiver, the phase interpolator selected phase is changed to a clock phase value equivalent to the cached clock phase value for communication with the receiver. Shifting by half a symbol requires two FPGA parallel clock cycles, which takes 5 ns in total, hence leading to a required 5-ns gap, followed by a further  $2/64$ th of a symbol shift at the beginning of the next packet. The  $2/64$ th of a symbol shift, if required, could alternatively be performed at the end of the previous packet. Shifting by greater than half a symbol is achieved by shifting the clock phase in the reverse direction by up to half a symbol. The gap could be reduced to subnanosecond by directly using the clock phase cached value to change the phase interpolator selected

phase, bypassing the PIPPM control circuit. This would remove the phase shift delay associated with the PIPPM control circuit, thereby enabling the phase of the transmitter clock to switch to any arbitrary phase in a subnanosecond time, which would fall within the inter-packet gap of 1 ns required to allow for optical switching time.

**Measurement of the transmitter-to-receiver clock phase offset.** For each clock phase update, the transmitter-to-receiver clock phase offset was measured by averaging the FPGA receiver raw CDR phase across the second and third packet sequences, and then averaging the clock phase offset of 32 packets separated by 1  $\mu$ s, taking a total of 2.08  $\mu$ s per update (including inter-packet gaps). The receiver CDR phase interpolator spans 128 evenly spaced phase steps split across two symbols (resolution of 64 phase steps per symbol). The raw CDR phase is output at a rate of 400 MHz, and this phase is equal to the clock phase shift applied to the receiver reference clock to keep it aligned with the bit edges in incoming data, which is achieved within the CDR by using a phase interpolator to measure the ratio between the number of times the clock edge occurs before and after the data edge in 64 bit intervals, and shifting the receiver reference clock such that this ratio equals 50%. A single packet is sent back to the transmitter node along the return link to update the clock phase offset at the transmitter once the clock phase offset has been calculated at the receiver. No oversampling is used. The built-in Xilinx GTY transceiver receiver-side continuous-time linear equalizer (CTLE) filter is used.

**Experimental emulation of the rate of change of temperature.** To evaluate the clock phase update rate required to maintain a CDR locking time under 16 symbols (under 625 ps) for a given rate of change of temperature across  $2 \times 2$ -km single-mode optical fibre, we introduced a 166-ps, voltage-controlled free-space optical delay line between the AWG and the clock input of node 1. We drove the fibre delay line with a sawtooth waveform to experimentally emulate the rate of change of phase that occurs for a given rate of change of temperature in SMF-28. The sawtooth waveform frequency was swept through a series of values equal to that theoretically experienced by a loose tube single-mode fibre of 4 km length with a temperature coefficient of delay of  $40.6 \text{ ps km}^{-1} \text{ } ^\circ\text{C}^{-1}$  at a series of different rates of change of temperature, as shown in Fig. 4b. A total of  $10^{11}$  bits per 16-bit bin were collected for each rate of change of temperature. To avoid the need to manually re-perform time synchronization of the two transmitters, we biased the switch into the crossbar state such that packets only from node 1 arrive at the receiver. To prevent the receiver CDR remembering the clock phase of incoming data between successive packets, we reset the receiver CDR phase within the 5-ns inter-packet gap.

**Clock jitter generation.** Jitter was applied to the electronic clock source shown in Fig. 1c by frequency modulating the 800-MHz clock source with a 1-MHz sinusoidal and white-noise voltage waveforms of different amplitudes. The peak-to-peak sinusoidal jitter applied to the clock was calculated after the clock output and before the MZM using a radiofrequency spectrum analyser by measuring the ratio of the power of the 800-MHz clock tone to the power of its  $\pm 1$ -MHz side tones. The r.m.s. jitter amplitude of the clock after frequency modulation with white Gaussian noise jitter was measured using a digital communications analyser, and  $10^{11}$  bits per 16-bit bin were collected for each applied jitter amplitude. To prevent the receiver CDR remembering the clock phase of incoming data between successive packets, we reset the receiver CDR phase within the 5-ns inter-packet gap.

**Estimation of network utilization in Fig. 2c.** To estimate the impact of the CDR locking time on the overall network utilization, we used an event-based network simulator that we developed in house, which we cross-validated against real data centre networks. The simulator models a network consisting of nodes interconnected by an optical switch. The line rate of each node to and from the optical switch was  $4 \times 25 \text{ Gb s}^{-1}$ . We generated a synthetic workload by randomly selecting the payload size for each packet from the distribution shown in Fig. 2b and by selecting sources and destinations with a uniform random distribution. As soon as a node finishes sending a packet, a new packet is generated and a new destination is selected. If a source-destination path is not already set up, before the packet payload can be received the optical switch needs to be reconfigured (optical switching time) and the receiver CDR has to lock onto the new incoming signal (CDR locking time). We set the optical switching time to 1 ns, based on recent advances in optical switching devices<sup>9-11</sup>, and we varied the CDR locking time from 0 to 400 ns (corresponding to a range from 0 to 1,000 symbols). To generate Fig. 2c, we recorded the network utilization (measured as the ratio of the time in which packet payloads were sent divided by the total simulation time) as we increased the CDR locking time.

