

# Countering detector manipulation attacks in quantum communication through detector self-testing

Lijiong Shen<sup>1</sup> and Christian Kurtsiefer<sup>1,2</sup>

<sup>1</sup>Centre for Quantum Technologies, National University of Singapore, 3 Science Drive 2, Singapore 117543

<sup>2</sup>Department of Physics, National University of Singapore, 2 Science Drive 3, Singapore 117551\*

(Dated: July 17, 2024)

In practical quantum key distribution systems, imperfect physical devices open security loopholes that challenge the core promise of this technology. Apart from various side channels, a vulnerability of single-photon detectors to blinding attacks has been one of the biggest concerns, and has been addressed both by technical means as well as advanced protocols. In this work, we present a countermeasure against such attacks based on self-testing of detectors to confirm their intended operation without relying on specific aspects of their inner working, and to reveal any manipulation attempts. We experimentally demonstrate this countermeasure with a typical InGaAs avalanche photodetector, but the scheme can be easily implemented with any single photon detector.

*Introduction* – Quantum key distribution (QKD) is a communication method that uses quantum states of light as a trusted courier such that any eavesdropping attempt in this information transmission is revealed as part of the underlying quantum physics of the measurement process on the states [1–3]. While the basic protocols are secure within their set of assumptions, practical QKD systems can exhibit vulnerabilities through imperfect implementation of the original protocol scenarios, through imperfect preparation and detection devices, or through side channels that leak information out of the supposedly safe perimeter of the two communication partners [4–6]. Families of such vulnerabilities have been identified and addressed through technical measures and advanced protocols. Examples are the photon number splitting attacks where single photons were approximated by faint coherent pulses [7, 8], Trojan horse attacks [3, 9], various timing attacks [10–12] and classes of information leakage into parasitic degrees of freedom.

A critical vulnerability of QKD systems is the detector blinding / fake state attack family on single-photon detectors [13]. This attack has been experimentally demonstrated to work for detectors based on avalanche photodiodes and superconducting nanowires [14–16], and allowed to completely recover a key generated by QKD without being noticed by the error detection step in a QKD implementation [17]. The attack is based on the fact that these single photon detectors can be blinded by a macroscopic light level into not giving any response, while an even stronger light pulse or a recovery event from a blinded state could create an output signal from the blinded detector that emulates a photon detection event [13] (see Fig. 1). This vulnerability can be exploited by carrying out an undetected man-in-the-middle attack, where an eavesdropper intercepts photon states carrying the information, measures the quantum state in a basis of his/her choice, and copies the measurement results into the pho-

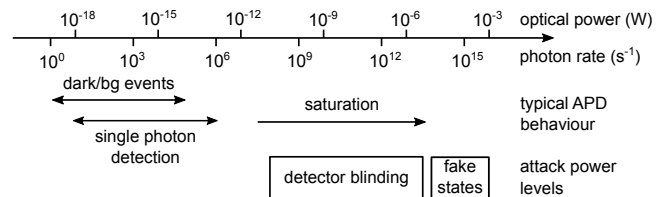


FIG. 1: Single photon avalanche photodiode properties underlying a blinding/fake state attack. At light levels below  $10^{-12}$  W, these devices respond with detection events that can be used to identify single photons. At higher power levels, they saturate and can eventually be brought into a blinded mode where they are not susceptible anymore to additional single photons. Very bright short pulses of light (“fake states”) can lead to a detector response that is indistinguishable from single photon response at low light levels. Photon rate/power level scaling is shown for a wavelength of 1300 nm.

ton detector of the legitimate receiver with macroscopic powers of light.

