

## A link layer protocol for quantum networks

Dahlberg, Axel; Skrzypczyk, Matthew; Coopmans, Tim; Wubben, Leon; Rozpdek, Filip; Pompili, Matteo; Stolk, Arian; Pawelczak, Przemyslaw; Knegjens, Rob; De Oliveira Filho, Julio

**DOI**

[10.1145/3341302.3342070](https://doi.org/10.1145/3341302.3342070)

**Publication date**

2019

**Document Version**

Accepted author manuscript

**Published in**

SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication

**Citation (APA)**

Dahlberg, A., Skrzypczyk, M., Coopmans, T., Wubben, L., Rozpdek, F., Pompili, M., Stolk, A., Pawelczak, P., Knegjens, R., De Oliveira Filho, J., Hanson, R., & Wehner, S. (2019). A link layer protocol for quantum networks. In *SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication* (pp. 159-173). (SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication). Association for Computing Machinery (ACM).  
<https://doi.org/10.1145/3341302.3342070>

**Important note**

To cite this publication, please use the final published version (if applicable).  
Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights.  
We will remove access to the work immediately and investigate your claim.

## A link layer protocol for quantum networks

Dahlberg, Axel; Skrzypczyk, Matthew; Coopmans, Tim; Wubben, Leon; Rozpdek, Filip; Pompili, Matteo; Stolk, Arian; Pawelczak, Przemyslaw; Knegjens, Rob; De Oliveira Filho, Julio

**DOI**

[10.1145/3341302.3342070](https://doi.org/10.1145/3341302.3342070)

**Publication date**

2019

**Document Version**

Accepted author manuscript

**Published in**

SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication

**Citation (APA)**

Dahlberg, A., Skrzypczyk, M., Coopmans, T., Wubben, L., Rozpdek, F., Pompili, M., Stolk, A., Pawelczak, P., Knegjens, R., De Oliveira Filho, J., Hanson, R., & Wehner, S. (2019). A link layer protocol for quantum networks. In *SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication* (pp. 159-173). (SIGCOMM 2019 - Proceedings of the 2019 Conference of the ACM Special Interest Group on Data Communication). Association for Computing Machinery (ACM).  
<https://doi.org/10.1145/3341302.3342070>

**Important note**

To cite this publication, please use the final published version (if applicable).  
Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights.  
We will remove access to the work immediately and investigate your claim.

# A Link Layer Protocol for Quantum Networks

Axel Dahlberg<sup>1,2</sup>, Matthew Skrzypczyk<sup>1,2</sup>, Tim Coopmans<sup>1,2</sup>, Leon Wubben<sup>1,2</sup>,  
Filip Rozpędek<sup>1,2</sup>, Matteo Pompili<sup>1,2</sup>, Arian Stolk<sup>1,2</sup>, Przemysław Pawełczak<sup>1</sup>,  
Robert Knegjens<sup>1</sup>, Julio de Oliveira Filho<sup>1</sup>, Ronald Hanson<sup>1,2</sup>, Stephanie Wehner<sup>1,2</sup>

<sup>1</sup>QuTech, Delft University of Technology and TNO

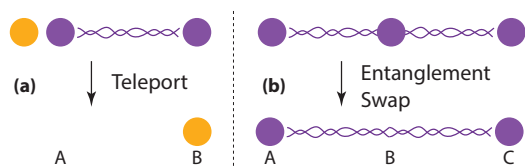
<sup>2</sup>Kavli Institute of Nanoscience, Delft University of Technology  
s.d.c.wehner@tudelft.nl

## ABSTRACT

Quantum communication brings radically new capabilities that are provably impossible to attain in any classical network. Here, we take the first step from a physics experiment to a fully fledged quantum internet system. We propose a functional allocation of a quantum network stack and construct the first physical and link layer protocols that turn ad-hoc physics experiments producing heralded entanglement between quantum processors into a well-defined and robust service. This lays the groundwork for designing and implementing scalable control and application protocols in platform-independent software. To design our protocol, we identify use cases, as well as fundamental and technological design considerations of quantum network hardware, illustrated by considering the state-of-the-art quantum processor platform available to us (Nitrogen-Vacancy (NV) centers in diamond). Using a purpose built discrete-event simulator for quantum networks, we examine the robustness and performance of our protocol using extensive simulations on a supercomputing cluster. We perform a full implementation of our protocol, where we successfully validate the physical simulation model against data gathered from the NV hardware. We first observe that our protocol is robust even in a regime of exaggerated losses of classical control messages with only little impact on the performance of the system. We proceed to study the performance of our protocols for 169 distinct simulation scenarios, including tradeoffs between traditional performance metrics such as throughput and the quality of entanglement. Finally, we initiate the study of quantum network scheduling strategies to optimize protocol performance for different use cases.

## 1 INTRODUCTION

Quantum communication enables the transmission of quantum bits (qubits) in order to achieve novel capabilities that are provably impossible using classical communication. As with any radically new technology, it is hard to predict all uses of a future Quantum Internet [62, 101], but several major applications have already been identified depending on



**Figure 1: Entanglement enables long-distance quantum communication: (a) once two qubits (purple/dark) are confirmed to be entangled (threaded links between qubits), a data qubit (yellow/light) can be sent deterministically using teleportation [11], consuming the entangled pair; (b) long-distance entanglement can be built from shorter segments: If node A is entangled with B (repeater), and B with C, then B can perform *entanglement swapping* [107] to create long-distance entanglement between the qubits at A and C.**

the stage of quantum network development [101], ranging from cryptography [10, 40], sensing and metrology [46, 63], distributed systems [9, 36], to secure quantum cloud computing [20, 43].

Qubits are fundamentally different from classical bits, which brings significant challenges both to the physical implementation of quantum networks, as well as the design of quantum network architectures. Qubits cannot be copied, ruling out signal amplification or repetition to overcome transmission losses to bridge great distances. Two qubits can share a special relation known as *entanglement*, even if these two qubits are stored at distant network nodes. Such entanglement is central not only to enable novel applications, but also provides a means to realize a quantum repeater, which enables quantum communication over long-distances (Figure 1).

At present, short-lived entanglement has been produced probabilistically over short distances ( $\approx 100$  km) on the ground by sending photons over standard telecom fiber (see e.g. [39, 55]), as well as from space over 1203km from a satellite [103]. Such systems can allow the realization of applications in the prepare-and-measure stage [101] of quantum

networks on point-to-point links, but cannot by themselves be concatenated to allow the transmission of qubits over longer distances.

In order to enable long-distance quantum communication and the execution of complex quantum applications, we would like to produce long-lived entanglement between two quantum nodes that are capable of storing and manipulating qubits. To do so efficiently (Section 3.1), we need to confirm entanglement generation by performing *heralded* entanglement generation. This means that there is a *heralding signal* that tells us if we have been successful in an attempt to generate entanglement.

The current world distance record for producing such entanglement is 1.3km, which has been achieved using a solid state platform known as Nitrogen-Vacancy (NV) centres in diamond [49]. Intuitively, this platform is a few qubit (as of now maximum 8 [21]) quantum computer capable of arbitrary quantum gates, with an optical interface for initialization, measurement and entanglement generation. Key capabilities of the NV platform have already been demonstrated, including qubit lifetimes of 1.46 s [3], entanglement production faster than it is lost [53], and using entanglement to teleport qubits between separated NV centres [78]. Other hardware platforms exist that are identical on an abstract level (quantum computer with an optical interface), and on which heralded long-lived entanglement generation has been demonstrated (e.g. Ion Traps [70], and Neutral Atoms [50]). Theoretical proposals and early stage demonstrations of individual components also exists for other physical platforms (e.g. quantum dots [35], rare earth ion-doped crystals [97], atomic gases [26, 57], and superconducting qubits [74]), but their performance is not yet good enough to generate entanglement faster than it is lost.

Up to now, the generation of long-lived entanglement has been the domain of highly sophisticated, but arguably ad-hoc physics experiments. We are now on the verge of seeing early stage quantum networks becoming a reality, entering a new phase of development which will require a joint effort across physics, computer science and engineering to overcome the many challenges in scaling such networks. In this paper, we take the first step from a physics experiment to a fully-fledged quantum communication *system*.

**Design considerations and use cases:** We identify general design considerations for quantum networks based on fundamental properties of entanglement, and technological limitations of near-term quantum hardware, illustrated with the example of our NV platform. For the first time, we identify systematic use cases, and employ them to guide the design of our stack and protocols.

**Functional allocation quantum network stack:** We propose a functional allocation of a quantum network stack, and define the service desired from its link layer to satisfy

| Application |                                 |
|-------------|---------------------------------|
| Transport   | Qubit transmission              |
| Network     | Long distance entanglement      |
| Link        | Robust entanglement generation  |
| Physical    | Attempt entanglement generation |

**Figure 2: Functional allocation in a quantum network stack. Entanglement is an inherent connection between quantum bits, which contrasts with classical networking where shared state is established at higher layers.**

use case requirements and design considerations. In analogy to classical networking, the link layer is responsible for producing entanglement between two nodes that share a direct physical connection (e.g. optical fiber).

**First physical and link layer entanglement generation protocols:** We proceed to construct the world’s first physical and link layer protocols that turn ad-hoc physics experiments producing heralded entanglement into a well defined service. This lays the groundwork for designing and implementing control and application protocols in platform independent software in order to build and scale quantum networks. At the physical layer, we focus primarily on the quantum hardware available to us (NV platform) but the same protocol could be realized directly using Ion Traps or Neutral Atoms, as well as —with small changes— other means of producing physical entanglement [89]. Our link layer protocol takes into account the intricacies of the NV platform, but is in itself already platform independent.

**Simulation validated against quantum hardware:** Using a purpose built discrete-event simulator for quantum networks, we examine the robustness and performance of our protocol using more than 169 scenarios totaling 94244h wall time and 707h simulated time on a supercomputing cluster. To this end, we perform a complete implementation of our protocols and let them use simulated quantum hardware and communication links. To illustrate their performance, we consider two concrete short and long-distance scenarios based on the NV platform: (1) LAB where the nodes *A* and *B* are 2m apart. Since this setup has already been realized, we can use it to compare the performance of the entanglement generation implemented on real quantum hardware against the simulation to validate its physical model, and (2) a planned implementation of QL2020 where *A* and *B* are in two European cities separated by  $\approx 25$  km over telecom fiber. Next, to investigate trade-offs between traditional performance metrics (e.g. throughput or latency) and genuinely

quantum ones (fidelity, Section 4.2), we take a first step in examining different quantum network scheduling strategies to optimize performance for different use cases.

## 2 RELATED WORK

At present there is no quantum network stack connected to quantum hardware, no link layer protocols have been defined to produce entanglement, and no quantum networks capable of end-to-end qubit transmission or entanglement production have been realized (see [101] and references therein).

A functional allocation of a stack for quantum repeaters and protocols controlling entanglement distillation (a procedure to increase the quality of entanglement) has been outlined in [6, 68, 69, 100], which is complementary to this work. This is very useful to ultimately realize entanglement distillation, even though no concrete control protocols or connection to a hardware system were yet given. We remark that here we do not draw layers from specific protocols like entanglement distillation, but focus on the service that these layers should provide (a layer protocol may of course choose distillation as a mean to realize requirements). An outline of a quantum network stack was also put forward in [80], including an appealing high level quantum information theory protocol transforming multi-partite entanglement. However, this high level protocol does not yet consider failure modes, hardware imperfections, nor the requirements on entanglement generation protocols and the impact of classical control. Plans to realize the physical layer of a quantum network from a systems view were put forward in [65], however development has taken a different route.

In the domain of single-use point-to-point links for quantum key distribution (QKD), software has been developed for trusted repeater networks [101] to make use of such key in e.g. VoIP [64]. However, these do not allow end-to-end transmission of qubits or generation of entanglement, and rely on trust in the intermediary nodes who can eavesdrop on the communication. Control using software defined networks (SDN) to assist trusted repeater nodes has been proposed, e.g. [81, 104]. These QKD-centric protocols however do not address control problems in true quantum networks aimed at end-to-end delivery of qubits, and the generation of long-lived entanglement.

In contrast, classical networking knows a vast literature on designing and analyzing network protocols. Some ideas can indeed be borrowed from classical networking such as scheduling methods, but fundamental properties of quantum entanglement (Section A.2), as well as technological considerations of quantum hardware capabilities (Section 4.5) call for new protocols and methods of network control and management. Naturally, there is a continuous flow of systems papers proposing new networking architectures, e.g. for SDN [17],

data center networks [48], content delivery networks [24] or cloud computing [106], to name a few. Yet, we are unaware of any system-level papers proposing a quantum network stack including protocols for concrete hardware implementations.

## 3 DESIGN CONSIDERATIONS FOR QUANTUM NETWORK ARCHITECTURES

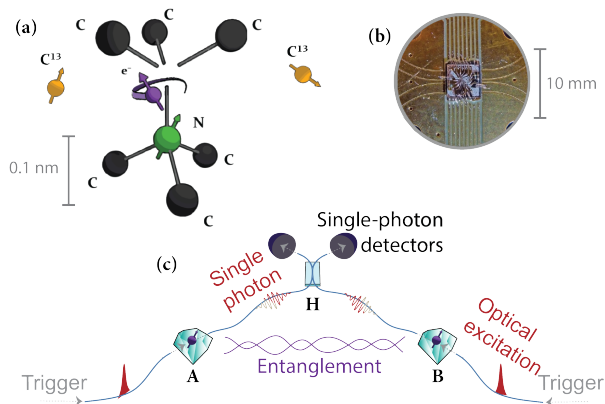
We first discuss design considerations of quantum networks themselves, followed by considerations specific to the physical and link layer (Section 4). These can be roughly subdivided into three categories: (i) fundamental considerations due to quantum entanglement, (ii) technological limitations of near-term quantum hardware, and (iii) requirements of quantum protocols themselves.

### 3.1 Qubits and Entanglement

We focus on properties of entanglement as relevant for usage and control (see Appendix and [76]). Teleportation [11] allows entanglement to be used to send qubits (see Figure 1). We will hence also call two entangled qubits an *entangled link* or *entangled pair*. Teleportation consumes the entangled link, and requires two additional classical bits to be transmitted per qubit teleported. Already at the level of qubit transmission we hence observe the need for a close integration between a quantum and classical communications. Specifically, we will need to match quantum data stored in quantum devices, with classical control information that is sent over a separate physical medium, akin to optical control plane architectures for classical optical networks [92]. To create long-distance entanglement, we can first attempt to produce short-distance entangled links, and then connect them to form longer distance ones [19, 72] via an operation known as entanglement swapping (see Figure 1). This procedure can be used iteratively to create entanglement along long chains, where we remark that the swapping operations can in principle be performed in parallel. From a resource perspective, we note that to store entanglement, both nodes need to store one qubit per entangled link. Proposals for enabling quantum communication by forward communication using quantum error correction also exist, which avoid entanglement swapping [71]. However, these have arguably much more stringent requirements in terms of hardware, putting them in a technologically more distant future: they require the ability to create entangled states consisting of a large number of photons (only 10 realized today [45]) and densely placed repeater stations performing near perfect operations [73].

Producing heralded entanglement does however allow long-distance quantum communication without the need to create entanglement consisting of many qubits. Here, the





**Figure 3: Heralded entanglement generation on the NV platform.** (a) NV centres are point defects in diamond with an electronic spin as a communication qubit (purple) and carbon-13 nuclear spins as memory qubits (yellow), realized in custom chips (b). (c) A trigger produces entanglement between the communication qubits of  $A$  and  $B$  (diamonds) and two qubits (photons) traveling over fiber to the heralding station  $H$ .  $H$  measures the photons by observing clicks in the left or right detector giving the *heralding signal*  $s$ : [failure] (none or both click), [success,  $|\Psi^+\rangle$ ] (left clicks), [success,  $|\Psi^-\rangle$ ] (right clicks). Success confirms one of two types of entangled pairs  $|\Psi^+\rangle$  or  $|\Psi^-\rangle$  (wiggly purple line).  $H$  sends  $s$  to  $A$  and  $B$  (not pictured).

heralding signal (see Figure 3) provides a confirmation that an entanglement generation attempt has succeeded. Such heralding allows long-distance quantum communication without exponential overheads [19], and without the need for more complex resources [8, 22]. Creating long-distance links between two controllable nodes by means of entanglement swapping (Section 3.2), and executing complex applications requires both nodes to know the state of their entangled links (which qubits belong to which entangled link, and who holds the other qubit of the entangled pair). As illustrated in Figure 1, remote nodes (“ $B$ ” in the figure) can change the state of such entangled links (“ $A$ ” and “ $C$ ” in the figure). Entanglement is an inherently connected element at the lowest physical level, whereas classical communications deal with unidirectional forward communication that abstracts the notion of a connection between a sender and receiver. In order to make use of entanglement for a quantum network special devices capable of producing entanglement and manipulating local qubits are required.

### 3.2 Quantum Network Devices

We focus on a high level summary of devices in a quantum network without delving into detailed physics (for more

details, see [7, 89, 101] and Section 4.4). Qubits can be sent optically through standard optical fiber using a variety of possible encodings, such as polarization [10, 67], time-bin [18], or absence and presence of a photon [22]. Such qubits can be emitted from the devices in quantum nodes [13, 14, 87], but in principle also transferred [59, 75, 87] from optical fiber into the node’s local quantum memory. Present day quantum memories have very limited lifetimes, making it highly desirable to avoid the exchange of additional control information before the entanglement can be used.

We distinguish two classes of quantum nodes. One, which we will call a *controllable quantum node*, offers the possibility to perform controllable quantum operations as well as storing qubits. Specifically, these nodes enable decision making, e.g. which nodes to connect by entanglement swapping. Such nodes can act as quantum repeaters and decision making routers in the network (e.g. NV platform or other quantum memories combined with auxiliary optics), and – if they support the execution of gates and measurements – function as *end nodes* [101] on which we run applications (e.g. NV centre in diamond or Ion Traps). Others, which we call *automated quantum nodes*, are typically only timing controlled, i.e. they perform the same preprogrammed action in each time step. Such nodes also support quantum operations and measurements, but only those necessary to perform their preprogrammed tasks. The latter is still very useful, for example, to establish entanglement along a chain of quantum repeaters performing the entanglement swapping operations [19, 72] (see again Figure 1). In Section 4.4 we give a concrete example of such a timing controlled node.

### 3.3 Use Cases

We distinguish four use cases of a quantum network: one related to producing long-distance entanglement, and three that come from application demands. Since no quantum network has been realized to date, we cannot gain insights from actual usage behavior. Instead we must resort to properties of application protocols known today. We desire flexibility to serve all use cases, including supporting multiple applications at the same time.

*Measure Directly (MD) Use Case:* The first application use case comes from application protocols that produce many ( $\geq 10^4$ ) pairs of entangled qubits sequentially, where both qubits are immediately measured to produce classical correlations. As such, no quantum memory is needed to store the entanglement and it is not necessary to produce all entangled pairs at the same time. It follows that applications making use of this use case may tolerate fluctuating delays in entanglement generation. Additionally, it is not essential to deliver error free correlations obtained from entanglement

to the application. Such applications will thus already anticipate error fluctuation across the many pairs. This contrasts with classical networking where errors are often corrected before the application layer. Examples of such applications are QKD [40], secure identification [33] and other two-party cryptographic protocols [4, 23, 34, 83, 102] at the prepare-and-measure network stage [101], and device-independent protocols at the entanglement network stage [101].

*Create and Keep (CK) Use Case:* The second application use case stems from protocols that require genuine entanglement, possibly even multiple entangled pairs to exist simultaneously. Here, we may wish to perform joint operations on multiple qubits, and perform quantum gates that depend on back and forth communication between two nodes while keeping the qubits in local quantum storage. While more applications can be realized with more qubits, this use case differs substantially in that we want to create relatively few (even just one) pairs, but want to store this entanglement. Since we typically want these pairs to be available at the same time, and memory lifetimes are short, we want to avoid delay between producing consecutive pairs, which is superficially similar to constraints in real time classical traffic. Also for CK, many applications can perform well with noisy entangled links and the amount of noise forms a performance metric (Section 4.2). Examples of such protocols lie in the domain of sensing [46], metrology [63], and distributed systems [9, 36] which lie in the quantum memory network stage and beyond [101].

*Send Qubit (SQ) Use Case:* While many application protocols known to date consume entanglement itself, some — such as distributed quantum computing applications — ask for the transmission of (unknown) qubits. This can be realized using teleportation over any distance as long as entanglement is confirmed between the sender and the receiver. For the link layer, this does not differ from CK, where we want to produce one entangled pair per qubit to be sent.

*Network Layer (NL) Use Case:* In analogy to the classical notion of a link layer, we take the link layer to refer to producing entanglement between neighboring nodes (see Section 3.4). The network layer will be responsible for producing entanglement between more distant ones. While usage behavior of quantum networks is unknown, it is expected (due to technological limitations) that routing decisions, i.e. how to form long-distance links from pairwise links, will not be entirely dynamic. One potential approach would be to first determine a path, and reserve it for some amount of time such that pairwise entanglement can be produced. Producing pairwise entanglement concurrently enables simultaneous entanglement swapping along the entire path with minimal delay to combat limited memory lifetimes. For this, the network layer needs to be capable of prioritizing entanglement production between neighboring nodes.

### 3.4 Network Stack

Based on these considerations, we propose an initial functional allocation of a quantum network stack (see Figure 2). In analogy to classical networking, we refer to the lowest element of the stack as the physical layer. This layer is realized by the actual quantum hardware devices and physical connections such as fibers. We take the physical layer to contain no decision making elements (comparable to e.g. ISP path tunneling architectures [77]) and keep no state about the production of entanglement (or the transmissions of qubits). The hardware at the physical layer is responsible for timing synchronization and other synchronization, such as laser phase stabilization [53], required to make attempts to produce heralded entanglement (Section 4.4). A typical realization of the physical layer involves two controllable quantum nodes linked by a (chain of) automated quantum nodes that attempt entanglement production in well defined time slots.

The task of the link layer is then to turn the physical layer making entanglement attempts into a robust entanglement generation service, that can produce entanglement between controllable quantum nodes connected by a (chain of) automated quantum nodes. Requests can be made by higher layers to the link layer to produce entanglement, where robust means that the link layer endows the physical system with additional guarantees: a request for entanglement generation will (eventually) be fulfilled or result in a time-out. This can be achieved by instructing the physical layer to perform many attempts to produce entanglement until success.

Built on top of the link layer rests the network layer, which is responsible for producing long-distance entanglement between nodes that are not connected directly or by automated quantum nodes. This may be achieved by means of entanglement swapping, using the link layer to generate entanglement between neighboring controllable nodes. In addition, it contains an entanglement manager that keeps track of entanglement in the network, and which may choose to pre-generate entanglement to service later requests from higher layers. It is possible that the network layer and entanglement manager may eventually be separated.

A transport layer takes responsibility for transmitting qubits deterministically (e.g. using teleportation). One may question why this warrants a separate layer, rather than a library. Use of a dedicated layer allows two nodes to pre-share entanglement that is used as applications of the system demand it. Here, entanglement is not assigned to one specific application (purpose ID, Section 4.1.1). This potentially increases the throughput of qubit transmission via teleportation as teleportation requires no additional connection negotiation, but only forward communication from a sender to the receiver. Implementing such functionality in a library

would incur delays in application behavior as entanglement would need to be generated on-demand rather than supplying it from pre-allocated resources.

## 4 LINK LAYER DESIGN CONSIDERATIONS

### 4.1 Desired Service

The link layer offers a robust entanglement creation service between a pair of controllable quantum nodes  $A$  and  $B$  that are connected by a quantum link, which may include automated nodes along the way. This service allows higher layers to operate independently of the underlying hardware platform, depending only on high-level parameters capturing the hardware capabilities.

*4.1.1 Requesting entanglement.* Our use cases bring specific requirements for such a service. Entanglement creation can be initiated at either  $A$  or  $B$  by a CREATE request from the higher layer with parameters:

- (1) *Remote node* with whom entanglement generation is desired if the node is connected directly to multiple others.
- (2) *Type of request* - create and keep (K), and create and measure (M). The first type of request (K) stores entanglement, addressing the use cases CK and NL (see Section 3.3. The second (M) leads to immediate measurement, supporting the use case MD. The reason for distinguishing these two scenarios is twofold: first, we will show later (Section 4.4 that a higher throughput can for some implementations be achieved for M than for K on the same system. Second, simple photonic quantum hardware without a quantum memory and sophisticated processing capabilities [90] only supports the M mode of operation.
- (3) *Number* of entangled pairs to be created. Allowing the higher layer to request several pairs at once can increase throughput by avoiding additional processing delays due to increased inter-layer communication (comparing to classical networks [96, Table 2]). It also helps the CK use case where an application actually needs several pairs concurrently.
- (4) *Atomic* is a flag that assists the requirement from the CK use case, by ensuring that the request must be fulfilled as one, and all pairs be made available at the same time.
- (5) *Consecutive* is a flag indicating an OK is returned for each pair made for a request (typical for NL use case). Otherwise, an OK is sent only when the entire request is completed (more common in application use cases).
- (6) *Waiting time*,  $t_{\max}$  can be used to indicate the maximum time that the higher layer is willing to wait for

completion of the request. This allows a general timeout to be set, and enables the NL and CK use case to specify strict requirements since the requested pairs may no longer be desirable if they are delivered too late.

- (7) A *purpose ID* can be specified which allows the higher layer to tag the entanglement for a specific purpose. For an application use case, this purpose ID may be considered analogous to a port number found in the TCP/IP stack and including it in the CREATE request allows both nodes to immediately pass the entanglement to the right application to proceed with processing at both ends without incurring further communication delays. A purpose ID is also useful to identify entanglement created by the NL use case for a specific long-distance path. We envision that an entanglement manager who may decide to pre-generate entanglement would use a special tag to indicate “ownership” of the requested pairs. For the NL use case for example, if the entanglement request does not correspond to a pre-agreed path, then the remote node may refuse to engage in entanglement generation. Finally, a purpose ID enables rejection of generation by the remote node based on scheduling or security policies, e.g. for NL when no path is agreed upon.
- (8) A *priority* that may be used by a scheduler. Here we use only three priorities (use cases NL, MD and CK).
- (9) Finally, we allow a specification of a purely quantum parameter (refer to Appendix A), the *desired minimum fidelity*,  $F_{\min}$ , of the entanglement [76]. Here, it is sufficient to note that the fidelity  $0 \leq F \leq 1$  measures the quality of entanglement, where a higher value of  $F$  means higher entanglement quality. The ideal target state has  $F = 1$ , while  $F \geq 1/2$  is often desirable [52].

The reason for allowing different  $F_{\min}$  instead of fixing one for each hardware platform is that the same platform can be used to produce higher or lower fidelity pairs, where a higher fidelity pair costs more time to prepare. Examples of this are the use of entanglement distillation [38, 60] where two lower quality pairs are combined into one higher quality one, and the bright state population  $\alpha$  (see Appendix D) used to generate entanglement.

*4.1.2 Response to entanglement requests.* If entanglement has been produced successfully, an OK message should be returned. In addition, the use cases specified in Section 3.3 desire several other pieces of information, which may also be tracked at higher layer:

- (1) An entanglement identifier  $Ent_{ID}$  unique in the network during the lifetime of the entanglement. This



allows both nodes to immediately process the entanglement without requiring an additional round of communication degrading the entanglement due to limited memory lifetimes.