## Data availability

Source data are provided with this paper.

Received: 12 November 2019; Accepted: 5 May 2020;

## References

- Singh, A. et al. Jupiter rising: a decade of Clos topologies and centralized control in Google's datacenter network. *ACM SIGCOMM Comput. Commun. Rev.* **45**, 183–197 (2015).
- Markov, I. Limits on fundamental limits to computation. *Nature* **512**, 147–154 (2014).
- Dorren, H. et al. Challenges for optically enabled high-radix switches for data center networks. *J. Lightwave Technol.* **33**, 1117–1125 (2015).
- Ghiasi, A. Large data centers interconnect bottlenecks. *Opt. Express* **23**, 2085–2090 (2015).
- Krishnamoorthy, A. et al. From chip to cloud: optical interconnects in engineered systems. *J. Lightwave Technol.* **35**, 3103–3115 (2017).
- Testa, F. & Pavesi, L. *Optical Switching in Next Generation Data Centers* (Springer, 2017).
- Ballani, H. et al. Bridging the last mile for optical switching in data centers. In *2018 Optical Fiber Communication Conference (OFC) W1C.3* (OSA, 2018).
- Clark, K. et al. Sub-nanosecond clock and data recovery in an optically-switched data centre network. In *2018 European Conference on Optical Communication (ECOC) 1–3* (IEEE, 2018).
- Chen, C. P. et al. Programmable dynamically-controlled silicon photonic switch fabric. *J. Lightwave Technol.* **34**, 2952–2958 (2016).
- Cheng, Q., Wonfor, A., Wei, J. L., Penty, R. V. & White, I. H. Low-energy, high-performance lossless 8×8 SOA switch. In *Proc. 2015 Optical Fiber Communication Conference (OFC) Th4E.6* (OSA, 2015).
- Shi, K. et al. System demonstration of nanosecond wavelength switching with burst-mode PAM4 transceiver. In *Proc. 2019 European Conference on Optical Communication (ECOC) 1–3* (IEEE, 2019).
- Xilinx. *DS893 (v1.12) Virtex UltraScale Architecture Data Sheet: DC and AC Switching Characteristics* (2019); [https://www.xilinx.com/support/documentation/data\\_sheets/ds893-virtex-ultrascale-data-sheet.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds893-virtex-ultrascale-data-sheet.pdf)
- Banu, M. & Dunlop, A. E. Clock recovery circuits with instantaneous locking. *Electron. Lett.* **28**, 2127–2130 (1992).
- Terada, J. et al. A 10.3 Gb/s burst-mode CDR using a  $\Delta\Sigma$  DAC. *IEEE J. Solid-State Circuits* **43**, 2921–2928 (2008).
- Rylyakov, A. et al. A 25 Gb/s burst-mode receiver for low latency photonic switch networks. *IEEE J. Solid-State Circuits* **50**, 3120–3132 (2015).
- Yang, Y., Wen, Y. J., Nirmalathas, A., Liu, H. F. & Novak, D. Optical clock recovery at line rates via injection locking of a long cavity Fabry–Pérot laser diode. *IEEE Photon. Technol. Lett.* **16**, 1561–1563 (2004).
- Ozkaya, I. et al. A 56 Gb/s burst-mode NRZ optical receiver with 6.8 ns power-on and CDR-lock time for adaptive optical links in 14 nm FinFET CMOS. In *Proc. 2018 IEEE International Solid-State Circuits Conference (ISSCC) 266–268* (IEEE, 2018).
- Bostica, B., Burzio, M., Gambini, P. & Zucchelli, L. In *Photonic Networks* (ed. Prati, G.) 362–376 (Springer, 1997).
- Timing and Synchronization Aspects in Packet Networks* G.8261 (ITU, 2008); <https://www.itu.int/rec/T-REC-G.8261-201908-I/en>
- Lee, K. S., Wang, H., Shrivastav, V. & Weatherspoon, H. Globally synchronized time via datacenter networks. In *Proc. 2016 ACM SIGCOMM Conference* 454–467 (ACM, 2016).
- Slavik, R. et al. Ultralow thermal sensitivity of phase and propagation delay in hollow core optical fibres. *Sci. Rep.* **5**, 15447 (2015).
- Data Center Networking Equipment—Issues and Best Practices* TC9.9 (ASHRAE, 2016); [https://tc0909.ashraetsc.org/documents/ASHRAE\\_TC0909\\_Power\\_White\\_Paper\\_22\\_June\\_2016\\_REVISIED.pdf](https://tc0909.ashraetsc.org/documents/ASHRAE_TC0909_Power_White_Paper_22_June_2016_REVISIED.pdf)
- Soleimani, S., Afzali-Kusha, A. & Forouzandeh, B. Temperature dependence of propagation delay characteristic in FinFET circuits. In *Proc. Int. Conference on Microelectronics (ICM) 276–279* (IEEE, 2008).
- Ma, L.-S., Jungner, P., Ye, J. & Hall, J. L. Delivering the same optical frequency at two places: accurate cancellation of phase noise introduced by an optical fiber or other time-varying path. *Opt. Lett.* **19**, 1777–1779 (1994).
- Lord, A., Blank, L. C., Boggis, J. M., Bryant, E. & Stallard, W. A. Theory of control mechanism for an optically time-division-multiplexed system. *Electron. Lett.* **24**, 2011–2012 (1988).
- Ellis, A. D., Widdowson, T., Phillips, I. D. & Pender, W. A. High speed OTDM networks employing electro-optic modulators. *Trans. Inst. Electron. Inf. Commun. Eng. Sect. E* **E81-C**, 1301–1308 (1998).
- Zhang, Q., Liu, V., Zeng, H. & Krishnamurthy, A. High-resolution measurement of data center microbursts. In *Proc. 2017 Internet Measurement Conference (IMC) 78–85* (ACM, 2017).
- IEEE Standard for Ethernet—Amendment 3: Physical Layer Specifications and Management Parameters for 40 Gb/s and 100 Gb/s Operation over Fiber Optic Cables* 802.3bm (IEEE, 2015); [https://standards.ieee.org/standard/802\\_3-2018.html](https://standards.ieee.org/standard/802_3-2018.html)
- Lipinski, M., Wlostowski, T., Serrano, J. & Alvarez, P. White rabbit: a PTP application for robust sub-nanosecond synchronization. In *Proc. 2011 IEEE International Symposium on Precision Clock Synchronization for Measurement, Control and Communication* (IEEE, 2011); <https://doi.org/10.1109/ISPCS.2011.6070148>
- Bousonville, M. et al. New phase stable optical fiber. In *Proc. Beam Instrumentation Workshop 2012 (BIW2012) 101–103* (JACoW, 2012).
- 40-Gigabit-Capable Passive Optical Networks 2 (NG-PON2)* G.989.2 (ITU, 2015); <https://www.itu.int/rec/T-REC-G.989.2-201902-I/en>
- Aijaz, A., Dohler, M., Aghvami, A. H., Friderikos, V. & Frodigh, M. Realizing the tactile internet: haptic communications over next generation 5G cellular networks. *IEEE Wirel. Commun.* **24**, 82–89 (2016).
- Takesue, H. et al. Quantum key distribution over a 40-dB channel loss using superconducting single-photon detectors. *Nat. Photon.* **1**, 343–348 (2007).
- Torres-Company, V. & Weiner, A. M. Optical frequency comb technology for ultra-broadband radio-frequency photonics. *Laser Photon. Rev.* **8**, 368–393 (2014).
- Kuo, B. P.-P., Myslivets, E., Alic, N. & Radic, S. Wavelength multicasting via frequency comb generation in a bandwidth-enhanced fiber optical parametric mixer. *J. Lightwave Technol.* **29**, 3515–3522 (2011).
- Zhang, M. et al. Broadband electro-optic frequency comb generation in a lithium niobate microring resonator. *Nature* **568**, 373–377 (2019).