Various countermeasures against the detector control attack have been suggested and implemented. One class of countermeasures addresses technical aspects of the detectors. Examples are using more than one detector or a multi-pixel detector for one measurement basis [18–21], including a watchdog detector for the blinding light [14, 22], effectively varying the detector efficiency at random timings [23, 24], and carefully monitoring the photocurrent or breakdown status of the detector [25, 26] to identify a detector manipulation. However, most of these countermeasures have operational drawbacks. For example, additional single photon detectors significantly increase the overall cost and complexity, and beam splitters in the receiver for watchdog detectors introduce additional optical losses. Varying the efficiency frequently to get enough statistics to identify the blinding attack could significantly affect the QKD bit rate, and changing the detector operation condition or monitoring its state increases the complexity of the electronic circuitry around the single photon detectors. Such countermea-

\*christian.kurtsiefer@gmail.com

sures may also introduce additional vulnerabilities that may be exploited in an arms race style [27].

An elegant countermeasure on the protocol level is provided by the so-called measurement-device independent quantum key distribution (MDI-QKD) [28], which further developed the idea of device-independent QKD where a photon pair source can be made public or even controlled by an eavesdropper [29] to a scenario where the detectors receiving single photons (or approximations thereof) can be public, or controlled by an eavesdropper. The scheme has been demonstrated experimentally several times by now [30–33]. It requires a pair of single photons (or weak coherent pulses) from two communication partners without a phase correlation to arrive within a coherence time on a Bell state analyzer, where single photon detection is carried out, and the result is published. This requires a matching of emission times and wavelengths of two spatially separated light sources with both communication partners.

The MDI-QKD approach counteracts any active or passive attack on single photon detectors, as their result need not to be private anymore. The communication partners can simply test if the detectors were performing single photon detection through a error detection process similar to the original QKD protocols.

In this work, we present a method of testing the proper operation of single photon detectors in a QKD scenario that does not require the synchronization of light sources like in the MDI-QKD approach, while also not touching the specific detector mechanism. It brings the idea of self-testing of quantum systems [34–36] to single photon detectors that can remain black boxes. We use a light emitter (LE) under control of a legitimate communication partner that is weakly coupled to its single photon detector for this self-testing. When the single photon detector is under a blinding attack, it is insensitive to low-intensity light fields used for quantum key distribution. Thus, by turning on the LE at times not predictable by an eavesdropper, “salt” optical detection events are generated in the detector when it operates normal, while it does not react to the test light when blinded. Complementary, the test light intensity can be raised to blinding levels of the photodetector, which thereby is desensitized to legitimate single photons. Registration of any detector events under self-blinding then suggests the presence of fake state events.

*Self-testing strategy* – In a generic QKD system, a transmitter generates photons containing quantum information in either polarization or time encoding, and sends them through an optical path (“quantum channel”) to a receiver. Therein, a measurement basis choice is made either through passive or active optical components, and the light arriving from the quantum channel is directed to single photon detectors. In a blinding/fake state attack, an eavesdropper measures a photon in the quantum channel, and copies the result into the corresponding photon detector of the legitimate receiver using blinding and fake state light levels. For detector testing, a light emit-

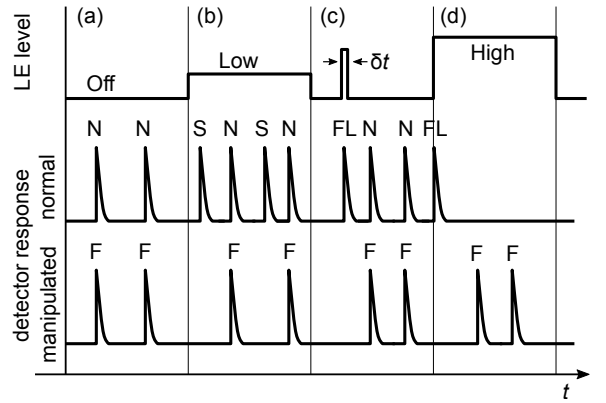


FIG. 2: Detector self testing. Top trace: light level of the light emitter LE, middle trace: normal detector response (no manipulation), lower trace: detector response under manipulation. Detector events are classified as normal (N), salt (S), “fake” (F), and flag (FL) signals. Segment (a) shows responses without self-testing, (b) with low LE power generating salt events, (c) with occasional test pulses at medium power, (d) with high LE power to self-blind the detector.

ter (LE) in the receiver is controlled by a random number generator and weakly coupled to the single photon detectors.