- (2) A qubit ID for  $K$ -type (create and keep) requests which identifies where the local qubit is in the quantum memory device.
- (3) The “Goodness”  $G$ , which for  $K$  requests is an estimate (Section B) of the fidelity — where  $G \geq F_{\min}$  should hold — and for  $M$  an estimate of the quantum bit error rate (QBER, see again Appendix A).
- (4) The measurement outcome for  $M$  type requests.
- (5) The time of entanglement creation.
- (6) The time the goodness parameter was established. The goodness may later be updated given fixed information about the underlying hardware platform.

Evidently, there are many possibilities of failure resulting in the return of error messages. This includes:

- Timeout when a request could not be fulfilled in a specific time frame (TIMEOUT).
- An immediate rejection of the request because the requested fidelity is not achievable in the given time frame (UNSUPP).
- The quantum storage is permanently (MEMEXCEEDED) or temporarily (OUTOFMEM) too small to simultaneously store all pairs of an atomic request.
- Refusal by the remote node to participate (DENIED).

Finally, we allow an EXPIRE message to be sent, indicating that the entanglement is no longer available. This in principle can be indicated by a quantum memory manager (see Appendix, Section 5.2.2) instead of the protocol, but we will that show that it allows for recovery from unlikely failures.

**4.1.3 Fixed hardware parameters.** Not included in these request or response messages are parameters that are fixed for the specific hardware platform, or change only very infrequently. As such, these may be obtained by high-level software by querying the low level system periodically, similarly to some classical network architectures (e.g. [66]). Such parameters include:

- The number of available qubits.
- The qubit memory lifetimes.
- Possible quantum operations.
- Attainable fidelities and generation time.
- The class of states that are produced.

The latter refers to the fact that more information than just the fidelity allows optimization at layers above the link layer.

## 4.2 Performance Metrics

Before designing any protocols that adhere to these requirements for entanglement generation, we consider the performance metrics that such protocols may wish to optimize. Standard metrics from networking also apply here, such as *throughput* (entangled pairs/s), and the *latency*. We distinguish between:

- (1) Latency per request (time between submission of a CREATE request and its successful completion at a requesting node).
- (2) Latency per pair (time between CREATE and OK at requesting node).
- (3) Latency per request per number of requested pairs (which we denote as the *scaled latency*).

Given requests may originate at both  $A$  and  $B$ , we also demand *fairness*, i.e., the metrics should be roughly independent of the origin of the request. Here, we also care about genuinely quantum quality metrics, specifically the fidelity  $F$  (at least  $F_{\min}$ ).

The non-quantum reader may wonder about the significance of  $F$ , and why we do not simply maximize throughput (e.g. [17, 91]) or minimize latency (e.g. [24, 37]). For instance, QKD (a MD use case as listed in Section 3.3), requires a minimum quantum bit error rate (QBER) between measurement outcomes at  $A$  and  $B$  (related to  $F$ , see Appendix A). A lower  $F$  results in a larger QBER, allowing less key to be produced per pair. We may thus achieve a higher throughput, but a lower number of key bits per second, or key generation may become impossible.

## 4.3 Error Detection

Link layer protocols for classical communication typically aim to correct or detect errors, e.g. using a CRC. In principle, there exists an exact analogy at the quantum level: We could use a checksum provided by a quantum error correcting code (QECC) [76, 94] to detect errors. This is technologically challenging and experimental implementations of QECC are in very early stages [29, 31, 86]. Yet, apart from technological limitations, future link layer protocols may not use quantum checksums due to different use case requirements: We typically only demand some minimum fidelity  $F_{\min}$  with high confidence that may also fluctuate slightly for pairs produced within a time window.

As we focus primarily on fidelity, we instead use a different mechanism: we intersperse test rounds during entanglement generation (for details, refer to Appendix B) to verify the quality of the link. Such test rounds are easy to produce without the need for complex gates or extra qubits. Evidently, there exists an exact analogy in the classical networking world, where we would transmit test bits to measure the

current quality of transmission, e.g. a direct analogy to network profiling [66, Section 4.3] to gain confidence that the non-test bits are also likely to be transmitted with roughly the same amount of error. Yet, there we typically care about correctness of a specific data item, rather than an enabling resources like entanglement.

#### 4.4 Physical Entanglement Generation

Let us now explain how heralded entanglement generation is actually performed between two controllable nodes  $A$  and  $B$  (see Appendix D for details). As an example, we focus on the hardware platform available to us (NV in diamond, Figure 3), but analogous implementations have been performed using remote Ion Traps [70] and Neutral Atoms [50].

Nodes  $A$  and  $B$  are few-qubit quantum processors, capable of storing and manipulating qubits. They are connected to an intermediate station called the *heralding station*  $H$  over optical fibers. This station is a much simpler automated node, built only from linear optical elements. Each node can have two types of qubits: *memory qubits* as a local memory, and *communication qubits* with an optical interface. To produce entanglement, a time synchronized trigger is used at both  $A$  and  $B$  to create entanglement between each communication qubit, and a corresponding traveling qubit (photon). These photons are sent from  $A$  and  $B$  to  $H$  over fiber. When both arrive at  $H$ ,  $H$  performs an automatic entanglement swapping operation which succeeds with some probability. Since  $H$  has no quantum memory, both photons must arrive at  $H$  at the same time to succeed. Success or failure is then transmitted back from  $H$  to the nodes  $A$  and  $B$  over a standard classical channel (e.g. 100Base-T). In the case of success, one of several entangled states may be produced, which can however be converted to one other using local quantum gates at  $A$  and  $B$ . After a generation attempt, the communication qubit may be moved to a memory qubit, in order to free the communication qubit to produce the next entangled pair. Many parameters influence the success and quality of this process, such as the quality of the qubits themselves, the probability of emission of a photon given a trigger signal, losses in fiber, and quality of the optical elements such as detectors used at  $H$  (Figure 3).

For further information on this process see [53]. For an overview on NV centres in diamond see [25]. Two different schemes for producing entanglement have been implemented, that differ in how the qubits are encoded into photons (time-bin [8], or presence/absence of a photon [22]). While physically different, both of these schemes fit into the framework of our physical and link layer protocols.

To evaluate the performance of the protocol (Section 6) and provide intuition of timings, we compare to data from the setup [53] which uses presence/absence of a photon as

encoding. A microwave pulse prepares the communication qubit depending on a parameter  $\alpha$ , followed by a laser pulse to trigger photon emission (total duration  $5.5\mu\text{s}$ ). A pair ( $|\Psi^+\rangle$  or  $|\Psi^-\rangle$ ) is successfully produced with fidelity  $F \approx 1 - \alpha$  (ignoring memory lifetimes and other errors, see Appendix D) with probability  $p_{\text{succ}} \approx 2\alpha p_{\text{det}}$ , where  $p_{\text{det}} \ll 1$  is the probability of emitting a photon followed by heralding success. The parameter  $\alpha$  thus allows a trade-off between the rate of generation ( $p_{\text{succ}}$ ), and the quality metric  $F$ . For  $K$  type requests, we may store the pair in the communication qubit, or move to a memory qubit (duration of  $1040\mu\text{s}$  for the qubit considered). The quality of this qubit degrades as we wait for  $H$  to reply. For  $M$  type requests, we may choose to measure immediately before receiving a reply (here readout  $3.7\mu\text{s}$ ). Important is the time of an attempt  $t_{\text{attempt}}$  (time preparing the communication qubit until receiving a reply from  $H$ , and completion of any post-processing such as moving to memory), and the maximum attempt rate  $r_{\text{attempt}}$  (maximum number of attempts that can be performed per second not including waiting for a reply from  $H$  or post-processing). The rate  $r_{\text{attempt}}$  can be larger than  $1/t_{\text{attempt}}$ : (1) for  $M$  the communication qubit is measured before receiving the reply from  $H$  and thus allows for multiple attempts to overlap and (2) for  $K$ , if the reply from  $H$  is failure, then no move to memory is done.

For performance evaluation we consider two physical setups as an example (see Appendix D) with additional parameters hereafter referred to as the LAB scenario and the QL2020 scenario. The LAB scenario already realized [53] with distance to the station 1 m from both  $A$  and  $B$  (communication delay to  $H$  negligible),  $p_{\text{succ}} \approx \alpha \cdot 10^{-3}$  ( $F$  vs.  $\alpha$ , Figure 8). For  $M$  requests, we act the same for LAB and QL2020 and always measure immediately before parsing the response from  $H$  to ease comparison (thus  $t_{\text{attempt}} = 1/r_{\text{attempt}} = 10.12\mu\text{s}$  which includes electron readout  $3.7\mu\text{s}$ , photon emission  $5.5\mu\text{s}$  and a 10 % extra delay to avoid race conditions). For  $K$  requests in LAB,  $t_{\text{attempt}} = 1045\mu\text{s}$  but  $1/r_{\text{attempt}} \approx 11\mu\text{s}$  as memory qubits need to be periodically initialized ( $330\mu\text{s}$  every  $3500\mu\text{s}$ ). The QL2020 scenario has not been realized and is based on a targeted implementation connecting two European cities by the end of 2020 ( $\approx 10\text{km}$  from  $A$  to  $H$  with a communication delay of  $48.4\mu\text{s}$  in fiber, and  $\approx 15\text{km}$  from  $B$  to  $H$  with a  $72.6\mu\text{s}$  delay). Frequency conversion of  $637\text{nm}$  to  $1588\text{nm}$  is performed on the photons emitted in our modeled NV centre while fiber losses at  $1588\text{nm}$  are taken to be  $0.5\text{ dB/km}$  (values for deployed QL2020 are fibers  $0.43\text{-}0.47\text{ dB/km}$ ). We model the use of optical cavities to enhance photon emission [15, 84] giving a probability of success  $p_{\text{succ}} \approx \alpha \cdot 10^{-3}$ .  $F$  is worse due to increased communication times from  $H$  (Figure 9). For QL2020  $t_{\text{attempt}} = 145\mu\text{s}$  for  $M$  (trigger, wait for reply from  $H$ ) and  $t_{\text{attempt}} = 1185\mu\text{s}$  for  $K$  (trigger, wait for reply from  $H$ , swap to carbon).

Maximum attempt rates are  $1/r_{\text{attempt}} = 10.120 \mu\text{s}$  (M) and  $1/r_{\text{attempt}} \approx 165 \mu\text{s}$  (K).

## 4.5 Hardware Considerations

Quantum hardware imposes design considerations for any link layer protocol based on top of such experiments.

*Trigger generation:* Entanglement can only be produced if both photons arrive at the heralding station at the same time. This means that the low level system requires tight timing control; such control (ns scale) is also required to keep the local qubits stable. This imposes hard real time constraints at the lowest level, with dedicated timing control (AWG) and software running on a dedicated microcontroller (Adwin ProII). When considering a functional allocation between the physical and link layer, this motivates taking all timing synchronization to happen at the physical layer. At this layer, we may then also timestamp classical messages traveling to and from  $H$ , to form an association between classical control information and entangled pairs.

*Scheduling and flow control:* Consequently, we make the link layer responsible for all higher level logic, including scheduling, while keeping the physical layer as simple as possible. An example of scheduling other than priorities, is flow control which controls the speed of generation, depending on the availability of memory on the remote node to store such entanglement.

Note that depending on the number of communication qubits, and parallelism of quantum operations that the platform allows, a node also has a global scheduler for the entire system and not only the actions of the link layer.

*Noise due to generation:* One may wonder why one does not continuously trigger entanglement generation locally whenever the node wants a pair, or why one does not continuously produce pairs and then this entanglement is either discarded or otherwise made directly available. In the NV system, triggering entanglement generation causes the memory qubits to degrade faster [58, 82]. As such we would like to achieve agreement between nodes to avoid triggering unless entanglement it is indeed desired.

This consideration also yields a security risk: if an attacker could trick a node into triggering entanglement generation, without a matching request on the other side, this would allow a rapid destruction of contents of the nodes' local quantum memory. For this reason, we want classical communication to be authenticated which can be achieved using standard methods.

*Memory allocation:* Decisions on which qubits to use for what purpose lies in the domain of higher level logic, where more information is available. We let such decisions be taken by a global quantum memory manager (QMM), which can assist the link layer to make a decision on which qubits to

employ. It can also translate logical qubit IDs into physical qubit IDs in case multiple qubits are used to redundantly form one logical storage qubit.

## 5 PROTOCOLS

We now present our protocols satisfying the requirements and considerations set forth in Sections 3 and 4. The entanglement generation protocol (EGP) at the link layer, uses the midpoint heralding protocol (MHP) at the physical layer. Classical communication is authenticated, and made reliable using standard methods (e.g. 802.1AE [54], authentication only).

### 5.1 Physical Layer MHP

Our MHP is a lightweight protocol built directly on top of physical implementations of the form of Section 4.4, supplementing them with some additional control information. With minor modifications this MHP can be adapted to other forms of heralded entanglement generation between controllable nodes, even using multiple automated middle nodes [47].

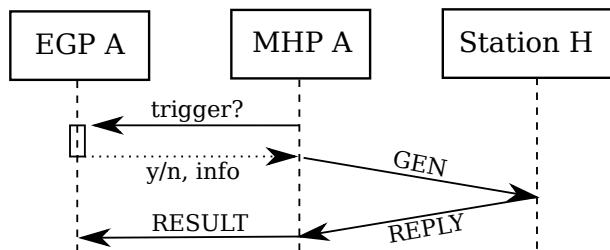
The MHP is meant to be implemented directly at the lowest level subject to tight timing constraints, which is why we let the MHP poll higher layer (Figure 4, the link layer EGP) at each timestep to determine whether entanglement generation is required, and keep no state. A batched operation is possible, should the delay incurred by the polling exceed the minimum time to make one entanglement generation attempt - *the MHP cycle* - and hence dominate the throughput. Upon polling, the higher layer may respond “no” in which case no attempt will be made or with “yes”, additionally providing parameters to use in the attempt. These parameters include the type of request (M, measure) or (K, store) passed on from the higher layer, for which the MHP takes the following actions.

*5.1.1 Protocol for Create and Keep (K).* The parameters given to the MHP with a “yes” response contain the following:

- An ID for the attempt that is forwarded to  $H$
- Generation parameters ( $\alpha$ , Section 4.4)
- The device qubits for storing the entanglement
- A sequence of operations to perform on the device memory <sup>1</sup>.

The higher layer may instruct the MHP to perform a gate on the communication qubit depending on the heralding signal from  $H$  allowing the conversion from the  $|\Psi^-\rangle$  state to the  $|\Psi^+\rangle$  state. Entanglement generation is then triggered at the start of the next time interval, including generation parameter  $\alpha$ , and a GEN message is sent to  $H$  which includes

<sup>1</sup>Less abstractly, by specifying microwave and laser pulse sequences controlling the chip (see Appendix).



**Figure 4: Timeline of the MHP polling higher layers to see if entanglement should be produced.**

a timestamp, and the given ID. The motivation for including the ID is to protect against errors in the classical control, for example losses.

The station  $H$  uses the timestamp to link the message to a detection window in which the accompanying photons arrived. Should messages from both nodes arrive, the midpoint verifies that the IDs transmitted with the GEN messages match, and checks the detection counts (Figure 3) from the corresponding detection window. The midpoint will then send a REPLY message indicating success or failure, and in the case of success which of the two states  $|\Psi^+\rangle$  and  $|\Psi^-\rangle$  was produced. The REPLY additionally contains a sequence number uniquely identifying the generated pair of entangled qubits chosen by  $H$ , which later enables the EGP to assign unique entanglement identifiers. This REPLY and the ID is forwarded to the link layer for post-processing. Note that the REPLY may be received many MHP cycles later, allowing the potential for emission multiplexing (Section 5.2).

**5.1.2 Protocol for Create and Measure (M).** Handling M type requests is very similar, differing only in two ways: Instead of performing a gate on the communication qubit, the “yes” message requests the MHP to perform a measurement on the communication qubit in a specified basis once the photon has been emitted, even before receiving the response from  $H$ . The outcome of the measurement and the REPLY are passed back to the EGP. In practice, the communication time from transmitting a GEN message to receiving a REPLY may currently exceed the duration of such a local measurement ( $3.7 \mu\text{s}$  vs. communication delay LAB 9.7 ns, and QL2020 145  $\mu\text{s}$ ), and the MHP may choose to perform the measurement only after a successful response is received.

## 5.2 Link Layer EGP

Here we present an implementation of a link layer protocol, dubbed EGP, satisfying the service requirements put forth in Section 4 (see Appendix E for details and message formats). We build up this protocol from different components:

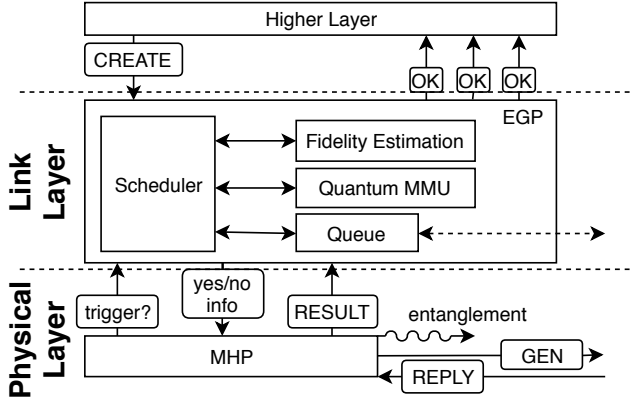
**5.2.1 Distributed Queue.** Both nodes that wish to establish entangled link(s) must trigger their MHP devices in a coordinated fashion (Section 4.4). To achieve such agreement, the EGP employs a distributed queue comprised of synchronized local queues at the controllable nodes. These local queues can be used to separate requests based on priority, where here we employ 3 for the different use cases (CK, NL, MD). Due to low errors in classical communication (estimated  $< 4 \times 10^{-8}$  on QL2020, see Appendix D.6.1), we let one node hold the master copy of the queue, and use a simple two-way handshake for enqueueing items, and a windowing mechanism to ensure fairness. Queue items include a *min\_time* that specifies the earliest possible time a request is deemed ready for processing by both nodes (depending on their distance). Specifying *min\_time* prevents either node from beginning entanglement generation in different timesteps.

**5.2.2 Quantum Memory Management (QMM).** The EGP uses the node’s QMM (Section 4.5) to determine which physical qubits to use for generating or storing entanglement.

**5.2.3 Fidelity estimation unit (FEU).** In order to provide information about the quality of entanglement, the EGP employs a fidelity estimation unit. This unit is given a desired quality parameter  $F_{\min}$ , and returns generation parameters (such as  $\alpha$ ) along with an estimated minimal completion time. Such a fidelity estimate can be calculated based on known hardware capabilities such as the quality of the memory and operations. To further improve this base estimate the EGP intersperses test rounds.

**5.2.4 Scheduler.** The EGP scheduler decides which request in the queue should be served next. In principle, any scheduling strategy is suitable as long as it is deterministic, ensuring that both nodes select the same request locally. This limits two-way communication, which adversely affects entanglement quality due to limited memory lifetimes.

**5.2.5 Protocol.** Figure 5 presents an architecture diagram visualizing the operation. The protocol begins when a higher layer at a controllable node issues a CREATE operation to the EGP specifying a desired number of entangled pairs along with  $F_{\min}$  and  $t_{\max}$ . Upon receipt of a request the EGP will query the FEU to obtain hardware parameters ( $\alpha$ ), and a minimum completion time. If this time is larger than  $t_{\max}$ , the EGP immediately rejects the request (UNSUPP). Should the request pass this evaluation, the local node will compute a fitting *min\_time* specifying the earliest MHP polling cycle the request may begin processing. The node then adds the request into the distributed queue shared by the nodes. This request may be rejected by the peer should the remote node have queue rules that do not accept the specified purpose ID. Then, the EGP locally rejects the request (DENIED).



**Figure 5: Flow diagram of the MHP and EGP operation. The EGP handles CREATE requests and schedules entanglement generation attempts are issued to the MHP. Replies from the midpoint are parsed and forwarded to the EGP from request management.**

The local scheduler selects the next request to be processed, given that there exists a ready one (as indicated by  $min\_time$ ). The QMM is then used to allocate qubits needed to fulfill the specified request type (create and keep  $K$  or create and measure  $M$ ). The EGP will then again ask the FEU to obtain a current parameter  $\alpha$  due to possible fluctuations in hardware parameters during the time spent in the queue. The scheduler then constructs a “yes” response to the MHP containing  $\alpha$  from the FEU, along with an ID containing the unique queue ID of the request in the distributed queue, and number of pairs already produced for the request. This response is then forwarded to the local MHP upon its next poll to the EGP. If no request is ready for processing, a “no” response is returned to the MHP. At this point the MHP behaves as described in the previous section and an attempt at generating entanglement is made.

Whenever a REPLY and ID is received from the MHP, the EGP uses the ID to match the REPLY to an outstanding request, and evaluates the REPLY for correctness. Should the attempt be successful, the number of pairs remaining to be generated for the request is decremented and an OK message is propagated to higher layers containing the information specified in Section 4.1.2, where the Goodness is obtained from the FEU.

In the Appendix, we consider a number of examples to illustrate decisions and possible pitfalls in the EGP. An example is the possibility of *emission multiplexing* [98]: The EGP can be polled by the MHP before receiving a response from the MHP for the previous cycle. This allows the choice to attempt entanglement generation multiple times in succession before receiving a reply from the midpoint, e.g., in

order to increase the throughput for the MD use case. Errors such as losses on the classical control link can lead to an inconsistency of state (of the distributed queue) at  $A$  and  $B$  from which we need to recover. Inconsistencies can also affect the higher layer. Since the probability of e.g. losses is extremely low, we choose not to perform additional two way discussion to further curb all inconsistencies at the link layer. Instead, the EGP can issue an EXPIRE message for an OK already issued if inconsistency is detected later, e.g. when the remote node never received an OK for this pair.

## 6 EVALUATION

We investigate the performance of our link layer protocol using a purpose built discrete event simulator for quantum networks (NetSquid [1], Python/C++) based on DynAA [41] (see Appendix C for details and more simulation results). Both the MHP and EGP are implemented in full in Python running on simulated nodes that have simulated versions of the quantum hardware components, fiber connections, etc. All simulations were performed on the supercomputer *Cartesius* at SURFsara [2], in a total of 2578 separate runs using a total of 94244 core hours, and 707 hours time in the simulation ( $\sim 250$  billion MHP cycles).

We conduct the following simulation runs:

- Long runs: To study robustness of our protocol, we simulate the 169 scenarios described below for an extended period of time. Each scenario was simulated twice for 120 wall time hours, simulating 502 – 13437 seconds. We present and analyze the data from these runs in sections 6.1, 6.2 and C.2.
- Short runs: We perform the following simulations for a shorter period of time (24 wall time hours, reaching 67 – 2356 simulated seconds):
  - Performance trade-offs: To study the trade-off between latency, throughput and fidelity we sweep the incoming request frequency and the requested minimum fidelity, see Figure 6.
  - Metric fluctuations: To be able to study the impact of different scheduling strategies on the performance metrics, we run 4 scenarios, 102 times each. The outcomes of these simulation runs are discussed in section 6.3.

To explore the performance at both short and long distances, the underlying hardware model is based on the LAB and QL2020 scenarios, where we validate the physical modeling of the simulation against data collected from the quantum hardware system of the LAB scenario already realized (Figure 8). For the quantum reader we note that while our simulations can also be used to predict behavior of physical implementations (such as QL2020), the focus here is on the performance and behavior of the link layer protocol.

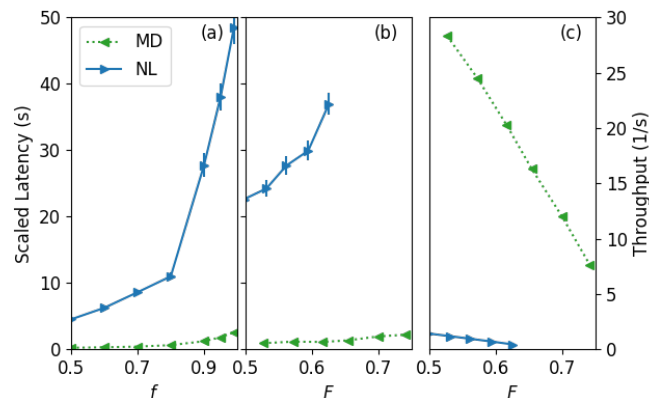
We structure the evaluation along the three different use cases (NL, CK, MD), leading to a total of 169 different simulation scenarios. First, we use different kinds of requests: (1) *NL* (K type request, consecutive flag, priority 1 (highest), store qubit in memory qubit), (2) *CK*, an application asking for one or more long-lived pairs (K type request, immediate return flag, priority 2 (high), store qubit in memory qubit) and (3) (*MD*) measuring directly (M type request, consecutive flag, priority 3 (lowest)). For an application such as QKD, one would not set the immediate return flag in practice for efficiency, but we do so here to ease comparison to the other two scenarios. Measurements in *MD* are performed in randomly chosen bases  $X$ ,  $Z$  and  $Y$  (see Appendix A).

In each MHP cycle, we randomly issue a new CREATE request for a random number of pairs  $k$  (max  $k_{\max}$ ), and random kind  $P \in \{NL, CK, MD\}$  with probability  $f_P \cdot p_{\text{succ}} / (E \cdot k)$ , where  $p_{\text{succ}}$  is the probability of one attempt succeeding (Section 4.4),  $f_P$  is a fraction determining the load in our system of kind  $P$ , and  $E$  is the expected number of MHP cycles to make one attempt ( $E = 1$  for MD and  $E \approx 1.1$  for NL/CK in LAB due to memory re-initialization and post-processing).  $E \approx 16$  for NL/CK in QL2020 due to classical communication delays with  $H$  ( $145\mu\text{s}$ ). In the long runs, we first study single kinds of requests (only one of MD/CK/NL), with  $f_P = 0.7$  (Low), 0.99 (HIGH) or 1.5 (ULTRA). For the long runs, we fix one target fidelity  $F_{\min} = 0.64$  to ease comparison. For each of the 3 kinds (MD/CK/NL), we examine (1)  $k_{\max} = 1$ , (2)  $k_{\max} = 3$ , and (3) only for MD,  $k_{\max} = 255$ . For ULTRA the number of outstanding requests intentionally grows until the queue is full (max 256), to study an overload of our system. To study fairness, we take 3 cases of CREATE origin for each single kind (MD/CK/NL) scenario: (1) all from  $A$  (master of the distributed queue), (2) all from  $B$ , (3)  $A$  or  $B$  are randomly chosen with equal probability. To examine scheduling, we additionally consider long runs with mixed kinds of requests (Appendix, e.g. Figure 7).