## Acknowledgements

We acknowledge financial support from Microsoft, Inphi Inc. and EPSRC grants EP/R041792/1 and EP/R035342/1 and Royal Society Paul Instrument Fund PIF/R1/180001. Eblana Photonics provided the lasers used in this work. We thank P. Watts for helpful discussion at the early stages of the work and E. Vohnhof for assistance in the generation of the figures.

## Author contributions

K.A.C., P.C., H.B., I.H., K.J., B.T., H.W. and T.G. conceived the concept of clock phase caching, which was later refined with help from K.S., D.C. and Z.L. K.A.C. and Z.L. conceived and constructed the experimental set-up. K.A.C. implemented clock phase caching in the experiment. K.A.C. designed and implemented all FPGA hardware code. K.A.C. led the experiment and collected all experimental results, supervised by Z.L., with support provided by P.C., H.B., B.T., P.B. and G.Z. P.C. and H.B. collected data-centre traffic and performed the optical switch network utilization analysis. K.A.C., Z.L., P.B., P.C. and H.B. wrote and revised the manuscript. All authors discussed the results and commented on the manuscript.

## Competing interests

A patent application, entitled 'Phase Caching for Fast Data Recovery', has been filed by Microsoft Technology Licensing, LLC with the US Patent and Trademark Office on 27 October 2017, on the technology described in this Article. This patent is currently pending (patent no. US20190132112A1).

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41928-020-0423-y>.

**Correspondence and requests for materials** should be addressed to K.A.C., P.C. or Z.L.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature Limited 2020