An unblinded single-photon detector generates events due to photons from the legitimate source or the background (labeled “N” in Fig. 2(a)). The brightness of the legitimate source, the transmission of the quantum channel, the efficiency of the single photon detectors, and the detector dark count rate determine the average number  $\bar{n}$  of the photon-detection events registered in a time interval  $T$ . An eavesdropper would choose a rate of “fake” detection events (labeled “F” in Fig. 2(a)) similar to normal QKD operation to prevent detecting the attack by monitoring photon detection statistics.

We illustrate three different examples of detector self-testing to detect detector manipulation attacks.

In the first one, the legitimate receiver switches occasionally the light emitter LE to a low light level for a test time interval  $T$  at a random timing unpredictable by an eavesdropper, while it is off for the rest of the time. In the test interval, an unblinded detector would see an increase of detector events above  $\bar{n}$  due to additional salt events (“S” in Fig. 2(b)). The legitimate receiver has complete control of the light emitter to make excess photon detection events statistically detectable in the probe interval  $T$ . A single photon detector under blinding attack would be insensitive to the low light levels of LE, so only detector events generated by positive detector manipulations like fake states would be registered (labeled “F” in Fig. 2(b)). A statistically significant presence of salt events in a time interval  $T$  would therefore allow to sense a negative detector manipulation e.g. through blinding. Note that the test interval  $T$  does not need to be distributed contiguously in time.

This leads to a second self-testing example, which turns on the light emitter for a short pulse time interval  $\delta t$  at a random timing and with a high enough energy (a few photons) to cause a detection event with almost unit probability in an unblinded single-photon detector. A blinded detector is again insensitive to such a short optical pulse as long as the light level is way below the fake state threshold. In this situation, detecting a single flag event can witness a non-blinded detector (see Fig. 2(c)).

The third self-testing example uses the light emitter in the receiver to locally blind the detector. The typical power necessary to blind an APD is on the order of a few nW, which can easily be accomplished by weakly coupling even faint light sources like LEDs. Detection events caused by single photons from the legitimate source will be suppressed by the local blinding light. In absence of a negative detector manipulation (e.g. detector blinding), the intense light at the onset of the self-blinding period will almost deterministically create a flag event in the detector, which then remains silent during the rest of the self-blinding interval (see Fig. 2(d)). However, any positive detector manipulation will overrule the local blinding, and cause a false detection event. Both the initial flag event and any possible later event can be easily checked. This method only requires a small number of registered events in a time interval  $T$  to discover both a negative and positive detector manipulation attack.

A detector event could also be triggered when the detector recovers from a (remote) blinding exposure [37]. Local blinding will suppress such “fake” detector events, so they may not get noticed by looking for signals under local blinding. However, in such a case, the flag event will also be suppressed. Therefore, a combination of checking for detection events during self-blinding and looking for a flag event is necessary to identify such an attack.

*Experimental results* – We demonstrate our countermeasure with a single-photon detector commonly used in quantum key distribution which is susceptible to manipulation attacks (see Fig. 3(a)). Light that simulates legitimate quantum signals and provides the larger power levels required for detector manipulation is generated by combining the output of a continuous wave (cw) laser diode (LD1) with light from a pulsed laser diode (LD2) on a fiber beam splitter (BS). The 2 ns long bright fake states from LD2 can be emitted upon detection events from an auxiliary avalanche photodetector (APD1) to emulate a credible (Poissonian) event distribution. On the receiver side, the light from the quantum channel passes through an interference filter (IF) before it is focused onto the main photodetector (APD2), a passively quenched InGaAs device (S-Fifteen Instruments IRSPD1) with a maximal count rate of  $5 \times 10^5 \text{ s}^{-1}$  and a dark count rate of  $7 \times 10^3 \text{ s}^{-1}$ . The light emitter (LE) for detector self-testing is a light emitting diode with a center wavelength of 940 nm (Vishay VSLY5940), which is reflected off the IF (acting as a dichroic beam splitter) onto APD2.