## 6.1 Robustness

To study robustness, we artificially increase the probability of losing classical control messages (100 Base T on QL2020 fiber  $< 4 \times 10^{-8}$ , Appendix D.6.1), which can lead to an inconsistency of state of the EGP but also at higher layers (Section 5.2). We ramp up loss probabilities up to  $10^{-4}$  (Appendix D.6.1) and observe our recovery mechanisms work to ensure stable execution in all cases (35 runs, 281 - 3973 s simulated time), with only small impact to the performance parameters (maximum relative differences<sup>2</sup> to the case of no losses, fidelity (0.005), throughput (0.027), latency (0.629), number of OKs (0.026) with no EXPIRE messages). We see a relatively large difference for latency, which may however

<sup>2</sup>Relative difference between  $m_1$  and  $m_2$  is  $|m_1 - m_2| / \max(|m_1|, |m_2|)$



**Figure 6: Performance trade-offs. (a) Scaled latency vs.  $f_P$  determining fraction of throughput (b) Scaled latency vs. fidelity  $F_{\min}$ . Demanding a higher  $F_{\min}$  lowers the probability of an attempt being successful (Section 4.4), meaning (c) throughput directly scales with  $F_{\min}$  (each point averaged over 40 short runs each 24 h, 93 – 2355 s simulated time, QL2020,  $k_{\max} = 3$ , for (b,c)  $f_P = 0.99$ ). Higher  $F_{\min}$  not satisfiable for NL in (b).**

be due to latency not reaching steady state in the course of this simulation ( $70 \times 70$  core hours).

## 6.2 Performance Metrics

We first consider runs including only a single kind of request (MD/CK/NL). In addition to the long runs, we conduct specific short runs examining the trade-off between latency and throughput for fixed target fidelity  $F_{\min}$  (Figure 6(a)), and the trade-off between latency (throughput) and the target fidelity in Figure 6(b) (Figure 6(c)).

Below we present the metrics extracted from the long runs with single kinds of requests:

*Fidelity:* As a benchmark, we began by recording the average fidelity  $F_{\text{avg}}$  in all 169 scenarios with fixed minimum fidelity. We observe that  $F_{\text{avg}}$  is independent of the other metrics but does depend on the distance, and whether we store or measure:  $0.745 < F_{\text{avg}} < 0.757$  NL/CK LAB,  $0.626 < F_{\text{avg}} < 0.653$  NL/CK QL2020,  $0.709 < F_{\text{avg}} < 0.779$  MD LAB,  $0.723 < F_{\text{avg}} < 0.767$  MD QL2020 (Fidelity MD extracted from QBER measurements, Appendix A). This is to be expected since (1) we fix one  $F_{\min}$  and (2) we consider an NV platform with only 1 available memory qubit (LAB).

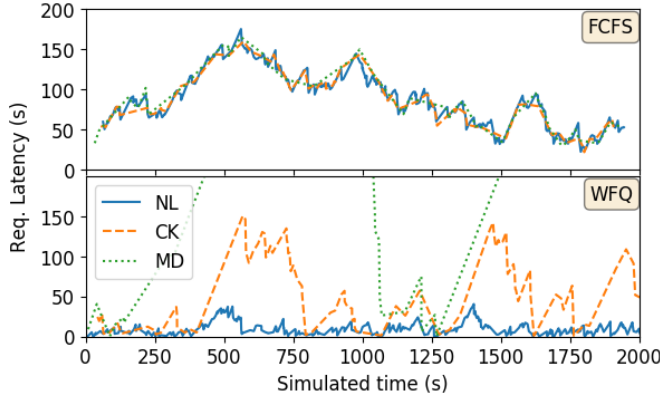
*Throughput:* All scenarios HIGH and ULTRA in LAB reach an average throughput  $th_{\text{avg}}$  (1/s) of  $6.05 < th_{\text{avg}} < 6.47$  NL/CK and  $6.51 < th_{\text{avg}} < 7.09$  for MD. It is expected that



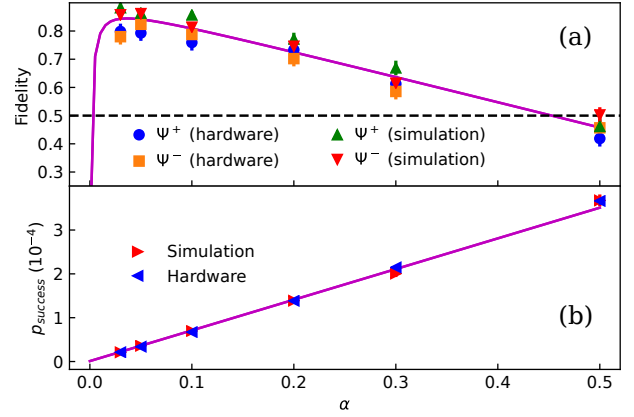
*MD* has higher throughput, since no memory qubit needs to be initialized. The time to move to memory ( $1040\mu\text{s}$ ) is less significant since many MHP cycles are needed to produce one pair, but we only move once. As expected for Low the throughput is slightly lower in all cases,  $4.44 < th_{\text{avg}} < 4.72$  *NL/CK*, and  $4.86 < th_{\text{avg}} < 5.22$  *MD*. For QL2020, the throughput for *NL/CK* is about 14 times lower, since we need to wait ( $145\mu\text{s}$ ) for a reply from *H* before MHP can make a new attempt.

*Latency*: The scaled latency highly depends on the incoming request frequency as the queue length causes higher latency. However, from running the same scenarios many (102) times for a shorter period (24 wall time hours), see section 6.3, we see that the average scaled latency fluctuates a lot, with a standard deviation of up to 6.6 s in some cases. For QL2020 with *NL* requests specifying 1-3 pairs from both nodes, we observe an average scaled latency of 10.97 s Low, 142.9 s HIGH and 521.5 s ULTRA. For *MD* requests, 0.544 s Low, 3.318 s HIGH and 32.34 s ULTRA. The longer scaled latency for *NL* is largely due to longer time to fulfill a request, and not that the queues are longer (average queue length for *NL*: 3.83 Low, 56.3 HIGH, 214 ULTRA), and for *MD*: 3.23 Low, 22.4 HIGH and 219 ULTRA).

*Fairness*: For 103 scenarios of the long runs (single kinds of requests (*MD/CK/NL*) randomly from *A* and *B*), we see only slight differences in fidelity, throughput or latency between requests from *A* and *B*. Maximum relative differences do not exceed: fidelity 0.033, throughput 0.100, latency 0.073, number of OKs 0.100 (which occurs for ULTRA).



**Figure 7: Request latency vs. time for two scheduling scenarios (long runs simulated 120 h wall time). As expected the max. latency for *NL* is decreased due to strict priority. In this scenario, there are more incoming *NL* requests ( $f_{NL} = 0.99 \cdot 4/5$ ,  $f_{CK} = 0.99 \cdot 1/5$  and  $f_{MD} = 0.99 \cdot 1/5$ ).**



**Figure 8: Validation against data from NV hardware (LAB scenario). Fidelity (a) and probability an attempt succeeds (b) in terms of  $\alpha$  (Section 4.4) shows good agreement between hardware and simulation points (each at least 300 pairs averaged, 5s–117s simulated time, 500k–10.000k attempts, 122 hours wall time). Theoretical model [53] as visual guide (solid line).**

### 6.3 Scheduling

We take a first step studying the effect of scheduling strategies on the performance when using mixed kinds of requests. Part of simulating the performance of a scheduling strategies can certainly be done without implementing all details of the physical entanglement generation. However, since we do simulate these details we can first confirm that different scheduling strategies below do not affect the average fidelity in these scenarios. Here, we examine two simple scheduling strategies: (1) first-come-first-serve (FCFS) and (2) a strategy where *NL* of priority 1 has a strict highest priority, and use a weighted fair queue (WFQ) for *CK* (priority 2) and *MD* (priority 3), where *CK* has 10 times the weight of *MD*. With these scheduling strategies, we simulate two different request patterns ((i) uniform and (ii) no *NL* more *MD*), 102 times over 24 wall time hours each and extract the performance metrics of throughput and scaled request latency, see table 1.

As expected we see a drastic decrease of the average scaled latency for *NL* when giving it strict priority: 10.3 s with FCFS and 3.5 s with WFQ. For *CK* there is similarly a decrease in average scaled latency, however smaller than for *NL*, 10.1 s and 6.5 s for FCFS and WFQ respectively. For *MD* the average scaled latency goes up in both cases when using WFQ instead of FCFS, by factors of 2.49 (uniform) and 1.28 (no *NL* more *MD*).

We observe that the throughput gets less affected by the scheduling strategy than the latency for these scenarios. The maximal difference between the throughput for FCFS and WFQ is by a factor of 1.16 (for *MD* in the scenario of no *NL*

**Table 1: Throughput (T) and scaled latency (SL) using scheduling strategies FCFS and WFQ for two request patterns: (i) with  $f_{NL} = f_{CK} = f_{MD} = 0.99 \cdot 1/3$ , i.e. a uniform load of the different priorities and (ii) with  $f_{NL} = 0$ ,  $f_{CK} = 0.99 \cdot 1/5$  and  $f_{MD} = 0.99 \cdot 4/5$ , i.e. no NL and more MD. The physical setup: QL2020 and number of pairs per request: 2 (NL), 2 (CK), and 10 (MD). Each value average over 102 short runs each 24 h, with standard error in parentheses.**

| T (1/s)   | NL            | CK            | MD            |
|-----------|---------------|---------------|---------------|
| (i) FCFS  | 0.146 (0.003) | 0.144 (0.003) | 2.464 (0.056) |
| (i) WFQ   | 0.154 (0.003) | 0.156 (0.003) | 2.130 (0.063) |
| (ii) FCFS | -             | 0.086 (0.003) | 5.912 (0.033) |
| (ii) WFQ  | -             | 0.096 (0.003) | 5.829 (0.049) |

| SL (s)    | NL             | CK             | MD            |
|-----------|----------------|----------------|---------------|
| (i) FCFS  | 10.272 (0.654) | 10.063 (0.631) | 1.740 (0.120) |
| (i) WFQ   | 3.520 (0.085)  | 6.548 (0.361)  | 4.331 (0.336) |
| (ii) FCFS | -              | 5.659 (0.313)  | 0.935 (0.062) |
| (ii) WFQ  | -              | 2.503 (0.100)  | 1.194 (0.093) |

and more MD). Furthermore, we see that the total throughput for all requests goes down from 2.75 (5.99) 1/s for FCFS to 2.44 (5.92) 1/s for WFQ in the case of uniform (no NL more MD).

## 7 CONCLUSION

Our top down inventory of design requirements, combined with a bottom up approach based on actual quantum hardware allowed us to take quantum networks a step further on the long path towards their large-scale realization. Our work readies QL2020, and paves the way towards the next step, a robust network layer control protocol. The link layer may now be used as a robust service without detailed knowledge of the physics of the devices. We expect that at the network layer, and when considering larger quantum memories, smart scheduling strategies will be important not only to combat memory lifetimes but also to coordinate actions of different nodes in time, calling for significant effort in computer science and engineering.

## ACKNOWLEDGEMENTS

We thank Kenneth Goodenough for comments in earlier drafts. This work was supported by ERC Starting Grant (Stephanie Wehner), ERC Consolidator Grant (Ronald Hanson), EU Flagship on Quantum Technologies, Quantum Internet Alliance, NWO VIDI (Stephanie Wehner), Marie Skłodowska-Curie actions – Nanoscale solidstate spin systems and emerging quantum technologies – Spin-NANO, grant agreement number 676108.

## REFERENCES

- [1] 2018. NetSQUID. <https://netsquid.org/>. (2018).
- [2] 2018. SURFsara. <https://userinfo.surfsara.nl/systems/cartesius>. (2018).
- [3] M. H. Aboeib, J. Cramer, M. A. Bakker, N. Kalb, M. Markham, D. J. Twitchen, and T. H. Taminiau. 2018. One-second coherence for a single electron spin coupled to a multi-qubit nuclear-spin environment. *Nature Communications* 9, 1 (Dec. 2018), 2552. <https://doi.org/10.1038/s41467-018-04916-z> arXiv:1801.01196
- [4] Dorit Aharonov, Amnon Ta-Shma, Umesh V Vazirani, and Andrew C Yao. 2000. Quantum bit escrow. In *Proceedings of the thirty-second annual ACM symposium on Theory of computing*. ACM, 705–714.
- [5] D. E. Amos. 1974. Computation of Modified Bessel Functions and Their Ratios. *Math. Comp.* 28, 125 (Jan. 1974), 239. <https://doi.org/10.2307/2005830>
- [6] L. Aparicio, R. Van Meter, and H. Esaki. 2011. Protocol design for quantum repeater networks. In *Proc. Asian Conference on Internet Engineering*. ACM, Pattaya, Thailand.
- [7] David D Awschalom, Ronald Hanson, Jörg Wrachtrup, and Brian B Zhou. 2018. Quantum technologies with optically interfaced solid-state spins. *Nature Photonics* 12, 9 (2018), 516.
- [8] Sean D Barrett and Pieter Kok. 2005. Efficient high-fidelity quantum computation using matter qubits and linear optics. *Physical Review A* 71, 6 (2005), 060310.
- [9] Michael Ben-Or and Avinatan Hassidim. 2005. Fast quantum Byzantine agreement. In *Proceedings of the thirty-seventh annual ACM symposium on Theory of computing*. ACM, 481–485.
- [10] Charles H Bennett and Gilles Brassard. 1984. Quantum Cryptography: Public Key Distribution, and Coin-Tossing. In *Proc. 1984 IEEE International Conference on Computers, Systems, and Signal Processing*. 175–179. <https://doi.org/10.1016/j.tcs.2011.08.039>
- [11] Charles H Bennett, Gilles Brassard, Claude Crépeau, Richard Jozsa, Asher Peres, and William K Wootters. 1993. Teleporting an unknown quantum state via dual classical and Einstein-Podolsky-Rosen channels. *Physical Review Letters* 70, 13 (1993), 1895.
- [12] Hannes Bernien. 2014. *Control, measurement and entanglement of remote quantum spin registers in diamond*. PhD Thesis. TU Delft. <https://doi.org/10.4233/uuid:75130c37-edb5-4a34-ac2f-c156d377ca55>
- [13] Hannes Bernien, Bas Hensen, Wolfgang Pfaff, Gerwin Koolstra, MS Blok, Lucio Robledo, TH Taminiau, Matthew Markham, DJ Twitchen, Lilian Childress, et al. 2013. Heralded entanglement between solid-state qubits separated by three metres. *Nature* 497, 7447 (2013), 86.
- [14] BB Blinov, DL Moehring, L-M Duan, and Chris Monroe. 2004. Observation of entanglement between a single trapped atom and a single photon. *Nature* 428, 6979 (2004), 153.
- [15] Stefan Bogdanović, Suzanne B van Dam, Cristian Bonato, Lisanne C Coenen, Anne-Marije J Zwerver, Bas Hensen, Madelaine SZ Liddy, Thomas Fink, Andreas Reiserer, Marko Lončar, et al. 2017. Design and low-temperature characterization of a tunable microcavity for diamond-based quantum networks. *Applied Physics Letters* 110, 17 (2017), 171103.
- [16] Agata M Branczyk. 2017. Hong-Ou-Mandel Interference. *arXiv preprint: 1711.00080* (2017).
- [17] Anat Bremler-Barr, Yotam Harchol, and David Hay. 2016. OpenBox: A Software-Defined Framework for Developing, Deploying, and Managing Network Functions. In *Proc. SIGCOMM*. ACM, Florianopolis, Brazil.
- [18] Jürgen Brendel, Nicolas Gisin, Wolfgang Tittel, and Hugo Zbinden. 1999. Pulsed energy-time entangled twin-photon source for quantum communication. *Physical Review Letters* 82, 12 (1999), 2594.

- [19] H-J Briegel, Wolfgang Dür, Juan I Cirac, and Peter Zoller. 1998. Quantum repeaters: the role of imperfect local operations in quantum communication. *Physical Review Letters* 81, 26 (1998), 5932.
- [20] Anne Broadbent, Joseph Fitzsimons, and Elham Kashefi. 2009. Universal blind quantum computation. In *Foundations of Computer Science, 2009. FOCS'09. 50th Annual IEEE Symposium on*. IEEE, 517–526.
- [21] M. H. Abobeih R. Berrevoets M. Degen M. A. Bakker R. F. L. Vermeulen M. Markham D. J. Twitchen T. H. Taminiau C. E. Bradley, J. Randall. in preparation. A 10-qubit solid-state spin register with quantum memory exceeding 10 seconds. (in preparation).
- [22] C Cabrillo, J Ignacio Cirac, P Garcia-Fernandez, and P Zoller. 1999. Creation of entangled states of distant atoms by interference. *Physical Review A* 59, 2 (1999), 1025.
- [23] André Chailloux and Iordanis Kerenidis. 2011. Optimal bounds for quantum bit commitment. In *Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on*. IEEE, 354–362.
- [24] Fangfei Chen, Ramesh K. Sitaraman, and Marcelo Torres. 2015. End-User Mapping: Next Generation Request Routing for Content Delivery. In *Proc. SIGCOMM*. ACM, London, UK.
- [25] Lilian Childress and Ronald Hanson. 2013. Diamond NV centers for quantum computing and quantum networks. *MRS Bulletin* 38, 2 (2013), 134–138. <https://doi.org/10.1557/mrs.2013.20>
- [26] Chin-Wen Chou, H De Riedmatten, D Felinto, SV Polyakov, SJ Van Enk, and H Jeff Kimble. 2005. Measurement-induced entanglement for excitation stored in remote atomic ensembles. *Nature* 438, 7069 (2005), 828.
- [27] Cisco. 2005. Calculating the Maximum Attenuation for Optical Fiber Links. (March 2005). <https://www.cisco.com/c/en/us/support/docs/optical-networking/ons-15454-sonet-multiservice-provisioning-platform-mspp/27042-max-att-27042.html>
- [28] NS-3 Consortium. 2019. NS-3 Documentation: ns3::WifiRemoteStationInfo Class Reference. (Jan. 2019). [https://www.nsnam.org/doxygen/classns3\\_1\\_1\\_wif%5Ci\\_remote\\_station\\_info.html](https://www.nsnam.org/doxygen/classns3_1_1_wif%5Ci_remote_station_info.html)
- [29] Antonio D Córcoles, Easwar Magesan, Srikanth J Srinivasan, Andrew W Cross, Matthias Steffen, Jay M Gambetta, and Jerry M Chow. 2015. Demonstration of a quantum error detection code using a square lattice of four superconducting qubits. *Nature communications* 6 (2015), 6979.
- [30] Eric Corndorf, Chuang Liang, Gregory S. Kanter, Prem Kumar, and Horace P. Yuen. 2004. Quantum-noise-protected data encryption for WDM fiber-optic networks. *SIGCOMM Comput. Commun. Rev.* 34, 5 (Oct. 2004), 21–30.
- [31] Julia Cramer, Norbert Kalb, M Adriaan Rol, Bas Hensen, Machiel S Blok, Matthew Markham, Daniel J Twitchen, Ronald Hanson, and Tim H Taminiau. 2016. Repeated quantum error correction on a continuously encoded qubit by real-time feedback. *Nature communications* 7 (2016), 11526.
- [32] J. Cramer, N. Kalb, M. A. Rol, B. Hensen, M. S. Blok, M. Markham, D. J. Twitchen, R. Hanson, and T. H. Taminiau. 2016. Repeated quantum error correction on a continuously encoded qubit by real-time feedback. *Nature Communications* (2016). <https://doi.org/10.1038/ncomms11526> arXiv:1508.01388
- [33] Ivan B Damgård, Serge Fehr, Louis Salvail, and Christian Schaffner. 2007. Secure identification and QKD in the bounded-quantum-storage model. In *Annual International Cryptology Conference*. Springer, 342–359.
- [34] Ivan B Damgård, Serge Fehr, Louis Salvail, and Christian Schaffner. 2008. Cryptography in the bounded-quantum-storage model. *SIAM J. Comput.* 37, 6 (2008), 1865–1890.
- [35] Aymeric Delteil, Zhe Sun, Wei-bo Gao, Emre Togan, Stefan Faelt, and Ataç Imamoğlu. 2016. Generation of heralded entanglement between distant hole spins. *Nature Physics* 12, 3 (March 2016), 218–223. <https://doi.org/10.1038/nphys3605> arXiv:1507.00465
- [36] Vasil S Denchev and Gopal Pandurangan. 2008. Distributed quantum computing: A new frontier in distributed systems or science fiction? *ACM SIGACT News* 39, 3 (2008), 77–95.
- [37] Fahad R. Dogar, Thomas Karagiannis, Hitesh Ballani, and Antony Rowstron. 2014. Decentralized Task-Aware Scheduling for Data Center Networks. In *Proc. SIGCOMM*. ACM, Chicago, IL, USA.
- [38] Wolfgang Dür and Hans J Briegel. 2007. Entanglement purification and quantum error correction. *Reports on Progress in Physics* 70, 8 (2007), 1381.
- [39] JF Dynes, H Takesue, ZL Yuan, AW Sharpe, K Harada, T Honjo, H Kamada, O Tadanaga, Y Nishida, M Asobe, et al. 2009. Efficient entanglement distribution over 200 kilometers. *Optics express* 17, 14 (2009), 11440–11449.
- [40] Artur K Ekert. 1991. Quantum cryptography based on Bell's theorem. *Physical Review Letters* 67, 6 (1991), 661.
- [41] Julio Filho, Zoltan Papp, Relja Djapic, and Job Oostveen. 2013. Model-based Design of Self-adapting Networked Signal Processing Systems. In *Proc. SASO*. IEEE, 41–50. <https://doi.org/10.1109/SASO.2013.16>
- [42] Julio Filho, Teus Vogel, and Jan de Gier. 2016. *Runtime Services and Tooling for Reconfiguration*. 69–92. [https://doi.org/10.1007/978-981-10-0715-6\\_3](https://doi.org/10.1007/978-981-10-0715-6_3)
- [43] Joseph F Fitzsimons and Elham Kashefi. 2017. Unconditionally verifiable blind quantum computation. *Physical Review A* 96, 1 (2017), 012303.
- [44] Toru Fujiwara, Tadao Kasami, and Shu Lin. 1989. Error detecting capabilities of the shortened Hamming codes adopted for error detection in IEEE Standard 802.3. *IEEE Trans. Commun.* 37, 9 (sep 1989), 986–989.
- [45] Wei-Bo Gao, Chao-Yang Lu, Xing-Can Yao, Ping Xu, Otfried Gühne, Alexander Goebel, Yu-Ao Chen, Cheng-Zhi Peng, Zeng-Bing Chen, and Jian-Wei Pan. 2010. Experimental demonstration of a hyper-entangled ten-qubit Schrödinger cat state. *Nature physics* 6, 5 (2010), 331.
- [46] Daniel Gottesman, Thomas Jennewein, and Sarah Croke. 2012. Longer-baseline telescopes using quantum repeaters. *Physical Review Letters* 109, 7 (2012), 070503.
- [47] Saikat Guha, Hari Krovi, Christopher A Fuchs, Zachary Dutton, Joshua A Slater, Christoph Simon, and Wolfgang Tittel. 2015. Rate-loss analysis of an efficient quantum repeater architecture. *Physical Review A* 92, 2 (2015), 022357.
- [48] Mark Handley, Costin Raiciu, Alexandru Agache, Andrei Voinescu, Andrew W. Moore, Gianni Antichi, and Marcin Wójcik. 2017. Re-architecting datacenter networks and stacks for low latency and high performance. In *Proc. SIGCOMM*. ACM, Los Angeles, CA, USA.
- [49] B. Hensen, H. Bernien, A. E. Dréau, A. Reiserer, N. Kalb, M. S. Blok, J. Ruitenberg, R. F. L. Vermeulen, R. N. Schouten, C. Abellán, W. Amaya, V. Pruneri, M. W. Mitchell, M. Markham, D. J. Twitchen, D. Elkouss, S. Wehner, T. H. Taminiau, and R. Hanson. 2015. Loophole-free Bell inequality violation using electron spins separated by 1.3 kilometres. *Nature* 526, 7575 (oct 2015), 682–686. <https://doi.org/10.1038/nature15759> arXiv:1508.05949
- [50] Julian Hofmann, Michael Krug, Norbert Ortégel, Lea Gérard, Markus Weber, Wenjamin Rosenfeld, and Harald Weinfurter. 2012. Heralded entanglement between widely separated atoms. *Science* 337, 6090 (2012), 72–75.
- [51] C.-K. Hong, Z.-Y. Ou, and L Mandel. 1987. Measurement of subpicosecond time intervals between two photons by interference. *Physical Review Letters* 59 (nov 1987), 2044–2046. <https://doi.org/10.1103/>

- PhysRevLett.59.2044
- [52] Michał Horodecki and Paweł Horodecki. 1999. Reduction criterion of separability and limits for a class of distillation protocols. *Physical Review A* 59, 6 (1999), 4206.
- [53] Peter C. Humphreys, Norbert Kalb, Jaco P.J. Morits, Raymond N. Schouten, Raymond F.L. Vermeulen, Daniel J. Twitchen, Matthew Markham, and Ronald Hanson. 2018. Deterministic delivery of remote entanglement on a quantum network. *Nature* (2018). <https://doi.org/10.1038/s41586-018-0200-5> arXiv:1712.07567
- [54] IEEE 802.1 working group. 2015. 802.1AE - Media Access Control (MAC) Security. (2015).
- [55] Takahiro Inagaki, Nobuyuki Matsuda, Osamu Tadanaga, Masaki Asobe, and Hiroki Takesue. 2013. Entanglement distribution over 300 km of fiber. *Optics express* 21, 20 (2013), 23241–23249.
- [56] Laura Bryony James. 2005. *Error Behaviour in Optical Networks*. Ph.D. Dissertation. University of Cambridge. <https://www.repository.cam.ac.uk/bitstream/handle/1810/265632/LauraJamesPhDdissertationErrorBehaviourInOpticalNetworks.pdf>
- [57] Brian Julsgaard, Alexander Kozhekin, and Eugene S Polzik. 2001. Experimental long-lived entanglement of two macroscopic objects. *Nature* 413, 6854 (2001), 400.
- [58] N. Kalb, P. C. Humphreys, J. J. Slim, and R. Hanson. 2018. Dephasing mechanisms of diamond-based nuclear-spin memories for quantum networks. (2018), 1–11. <https://doi.org/10.1103/PhysRevA.97.062330> arXiv:1802.05996
- [59] Norbert Kalb, Andreas Reiserer, Stephan Ritter, and Gerhard Rempe. 2015. Heralded storage of a photonic quantum bit in a single atom. *Physical Review Letters* 114, 22 (2015), 220501.
- [60] N. Kalb, A. A. Reiserer, P. C. Humphreys, J. J. W. Bakermans, S. J. Kamerling, N. H. Nickerson, S. C. Benjamin, D. J. Twitchen, M. Markham, and R. Hanson. 2017. Entanglement distillation between solid-state quantum network nodes. *Science* 356, 6341 (jun 2017), 928–932. <https://doi.org/10.1126/science.aan0070> arXiv:1703.03244
- [61] Globix Vertriebs GmbH & Co. KG. 2013. SFP-ZX 1.25 Gb/s Single-Mode SFP Transceiver 1000BASE-ZX. (Oct. 2013). <http://www.globix-it.de/datasheets/SFP-ZX.pdf>
- [62] H Jeff Kimble. 2008. The quantum internet. *Nature* 453, 7198 (2008), 1023.
- [63] Peter Komar, Eric M Kessler, Michael Bishof, Liang Jiang, Anders S Sørensen, Jun Ye, and Mikhail D Lukin. 2014. A quantum network of clocks. *Nature Physics* 10, 8 (2014), 582.
- [64] Bo Liu, Baokang Zhao, Ziling Wei, Chunqing Wu, Jinshu Su, Wanrong Yu, Fei Wang, and Shihai Sun. 2013. Qphone: A Quantum Security VoIP Phone. In *Proc. SIGCOMM*. ACM, Hong Kong.
- [65] Seth Lloyd, Jeffrey H. Shapiro, Franco N. C. Wong, Prem Kumar, Selim M. Shahriar, and Horace P. Yuen. 2004. Infrastructure for the Quantum Internet. *SIGCOMM Comput. Commun. Rev.* 34, 5 (Oct. 2004), 9–20.
- [66] Ilias Marinos, Robert N. M. Watson, and Mark Handley. 2014. Network Stack Specialization for Performance. In *Proc. SIGCOMM*. ACM, Chicago, IL, USA.
- [67] Klaus Mattle, Harald Weinfurter, Paul G Kwiat, and Anton Zeilinger. 1996. Dense coding in experimental quantum communication. *Physical Review Letters* 76, 25 (1996), 4656.
- [68] R. Van Meter. 2012. Quantum networking and internetworking. *IEEE Network* 26, 4 (2012), 59–64.
- [69] Rodney Van Meter and Joe Touch. 2013. Designing quantum repeater networks. *IEEE Comm. Mag.* 51, 8 (Aug. 2013).
- [70] DL Moehring, P Maunz, S Olmschenk, KC Younge, DN Matsukevich, L-M Duan, and C Monroe. 2007. Entanglement of single-atom quantum bits at a distance. *Nature* 449, 7158 (2007), 68.
- [71] WJ Munro, AM Stephens, SJ Devitt, KA Harrison, and Kae Nemoto. 2012. Quantum communication without the necessity of quantum memories. *Nature Photonics* 6, 11 (2012), 777.
- [72] William J Munro, Koji Azuma, Kiyoshi Tamaki, and Kae Nemoto. 2015. Inside quantum repeaters. *IEEE Journal of Selected Topics in Quantum Electronics* 21, 3 (2015), 78–90.
- [73] Sreraman Muralidharan, Jungsang Kim, Norbert Lütkenhaus, Mikhail D Lukin, and Liang Jiang. 2014. Ultrafast and fault-tolerant quantum communication across long distances. *Physical Review Letters* 112, 25 (2014), 250501.
- [74] A Narla, S Shankar, M Hatridge, Zaki Leghtas, KM Sliwa, E Zalys-Geller, SO Mundhada, W Pfaff, L Frunzio, RJ Schoelkopf, et al. 2016. Robust concurrent remote entanglement between two superconducting qubits. *Physical Review X* 6, 3 (2016), 031036.
- [75] Kae Nemoto, Michael Trupke, Simon J Devitt, Burkhard Scharfenberger, Kathrin Buczak, Jörg Schmiedmayer, and William J Munro. 2016. Photonic Quantum Networks formed from NV- centers. *Scientific reports* 6 (2016), 26284.
- [76] Michael A. Nielsen and Isaac L. Chuang. 2010. *Quantum Computation and Quantum Information* (10th ed.). Cambridge University Press, Cambridge. <https://doi.org/10.1017/CBO9780511976667>
- [77] Simon Peter, Umar Javed, Qiao Zhang, Doug Woos, Thomas Anderson, and Arvind Krishnamurthy. 2014. One Tunnel is (Often) Enough. In *Proc. SIGCOMM*. ACM, Chicago, IL, USA.
- [78] W. Pfaff, B. J. Hensen, H. Bernien, S. B. van Dam, M. S. Blok, T. H. Taminiau, M. J. Tiggelman, R. N. Schouten, M. Markham, D. J. Twitchen, and R. Hanson. 2014. Unconditional quantum teleportation between distant solid-state quantum bits. *Science* 345, 6196 (aug 2014), 532–535. <https://doi.org/10.1126/science.1253512> arXiv:1404.4369
- [79] C. Pfister, M. Adriaan Rol, A. Mantri, M. Tomamichel, and S. Wehner. 2018. Capacity estimation and verification of quantum channels with arbitrarily correlated errors. 9, 27 (2018).
- [80] A Pirker and W Dür. 2018. A quantum network stack and protocols for reliable entanglement-based networks. (Oct. 2018). arXiv:1810.03556 <https://arxiv.org/pdf/1810.03556.pdf> <http://arxiv.org/abs/1810.03556>
- [81] J. G. Rarity, M. Thompson, A. Lord, R. Nejabati, and D. Simeonidou. 2017. Secure NFV Orchestration Over an SDN-Controlled Optical Network With Time-Shared Quantum Key Distribution Resources. *Journal of Lightwave Technology* 35, 8 (2017), 1357–1362.
- [82] Andreas Reiserer, Norbert Kalb, Machiel S. Blok, Koen J.M. van Bemmelen, Tim H. Taminiau, Ronald Hanson, Daniel J. Twitchen, and Matthew Markham. 2016. Robust Quantum-Network Memory Using Decoherence-Protected Subspaces of Nuclear Spins. *Physical Review X* 6, 2 (jun 2016), 021040. <https://doi.org/10.1103/PhysRevX.6.021040> arXiv:1603.01602
- [83] Jérémy Ribeiro and Frédéric Grosshans. 2015. A tight lower bound for the bb84-states quantum-position-verification protocol. *arXiv preprint arXiv:1504.07171* (2015).
- [84] Daniel Riedel, Immo Söllner, Brendan J Shields, Sebastian Starosielec, Patrick Appel, Elke Neu, Patrick Maletinsky, and Richard J Warburton. 2017. Deterministic enhancement of coherent photon generation from a nitrogen-vacancy center in ultrapure diamond. *Physical Review X* 7, 3 (2017), 031040.
- [85] Daniel Riedel, Immo Söllner, Brendan J. Shields, Sebastian Starosielec, Patrick Appel, Elke Neu, Patrick Maletinsky, and Richard J. Warburton. 2017. Deterministic Enhancement of Coherent Photon Generation from a Nitrogen-Vacancy Center in Ultrapure Diamond. *Physical Review X* 7, 3 (sep 2017), 031040. <https://doi.org/10.1103/PhysRevX.7.031040>
- [86] Diego Riste, Stefano Poletto, M-Z Huang, Alessandro Bruno, Visa Vesterinen, O-P Saira, and Leonardo DiCarlo. 2015. Detecting bit-flip errors in a logical qubit using stabilizer measurements. *Nature*

- communications* 6 (2015), 6983.
- [87] Stephan Ritter, Christian Nölleke, Carolin Hahn, Andreas Reiserer, Andreas Neuzner, Manuel Uphoff, Martin Mücke, Eden Figueroa, Joerg Bochmann, and Gerhard Rempe. 2012. An elementary quantum network of single atoms in optical cavities. *Nature* 484, 7393 (2012), 195.
- [88] Filip Rozpędek, Raja Yehia, Kenneth Goodenough, Maximilian Ruf, Peter C. Humphreys, Ronald Hanson, Stephanie Wehner, and David Elkouss. 2018. Near-term quantum repeater experiments with NV centers: overcoming the limitations of direct transmission. (sep 2018), 1–35. arXiv:1809.00364 <http://arxiv.org/abs/1809.00364>
- [89] Nicolas Sangouard, Christoph Simon, Hugues De Riedmatten, and Nicolas Gisin. 2011. Quantum repeaters based on atomic ensembles and linear optics. *Reviews of Modern Physics* 83, 1 (2011), 33.
- [90] Valerio Scarani, Helle Bechmann-Pasquinucci, Nicolas J Cerf, Miloslav Dušek, Norbert Lütkenhaus, and Momtchil Peev. 2009. The security of practical quantum key distribution. *Reviews of modern physics* 81, 3 (2009), 1301.
- [91] Rachee Singh, Manya Ghobadi, and Klaus-Tycho Foerster. 2018. RAD-WAN: Rate Adaptive Wide Area Network. In *Proc. SIGCOMM*. ACM, Budapest, Hungary.
- [92] John Strand, Angela L. Chiu, and Robert Tkach. 2001. Issues For Routing In The Optical Layer. *IEEE Comm. Mag.* 39, 2 (2001), 81–87.
- [93] T. H. Taminiau, J. Cramer, T. van der Sar, V. V. Dobrovitski, and R. Hanson. 2014. Universal control and error correction in multi-qubit spin registers in diamond. *Nature Nanotechnology* 9, 3 (mar 2014), 171–176. <https://doi.org/10.1038/nnano.2014.2> arXiv:1309.5452
- [94] Barbara M Terhal. 2015. Quantum error correction for quantum memories. *Reviews of Modern Physics* 87, 2 (2015), 307.
- [95] M. Tomamichel, C.Lim, N. Gisin, and R. Renner. 2012. Tight finite-key analysis for quantum cryptography. *Nature Communications* 3, 634 (2012).
- [96] unpeng James Liu, Peter Xiang Gao, Bernard Wong, and Srinivasan Keshav. 2014. Quartz: A New Design Element for Low-Latency DCN. In *Proc. SIGCOMM*. ACM, Chicago, IL, USA.
- [97] Raju Valivarthi, Marcelli Grima Puigibert, Qiang Zhou, Gabriel H. Aguilar, Varun B. Verma, Francesco Marsili, Matthew D. Shaw, Sae Woo Nam, Daniel Oblak, and Wolfgang Tittel. 2016. Quantum teleportation across a metropolitan fibre network. *Nature Photonics* 10, 10 (oct 2016), 676–680. <https://doi.org/10.1038/nphoton.2016.180> arXiv:1605.08814
- [98] Suzanne B van Dam, Peter C Humphreys, Filip Rozpędek, Stephanie Wehner, and Ronald Hanson. 2017. Multiplexed entanglement generation over quantum networks using multi-qubit nodes. *Quantum Science and Technology* 2, 3 (2017), 034002.
- [99] C. J. van Leeuwen, J. M. de Gier, Julio A. de Oliveira Filho, and Zoltán Papp. 2014. Model-Based Architecture Optimization for Self-Adaptive Networked Signal Processing Systems. In *Eighth IEEE International Conference on Self-Adaptive and Self-Organizing Systems, SASO 2014, London, United Kingdom, September 8-12, 2014*. London, United Kingdom, 187–188. <https://doi.org/10.1109/SASO.2014.37>
- [100] Rodney Van Meter, T.D. Ladd, W.J. Munro, and Kae Nemoto. 2009. System Design for a Long-Line Quantum Repeater. *IEEE/ACM Transactions on Networking* 17, 3 (June 2009), 1002–1013. <https://doi.org/10.1109/TNET.2008.927260> arXiv:0705.4128
- [101] Stephanie Wehner, David Elkouss, and Ronald Hanson. 2018. Quantum internet: A vision for the road ahead. *Science* 362, 6412 (oct 2018), eaam9288. <https://doi.org/10.1126/science.aam9288>
- [102] Stephanie Wehner, Christian Schaffner, and Barbara M Terhal. 2008. Cryptography from noisy storage. *Physical Review Letters* 100, 22 (2008), 220502.
- [103] Juan Yin, Yuan Cao, Yu-Huai Li, Sheng-Kai Liao, Liang Zhang, Ji-Gang Ren, Wen-Qi Cai, Wei-Yue Liu, Bo Li, Hui Dai, Guang-Bing Li, Qi-Ming Lu, Yun-Hong Gong, Yu Xu, Shuang-Lin Li, Feng-Zhi Li, Ya-Yun Yin, Zi-Qing Jiang, Ming Li, Jian-Jun Jia, Ge Ren, Dong He, Yi-Lin Zhou, Xiao-Xiang Zhang, Na Wang, Xiang Chang, Zhen-Cai Zhu, Nai-Le Liu, Yu-Ao Chen, Chao-Yang Lu, Rong Shu, Cheng-Zhi Peng, Jian-Yu Wang, and Jian-Wei Pan. 2017. Satellite-based entanglement distribution over 1200 kilometers. *Science* 356, 6343 (jun 2017), 1140–1144. <https://doi.org/10.1126/science.aan3211> arXiv:1707.01339[quant-ph]
- [104] Wanrong Yu, Baokang Zhao, and Zhe Yan. 2018. Software defined quantum key distribution network. *2017 3rd IEEE International Conference on Computer and Communications, ICC3 2017 2018-Janua* (2018), 1293–1297. <https://doi.org/10.1109/CompComm.2017.8322751>
- [105] Sebastian Zaske, Andreas Lenhard, Christian A. Keßler, Jan Kettler, Christian Hepp, Carsten Arend, Roland Albrecht, Wolfgang-Michael Schulz, Michael Jetter, Peter Michler, and Christoph Becher. 2012. Visible-to-Telecom Quantum Frequency Conversion of Light from a Single Quantum Emitter. *Physical Review Letters* 109, 14 (oct 2012), 147404. <https://doi.org/10.1103/PhysRevLett.109.147404>
- [106] Liang Zheng, Carlee Joe-Wong, Chee Wei Tan, Mung Chiang, and Xinyu Wang. 2015. How to Bid the Cloud. In *Proc. SIGCOMM*. ACM, London, UK.
- [107] Marek Zukowski, Anton Zeilinger, Michael A Horne, and Arthur K Ekert. 1993. “Event-ready-detectors” Bell experiment via entanglement swapping. *Physical Review Letters* 71 (1993), 4287–4290.

## APPENDIX

In this Appendix, we provide further background and details, as well as more in-depth simulation results.

- In Section A, we provide a very brief introduction to quantum information, including concepts like entanglement, fidelity, QBER and provide an intuition on how the fidelity is related to memory lifetimes.
- In Section B, we explain how to estimate the fidelity by interspersing test rounds.
- In Section C, we provide further simulation results to illustrate protocol performance and further validate our simulation against hardware.
- In Section D, we provide further details of the NV platform as relevant for the design considerations of the proposed protocol. Furthermore, we provide details the physical modeling as well as numerical methods used to conduct the simulation.
- In Section E, we provide a complete description of the proposed protocol.

### A QUANTUM PRELUDE

This section provides a very short introduction to quantum information and in particular quantum states, entanglement, fidelity and decoherence. For a deeper introduction to quantum information, see for example [76].

#### A.1 Qubits and states

A quantum bit (*qubit*) is a two-level system, where the two levels are usually denoted  $|0\rangle$  and  $|1\rangle$  respectively (“ket”-notation) and called the *basis states* of the qubit. These levels can for example be two energy levels of an electron spin or - when considering transmitting qubits - vertical and horizontal polarization of a photon, presence or absence of a photon, or a time-bin of early and late. Compared to a “classical” bit  $|0\rangle$  or  $|1\rangle$ , a qubit can be in *superpositions* thereof. Mathematically, a state  $|\phi\rangle$  of a qubit is written as

$$|\phi\rangle = \alpha |0\rangle + \beta |1\rangle \quad (1)$$

where  $\alpha$  and  $\beta$  are arbitrary complex numbers with the constraint that  $|\alpha|^2 + |\beta|^2 = 1$ , and

$$|0\rangle = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, |1\rangle = \begin{pmatrix} 0 \\ 1 \end{pmatrix}. \quad (2)$$

Note that  $|0\rangle$  and  $|1\rangle$  form a basis for  $\mathbb{C}^2$ . Some common states are

$$|X, 0\rangle = \frac{1}{\sqrt{2}}(|0\rangle + |1\rangle) \quad |X, 1\rangle = \frac{1}{\sqrt{2}}(|0\rangle - |1\rangle) \quad (3)$$

$$|Y, 0\rangle = \frac{1}{\sqrt{2}}(|0\rangle + i|1\rangle) \quad |Y, 1\rangle = \frac{1}{\sqrt{2}}(|0\rangle - i|1\rangle) \quad (4)$$

$$|Z, 0\rangle = |0\rangle \quad |Z, 1\rangle = |1\rangle \quad (5)$$

corresponding to a ‘0’ or ‘1’ in the three different bases labeled  $X$ ,  $Y$ , and  $Z$ . The label  $Z$  also refers to the *standard basis*. We also use  $\langle\phi| = (|\phi\rangle^*)^T$  to denote the conjugate transpose of  $|\phi\rangle$ .

Measuring a qubit in the standard ( $Z$ ) basis ( $|0\rangle$ ,  $|1\rangle$ ), gives measurement outcomes ‘0’ (i.e.  $|0\rangle$ ) or ‘1’ (i.e.  $|1\rangle$ ). Measuring a qubit which is in the state  $|\phi\rangle$  as in equation 1 in the *standard basis*, yields the outcomes 0 or 1 with the following probabilities

$$P[\text{"measuring 0"}|Z\text{-basis}] = |\alpha|^2, \quad P[\text{"measuring 1"}|Z\text{-basis}] = |\beta|^2. \quad (6)$$

which is why  $|\phi\rangle$  needs to be normalized. Measuring a qubit in the standard basis collapses it to  $|0\rangle$  or  $|1\rangle$ . Measuring a qubit in the  $X$ - or  $Y$ -basis yields outcomes with probabilities

$$P[\text{"measuring 0"}|X/Y\text{-basis}] = |\langle X/Y, 0|\psi\rangle|^2, \quad P[\text{"measuring 1"}|X/Y\text{-basis}] = |\langle X/Y, 1|\psi\rangle|^2 \quad (7)$$

where  $\langle\cdot|\cdot\rangle$  is the inner product.

#### A.2 Entangled states

If qubit  $A$  is in a state  $|\phi_1\rangle$  and qubit  $B$  is in the state  $|\phi_2\rangle$ , then their joint state (at possibly remote nodes) is given by the *tensor product* of the individual states  $|\phi_1\rangle_A$  and  $|\phi_2\rangle_B$ , i.e. as

$$|\text{"joint state"}\rangle = |\phi_1\rangle_A \otimes |\phi_2\rangle_B. \quad (8)$$



Importantly, for the discussion here is that not all joint states can be factorized into single qubit states  $|\phi_1\rangle_A$  and  $|\phi_2\rangle_B$  in this way. These are called *entangled* states. For example, consider the state

$$|\Phi^+\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A \otimes |0\rangle_B + |1\rangle_A \otimes |1\rangle_B), \quad (9)$$

which is a superposition of (1) both qubits being in the state  $|0\rangle$  and (2) both qubits being in the state  $|1\rangle$ . This is an entangled state, i.e., it cannot be factorized into two individual states, giving rise to genuinely quantum correlations between  $A$  and  $B$  that have no classical analogue. The state  $|\Phi^+\rangle$  is one of the so called *Bell states*. These are entangled states, where the other three are given as

$$|\Phi^-\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A \otimes |0\rangle_B - |1\rangle_A \otimes |1\rangle_B), \quad (10)$$

$$|\Psi^+\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A \otimes |1\rangle_B + |1\rangle_A \otimes |0\rangle_B), \quad (11)$$

$$|\Psi^-\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A \otimes |1\rangle_B - |1\rangle_A \otimes |0\rangle_B). \quad (12)$$

Measurement outcomes of measuring the two qubits in any of the Bell-states in the bases  $X$ ,  $Z$  and  $Y$  are either perfectly correlated or perfectly anti-correlated. For example, for  $|\Phi^+\rangle$  the measurement outcomes are perfectly correlated in the  $X$  and  $Z$  bases but perfectly anti-correlated in the  $Y$  basis. On the other hand, for  $|\Psi^-\rangle$  the measurement outcomes are perfectly anti-correlated in all three bases.

Relevant to understand the generation of entanglement is that all the Bell-states can be transformed to one another by only applying *local* quantum gates to one of the qubits. Two useful gates are the bit flip  $X|x\rangle = |x+1 \pmod{2}\rangle$  and phase flip  $Z|x\rangle = (-1)^x|x\rangle$ . Applying them on qubit  $A$  (at node  $A$  only) allows one to transform:

$$|\Phi^-\rangle = Z^A|\Phi^+\rangle, \quad |\Psi^+\rangle = X^A|\Phi^+\rangle, \quad |\Psi^-\rangle = Z^AX^A|\Phi^+\rangle, \quad (13)$$

where we added the index  $A$  to emphasize the gates are applied to qubit  $A$ . We could also apply such gates to qubit  $B$  to have the same effect.

In the heralded entanglement generation (Section 4.4), we can obtain either failure, or else success. In the case of success, an additional bit indicates whether we produced the state  $|\Psi^+\rangle$  or  $|\Psi^-\rangle$ . From equation (13), these two states can be transformed between each other by simply applying a  $Z$ -gate to one of the qubits.

### A.3 Fidelity and QBER

In any real implementation of a quantum network, the generated entangled states will always differ from the perfect Bell states above due to noise in the system. When writing noisy states, it is convenient to express the state as a *density matrix*. For a perfectly prepared state  $|\Psi^-\rangle$ , the density matrix is  $\rho = |\Psi^-\rangle\langle\Psi^-|$ . This allows one to express noise. For example, the analogue of applying a classical bit flip error  $X$  with some probability  $p_{\text{err}}$  can be written as

$$\rho_{\text{noisy}} = (1 - p_{\text{err}})\rho + p_{\text{err}}X\rho X. \quad (14)$$

The *fidelity*  $F$  measures how close a realized state  $\rho$  is to an ideal target state  $|\Psi^-\rangle$ . The fidelity of a state  $\rho$  with the target state  $|\Psi^-\rangle$  can be written as

$$F[|\Psi^-\rangle] = \langle\Psi^-|\rho|\Psi^-\rangle \quad (15)$$

where  $F = 1$  if  $\rho$  is identical to the target state. We have  $0 \leq F \leq 1$ , where a larger value of  $F$  means we are closer to the target state.

It is important to note that one cannot measure the fidelity of a single instance of a quantum state. However, if we produce the same state many times in succession, we can estimate its fidelity. One way to do this is to measure the qubit-error-rate (*QBER*). Consider  $|\Psi^-\rangle$  above and recall that measurement outcomes in the  $X$ ,  $Z$  and  $Y$  bases are always perfectly anti-correlated in this case. I.e. we always get different measurement outcome for qubit  $A$  and for qubit  $B$ . In case the state is noisy, this is no longer the case. For a fixed basis (say  $Z$ ) the QBER (here  $\text{QBER}_Z$ ) is the probability of receiving equal<sup>3</sup> measurement outcomes,

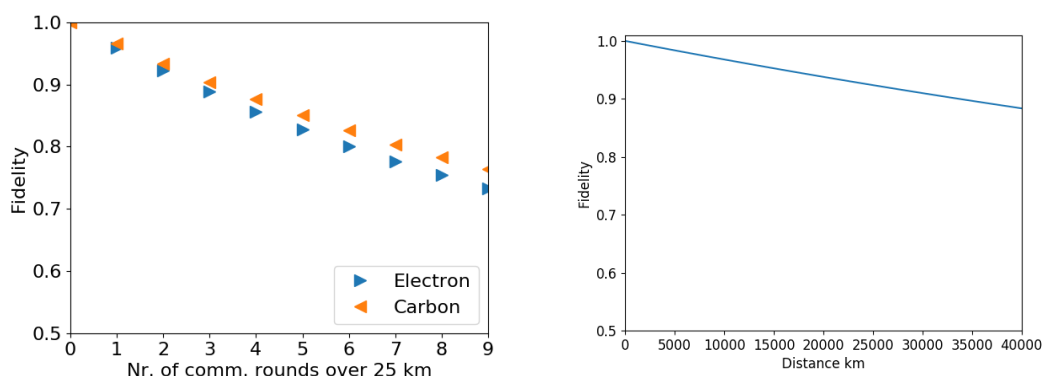
<sup>3</sup>QBER for the other Bell states is defined in a similar manner, taking into account that measurement outcomes are always equal in some bases for the other ideal Bell states.

when measuring qubit  $A$  and qubit  $B$  in the  $Z$  basis. Similarly, we can define  $\text{QBER}_X$  and  $\text{QBER}_Y$  for measurements in the  $X$  and  $Y$  bases. One can show that the fidelity and QBER of the Bell state state  $|\Psi^-\rangle$  are related as

$$F[|\Psi^-\rangle] = 1 - \frac{\text{QBER}_X + \text{QBER}_Y + \text{QBER}_Z}{2}. \quad (16)$$

## A.4 Decoherence

Quantum memories are inherently noisy and the amount of noise a qubit experiences depends on how long it stays in the memory. How long a qubit state is preserved is usually captured by the two numbers  $T_1$  (energy/thermal relaxation time) and  $T_2$  (dephasing time) of the qubit [76], as well as free-induction decay  $T_2^*$  (see e.g. [60]). Here, we focus on our performance metric (the fidelity) and illustrate in Figure 9a how it behaves as a function of time, where to highlight the actual effect of limited memory lifetimes we show the timescales in terms kms in fiber, where  $c = 206753$  km/s is the speed of light in fiber. Figure 9 shows the fidelity of an entangled state stored in two electron states with a coherence time ( $T_2$ ) of 1.46 s as a function of the time it takes to communicate over a certain distance.



(a) Reduction in fidelity  $F$  when storing a perfect entangled state  $|\Psi^+\rangle$  in the communication (blue, left triangles) and memory (orange, right triangles) qubit in terms of the number of communication rounds in the QL2020 scenario (over 25 km). Noise parameters listed in Table 6 ( $T_1 = 2.68$  ms and  $T_2 = 1.00$  ms for communication and  $T_1 = \infty$  and  $T_2 = 3.5$  ms for memory qubit).

(b) Illustration of an improved communication qubit by means of dynamical decoupling with  $T_2 = 1.46$  s ( $T_1 = \infty$ ). If such a qubit was used in an NV platform connected to a network, the qubit could be kept alive while waiting for classical control communication over long distances.

Figure 9

## B TESTING

We now explain the test used to gain confidence in transmission quality, specifically to estimate the quality parameter fidelity  $F$  used in the FEU (Section 5.2.3). The following is a standard procedure in quantum information to estimate the fidelity  $F = F[|\Psi^-\rangle]$  of a state  $\rho$  to the entangled target state  $|\Psi^-\rangle$ . We emphasize that it is not possible to measure  $F$  from a single copy of the state  $\rho$ . The matrices  $\rho$  are a mathematical description of an underlying quantum system, and not a matrix that one can read or access like classical information.

We first describe the standard procedure in the way that it is normally used. We then outline how this protocol can be extended to the case of interest here, and why we can draw conclusions even in real world scenarios in which we can experience arbitrary correlated errors.

Let us first assume a simpler scenario, in which  $n$  identical noisy entangled states  $\rho$  are produced in succession and we want to estimate  $F$ . We remark that when using imperfect quantum devices it is evidently a highly idealized situation that all states  $\rho$  are exactly identical. We can see from Eq. (16) that we can express  $F$  in terms of the quantum bit error rates  $\text{QBER}_X$ ,  $\text{QBER}_Z$ , and  $\text{QBER}_Y$ , which immediately suggests a protocol: specifically, we will estimate the QBERs in bases  $X$ ,  $Z$  and  $Y$  to obtain  $F$ . We sketch such a protocol in a specific way to build intuition for the more general procedure below:

## A Link Layer Protocol for Quantum Networks

- Node  $A$  randomly chooses an  $n$  element string  $r = r_1, \dots, r_n \in \{X, Z, Y\}$  and sends it to Node  $B$ .
- Nodes  $A$  and  $B$  now perform the following procedure for  $1 \leq j \leq n$  rounds:
  - Node  $A$  produces one entangled pair  $\rho$  with Node  $B$ .
  - Nodes  $A$  and  $B$  both measure their respective qubits in the basis  $r_j$  and record outcomes  $x_j^A$  (Node  $A$ ) and  $x_j^B$  (Node  $B$ ) respectively.
- Node  $B$  ( $A$ ) transmits the outcome string  $x^B = x_1^B, \dots, x_n^B$  ( $x^A = x_1^A, \dots, x_n^A$ ) to Node  $A$  ( $B$ ).
- Both nodes estimate the error rates

$$QBER_Z \approx \frac{\#\{j \mid x_j^A = x_j^B, r_j = Z\}}{\#\{j \mid r_j = Z\}}, \quad (17)$$

$$QBER_X \approx \frac{\#\{j \mid x_j^A = x_j^B, r_j = X\}}{\#\{j \mid r_j = X\}}, \quad (18)$$

$$QBER_Y \approx \frac{\#\{j \mid x_j^A = x_j^B, r_j = Y\}}{\#\{j \mid r_j = Y\}}, \quad (19)$$

$$(20)$$

where  $\#\{j \mid \text{condition}\}$  is the number of indices  $1 \leq j \leq n$  satisfying the stated condition.