For the demonstration, we consider an event rate of  $\approx 5 \times 10^4 \text{ s}^{-1}$  at APD2, which is about an order of mag-

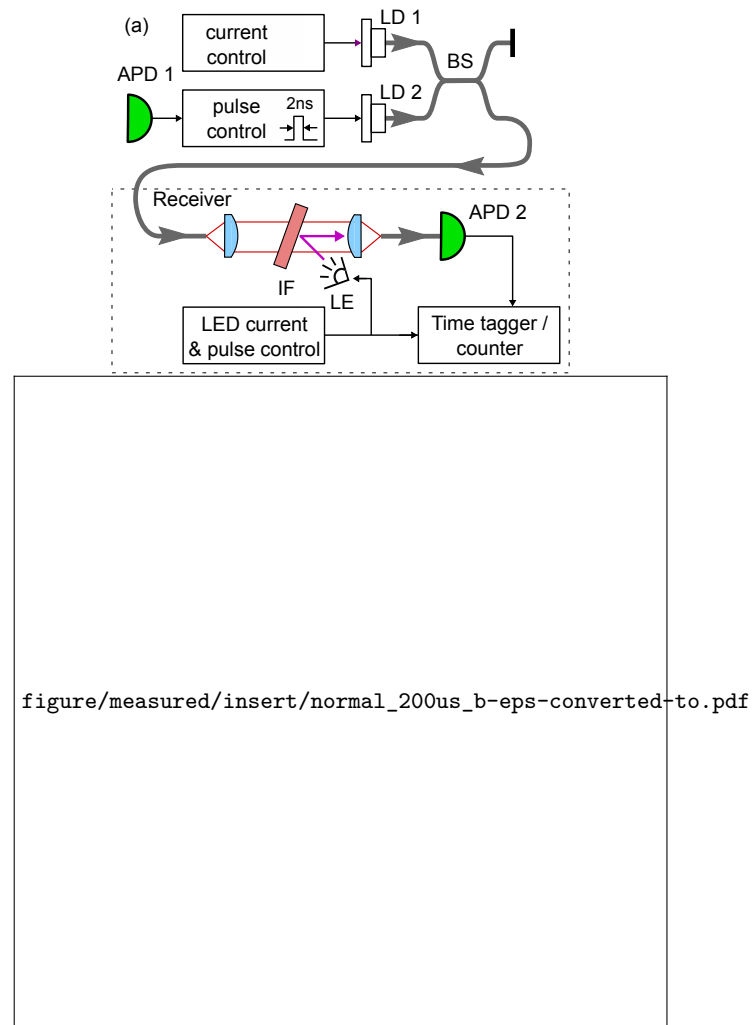


FIG. 3: (a) Setup to demonstrate detector self-testing. Light from a cw laser diode (LD1) and pulsed laser diode (LD2), both around 1310 nm, is combined in a fiber beamsplitter (BS) to simulate different illumination scenarios. Besides the single photon InGaAs detector APD2, the receiver contains an LED (940 nm) as a light emitter (LE) for local testing of APD2. An interference (IF) filter prevents leakage of LE light out of the receiver. (b) Distribution of photodetection events in a time window of  $T = 200 \mu\text{s}$  under “normal” operation under illumination of the detector with a low power level from LD1.

nitude below the maximal detection rate to not reduce the detector efficiency significantly. Figure 3(b) shows a histogram of detection events in a time interval of  $T = 200 \mu\text{s}$  generated by choosing an appropriate light level of LD1. The result with a mean photodetection number  $\bar{n} \approx 10$  differs slightly from a Poisson distribution since the detector has an after-pulse possibility of about 40%. To implement a detector manipulation with the same event characteristic, we elevate the optical output power of LD1 to 500 pW, the minimal power to completely blind detector APD2. Fake states that emulate photodetection events in APD2 are generated with