Using Eq. (16) then yields an estimate of  $F$ .

Before we continue it may be instructive to compare the procedure above to the classical world. Evidently, classically, one way to gain confidence in a channels ability to transmit classical bits would be rather similar: Instead of preparing states  $\rho$ , we choose  $n$  random bits and send them. In the end, we estimate the error rate. Translated to the quantum setting, we would be preparing random bits  $|0\rangle$  and  $|1\rangle$  at node  $A$ , and sending them to node  $B$  which measures them in the  $Z$  bases to obtain an estimate of the bit error rate, similarly to  $QBER_Z$ . Such an estimate can give us confidence that also future bits are likely to be transmitted with roughly the same amount of errors as the test bits. This of course does *not* allow the same level of confidence as error detection in the quality of transmission. Specifically, a CRC is a check for a *specific* piece of data (e.g. one frame in 100 Base T), whereas such a test only yields a confidence in transmission quality.

Creating an analogous quantum CRC is theoretically possible by using a quantum error correcting code [76], but technologically highly challenging and highly infeasible for many years to come. Yet, we remark that also in a future in which such methods would become feasible we may not want to employ them because the requirements of our use cases are different. Since many protocols for our use cases are probabilistic, or make many pairs (especially NL and MD use cases), we often do not require more confidence on the exact quality of a single pair. Indeed, we can pass errors all the way up to the application level (such as for example in QKD [10]), where errors are then corrected using classical instead of quantum error correction. In such cases, fluctuations in quality are indeed expected at the application level. Here, using fast and easy to produce test rounds may remain preferable over more time consuming quantum CRCs.

The protocol above is limited in two ways: (1) all states were assumed to be both identical and independent from each other. I.e., there are no memory effects in the noise. Such memory effects are non-trivial in the quantum regime since they may inadvertently create (some amount of) entanglement not only between  $A$  and  $B$ , but between subsequent pairs produced. (2) We measured all  $n$  rounds, consuming all entangled pairs. Instead, we would like a protocol in which only test rounds are interspersed, and we can draw an inference about the pairs which we did *not* measure. Again, the possibility of quantum correlated noise between subsequent rounds makes this non-trivial.

To achieve this, we use a slight variant of the above as in [79]. Precise statistical statements are relatively straightforward - but very lengthy - to obtain using the techniques in [79, 95] and are out of scope of this paper. Here, we focus on the practical protocol and intuition without the need for mathematical tools from quantum information, which is a direct extension of the one above:

- Nodes  $A$  and  $B$  agree on a sampling window  $N$ .
- Nodes  $A$  and  $B$  randomly pick an  $N$  bit string  $t = t_1, \dots, t_N$  where  $\Pr[t_j = 1] = q$  for some parameter  $q$  determining the frequency of using test rounds.  $A$  and  $B$  periodically refresh  $t$  as needed.
- Nodes  $A$  and  $B$  randomly pick an  $N$  element basis string  $r = r_1, \dots, r_N \in \{X, Z, Y\}$ .  $A$  and  $B$  periodically refresh  $r$  as needed.

- The EGP uses  $t$  to determine when to intersperse a test round. When producing the  $j$ -th response to the MHP, the EGP checks whether  $t_j = 1$ . If so, it uses a standard test response instead to attempt to produce a test pair  $\rho$ , and takes as the measurement basis the next available in the random basis string  $r$ .
- $A$  and  $B$  record their measurement outcomes.
- $A$  and  $B$  estimate  $\text{QBER}_X, \text{QBER}_Z, \text{QBER}_Y$  over the past  $N$  rounds of producing entanglement (tested and untested 'data' rounds)

The key insight in the analysis of this procedure is that we can (with some amount of confidence depending on  $N$  and  $q$ ) use the QBER measured on the test rounds to determine the QBER on the untested - i.e. data - rounds [95, Inequality 1.3]. Using Eq. (16), then again allows one to draw conclusions about the average fidelity of the untested rounds to inform the FEU.

## C SIMULATION AND MODELING

We here provide additional simulation results, and further verification against the quantum hardware.

For our simulations we make use of a purpose built discrete event simulator for quantum networks: NetSquid<sup>4</sup>. By utilizing the discrete event paradigm NetSquid is capable of efficiently simulating the transmission and decay of quantum information in combination with the complex and stochastic nature of communication protocols. NetSquid can simulate both arbitrary quantum operations and Clifford-only gates, the former allowing for a precise simulation of small networks, while the latter allowing for networks containing thousands of nodes and qubits to be studied. Complete libraries of base classes enable users to simulate protocols and model physical devices at different levels of abstraction; for instance, (quantum) channels with modular noise, loss and delay models, or quantum processing devices with configurable gate topologies. NetSquid thus provides an ideal tool to validate network design choices and verify the performance of quantum network protocols in a physically-realistic setting.

The core simulation engine used by NetSquid is based on DynAA[41, 42, 99], a computer-aided analysis and design tool for the development of large, distributed, adaptive, and networked systems. It combines the best of network and system simulation technologies in a discrete-event modeling framework. DynAA provides a simple, yet powerful language able to describe large and complex system architectures, and innovative constructs to express adaptation mechanisms of the system, such as dynamic parameterization, and functional and architectural reconfiguration. A DynAA model can be simulated and/or analyzed to reveal system wide performance indicators, such as – but not limited to – throughput, power consumption, connectivity, reliability, and availability.

### C.1 Validation of simulation

We compare our simulation model against further data gathered from the NV platform LABscenario. Node A rotates its qubit over the Z-axis of the Bloch sphere by a fixed angle, followed by measuring its communication qubit in a basis ( $X, Y$  or  $Z$ ) that the nodes agreed upon beforehand. Node B only performs the measurement on its communication qubit, in the same pre-agreed basis. Regardless of the signal from the heralding station, both nodes initialize their qubit in  $|0\rangle$  before the start of the next round.

We compute the correlations of the measurement outcomes ( $m_A, m_B = \pm 1$ ) as shown in Figure 10 using

$$\Pr(m_A \neq m_B) = \frac{1 - \langle \mathcal{B} \otimes \mathcal{B} \rangle}{2}$$

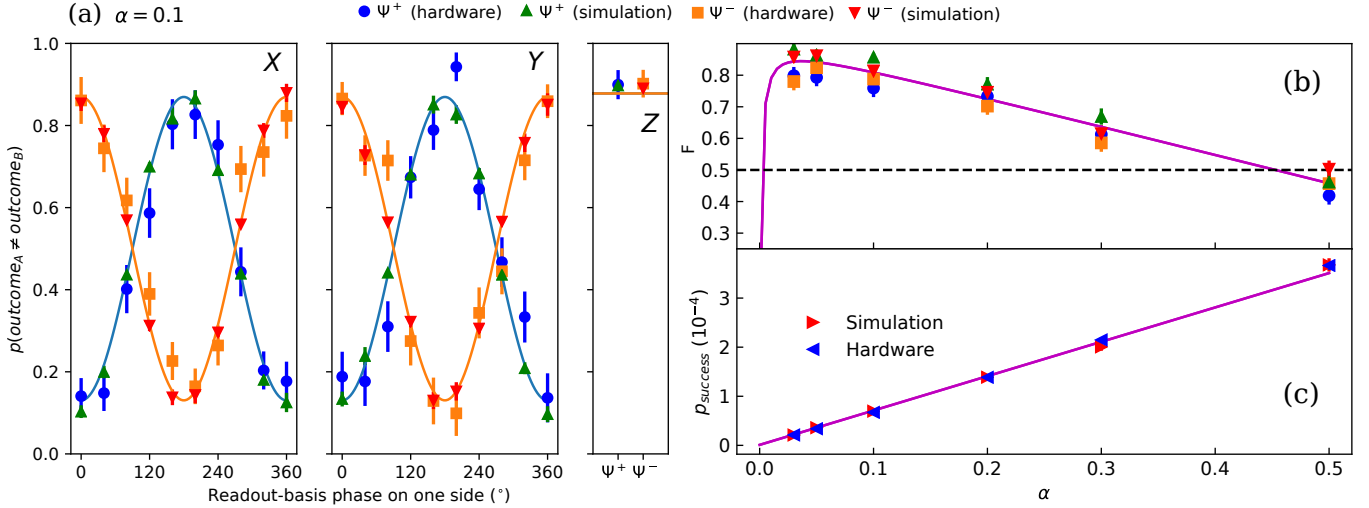
where  $\langle \mathcal{B} \otimes \mathcal{B} \rangle$  is the expectation value of the product of joint measurement outcomes  $m_A \cdot m_B$  with  $\mathcal{B} \in \{X, Y, Z\}$  the measurement basis after rotation.

The fidelity with the target state  $|\Psi^\pm\rangle$  (where  $\pm$  denotes the heralding detector) can be expressed as a function of the correlations as

$$\frac{1}{4} [1 \pm \langle X \otimes X \rangle \pm \langle Y \otimes Y \rangle - \langle Z \otimes Z \rangle].$$

Assuming independence between the different rounds, propagation of standard deviations can be computed using standard techniques.

<sup>4</sup>NetSquid is an acronym for Network Simulator for Quantum Information using Discrete events.



**Figure 10: Comparison of simulation results with data from NV hardware from [53] (LAB scenario), showing good agreement. (a) Probability of success that the two nodes' measurements in basis  $X/Y/Z$  on the state after a one-sided  $Z$ -rotation are unequal, at  $\alpha = 0.1$  and (b) fidelities, both computed from correlations in the measurement outcomes (see text). (c) Probability that a single generation attempt succeeds, which is computed as  $1/\bar{N}$  where  $\bar{N}$  is the average number of runs up to and including successful heralding of entanglement. Solid line is the theoretical model from [53]. Error bars indicate 1 standard deviation. The simulation data was extracted by running our model implemented in NetSquid on the supercomputer *Cartesius* at SURFsara[2] for 122 hours of wall clock time using 63 cores. A single data point is the average over at least (a) 100 pairs, (b) 300 pairs and (c) 600 pairs, which took between 500k (for  $\alpha = 0.5$ ) and 10.000k (for  $\alpha = 0.03$ ) entanglement generation attempts, with elapsed simulated time between 5 and 117 seconds.**

## C.2 Simulation data

In this section we present further results from the simulations of our proposed link layer protocol. Simulation data will be made available upon request. In total 1618 simulation runs were performed:  $2 \times 169$  long runs (120 h wall time each) with 169 scenarios and 1280 shorter runs (24 h wall time each) with varying request load and minimal requested fidelity. Out of the 169 scenarios used in the long runs,  $2 \times 5 \times 3$  concerned scenarios where the entanglement generation requests where a mix of the priorities  $NL$ ,  $CK$  and  $MD$ . For these mixed scenarios we considered (1) two physical setups, LAB and QL2020, (2) five usage patterns (described below) and (3) three different schedulers, FCFS, LOWERWFQ and HIGHERWFQ.

We implement different usage patterns of the link layer by, in every MHP cycle, issuing a new CREATE request for a random number of pairs  $k$  (max  $k_{\max}$ ) with probability  $f \cdot p_{\text{succ}} / (E \cdot k)$ , where  $p_{\text{succ}}$  is the probability of an attempt being successful,  $f$  is a fraction determining load of our system and  $E$  is the expected number of MHP cycles to make one attempt. For LAB(QL2020)  $E = 1$  ( $E = 1$ ) for M requests and  $E \approx 1.1$  ( $E \approx 16$ ) for K requests. We consider five different use patterns with  $f$  and  $k_{\max}$  defined in table 2.

**Table 2**

| Usage pattern | $NL$                               | $CK$                               | $MD$                                 |
|---------------|------------------------------------|------------------------------------|--------------------------------------|
| UNIFORM       | $f = 0.99 \cdot 1/3, k_{\max} = 1$ | $f = 0.99 \cdot 1/3, k_{\max} = 1$ | $f = 0.99 \cdot 1/3, k_{\max} = 1$   |
| MORENL        | $f = 0.99 \cdot 4/6, k_{\max} = 3$ | $f = 0.99 \cdot 1/6, k_{\max} = 3$ | $f = 0.99 \cdot 1/6, k_{\max} = 256$ |
| MORECK        | $f = 0.99 \cdot 1/6, k_{\max} = 3$ | $f = 0.99 \cdot 4/6, k_{\max} = 3$ | $f = 0.99 \cdot 1/6, k_{\max} = 256$ |
| MOREMD        | $f = 0.99 \cdot 1/6, k_{\max} = 3$ | $f = 0.99 \cdot 1/6, k_{\max} = 3$ | $f = 0.99 \cdot 4/6, k_{\max} = 256$ |
| NONLMOECK     | $f = 0, k_{\max} = 3$              | $f = 0.99 \cdot 4/5, k_{\max} = 3$ | $f = 0.99 \cdot 1/5, k_{\max} = 256$ |
| NONLMOEMD     | $f = 0, k_{\max} = 3$              | $f = 0.99 \cdot 1/5, k_{\max} = 3$ | $f = 0.99 \cdot 4/5, k_{\max} = 256$ |

We make use of the following three scheduling strategies:

- FCFS: First-come-first-serve with a single queue.
- LOWERWFQ: NL are always service first (strict priority) and a weighted fair queue (WFQ) is used between CK (weight 2) and MD (weight 1).
- HIGHERWFQ: NL are always service first (strict priority) and a weighted fair queue (WFQ) is used between CK (weight 10) and MD (weight 1).

Figures 11-16 show scaled latencies and request latencies as functions of simulated time for 20 of the first long simulation runs using the different scenarios of mixed request types. Furthermore, Figures 17-22 show the throughput as a function of simulated time for the same runs. We also ran a second batch with the same scenarios which produced similar plots. When using FCFS the request latency for the different requests are highly correlated, which is to be expected since all requests are in the same queue. On the other hand the scaled latency for the different request priorities, in particular MD, deviates from the others, which is due to the varying number of number of pairs in the requests. From Figures 17-22 one can observe that the type of scheduler, at least for these simulated scenarios, have a relative small affect on the throughput. In table 3 and 4 the average throughputs, scaled latency and request latencies are collected for the same 20 simulation runs.