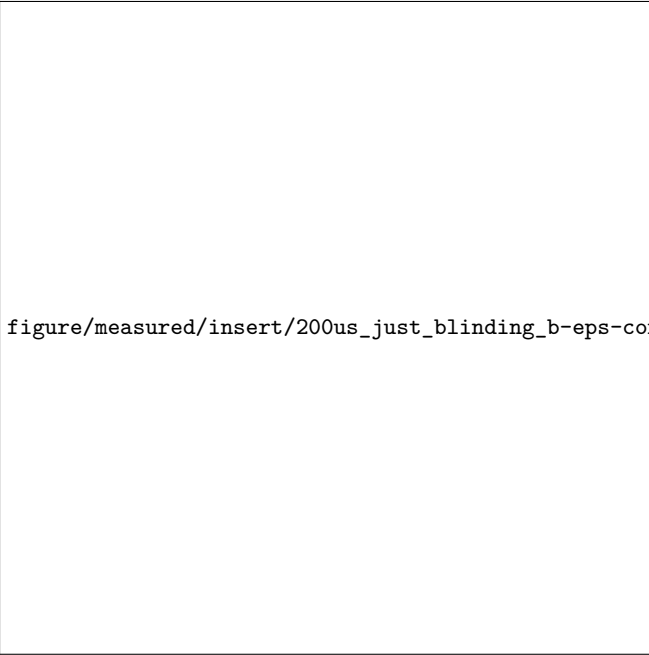


FIG. 4: Distribution of detector events in the presence of self-seeding light in a test interval of  $T = 200 \mu\text{s}$  for a normally operating, and a manipulated detector. The manipulated detector shows a similar distribution as the one in Fig. 3(b), while the normally operating detector shows a distinctly higher event number. Error bars indicate Poissonian standard deviations resulting from 7432 and 7686 test runs for a normal detector and a manipulated detector, respectively.

optical pulses through LD2 with a peak power of  $3 \mu\text{W}$ .

To demonstrate the first example of detector self-testing, we turn on the light emitter LE in the test interval  $T$  both for a normally operating and a manipulated detector. The resulting detection event distributions are shown in Fig. 4. For a normally operating detector, the observed APD2 events in the test interval increase significantly to a mean of about  $\bar{n}_{T1} \approx 100$ , while for a manipulated detector, the distribution is similar to the “normal” distribution with  $\bar{n}_N \approx 10$  in Fig. 3(b). With a threshold at  $n = 50$ , the two distributions can be easily distinguished, and a detector manipulation attempt (specifically: the presence of a blinding light level) easily identified in a single measurement interval  $T$ ; in the experiment, the unmanipulated detector never showed less than 78 events, while the manipulated showed never more than 30 events.

The necessary time to detect a manipulated detector can be shortened even further with the second example of self-testing. We demonstrate this by driving the light emitter LE to emit  $\delta t = 25 \text{ ns}$  long pulses, and increasing the coupling to the detector APD2 compared to the previous example. Figure 5 shows the probability of registering a signal from APD2 as a function of the time  $\Delta t$  after the start of the self-testing pulse. A non-manipulated detector shows an overall detector response