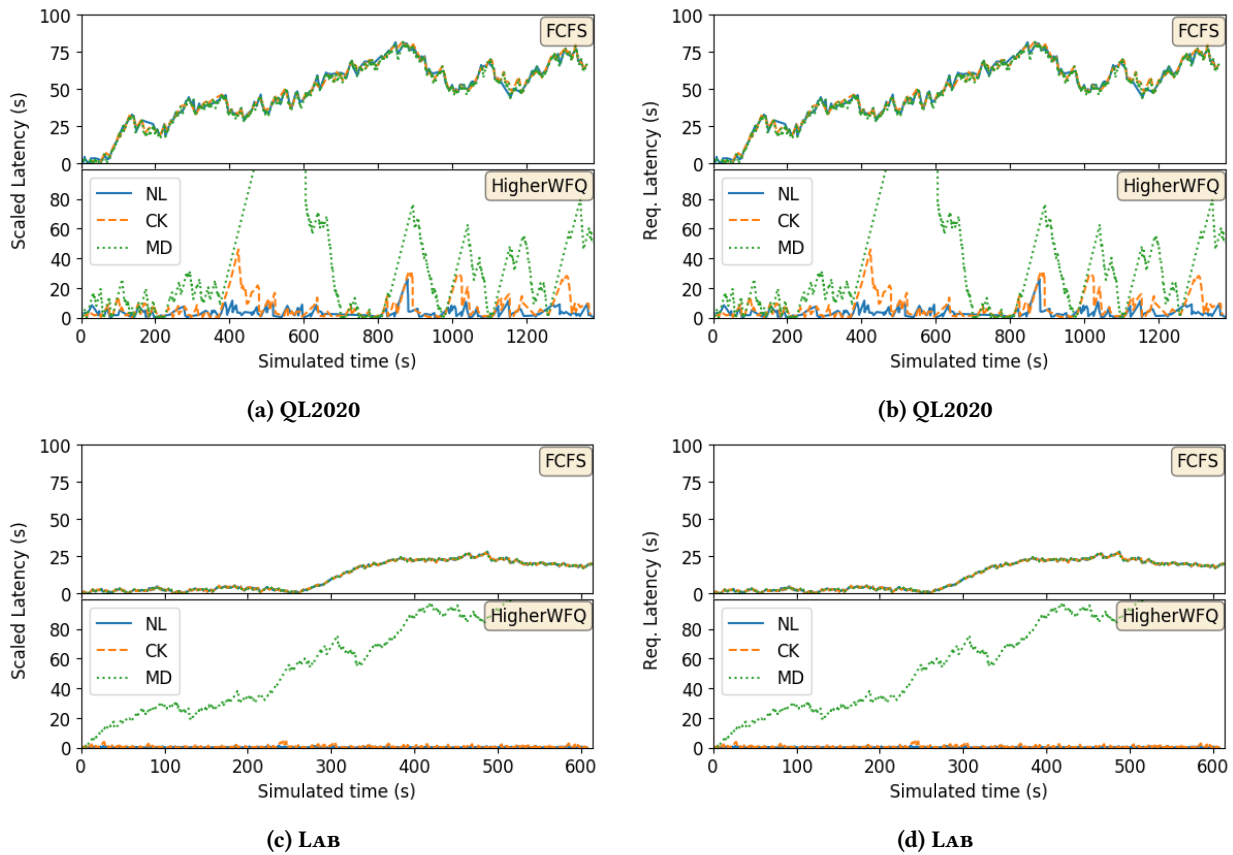


Figure 11: Latencies for UNIFORM



# A Link Layer Protocol for Quantum Networks

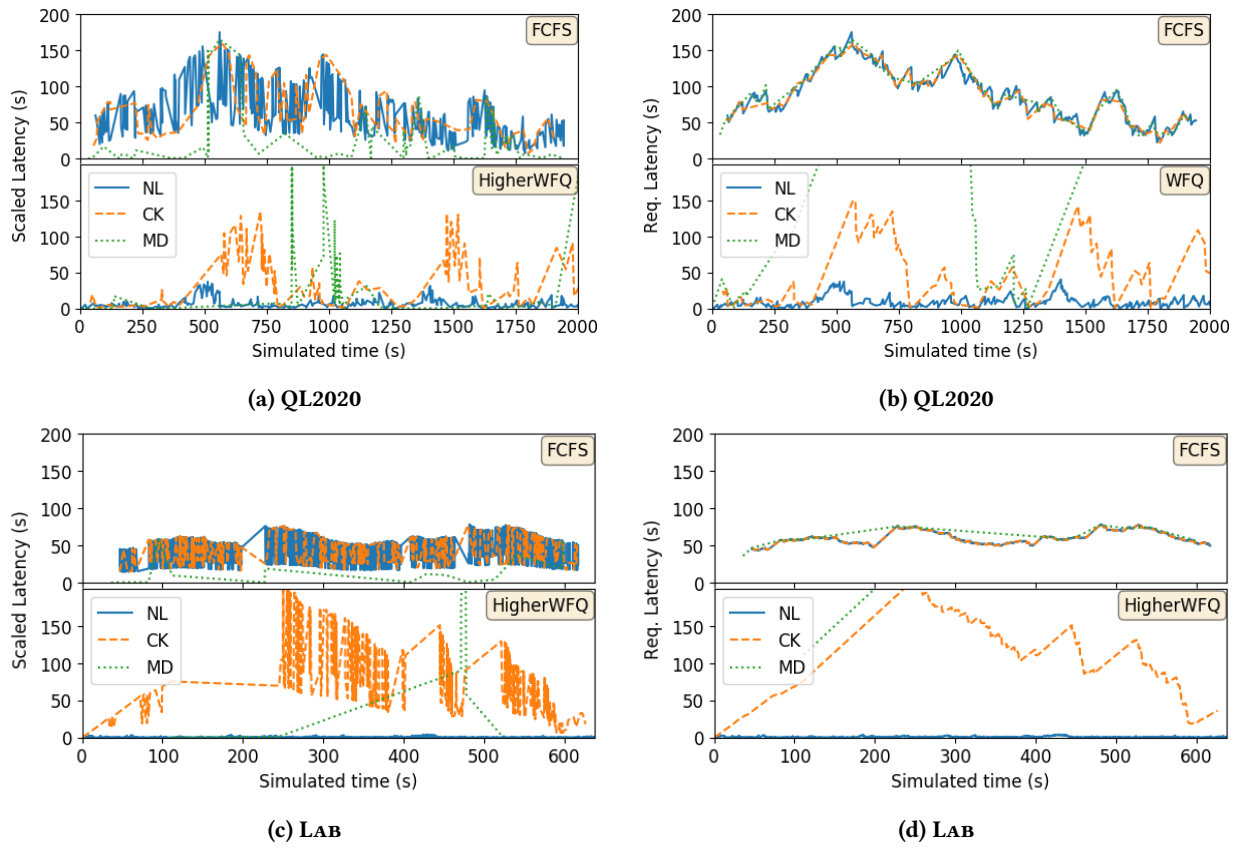


Figure 12: Latencies for MORENL

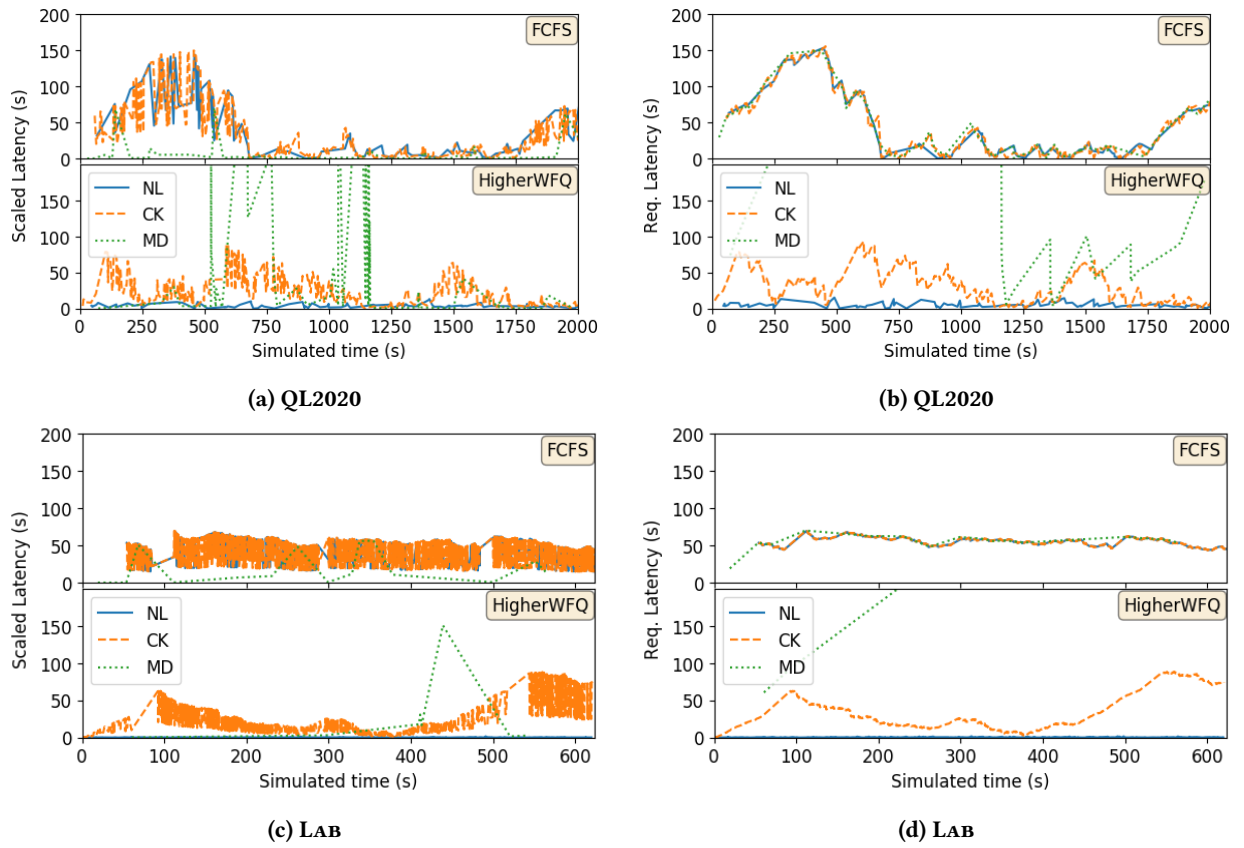


Figure 13: Latencies for MORECK

# A Link Layer Protocol for Quantum Networks

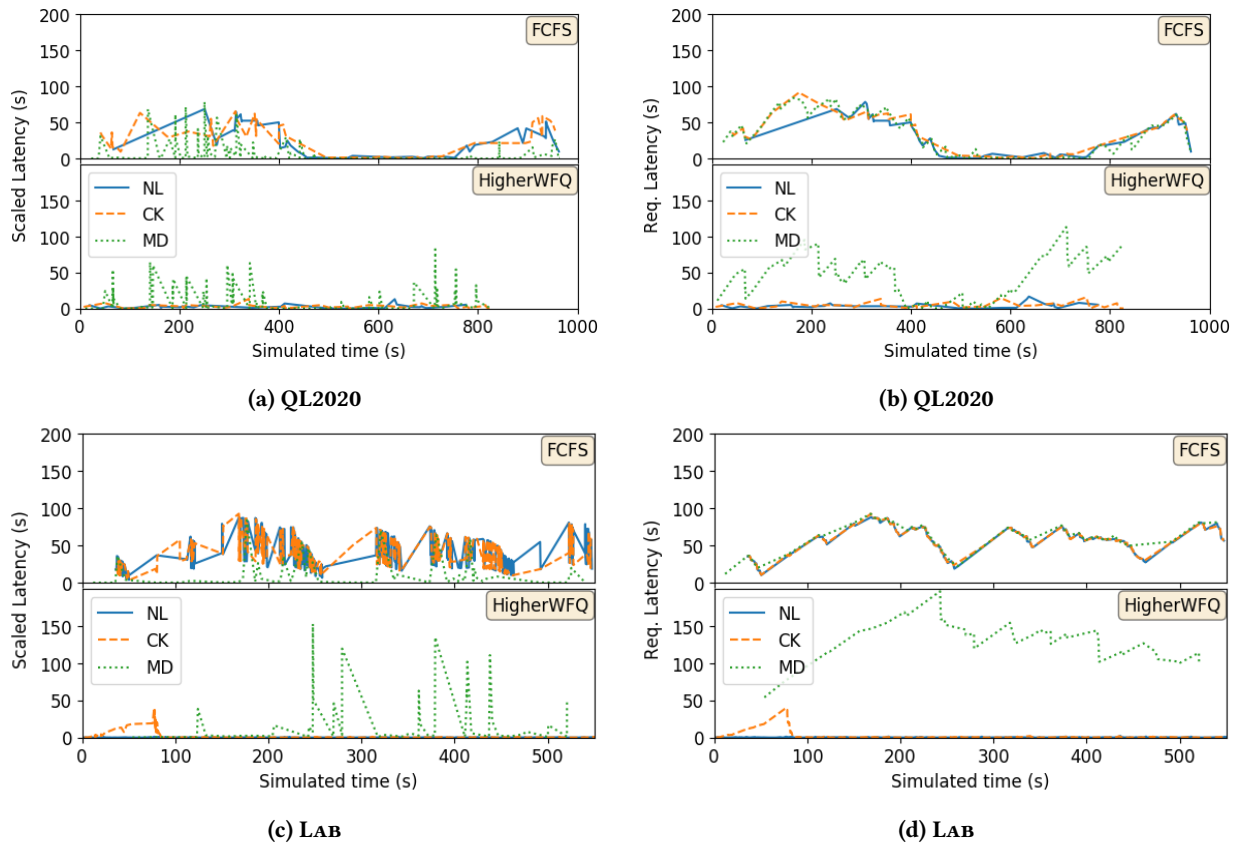


Figure 14: Latencies for MOREMD

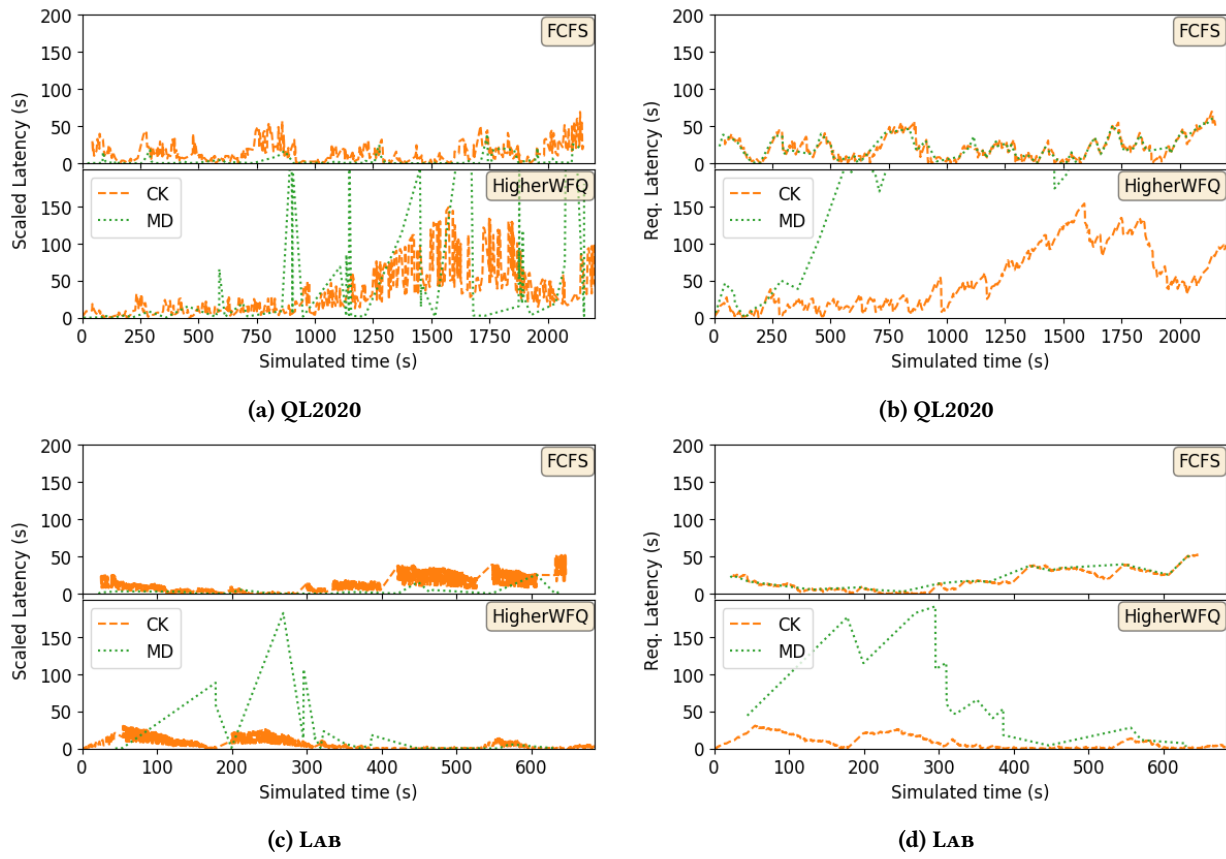


Figure 15: Latencies for NonLMORECK

## A Link Layer Protocol for Quantum Networks

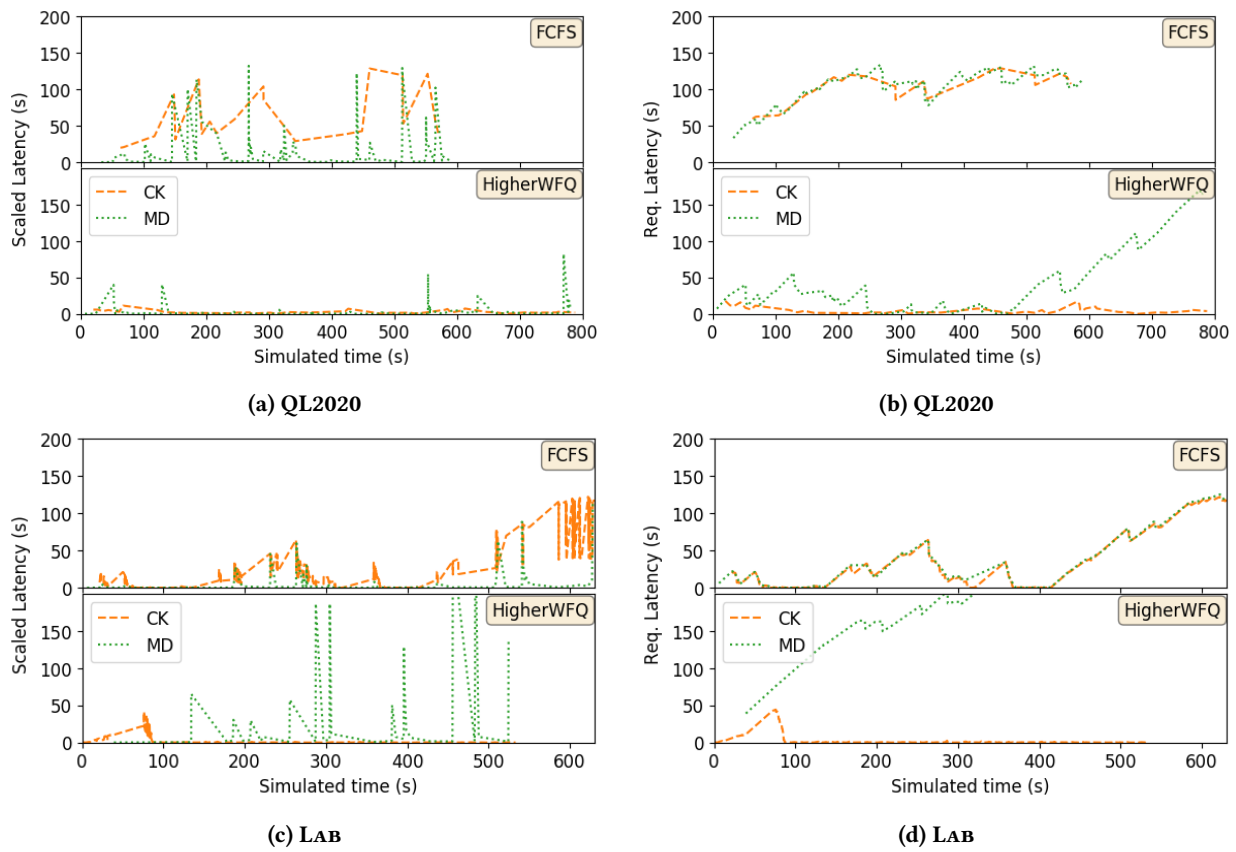


Figure 16: Latencies for NonLMOREMD

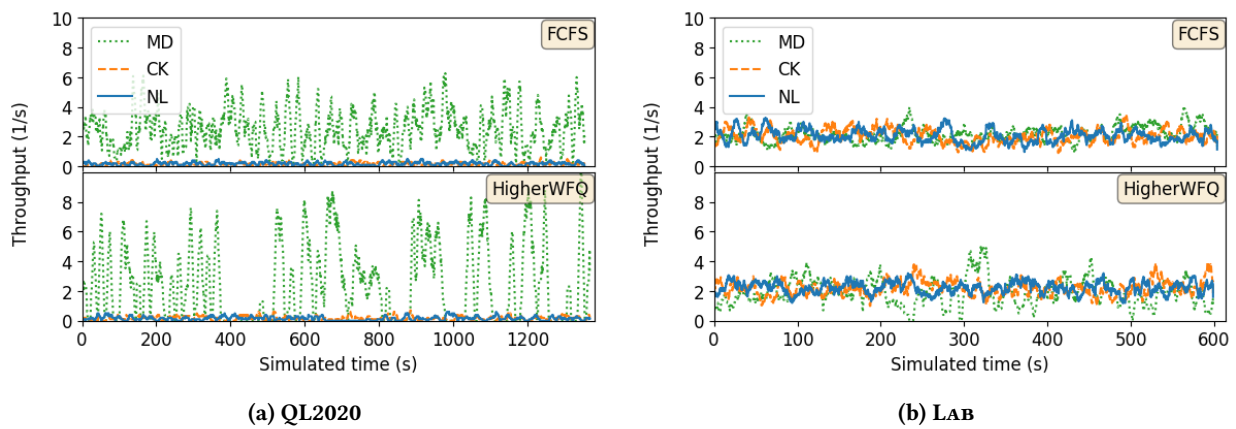
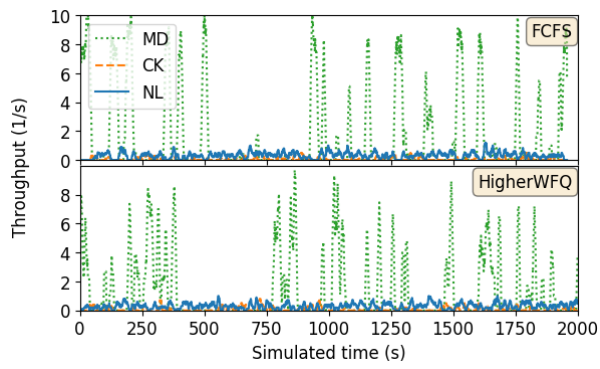
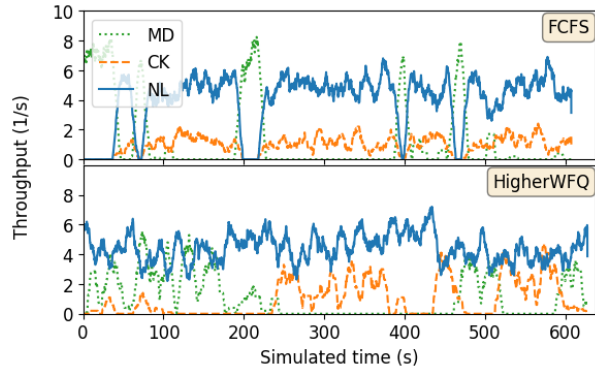


Figure 17: Throughputs for UNIFORM

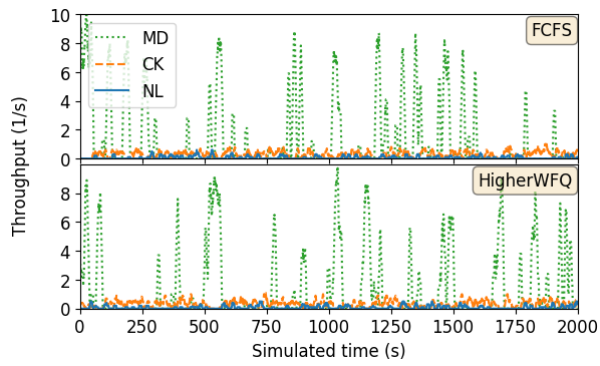


(a) QL2020

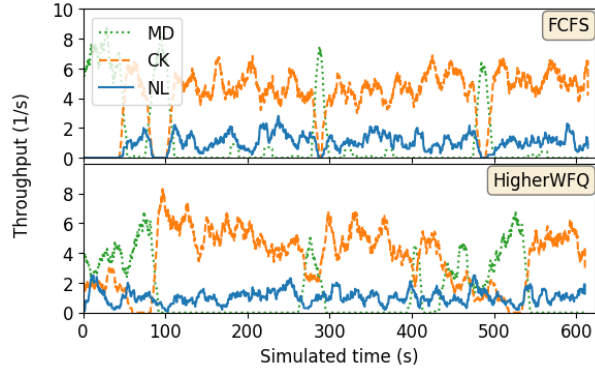


(b) LAB

Figure 18: Throughputs for MORENL

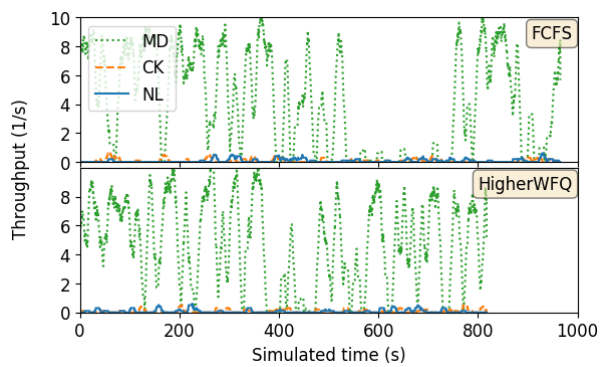


(a) QL2020

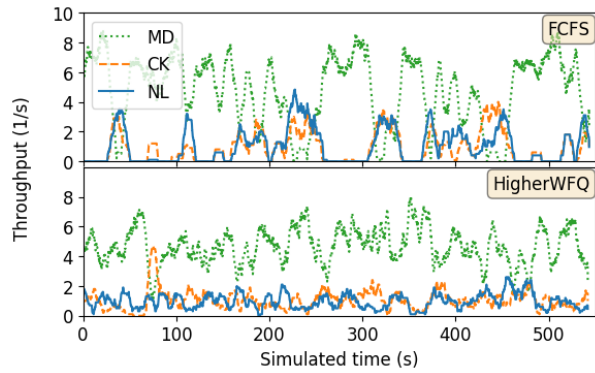


(b) LAB

Figure 19: Throughputs for MORECK



(a) QL2020



(b) LAB

Figure 20: Throughputs for MOREMD

# A Link Layer Protocol for Quantum Networks

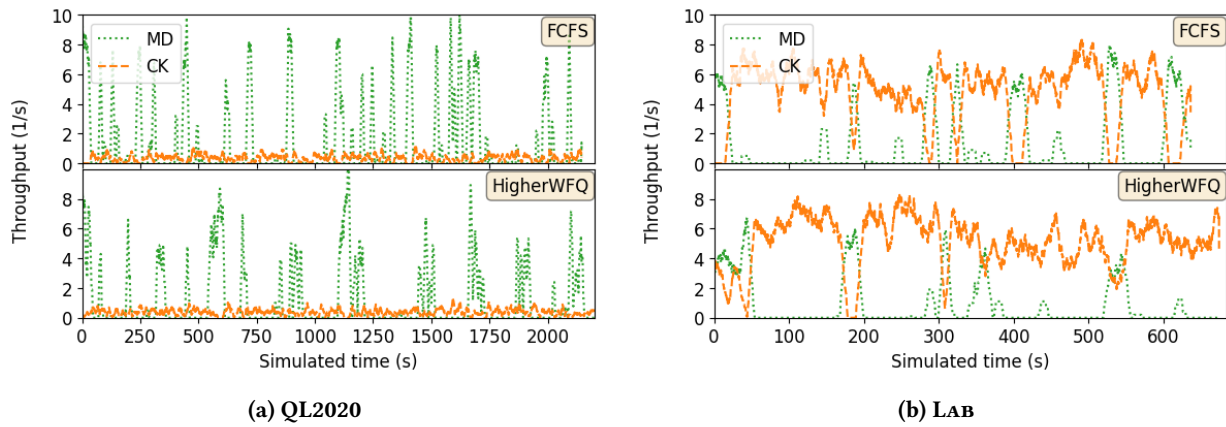


Figure 21: Throughputs for NonLMORECK

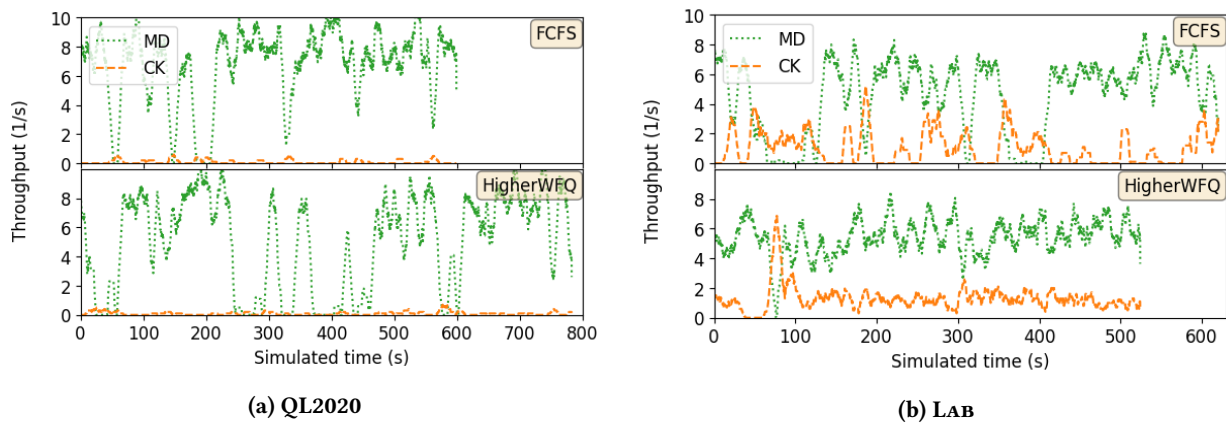


Figure 22: Throughputs for NonLMOREMD



**Table 3: Average throughputs for 20 simulation scenarios, as defined above, of requests with mixed priorities. Each value is computed as  $(\#pairs/t_{sim\_max})$  where  $t_{sim\_max}$  is the reached simulated time (93 - 2355 s) of a single simulation run (24 h wall time).**

| Scenario                    | Throughput_NL (1/s) | Throughput_CK (1/s) | Throughput_MD (1/s) |
|-----------------------------|---------------------|---------------------|---------------------|
| LAB_MORECK_FCFS             | 0.976               | 4.126               | 1.187               |
| LAB_MORECK_HIGHERWFQ        | 1.046               | 3.719               | 1.408               |
| LAB_MOREMD_FCFS             | 1.025               | 0.905               | 4.771               |
| LAB_MOREMD_HIGHERWFQ        | 0.981               | 1.058               | 4.659               |
| LAB_MORENL_FCFS             | 3.975               | 0.950               | 1.375               |
| LAB_MORENL_HIGHERWFQ        | 4.447               | 0.986               | 1.117               |
| LAB_NoNLMORECK_FCFS         | -                   | 4.696               | 1.366               |
| LAB_NoNLMORECK_HIGHERWFQ    | -                   | 5.101               | 0.916               |
| LAB_NoNLMOREMD_FCFS         | -                   | 1.044               | 4.600               |
| LAB_NoNLMOREMD_HIGHERWFQ    | -                   | 1.300               | 5.408               |
| LAB_UNIFORM_FCFS            | 2.066               | 2.035               | 2.170               |
| LAB_UNIFORM_HIGHERWFQ       | 2.210               | 2.186               | 1.908               |
| QL2020_MORECK_FCFS          | 0.064               | 0.302               | 1.398               |
| QL2020_MORECK_HIGHERWFQ     | 0.078               | 0.329               | 1.146               |
| QL2020_MOREMD_FCFS          | 0.075               | 0.078               | 4.139               |
| QL2020_MOREMD_HIGHERWFQ     | 0.066               | 0.073               | 4.793               |
| QL2020_MORENL_FCFS          | 0.312               | 0.066               | 1.667               |
| QL2020_MORENL_HIGHERWFQ     | 0.292               | 0.084               | 1.374               |
| QL2020_NoNLMORECK_FCFS      | -                   | 0.355               | 1.480               |
| QL2020_NoNLMORECK_HIGHERWFQ | -                   | 0.374               | 1.180               |
| QL2020_NoNLMOREMD_FCFS      | -                   | 0.084               | 6.756               |
| QL2020_NoNLMOREMD_HIGHERWFQ | -                   | 0.091               | 5.036               |
| QL2020_UNIFORM_FCFS         | 0.175               | 0.143               | 2.538               |
| QL2020_UNIFORM_HIGHERWFQ    | 0.154               | 0.166               | 2.483               |

**Table 4: Average scaled latencies (SL) and request latencies (RL) for 20 simulation scenarios, as defined above, of requests with mixed priorities of a single simulation run (93 - 2355 s simulated time and 24 h wall time). Values in parenthesis are estimates of standard errors, computed as  $s_n/\sqrt{n}$  where  $s_n$  is the sample standard deviation and  $n$  is the number of data points used for averaging.**

| Scenario                    | SL_NL (s)    | SL_CK (s)    | SL_MD (s)     | RL_NL (s)    | RL_CK (s)     | RL_MD (s)      |
|-----------------------------|--------------|--------------|---------------|--------------|---------------|----------------|
| LAB_MORECK_FCFS             | 40.18 (0.90) | 41.09 (0.42) | 19.64 (5.52)  | 55.50 (0.30) | 55.58 (0.14)  | 55.11 (2.83)   |
| LAB_MORECK_HIGHERWFQ        | 0.30 (0.01)  | 25.63 (0.61) | 24.95 (15.20) | 0.46 (0.02)  | 33.88 (0.66)  | 264.28 (35.75) |
| LAB_MOREMD_FCFS             | 40.25 (1.15) | 41.70 (1.28) | 13.87 (2.91)  | 55.63 (1.00) | 57.24 (1.00)  | 62.55 (2.43)   |
| LAB_MOREMD_HIGHERWFQ        | 0.23 (0.01)  | 2.09 (0.29)  | 27.53 (6.00)  | 0.36 (0.02)  | 2.65 (0.35)   | 129.30 (3.51)  |
| LAB_MORENL_FCFS             | 45.59 (0.47) | 46.44 (0.95) | 14.70 (4.65)  | 60.36 (0.21) | 60.59 (0.42)  | 61.80 (2.76)   |
| LAB_MORENL_HIGHERWFQ        | 0.69 (0.02)  | 83.79 (2.64) | 98.34 (46.15) | 0.97 (0.02)  | 114.89 (2.57) | 299.28 (37.25) |
| LAB_NoNLMORECK_FCFS         | -            | 13.05 (0.27) | 3.55 (1.36)   | -            | 17.58 (0.30)  | 23.04 (3.21)   |
| LAB_NoNLMORECK_HIGHERWFQ    | -            | 6.70 (0.16)  | 26.14 (9.05)  | -            | 9.04 (0.19)   | 76.02 (12.51)  |
| LAB_NoNLMOREMD_FCFS         | -            | 23.45 (1.65) | 10.97 (2.51)  | -            | 30.86 (1.94)  | 39.33 (4.09)   |
| LAB_NoNLMOREMD_HIGHERWFQ    | -            | 2.26 (0.31)  | 44.33 (11.92) | -            | 3.34 (0.43)   | 204.78 (9.54)  |
| LAB_UNIFORM_FCFS            | 11.41 (0.27) | 11.46 (0.27) | 12.38 (0.26)  | 11.41 (0.27) | 11.46 (0.27)  | 12.38 (0.26)   |
| LAB_UNIFORM_HIGHERWFQ       | 0.35 (0.01)  | 0.73 (0.02)  | 61.19 (0.97)  | 0.35 (0.01)  | 0.73 (0.02)   | 61.19 (0.97)   |
| QL2020_MORECK_FCFS          | 40.65 (4.43) | 37.46 (1.93) | 12.22 (3.82)  | 52.20 (5.21) | 51.34 (2.33)  | 43.41 (5.79)   |
| QL2020_MORECK_HIGHERWFQ     | 4.11 (0.24)  | 26.66 (0.95) | 76.29 (16.50) | 5.72 (0.34)  | 35.91 (1.05)  | 238.09 (21.00) |
| QL2020_MOREMD_FCFS          | 25.45 (3.03) | 28.34 (3.01) | 9.30 (1.63)   | 32.94 (3.43) | 38.90 (3.65)  | 37.97 (2.67)   |
| QL2020_MOREMD_HIGHERWFQ     | 2.62 (0.42)  | 3.04 (0.42)  | 14.44 (1.92)  | 3.79 (0.60)  | 4.51 (0.62)   | 47.64 (2.75)   |
| QL2020_MORENL_FCFS          | 65.92 (1.93) | 64.05 (4.18) | 26.16 (5.84)  | 89.83 (1.79) | 85.99 (3.70)  | 85.98 (5.24)   |
| QL2020_MORENL_HIGHERWFQ     | 7.04 (0.35)  | 45.87 (3.73) | 50.55 (12.63) | 9.78 (0.41)  | 59.92 (4.15)  | 236.03 (21.62) |
| QL2020_NoNLMORECK_FCFS      | -            | 15.98 (0.64) | 4.63 (0.87)   | -            | 21.90 (0.72)  | 21.27 (1.63)   |
| QL2020_NoNLMORECK_HIGHERWFQ | -            | 39.31 (1.64) | 70.64 (16.03) | -            | 55.94 (1.93)  | 277.08 (19.97) |
| QL2020_NoNLMOREMD_FCFS      | -            | 60.12 (6.95) | 22.28 (4.04)  | -            | 102.23 (4.17) | 104.88 (2.69)  |
| QL2020_NoNLMOREMD_HIGHERWFQ | -            | 3.01 (0.35)  | 6.21 (1.58)   | -            | 5.18 (0.70)   | 38.91 (4.89)   |
| QL2020_UNIFORM_FCFS         | 49.13 (1.29) | 50.85 (1.38) | 47.39 (0.33)  | 49.13 (1.29) | 50.85 (1.38)  | 47.39 (0.33)   |
| QL2020_UNIFORM_HIGHERWFQ    | 4.16 (0.26)  | 8.22 (0.58)  | 34.39 (0.61)  | 4.16 (0.26)  | 8.22 (0.58)   | 34.39 (0.61)   |

**Table 5: Relative difference (Rel. Diff.) for the metrics: fidelity (Fid.), throughput (Throughp.) scaled latency (Laten.) and number of generated entangled pairs (Nr. pairs), between two identical scenarios except that the probability of losing a classical message ( $p_{\text{loss}}$ ) is zero for one and between  $10^{-10}$  and  $10^{-4}$  for the other. The relative difference is maximized over three simulation runs (281 - 3973 s simulated time and 70 h wall time each) with requests of priority either *NL*, *CK* or *MD* ( $f = 0.99$ ,  $k_{\text{max}} = 3$ ), with equal  $p_{\text{loss}}$ .**

| $p_{\text{loss}}$ | Max Rel. Diff. Fid. | Max Rel. Diff. Throughp. | Max Rel. Diff. Laten. | Max Rel. Diff. Nr pairs |
|-------------------|---------------------|--------------------------|-----------------------|-------------------------|
| $10^{-4}$         | 0.005               | 0.027                    | 0.629                 | 0.026                   |
| $10^{-5}$         | 0.004               | 0.012                    | 0.469                 | 0.008                   |
| $10^{-6}$         | 0.016               | 0.037                    | 0.332                 | 0.047                   |
| $10^{-7}$         | 0.040               | 0.026                    | 0.576                 | 0.020                   |
| $10^{-8}$         | 0.007               | 0.023                    | 0.623                 | 0.020                   |
| $10^{-9}$         | 0.004               | 0.026                    | 0.338                 | 0.021                   |
| $10^{-10}$        | 0.018               | 0.075                    | 0.742                 | 0.077                   |

## D UNDER THE HOOD

We now provide some more details on the simulation, numerical methods and the underlying NV platform. We remark that physical models for different parts of the NV platform are well established and validated [12]. The validation performed here is thus only about how the combined system performing entanglement generation matched with our simulation.