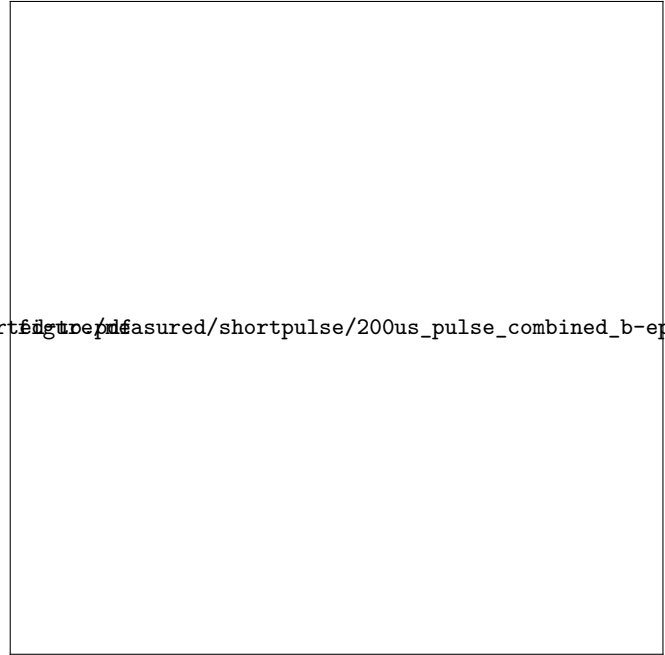


FIG. 5: Detector event probability for a  $25 \text{ ns}$  long bright pulse of the self-testing light emitter LE for a manipulated and normal detector vs the time difference  $\Delta t$  between detector event and a self testing pulse edge. A non-manipulated detector reacts with an event with high probability within less than  $60 \text{ ns}$ . Optical and electrical delays shift the detector response away from  $\Delta t = 0$ , and error bars indicate Poissonian standard deviations resulting from 12542 and 12380 test runs for normal detector and manipulated detector, respectively.

probability  $p_1 = 93.4\%$  within  $60 \text{ ns}$  (11720 photon detection events out of 12542 optical pulse). This number does not reach  $100\%$ , as the detector may have been in a recovery state from a previous detection event. For a manipulated detector, i.e., in presence of both detector blinding and fake states, we find an integral detector event probability  $p_2 = 0.3\%$  (36 out of 12380 test pulses). These events were caused by fake states, not by light from the LE. Detector manipulation (specifically, the detector blinding) can therefore be identified with a few short test pulses to a very high statistical significance.

To demonstrate the third example of detector self-testing, we increased the optical power of LE on detector APD2 to a level that it could reliably blind the detector. Figure 6 shows both a distribution of detection events in a test interval  $T = 200 \mu\text{s}$ , taken  $60 \text{ ns}$  after the onset of light emission by LE. The un-manipulated detector is insensitive to single photons in this interval; we observed only 8 events in 7608 test runs (likely due to electrical noise), while a manipulated detector still reported events due to fake states present at the input; we observed 7655 out of 7658 events (with the missing events compatible with statistics). The onset of the test light emission triggered a detector reaction within the first  $60 \text{ ns}$  with a probability  $p_1 = 97.6\%$  (7426 detector events out of 7608 test runs, see inset of Fig. 6) for a non-manipulated detec-

figure/measured/selfb/selfb\_200us\_combined\_b-eps-converted-to.pdf

FIG. 6: Detector event distribution in a test interval  $T = 200 \mu\text{s}$  in the presence of self-blinding light for a normal and manipulated detector, registered 60 ns *after* the onset of the self-blinding light. A manipulated detector still reports events due to fake states. Inset: probability of a detector event in the first 60 ns after switching on the self-blinding light. This scheme allows to detect the presence of both a blinding and fake state detector manipulation.

tor, while the probability of an onset event was  $p_2 = 0.2\%$  (17 out of 7658 runs) for a manipulated detector caused by fake states. A local light emitter that is able to self-

blind the detector is thus able to reveal the presence of both blinding and the fake state in a detector manipulation attempt.

*Conclusion* – We demonstrated self-testing of single photon detectors that can reliably reveal manipulation attacks. The self-testing strategy relies on a light source near the detector under possible external manipulation, and is able to detect both negative manipulations (i.e. suppression of single photon detections) and positive manipulations (i.e., generating detector events that are not caused by single photon detections) in a relatively short time with a high statistical significance. The detector self-testing makes no assumptions on the nature of the manipulation attack of the detector, and thus also covers manipulations that are not of the known nature like detector blinding and fake states. As the self-testing can be accomplished by a relatively simple light source (as long as this is outside the control and knowledge of an adversary), this scheme can address one of the most significant hardware vulnerabilities of QKD systems in a significantly simpler way as compared to device-independent or measurement-device independent approaches, and may even be a suitable to retrofit existing QKD systems to make them resilient against detector manipulation attacks.

### Acknowledgements

This work was supported by the National Research Foundation (NRF) Singapore (Grant: QEP-P1), and through the Research Centres of Excellence programme supported by NRF Singapore and the Ministry of Education, Singapore.

- 
- [1] C. H. Bennett and G. Brassard, Proceedings of IEEE International Conference on Computers, Systems and Signal Processing , 175 (1984).
  - [2] A. K. Ekert, Phys. Rev. Lett. **67**, 661 (1991).
  - [3] N. Gisin, G. Ribordy, W. Tittel, and H. Zbinden, Rev. Mod. Phys. **74**, 145 (2002).
  - [4] V. Scarani, H. Bechmann-Pasquinucci, N. J. Cerf, M. Dušek, N. Lütkenhaus, and M. Peev, Rev. Mod. Phys. **81**, 1301 (2009).
  - [5] V. Scarani and C. Kurtsiefer, Theor. Comput. Sci. **560**, 27 (2014).
  - [6] F. Xu, X. Ma, Q. Zhang, H. K. Lo, and J. W. Pan, Rev. Mod. Phys. **92**, 025002 (2020).
  - [7] G. Brassard, N. Lütkenhaus, T. Mor, and B. C. Sanders, Phys. Rev. Lett. **85**, 1330 (2000).
  - [8] N. Lütkenhaus, Phys. Rev. A **61**, 10 (2000).
  - [9] A. Vakhitov, V. Makarov, and D. R. Hjelle, J. Mod. Opt. **48**, 2023 (2001).
  - [10] B. Qi, C. H. F. Fung, H. K. Lo, and X. Ma, Quantum Inf. Comput. **7**, 73 (2007).
  - [11] A. Lamas-Linares and C. Kurtsiefer, Opt. Express **15**, 9388 (2007).
  - [12] Y. Zhao, C. H. F. Fung, B. Qi, C. Chen, and H. K. Lo, Phys. Rev. A **78**, 042333 (2008).
  - [13] V. Makarov, New J. Phys. **11**, 065003 (2009).
  - [14] L. Lydersen, C. Wiechers, C. Wittmann, D. Elser, J. Skaar, and V. Makarov, **4**, 686 (2010).
  - [15] L. Lydersen, M. K. Akhlaghi, A. H. Majedi, J. Skaar, and V. Makarov, New J. Phys. **13**, 113042 (2011).
  - [16] G. Goltsman, M. Elezov, R. Ozhegov, V. Makarov, G. Goltsman, V. Makarov, V. Makarov, V. Makarov, and V. Makarov, Opt. Express, Vol. 27, Issue 21, pp. 30979-30988 **27**, 30979 (2019).
  - [17] I. Gerhardt, Q. Liu, A. A. Lamas-Linares, J. Skaar, C. Kurtsiefer, and V. Makarov, Nat. Commun. **2**, 1 (2011).
  - [18] T. Honjo, M. Fujiwara, K. Shimizu, K. Tamaki, S. Miki, T. Yamashita, H. Terai, Z. Wang, and M. Sasaki, Opt. Express **21**, 2667 (2013).
  - [19] T. Ferreira Da Silva, G. C. Do Amaral, G. B. Xavier, G. P. Temporao, and J. P. Von Der Weid, IEEE J. Sel. Top. Quantum Electron. **21**, 159 (2015).
  - [20] J. Wang, H. Wang, X. Qin, Z. Wei, and Z. Zhang, Eur. Phys. J. D **70**, 5 (2016).