### D.1 The simulated network

To understand our simulation we perform a full implementation of the MHP and EGP, running on simulated quantum hardware. To achieve this, we start by defining basic components in NetSquid, which, inspired by NS-3, is entirely modular and can be used to construct complex simulation scenarios by combining component elements. The components in the simulation are as follows, and our simulation could easily be configured to examine the performance of our protocol on other underlying hardware platforms such as Ion Traps.

- A `QuantumProcessingDevice`, which is a general component we defined in NetSquid. Abstractly, such a `QuantumProcessingDevice` is described by the following:
  - (1) A number of communication and memory qubits. Each such qubit is associated with a noise-model that describes how quantum information decoheres over time when kept in the memory itself. Concrete parameters for the NV platform are given in Section D.2.
  - (2) Possible one or two-qubit quantum gates to be performed on each (pair of) qubit(s), the time required to execute the gate, as well as a noise-model associated with each such gate that may differ from qubit to qubit. For the NV platform we will only need to consider the gates given in Section D.2.2.
  - (3) With each communication qubit we associate a trigger that allows the generation of entanglement between this communication qubit and a traveling qubit (a photon). Such an operation only succeeds with some probability of success, requires a certain amount of time, and can also be noisy. For the NV platform, we explain this in Section D.4.
  - (4) Readout - i.e. measurement of a qubit. This takes a certain amount of time, and does again carry a noise-model. For NV, we explain this in Section D.3.4.
- A `FiberConnection`, which is a general NetSquid component that allows us to simulate optical fibers, including photon loss per km.
- A `HeraldingStation`, which automatically measures incoming photons in a certain time interval. This process is again subject to several possible errors explained in Section D.5.
- `ClassicalConnection`, which we use to transmit classical messages allowing us to simulate, for example, losses on the channel. Section D.6.1 explains the model considered here.
- `Node`, which includes a `QuantumProcessingDevice`, and enjoys fiber connections with the `HeraldingStation` or other nodes. Each `Node` can run programs in the same way that they could be run on e.g. a computer or microcontroller, allowing these Programs to make use of - for example - the `QuantumProcessingDevice`. This allows us to perform a full implementation (here in Python) of the MHP and EGP including all subsystems in the same way as the program will later run on an actual microcontroller (however, in C).

We briefly review properties of the nitrogen-vacancy (NV) centre in diamond [12] and how they affect the design and performance of our protocol. We will also highlight how this can be modeled in simulation. Important to the design and performance of our protocol is how long operations on qubits stored in the NV-center take. Additionally, the coherence time, i.e. how long a qubit can be stored, has an impact on our protocol. We summarize typical values for the noise-level and execution time of the allowed operations of a NV-center together with coherence times. These are the values we used in our simulation. Note however that these values can vary significantly between samples.

### D.2 Qubits on the NV Platform

A NV centre is formed by replacing a carbon atom in a diamond lattice with a nitrogen atom and removing a neighboring carbon (vacancy). This structure traps electrons which together form a spin-1 system. Two of the levels of the collective spin-1 state can be used as a communication qubit in a quantum network. Around the NV centre there is also a natural abundance of carbon-13 atoms which interact with the communication qubit (electron spin). The surrounding carbon spins can be addressed using the communication qubit and can thus be used as memory qubits. We here consider a situation with only one communication qubit, and one memory qubit.

|  | (Unsqured) fidelity                         | Duration/time | Experimentally realized   |
|--|---|---------------|---|
| Electron $T_1$                               | -   | 2.86 ms       | > 1h[3]   |
| Electron $T_2^*$                             | -   | 1.00 ms       | 1.46 s[3]   |
| Carbon $T_1$                                 | -   | $\infty$      | > 6m [21]   |
| Carbon $T_2^*$                               | -   | 3.5 ms        | $\approx$ 10ms [21]   |
| Electron single-qubit gate                   | 1.0   | 5 ns          | > 0.995 (100 ns) [60]   |
| E-C controlled- $\sqrt{X}$ -gate (E=control) | 0.992                                       | 500 $\mu$ s   | 0.992 (500-1000 $\mu$ s) fig 2 in [60]                          |
| Carbon Rot-Z-gate                            | 0.999                                       | 20 $\mu$ s    | 1.0 (20 $\mu$ s) [93]   |
| Electron initialization in $ 0\rangle$       | 0.95  | 2 $\mu$ s     | 0.99 (2 $\mu$ s) [82]   |
| Carbon initialization in $ 0\rangle$         | 0.95  | 310 $\mu$ s   | 0.95 (300 $\mu$ s) [32]   |
| Electron readout                             | 0.95 ( $ 0\rangle$ ), 0.995 ( $ 1\rangle$ ) | 3.7 $\mu$ s   | 0.95 ( $ 0\rangle$ ), 0.995 ( $ 1\rangle$ ) (3-10 $\mu$ s) [53] |

**Table 6: Gates and coherence times used in simulation. Values used in the simulation corresponding to LAB. We remark that since these are custom chips, no two are exactly identical. Individual values have since seen significant improvements (Experimentally realized), but not been realized simultaneously for producing entanglement that would allow a direct comparison to simulation. We have thus focused in simulation only what enables a comparison to data gathered from entanglement generation on hardware.**

*D.2.1 Noise model - Free evolution of the electronic and nuclear spins.* Noise in experimental implementations is described in terms of  $T_1$ ,  $T_2$ ,  $T_2^*$  times, where Section A.4 serves to provide intuition on how our quantity of interest - the fidelity to the maximally entangled target state  $|\Psi^+\rangle$  - depends on their values. Table 6 lists values used in simulation (reflecting LAB), and state of the art for the communication qubit (Electron), and memory qubit (Carbon).

### D.2.2 Quantum gates.

*Procedure and parallelism constraints.* Quantum gates can be realized by applying microwave pulses. Of specific interest that affects the throughput is the duration of such operations given in Table 6. While not absolutely necessary for the understanding of the simulation, we briefly sketch how operations are performed also on the carbons to illustrate one feature of this system that is relevant for the performance of our protocols - namely the parallelism of the allowed gate operations. We remark that operations on the carbon spins are performed using the following pulse sequence

$$(\tau - \pi - 2\tau - \pi - \tau)^{N/2}, \quad (21)$$

where  $\pi$  is a microwave- $\pi$ -pulse on the electron spin,  $2\tau$  is the time between the pulses and  $N$  is the total number of pulses. The target carbon spin can be chosen by picking  $\tau$  such that it is precisely in resonance with the target carbon spin's hyperfine interaction with the electron spin. If the electron spin is in the state  $|0\rangle$  ( $|1\rangle$ ) the target carbon spin will rotate around the  $X$ -axis of the Bloch sphere in the positive (negative) direction, with an angle  $\theta$  which depends on the total number of pulses  $N$ . This means that one can perform quantum gates on the carbon that are controlled by the state of the electron spin. The effective unitary operation (E=control) on the electron and the target carbon spin is then given as

$$\begin{pmatrix} R_X(\theta) & 0 \\ 0 & R_X(-\theta) \end{pmatrix}, \quad (22)$$

where  $R_X(\theta) = \exp(i\theta/2X)$  denotes a rotation around the  $X$ -axis of an angle  $\theta$ . Not only does the pulse sequence (21) manipulate the carbon spin, but it also decouples the electron from its environment, thereby prolonging its coherence time and is thus also called *dynamical decoupling*, allowing longer memory lifetimes (Figure 9a). We thus see a limit to the amount of parallelism when operating on the carbon and the electron spin.

Other quantum gates are however simpler: Since the carbon continuously precess around the  $Z$  axis of the Bloch sphere, rotations around  $Z$  (Carbon Rot-Z) are simply done by waiting a correct amount of time. Thus, also controlled rotations around other axes than  $X$  can be performed by correctly interspersed waiting times during the pulse sequence above.

### D.3 Gates and their noise

In this section we collect parameters for noise and delays of gates used in our simulation. Table 6 summarize the possible gates that can be performed on the electron and carbon spins in the NV system, together with decoherence times. Section D.4 describes how the noise occurring from entanglement generation attempts is modeled.

Here the E-C controlled- $X$  rotates the carbon spin around the  $X$ -axis in the positive (negative) direction if the electron is in the  $|0\rangle$  ( $|1\rangle$ ) state. Furthermore, note that there is an asymmetry in reading out the  $|0\rangle$ -state and the  $|1\rangle$ -state of the electron.

*D.3.1 Modeling noisy operations.* Noise in gate operations is modeled by applying noise after a perfect gate (a standard method):

$$U_{\text{noisy}}(\rho) = \mathcal{N}_{\text{dephas}}^f \circ U_{\text{perfect}}(\rho),$$

where

$$\mathcal{N}_{\text{dephas}}^p : \rho \mapsto f\rho + (1-f)\mathcal{Z}\rho\mathcal{Z}$$

is the dephasing channel in  $\mathcal{Z}$  and  $f$  is the gate fidelity as given in table 6. States are initialized as  $\mathcal{N}_{\text{depol}}^p(|0\rangle\langle 0|)$ , where

$$\mathcal{N}_{\text{depol}}^p : \rho \mapsto f\rho + \frac{1-f}{3} [X\rho X + Y\rho Y + Z\rho Z]$$

denotes the depolarization channel by.

*D.3.2 How the electron spin is initialized.* Initialization of the electron spin means setting the state to  $|0\rangle$  [82]. Initialization of the electron spin is done by performing *optical pumping*, in which light shines onto the electron, thereby bringing its quantum state in a higher energy level, given it was in  $|1\rangle$ , after which it falls back to either  $|0\rangle$  or  $|1\rangle$  with a given probability. If the electron falls down to the state  $|0\rangle$  it will stay there, thus after many repetitions of this process, the electron is with high confidence in the state  $|0\rangle$ . For our discussion here, it will be relevant to remark that this operation takes time (Table 6), and we will need to perform it repeatedly as the first step in producing entanglement.

*D.3.3 Moving a qubit to memory.* When moving a qubit from the communication qubit to the memory qubit, the memory qubit needs to already be initialized. Initialization of the carbon is done by effectively swapping the  $|0\rangle$  state from the electron to the carbon and cannot therefore be performed while having an entangled state in the electron. For this reason, initialization of the carbon (310  $\mu\text{s}$ ) needs to be performed before a photon is emitted from the electron during an entanglement generation attempt. However, it is not necessary to re-initialize the carbon before every entanglement generation attempt but simply periodically depending on the coherence time. In our simulation we assumed  $T_1$  to be 3500  $\mu\text{s}$  and thus re-initialize the carbon every 3500  $\mu\text{s}$  (every 350th MHP cycle).

Swapping a state in the electron to the carbon can be done by 2 E-C controlled- $\sqrt{X}$ -gates and single qubit gates (total time 1040  $\mu\text{s}$ ) [60].

*D.3.4 How a measurement (readout) is performed.*

*Readout of the communication qubit.* First, we are again interested in the time to perform this operation given in Table 6 which will be relevant in the MD use case. Evidently, also a readout can be noisy, where we here remark that the noise is asymmetric in that the probability of incorrectly obtaining measurement outcome 0 is much lower than incorrectly getting outcome 1.

*Reading out a memory qubit.* We again remark that next to timing constraints (Table 6), we have limited parallelism on the current NV platform, since we need the electron spin to readout the memory qubit. Reading out the nuclear spin is done by performing the following steps:

- (1) initialize the electron spin,
- (2) apply an effective controlled NOT operation with the nuclear spin as control (consisting of one E-C controlled- $\sqrt{X}$ -gate and single-qubit gates),
- (3) measure (readout) the electron spin.

The reason why a controlled NOT is sufficient, rather than a full swap, is the following: If the nuclear spin is in state  $\alpha|0\rangle + \beta|1\rangle$ , then after the CNOT, the combined state is  $\alpha|00\rangle_{EC} + \beta|11\rangle_{EC}$ . The reduced state [76] of the electron is then  $|\alpha|^2|0\rangle\langle 0| + |\beta|^2|1\rangle\langle 1|$ , so measuring in the standard basis yields the same statistics as measuring  $\alpha|0\rangle + \beta|1\rangle$  in the same basis.

*Readout noise.* Readout is modeled by performing a POVM measurement with the following Kraus operators (see [76] for definition)

$$M_0 = \begin{pmatrix} \sqrt{f_0} & 0 \\ 0 & \sqrt{1-f_1} \end{pmatrix}, \quad M_1 = \begin{pmatrix} \sqrt{1-f_0} & 0 \\ 0 & \sqrt{f_1} \end{pmatrix} \quad (23)$$

where  $f_0$  ( $f_1$ ) is the readout fidelity of the  $|0\rangle$ -state ( $|1\rangle$ ) as given in table 6

#### D.4 Physical Entanglement Generation and Noise

We here consider the single-click scheme of LAB. To understand timing and quality, as well as parameter choices let us give a high-level overview of the single-click scheme: A microwave pulse is used to prepare the communication qubit in the state  $\sqrt{\alpha}|0\rangle + \sqrt{1-\alpha}|1\rangle$  (max.  $5.5\mu\text{s}$  for  $A$  and  $B$ ), where  $|0\rangle$  is also called the bright state, and  $\alpha$  the bright state population. A resonant laser pulse is then used to cause emission of a photon, if the state was in the bright state, preparing the joint state of the communication qubit ( $C$ ) and an emitted photon ( $P$ ) in the state  $\sqrt{\alpha}|0\rangle_C|1\rangle_P + \sqrt{1-\alpha}|1\rangle_C|0\rangle_P$ , where  $|0\rangle_P$  ( $|1\rangle$ ) denotes the absence (presence) of a photon. This process succeeds with probability  $p_{\text{em}} \approx 0.03$  without Purcell enhancement using optical cavities. To ensure phase-stabilization only one laser may be used for both nodes, combined with local shutters to allow node control. We remark that local control at the node is still desirable at a distance due to aligning with other operations such as performing gates, to keep qubits stable. The heralding station interferes both incoming photons on a beam splitter, thereby performing a probabilistic entanglement swap. Intuitively, this can be understood as a measurement of the incoming qubits in the Bell basis, where we can only obtain outcomes  $|\Psi^+\rangle$ ,  $|\Psi^-\rangle$  or “other”. Since “other” is more than one possible state, we declare failure in this case.

Depending on the clicks observed in the detectors, we have projected the state of the communication qubits at  $A$  and  $B$  in the state  $|\Psi^+\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A|1\rangle_B + |1\rangle_A|0\rangle_B)$  (only left detector clicks),  $|\Psi^-\rangle = \frac{1}{\sqrt{2}}(|0\rangle_A|1\rangle_B - |1\rangle_A|0\rangle_B)$  (only right detector clicks), or we declare failure (either none or both of the detectors click). Success occurs with probability  $p_{\text{succ}} \approx 2\alpha p_{\text{det}}$ , where  $p_{\text{det}}$  is the probability of detecting an emitted photon.

During entanglement generation, a variety of noise processes occur:

- (1) Dephasing of the nuclear spins (memory qubits) due to resetting the electron during entanglement generation attempts.
- (2) Effective dephasing of the photon state due to uncertainty in the phase between the paths the photon travel to the beam-splitter.
- (3) Effective dephasing of the photons state due to non-zero probability of emitting two photons.
- (4) Effective amplitude damping of the photon state due to coherent emission, i.e. the photon is in a super-position of being emitted at different times.
- (5) Effective amplitude damping due to collection losses from non-unity probability of emitting the photon in the zero-phonon line and non-unity collection efficiency into the fiber.
- (6) Effective amplitude damping due to losses in fiber.
- (7) Non-perfect beam-splitter measurement at the heralding stations due to photons not being perfectly indistinguishable.
- (8) Errors in the classical outcome from the detectors due to non-unity detection efficiency and dark counts.

*D.4.1 Dephasing mechanism on nuclear spins during entanglement attempts.* Between every entanglement attempt, the electron spin (communication qubit) needs to be reset. The dominant source of noise on the nuclear spins (memory qubits) during the entanglement attempts is due to this re-initialization of the electron spin, as described in [58]. We model the noise on the nuclear spins as a fixed amount of dephasing noise

$$\mathcal{D}_{p_d}(\rho_n) = (1-p_d)\rho_n + p_d Z \rho_n Z \quad (24)$$

for every entanglement attempt. The dephasing parameter depends on: the bright state population  $\alpha$ , the coupling strength  $\Delta\omega$  and a decay constant  $\tau_d$  as follows

$$p_d = \frac{\alpha}{2} (1 - \exp(-\Delta\omega^2 \tau_d^2 / 2)), \quad (25)$$

see equation (2) in [58]. If the length of the Bloch vector in the equatorial plane of the state in the nuclear spin is before the entanglement attempts  $r_{XY}$ , then after  $N$  attempts the length will be

$$(1-p_d)^N r_{XY}. \quad (26)$$

The bright state population can be chosen per experiment but the coupling strength and decay constant are fixed but vary between different nuclear spins. The decay constant can also vary by performing different microwave control pulses of the

electron spin. As an example of these parameters, for the nuclear spin  $C_1$  and the standard microwave control pulses, the coupling strength is  $\Delta\omega = 2\pi \times 377$  kHz and the decay constant is  $\tau_d = 82$  ns, see [58]. In the simulations we use these values for the coupling strength and decay constant.

*D.4.2 Phase uncertainty for photon paths.* There is uncertainty in the phase between the paths the photon travels from the nodes to the beam-splitter, due to for example uncontrolled stretching of the fiber. If this phase difference is  $\Delta\phi$  then the state after a successful measurement at the heralding station (conditioned on there being only one photon) is

$$|0\rangle_{e_A} |1\rangle_{e_B} \pm e^{i\Delta\phi} |1\rangle_{e_A} |0\rangle_{e_B}. \quad (27)$$

where  $e_A$  is the electron spin at node  $A$  and  $e_B$  is the electron spin at node  $B$ .

We model this by performing dephasing noise to both qubits encoding the presence/absence of photon from  $A$  and  $B$ . As shown in [88], if we chose the dephasing parameter to be

$$p_d = \left(1 - \frac{I_1(\sigma(\phi)^{-2})}{I_0(\sigma(\phi)^{-2})}\right)/2 \quad (28)$$

then the standard deviation of the phase  $\phi$  in the state between the electron ( $e_A/e_B$ ) and the photon ( $p_A/p_B$ )

$$|0\rangle_{e_A} |1\rangle_{p_A} \pm e^{i\phi} |1\rangle_{e_A} |0\rangle_{p_A} \quad (29)$$

is precisely  $\sigma(\phi)$ . To efficiently compute quotients of modified Bessel functions we implemented the algorithm described in [5]. Note that the variance of the phase in equation (27) is twice the variance of the phase in equation (29). In experiments the standard deviation of the phase for the state between the electrons, i.e. (27), has in [53] shown to be  $14.3^\circ$ . Thus, we set the standard deviation in equation (29) to be  $14.3^\circ/\sqrt{2}$ .

*D.4.3 Two-photon emission.* There is a probability that two photons are emitted from the electron during the entanglement generation attempt at a node. For the physical setup we assume the probability of there being two photons emitted, given that at least one is emitted, to be approximately 4% [53]. As argued in [53], the two-photon event can effectively be modeled as applying dephasing noise to the electron qubit which is part of the generated entanglement.

*D.4.4 Coherent emission, superposition of time-modes.* The emission of the photon is a coherent process and the photon is actually in a super-position of being emitted at different times. As shown in [88], the effect of a finite detection window can be seen as effective amplitude damping noise to the qubit encoding the presence/absence of a photon. The amplitude damping parameter is then given by

$$p_a = \exp(-t_w/\tau_e), \quad (30)$$

where  $t_w$  is the detection time window and  $\tau_e$  is the characteristic time of the NV emission (12 ns without cavity [85] and 6.48 ns with cavity [88]).

*D.4.5 Collection losses.* We model non-unity collection efficiencies by effective amplitude damping noise on the qubit encoding the presence/absence of a photon. The amplitude damping parameter is given by

$$p_a = (1 - p_{\text{zero\_phonon}} \cdot p_{\text{collection}}), \quad (31)$$

where  $p_{\text{zero\_phonon}}$  is the probability of emitting a photon in the zero phonon line (3% without cavity and 46% with cavity [85]) and  $p_{\text{collection}}$  is the probability of collection the photon into the fiber. From [53] we know that the total detection efficiency of the system is  $4 \cdot 10^{-4}$ , which can be decomposed as

$$p_{\text{total}} = p_{\text{zero\_phonon}} \cdot p_{\text{collection}} \cdot p_{\text{transmission}} \cdot p_{\text{detection}}, \quad (32)$$

where  $p_{\text{transmission}}$  is the probability that the photon is not lost during transmission in the fiber and  $p_{\text{detection}}$  is the probability that the detector clicks, given that there was a photon. Using equation (32) we find that  $p_{\text{collection}} = 0.014$  given the numbers in [53]. Frequency conversion succeeds with probability 30% [105], so in this case we use  $p_{\text{collection}} = 0.014 \cdot 0.3$ .

*D.4.6 Transmission losses.* Since the qubit sent from the node to the heralding station is encoded in the presence/absence of a photon, the losses during transmission over fiber are modeled as amplitude damping noise. We use an amplitude damping parameter  $p_{\text{t\_loss}}$  given as

$$p_{\text{t\_loss}} = 1 - 10^{-L \cdot \gamma/10}, \quad (33)$$

where  $L$  is the length of the fiber (in km) and  $\gamma$  is assumed to be 5 dB/km without frequency conversion and 0.5 dB/km with frequency conversion.



*D.4.7 Distinguishable photons.* Entanglement is generated between the electrons of the two nodes since the beam-splitter in the heralding station effectively deletes the information of which node a detected photon came from. This information is only perfectly detected if the photons emitted from the nodes are completely indistinguishable. In reality however, the photons properties (spectral, temporal etc.) can be slightly different and they are therefore not completely indistinguishable. In section D.4.7 we derive effective measurement operators of a beam-splitter measurement, taking photon indistinguishability into account, which we make use of in our simulation. For the physical setup we simulate, the overlap (visibility) of the photons coming from the nodes is approximately 0.9 [53], i.e.  $|\mu|^2 = 0.9$  where  $\mu$  is defined in equation (66).

*D.4.8 Detection losses and dark counts.* Detection losses and dark counts are modeled by probabilistically changing the ideal classical outcome from the detectors at the heralding station. For each detector, if the ideal detector clicked the noisy detector also clicks with probability  $p_{\text{detection}}$  and otherwise not. In the simulations we use  $p_{\text{detection}} = 0.8$ , as measured in [49].

If the ideal detector did not click the noisy detector does click with probability  $p_{\text{dark}}$ . The parameter used for the dark count is the dark count rate  $\lambda_{\text{dark}} = 20$  per second [53]. The dark counts follow a Poisson distribution and we have that

$$p_{\text{dark}} = 1 - \exp(-t_w \cdot \lambda_{\text{dark}}), \quad (34)$$

where  $t_w$  is the detection time window.

## D.5 Heralding station

Let us now consider the measurement at the Heralding Station in more detail in order to understand its error models.

*D.5.1 Distinguishable photons.* We here describe how we model a beam-splitter measurement of two photons which are not perfectly indistinguishable. This is relevant for many heralding entanglement generation schemes, since if photons are distinguishable the beam-splitter will not erase the information of where the photons came from. Two perfectly indistinguishable photons incident on a beam-splitter will always go to the same output arm, as captured by the Hong-Ou-Mandel effect [51]. However, if the photons are distinguishable they do not necessarily go to the same output arm, which can be detected in experiment. For a given setup, let's denote the probability that two incident photons on the beam-splitter go to different output arms as  $\chi$ .

We will in this section derive the effective POVM and Kraus operators correspond to detecting photons at the ends of the output arms of the beam-splitter in terms of  $\chi$ , under the assumptions described below and using ideas from the paper [16] where  $\chi$  is computed.

*D.5.2 Model.* Assume that there is a 50:50 beam-splitter with input arms  $a$  and  $b$  and output arms  $c$  and  $d$ . At the end of the output arms there are photon detectors that can click. We will assume that the detectors have a flat frequency response and at first that the detectors can count photons, i.e. there are different measurement outcomes for there being one or two photons incident on a detector. However we will show below how one can easily consider detectors which cannot count photons from the analysis in this note.

*Photons.* In many simulations we model the presence or absence of a photon as a two-level system, i.e. a qubit  $\alpha |0\rangle + \beta |1\rangle$ , where  $|0\rangle$  means no photon and  $|1\rangle$  one photon. We would then describe the state before the beam-splitter as a state living in a 2-qubit Hilbert space spanned by the following four basis vectors:

$$|00\rangle_{lr}, \quad |01\rangle_{lr}, \quad |10\rangle_{lr}, \quad |11\rangle_{lr} \quad (35)$$

describing 0 photons, photon on the right, photon on the left and two photons. Here  $l$  and  $r$  corresponds to arm  $a$  and  $b$  of the beam-splitter, but we distinguish these since we will denote  $a$  (and  $b$ ) as the infinite dimensional Hilbert space describing the spectral property of the photon. Note that we assume that there are never more than one photon per arm.