- [21] G. Gras, D. Rusca, H. Zbinden, and F. Bussi eres, *Phys. Rev. Appl.* **15**, 034052 (2021).
- [22] A. R. Dixon, J. F. Dynes, M. Lucamarini, B. Fr ohlich, A. W. Sharpe, A. Plews, W. Tam, Z. L. Yuan, Y. Tanizawa, H. Sato, S. Kawamura, M. Fujiwara, M. Sasaki, and A. J. Shields, *Sci. Rep.* **7**, 1 (2017).
- [23] C. C. W. Lim, N. Walenta, M. Legr e, N. Gisin, and H. Zbinden, *IEEE J. Sel. Top. Quantum Electron.* **21**, 192 (2015).
- [24] Y.-J. Qian, D.-Y. He, S. Wang, W. Chen, Z.-Q. Yin, G.-C. Guo, and Z.-F. Han, *Optica* **6**, 1178 (2019).
- [25] Z. L. Yuan, J. F. Dynes, and A. J. Shields, “Avoiding the blinding attack in QKD,” (2010).
- [26] Z. L. Yuan, J. F. Dynes, and A. J. Shields, *Appl. Phys. Lett.* **98**, 231104 (2011).
- [27] A. Huang, S. Sajeed, P. Chaiwongkhot, M. Soucarros, M. Legr e, and V. Makarov, *IEEE Journal of Quantum Electronics* **52** (2016), 10.1109/JQE.2016.2611443.
- [28] H.-K. Lo, M. Curty, and B. Qi, *Phys. Rev. Lett.* **108**, 130503 (2012).
- [29] S. Pironio, A. Ac ın, N. Brunner, N. Gisin, S. Massar, and V. Scarani, **11**, 045021 (2009).
- [30] Y. Liu, T.-Y. Chen, L.-J. Wang, H. Liang, G.-L. Shentu, J. Wang, K. Cui, H.-L. Yin, N.-L. Liu, L. Li, X. Ma, J. S. Pelc, M. M. Fejer, C.-Z. Peng, Q. Zhang, and J.-W. Pan, *Physical Review Letters* **111**, 130502 (2013).
- [31] Y.-L. Tang, H.-L. Yin, S.-J. Chen, Y. Liu, W.-J. Zhang, X. Jiang, L. Zhang, J. Wang, L.-X. You, J.-Y. Guan, D.-X. Yang, Z. Wang, H. Liang, Z. Zhang, N. Zhou, X. Ma, T.-Y. Chen, Q. Zhang, and J.-W. Pan, *Physical Review Letters* **113**, 190501 (2014).
- [32] H.-L. Yin, T.-Y. Chen, Z.-W. Yu, H. Liu, L.-X. You, Y.-H. Zhou, S.-J. Chen, Y. Mao, M.-Q. Huang, W.-J. Zhang, H. Chen, M. J. Li, D. Nolan, F. Zhou, X. Jiang, Z. Wang, Q. Zhang, X.-B. Wang, and J.-W. Pan, *Physical Review Letters* **117**, 190501 (2016).
- [33] H. Liu, W. Wang, K. Wei, X.-T. Fang, L. Li, N.-L. Liu, H. Liang, S.-J. Zhang, W. Zhang, H. Li, L. You, Z. Wang, H.-K. Lo, T.-Y. Chen, F. Xu, and J.-W. Pan, *Physical Review Letters* **122**, 160501 (2019).
- [34] D. Mayers and A. C.-C. Yao, *Proceedings 39th Annual Symposium on Foundations of Computer Science (Cat. No.98CB36280)*, 503 (1998).
- [35] W. van Dam, F. Magniez, M. Mosca, and M. Santha, in *Proceedings of the Thirty-Second Annual ACM Symposium on Theory of Computing*, STOC ’00 (Association for Computing Machinery, New York, NY, USA, 2000) p. 688696.
- [36] I.  upic and J. Bowles, *Quantum* **4**, 337 (2020).
- [37] M. G. Tanner, R. H. Hadfield, and V. Makarov, *Opt. Express*, Vol. 22, Issue 6, pp. 6734-6748 **22**, 6734 (2014).