Describing the presence and absence of a photon as a qubit masks the fact that a photon can have many other degrees of freedom, such as polarization, spectral and temporal properties. We will here focus on spectral and temporal properties and will therefore model a photon in arm  $a$  with a spectral amplitude function  $\phi$  as the state

$$\int d\omega \phi(\omega) a^\dagger(\omega) |0\rangle_a, \quad (36)$$

where  $a^\dagger(\omega)$  is the creation operator of a photon in arm  $a$  of frequency  $\omega$  and  $|0\rangle_a$  is the vacuum and  $\phi$  is normalized such that

$$\int d\omega |\phi(\omega)|^2 = 1. \quad (37)$$

Furthermore, the state of arm  $b$  will be described by a spectral amplitude function  $\psi$  as

$$\int d\omega \psi(\omega) b^\dagger(\omega) |0\rangle_b. \quad (38)$$

Two photons arriving at the beam-splitter can have different spectral properties, captured by  $\phi$  and  $\psi$  being different. We will also include a possible temporal shift  $\tau$  between the arrival times of the two photons. As described in equation (16) of [16], a temporal shift of a photon in arm  $b$  induces the following action on the creation operators

$$b^\dagger(\omega) \rightarrow b^\dagger(\omega) e^{-i\omega\tau}. \quad (39)$$

*Beam-splitter.* The 50:50 beam-splitter acts on the creation operators in the following way:

$$a^\dagger(\omega) \rightarrow \frac{1}{\sqrt{2}}(c^\dagger(\omega) + d^\dagger(\omega)) \quad (40)$$

$$b^\dagger(\omega) \rightarrow \frac{1}{\sqrt{2}}(c^\dagger(\omega) - d^\dagger(\omega)). \quad (41)$$

Thus the state of a photon described as in equation (36), i.e. one photon in the input arm  $a$ , will after the beam-splitter become

$$|\phi\rangle_{cd} = \frac{1}{\sqrt{2}} \int d\omega \phi(\omega) (c^\dagger(\omega) + d^\dagger(\omega)) |0\rangle_{cd}. \quad (42)$$

Furthermore, the three other cases of no photon, one photon in the input arm  $b$  and one photon in each input arm becomes after the beam-splitter:

$$|0\rangle_{cd} \quad (43)$$

$$|\psi\rangle_{cd} = \frac{1}{\sqrt{2}} \int d\omega \psi(\omega) e^{-i\omega\tau} (c^\dagger(\omega) - d^\dagger(\omega)) |0\rangle_{cd} \quad (44)$$

$$|\phi, \psi\rangle_{cd} = \frac{1}{2} \int d\omega_1 \int d\omega_2 \phi(\omega_1) \psi(\omega_2) e^{-i\omega_1\tau} (c^\dagger(\omega_1) + d^\dagger(\omega_1)) (c^\dagger(\omega_2) - d^\dagger(\omega_2)) |0\rangle_{cd}. \quad (45)$$

Where the states  $|0\rangle_{cd}$ ,  $|\phi\rangle_{cd}$ ,  $|\psi\rangle_{cd}$  and  $|\phi, \psi\rangle_{cd}$  should be thought of as the corresponding states to the states in equation (35). Below, we will in fact formally define an isometry between these two Hilbert spaces.

*Detectors.* As mentioned we assume that the detectors have a flat frequency response. The event that the detector in arm  $c$  detected one photon can then be described by the projector

$$P_{1,0} = \int d\omega c^\dagger(\omega) |0\rangle\langle 0|_{cd} c(\omega) \quad (46)$$

Since we assume that there is maximally one photon arriving at each input arm of the beam-splitter the only other possible measurement outcomes are described by the following projectors:

$$P_{0,0} = |0\rangle\langle 0|_{cd} \quad (47)$$

$$P_{0,1} = \int d\omega d^\dagger(\omega) |0\rangle\langle 0|_{cd} d(\omega) \quad (48)$$

$$P_{1,1} = P_{1,0} \otimes P_{0,1} = \int d\omega_1 \int d\omega_2 c^\dagger(\omega_1) d^\dagger(\omega_2) |0\rangle\langle 0|_{cd} c(\omega_1) d(\omega_2) \quad (49)$$

$$P_{2,0} = \frac{1}{2} \int d\omega_1 \int d\omega_2 c^\dagger(\omega_1) c^\dagger(\omega_2) |0\rangle\langle 0|_{cd} c(\omega_1) c(\omega_2) \quad (50)$$

$$P_{0,2} = \frac{1}{2} \int d\omega_1 \int d\omega_2 d^\dagger(\omega_1) d^\dagger(\omega_2) |0\rangle\langle 0|_{cd} d(\omega_1) d(\omega_2) \quad (51)$$

where  $P_{0,0}$  corresponds to no photon,  $P_{0,1}$  one photon in arm  $d$ ,  $P_{1,1}$  one photon in each arm,  $P_{2,0}$  two photons in arm  $c$  and  $P_{0,2}$  two photons in arm  $d$ . Note that the factors of  $\frac{1}{2}$  are needed for  $P_{20}$  such that  $P_{20}^2 = P_{20}$  and similarly with  $P_{02}$ .

*Deriving effective POVM on presence/absence description.* The goal of this note is to derive the effective POVM on the Hilbert space  $lr$ , spanned by vectors in equation (35), induced by the projective measurements in equations (46)-(51) on the infinite-dimensional Hilbert space  $cd$ . To do this we will first define an isometry  $U_{lr \rightarrow cd}$  from the Hilbert space  $lr$  to  $cd$ , using the states in equation (35) and equations (42)-(45). This isometry will have the following action on the basis states of  $lr$ :

$$|00\rangle_{lr} \rightarrow |0\rangle_{cd} \quad (52)$$

$$|01\rangle_{lr} \rightarrow |\psi\rangle_{cd} \quad (53)$$

$$|10\rangle_{lr} \rightarrow |\phi\rangle_{cd} \quad (54)$$

$$|11\rangle_{lr} \rightarrow |\phi, \psi\rangle_{cd} \quad (55)$$

and will therefore be given as

$$U_{lr \rightarrow cd} = |0\rangle_{cd} \langle 00|_{lr} + |\psi\rangle_{cd} \langle 01|_{lr} + |\phi\rangle_{cd} \langle 10|_{lr} + |\phi, \psi\rangle_{cd} \langle 11|_{lr}. \quad (56)$$

One can easily check that the states  $|0\rangle_{cd}$ ,  $|\phi\rangle_{cd}$ ,  $|\psi\rangle_{cd}$  and  $|\phi, \psi\rangle_{cd}$  are mutually orthogonal and that  $U_{lr \rightarrow cd}$  is therefore indeed an isometry, i.e.

$$(U_{lr \rightarrow cd})^\dagger U_{lr \rightarrow cd} = \mathbb{1}_{lr}. \quad (57)$$

Lets assume that  $|\Phi\rangle_{lr}$  is a state in  $lr$  and we wish to compute the probability of receiving a measurement outcome corresponding to the projector  $P \in \{P_{00}, P_{10}, P_{01}, P_{11}, P_{20}, P_{02}\}$  for the state  $U_{lr \rightarrow cd} |\Phi\rangle_{lr}$ . Using Born's rule we find that this probability is given as

$$\langle \Phi|_{lr} (U_{lr \rightarrow cd})^\dagger P U_{lr \rightarrow cd} |\Phi\rangle_{lr} = \text{tr}[(U_{lr \rightarrow cd})^\dagger P U_{lr \rightarrow cd} |\Phi\rangle\langle \Phi|_{lr}]. \quad (58)$$

From the above equation we find that the effective POVM on  $lr$  is given as

$$\{(U_{lr \rightarrow cd})^\dagger P U_{lr \rightarrow cd} : P \in \{P_{00}, P_{10}, P_{01}, P_{11}, P_{20}, P_{02}\}\}. \quad (59)$$

whose elements we will denote as  $M_{00}$ ,  $M_{10}$ ,  $M_{01}$ ,  $M_{11}$ ,  $M_{20}$  and  $M_{02}$ . In section D.5.2 we compute what these POVM-elements are and find a choice of Kraus operators in section D.5.3 for both the case then the detector can count photons and when it cannot.

*Effective POVMs.* Here we compute the POVM-elements in equation (59) one-by-one.

$M_{11}$ :

Let's start with  $M_{11}$  since this will allow us to relate these POVM-elements to  $\chi$ , i.e. the probability that both detectors click, given that there were one photon in each input arm. The operator  $P_{11}$  only has non-zero overlap with the term  $|\phi, \psi\rangle_{cd} \langle 11|_{lr}$  of  $U_{lr \rightarrow cd}$  and is therefore given as

$$M_{11} = (U_{lr \rightarrow cd})^\dagger P_{11} U_{lr \rightarrow cd} = |11\rangle_{lr} \langle \phi, \psi|_{cd} P_{11} |\phi, \psi\rangle_{cd} \langle 11|_{lr}. \quad (60)$$

Lets evaluate the factor  $\langle \phi, \psi|_{cd} P_{11} |\phi, \psi\rangle_{cd}$ . Using equation (45) and equation (49) we find that the above expression evaluates to

$$\begin{aligned} \langle \phi, \psi|_{cd} P_{11} |\phi, \psi\rangle_{cd} &= \frac{1}{2} \int d\omega_1 \int d\omega_2 \phi^*(\omega_1) \psi^*(\omega_2) e^{i\omega_2 \tau} \langle 0|_{cd} (c(\omega_1) + d(\omega_1))(c(\omega_2) - d(\omega_2)) \\ &\quad \times \int d\omega_3 \int d\omega_4 c^\dagger(\omega_3) d^\dagger(\omega_4) |0\rangle\langle 0|_{cd} c(\omega_3) d(\omega_4) \\ &\quad \times \frac{1}{2} \int d\omega_5 \int d\omega_6 (c^\dagger(\omega_5) + d^\dagger(\omega_5))(c^\dagger(\omega_6) - d^\dagger(\omega_6)) |0\rangle_{cd} \phi(\omega_5) \psi(\omega_6) e^{-i\omega_6 \tau} \\ &= \frac{1}{4} \int d\omega_1 \int d\omega_2 \int d\omega_3 \int d\omega_4 \int d\omega_5 \int d\omega_6 \phi^*(\omega_1) \psi^*(\omega_2) \phi(\omega_5) \psi(\omega_6) e^{i\omega_2 \tau} e^{-i\omega_6 \tau} \left( \right. \\ &\quad + \delta(\omega_2 - \omega_3) \delta(\omega_3 - \omega_6) \delta(\omega_1 - \omega_4) \delta(\omega_4 - \omega_5) \\ &\quad - \delta(\omega_2 - \omega_3) \delta(\omega_3 - \omega_5) \delta(\omega_1 - \omega_4) \delta(\omega_4 - \omega_6) \\ &\quad - \delta(\omega_1 - \omega_3) \delta(\omega_3 - \omega_6) \delta(\omega_2 - \omega_4) \delta(\omega_4 - \omega_5) \\ &\quad \left. + \delta(\omega_1 - \omega_3) \delta(\omega_3 - \omega_5) \delta(\omega_2 - \omega_4) \delta(\omega_4 - \omega_6) \right) \end{aligned} \quad (61)$$

where we used the fact that  $\langle 0|_{cd} c(\omega_1)c^\dagger(\omega_2)|0\rangle_{cd} = \delta(\omega_1 - \omega_2)$  and similarly for arm  $d$ . Using that

$$\int d\omega_2 f(\omega_2)\delta(\omega_1 - \omega_2) = f(\omega_1) \quad (63)$$

we find that equation (62) evaluates to

$$\frac{1}{2} \int d\omega_1 \int d\omega_2 (\phi^*(\omega_1)\psi^*(\omega_2)\phi(\omega_1)\psi(\omega_2)e^{i\omega_2\tau}e^{-i\omega_2\tau} - \phi^*(\omega_1)\psi^*(\omega_2)\phi(\omega_2)\psi(\omega_1)e^{i\omega_2\tau}e^{-i\omega_1\tau}) \quad (64)$$

Finally using equation (37) we find that  $M_{11}$  evaluates to

$$M_{11} = \frac{1}{2}(1 - |\mu|^2) |11\rangle\langle 11|_{lr} \quad (65)$$

where

$$\mu = \int d\omega \phi^*(\omega)\psi(\omega)e^{-i\omega\tau}. \quad (66)$$

From equation (65) we can relate  $|\mu|$  to  $\chi$  as

$$\chi = \frac{1}{2}(1 - |\mu|^2). \quad (67)$$

$M_{20}$ :

The operator  $P_{20}$  only has non-zero overlap with the term  $|\phi, \psi\rangle_{cd} \langle 11|_{lr}$  of  $U_{lr \rightarrow cd}$  and is therefore given as

$$M_{11} = (U_{lr \rightarrow cd})^\dagger P_{20} U_{lr \rightarrow cd} = |11\rangle_{lr} \langle \phi, \psi|_{cd} P_{20} |\phi, \psi\rangle_{cd} \langle 11|_{lr}. \quad (68)$$

Lets evaluate the factor  $\langle \phi, \psi|_{cd} P_{20} |\phi, \psi\rangle_{cd}$ . Using equation (45) and equation (50) we find that the above expression evaluates to

$$\begin{aligned} \langle \phi, \psi|_{cd} P_{20} |\phi, \psi\rangle_{cd} &= \frac{1}{2} \int d\omega_1 \int d\omega_2 \phi^*(\omega_1)\psi^*(\omega_2)e^{i\omega_2\tau} \langle 0|_{cd} (c(\omega_1) + d(\omega_1))(c(\omega_2) - d(\omega_2)) \\ &\quad \times \frac{1}{2} \int d\omega_3 \int d\omega_4 c^\dagger(\omega_3)c^\dagger(\omega_4)|0\rangle\langle 0|_{cd} c(\omega_3)c(\omega_4) \\ &\quad \times \frac{1}{2} \int d\omega_5 \int d\omega_6 (c^\dagger(\omega_5) + d^\dagger(\omega_5))(c^\dagger(\omega_6) - d^\dagger(\omega_6))|0\rangle_{cd} \phi(\omega_5)\psi(\omega_6)e^{-i\omega_6\tau} \\ &= \frac{1}{8} \int d\omega_1 \int d\omega_2 \int d\omega_3 \int d\omega_4 \int d\omega_5 \int d\omega_6 \phi^*(\omega_1)\psi^*(\omega_2)\phi(\omega_5)\psi(\omega_6)e^{i\omega_2\tau}e^{-i\omega_6\tau} \\ &\quad \times \left( \delta(\omega_1 - \omega_4)\delta(\omega_2 - \omega_3) + \delta(\omega_1 - \omega_3)\delta(\omega_2 - \omega_4) \right) \\ &\quad \times \left( \delta(\omega_3 - \omega_6)\delta(\omega_4 - \omega_5) + \delta(\omega_3 - \omega_5)\delta(\omega_4 - \omega_6) \right) \end{aligned} \quad (69)$$

where we used the fact that  $\langle 0|_{cd} c(\omega_1)c(\omega_2)c^\dagger(\omega_3)c^\dagger(\omega_4)|0\rangle_{cd} = \delta(\omega_1 - \omega_3)\delta(\omega_2 - \omega_4) + \delta(\omega_2 - \omega_3)\delta(\omega_1 - \omega_4)$ . Then similarly to  $M_{11}$  we find that equation (70) evaluates to

$$\langle \phi, \psi|_{cd} P_{20} |\phi, \psi\rangle_{cd} = \frac{1}{4}(1 + |\mu|^2) \quad (71)$$

and we thus find  $M_{20}$  to be

$$M_{20} = \frac{1}{4}(1 + |\mu|^2) |11\rangle\langle 11|_{lr}. \quad (72)$$

$M_{20}$ :

Similarly to  $M_{20}$  we find that  $M_{02}$  evaluates to

$$M_{02} = \frac{1}{2}(1 + |\mu|^2) |11\rangle\langle 11|_{lr}. \quad (73)$$

$M_{10}$ :

The operator  $P_{10}$  only has non-zero overlap with the terms  $|\phi\rangle_{cd} \langle 10|_{lr}$  and  $|\psi\rangle_{cd} \langle 01|_{lr}$  of  $U_{lr \rightarrow cd}$  and is therefore given as

$$M_{10} = (U_{lr \rightarrow cd})^\dagger P_{10} U_{lr \rightarrow cd} = \left( |10\rangle_{lr} \langle \phi|_{cd} + |01\rangle_{lr} \langle \psi|_{cd} \right) P_{10} \left( |\phi\rangle_{cd} \langle 10|_{lr} + |\psi\rangle_{cd} \langle 01|_{lr} \right). \quad (74)$$

## A Link Layer Protocol for Quantum Networks

Lets evaluate the factors  $\langle \phi|_{cd} P_{10} |\phi\rangle_{cd}$ ,  $\langle \psi|_{cd} P_{10} |\psi\rangle_{cd}$ ,  $\langle \phi|_{cd} P_{10} |\psi\rangle_{cd}$  and  $\langle \psi|_{cd} P_{10} |\phi\rangle_{cd}$  one-by-one. First we have that:

$$\begin{aligned}
 \langle \phi|_{cd} P_{10} |\phi\rangle_{cd} &= \frac{1}{\sqrt{2}} \int d\omega_1 \phi^*(\omega_1) \langle 0|_{cd} (c(\omega_1) + d(\omega_1)) \\
 &\quad \times \int d\omega_2 c^\dagger(\omega_2) |0\rangle\langle 0|_{cd} c(\omega_2) \\
 &\quad \times \frac{1}{\sqrt{2}} \int d\omega_3 (c^\dagger(\omega_3) + d^\dagger(\omega_3)) |0\rangle_{cd} \phi(\omega_3) \\
 &= \frac{1}{2} \int d\omega_1 \int d\omega_2 \int d\omega_3 \phi^*(\omega_1) \phi(\omega_3) \delta(\omega_1 - \omega_2) \delta(\omega_2 - \omega_3) \\
 &= \frac{1}{2} \int d\omega |\phi(\omega)|^2 \\
 &= \frac{1}{2}.
 \end{aligned} \tag{75}$$

and similarly that

$$\langle \psi|_{cd} P_{10} |\psi\rangle_{cd} = \frac{1}{2}. \tag{76}$$

Furthermore, we find that

$$\begin{aligned}
 \langle \phi|_{cd} P_{10} |\psi\rangle_{cd} &= \frac{1}{\sqrt{2}} \int d\omega_1 \phi^*(\omega_1) \langle 0|_{cd} (c(\omega_1) + d(\omega_1)) \\
 &\quad \times \int d\omega_2 c^\dagger(\omega_2) |0\rangle\langle 0|_{cd} c(\omega_2) \\
 &\quad \times \frac{1}{\sqrt{2}} \int d\omega_3 (c^\dagger(\omega_3) - d^\dagger(\omega_3)) |0\rangle_{cd} \psi(\omega_3) e^{-i\omega_3 \tau} \\
 &= \frac{1}{2} \int d\omega_1 \int d\omega_2 \int d\omega_3 \phi^*(\omega_1) \psi(\omega_3) e^{-i\omega_3 \tau} \delta(\omega_1 - \omega_2) \delta(\omega_2 - \omega_3) \\
 &= \frac{1}{2} \int d\omega \phi^*(\omega) \psi(\omega) e^{-i\omega \tau} \\
 &= \frac{1}{2} \mu.
 \end{aligned} \tag{77}$$

where  $\mu$  is defined in equation (66). One easily then finds that

$$\langle \psi|_{cd} P_{10} |\phi\rangle_{cd} = (\langle \phi|_{cd} P_{10} |\psi\rangle_{cd})^* = \frac{1}{2} \mu^*. \tag{78}$$

Combining the above results, we find that  $M_{10}$  is given as

$$M_{10} = \frac{1}{2} \left( |10\rangle\langle 10|_{lr} + |01\rangle\langle 01|_{lr} + \mu |10\rangle\langle 01|_{lr} + \mu^* |01\rangle\langle 10| \right) \tag{79}$$

$M_{10}$ :

Similarly to  $M_{10}$  one finds that  $M_{01}$  evaluates to

$$M_{01} = \frac{1}{2} \left( |10\rangle\langle 10|_{lr} + |01\rangle\langle 01|_{lr} - \mu |10\rangle\langle 01|_{lr} - \mu^* |01\rangle\langle 10| \right) \tag{80}$$

$M_{00}$ :

Its easy to see that

$$M_{00} = |00\rangle\langle 00|_{lr}. \tag{81}$$

POVM for photon-counter detectors. To summarize we found that the POVM-elements are given as

$$M_{00} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (82)$$

$$M_{10} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & \mu & 0 \\ 0 & \mu^* & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (83)$$

$$M_{01} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & -\mu & 0 \\ 0 & -\mu^* & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (84)$$

$$M_{11} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 - |\mu|^2 \end{pmatrix} \quad (85)$$

$$M_{20} = \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 + |\mu|^2 \end{pmatrix} \quad (86)$$

$$M_{02} = \frac{1}{4} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 + |\mu|^2 \end{pmatrix} \quad (87)$$

where the rows and columns of the above matrices are ordered as  $|00\rangle\langle 00|_{lr}$ ,  $|10\rangle\langle 10|_{lr}$ ,  $|01\rangle\langle 01|_{lr}$ ,  $|11\rangle\langle 11|_{lr}$  and  $\mu$  is given as

$$\mu = \int d\omega \phi^*(\omega)\psi(\omega)e^{-i\omega\tau}. \quad (88)$$

and is related to the probability that both detectors click, given that there were one photon in each input arm  $\chi$  as

$$\chi = \frac{1}{2}(1 - |\mu|^2). \quad (89)$$

*POVM for non-photon-counter detectors.* If the detectors used cannot distinguish between one and two photons we can simply add the POVM elements  $M_{10}$  and  $M_{20}$  to get a new POVM given as

$$\tilde{M}_{00} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (90)$$

$$\tilde{M}_{10} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & \mu & 0 \\ 0 & \mu^* & 1 & 0 \\ 0 & 0 & 0 & (1 + |\mu|^2)/2 \end{pmatrix} \quad (91)$$

$$\tilde{M}_{01} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & -\mu & 0 \\ 0 & -\mu^* & 1 & 0 \\ 0 & 0 & 0 & (1 + |\mu|^2)/2 \end{pmatrix} \quad (92)$$

$$\tilde{M}_{11} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 - |\mu|^2 \end{pmatrix} \quad (93)$$

*D.5.3 Effective Kraus operators.* Given the POVMs in equation (82)-(87) and equation (90)-(93) one can choose corresponding Kraus operators for these measurements by taking the matrix square root of the corresponding POVM-elements. Assuming that  $\mu$  is real one finds a set of Kraus operators of the POVM  $\{\tilde{M}_{00}, \tilde{M}_{10}, \tilde{M}_{01}, \tilde{M}_{11}\}$  to be

$$\tilde{E}_{00} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (94)$$

$$\tilde{E}_{10} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & (\sqrt{1+\mu} + \sqrt{1-\mu})/\sqrt{2} & (\sqrt{1+\mu} - \sqrt{1-\mu})/\sqrt{2} & 0 \\ 0 & (\sqrt{1+\mu} - \sqrt{1-\mu})/\sqrt{2} & (\sqrt{1+\mu} + \sqrt{1-\mu})/\sqrt{2} & 0 \\ 0 & 0 & 0 & \sqrt{1+|\mu|^2} \end{pmatrix} \quad (95)$$

$$\tilde{E}_{01} = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & (\sqrt{1+\mu} + \sqrt{1-\mu})/\sqrt{2} & (\sqrt{1-\mu} - \sqrt{1+\mu})/\sqrt{2} & 0 \\ 0 & (\sqrt{1-\mu} - \sqrt{1+\mu})/\sqrt{2} & (\sqrt{1+\mu} + \sqrt{1-\mu})/\sqrt{2} & 0 \\ 0 & 0 & 0 & \sqrt{1+|\mu|^2} \end{pmatrix} \quad (96)$$

$$\tilde{E}_{11} = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sqrt{1-|\mu|^2} \end{pmatrix}. \quad (97)$$

## D.6 Classical communication

*D.6.1 Optical Link Error Model.* We claimed that we highly inflated losses in the simulation to stress test our protocol. We now consider more realistic values for such errors by considering a realistic *packet-level* error model for the non-quantum optical link. For this we have assumed that two quantum internet end nodes are connected by a legacy 1000BASE-ZX single-mode 1550 nm wavelength Gigabit Ethernet link. The reason for choosing 1000BASE-ZX interface is (i) its achievable long-distance transmission at least up to 70 km with no dependency on optical repeaters and (ii) decades of its successful deployment within magnitude of networks worldwide.



To be conservative, our optical Gigabit Ethernet model assumes a typical worst-case optical link budget (0.5 dB/km attenuation<sup>5</sup>, 0.7 dB/connector loss, 0.1 dB/splice/(joint) loss, and 3 dB safety margin) [27]. We also assume a typical worst-case  $-1$  dBm optical transmission power and  $-24$  dB receiver sensitivity of a 1000BASE-ZX small form-factor hot pluggable transceiver, see e.g. [61]. For a maximum realism of link error over an optical link we model a *IEEE 802.3 frame errors*, instead of modeling individual bit errors of every message sent across the network. The latter would require a software implementation of a complete modulation and coding layer of IEEE 802.3 which is beyond the scope of this work. Using measurement trace-driven packet-level Gigabit Ethernet frame error data from [56, Table 6.1] we have mapped the received SNR per transmitter/receiver distance to the respective frame error probability, which was then applied to every classical message sent over an optical link between quantum end nodes. SNR values that were not represented in the measurements of [56] have been linearly interpolated. We have not distinguished between the lengths of each classical message as the model of [56] has aggregated over all messages captured over a measured campus Ethernet link (cf. [56, Fig. 6.1]). We note that our modeling approach is equivalent to the frame error models applied in e.g. NS-3 [28] for WiFi frame errors.

For two example long-distance Quantum Internet typologies (node-to-node distance of 15 km and 20 km, respectively) we have ended up in a perfect frame error probability, with the assumption that amount of splices is zero<sup>6</sup> (we only start to observe frame errors only at transmitter/receiver distance exceeding 40 km for the above model variables, with a very narrow transition error between no frame error rate and disconnected interface, i.e. frame error rate of one). Even when we increase the number of splices to an exaggerated level, say 30 splices for a 15 km interface (with 0.3 dB loss/splice), we still observe a very low frame error probability of  $4 \times 10^{-8}$ . Therefore, to test the effect of frame errors on the non-quantum optical link on the Quantum Internet protocol stack—in the cases of extreme frame loss—we have increased the value of frame error to  $10^{-4}$  (and tested frame error rate to up to  $10^{-10}$ —an error rate level of a 20 km link with 21 splices—in steps of  $10^{-1}$ ). If our protocol would work in such a high (but unrealistic) condition then it would also work on a realistic low-error optical link.

*D.6.2 Optical Link CRC Error Model.* Additionally, we have investigated a non-zero probability of CRC not being able to detect a frame error. Assuming the same optical link type (e.g. 1000BASE-ZX) we have used a model of [44] to calculate the respective probability of not detecting a CRC frame error within a IEEE 802.3 frame. For this we have mapped the transmitter/receiver distance to the respective SNR (the same way as described in Section D.6.1). Then we mapped the SNR to the respective BER using [56, Table 6.3] (performing the same process of interpolating SNR between the points not measured by [56] as for the optical link error model, see again Section D.6.1) and then using [44, Fig. 1] mapped this resulting SNR to the respective probability of undetected error. We have assumed a worst case scenario of the longest IEEE 802.3 frame (i.e.  $n = 12144$  bits, that is a maximum MTU). Again, for any of the two Quantum Internet lengths mentioned above, we do not find any CRC errors. At the highly-spliced case, considered in Section D.6.1, we obtain an *extremely low* CRC error rate of  $1.4 \times 10^{-23}$ . Therefore such errors were decided to be ignored in our implementation. Another reason for not considering these errors: it would require a full implementation of en- and decoding of classical frames which outside the scope of this work.

## E PROTOCOLS

Here we give details of the implementation of the physical and link layer protocols in our simulations. Python implementation is available upon request.

### E.1 Distributed Queue Protocol

In order to track the individual applications that the entangled qubits belong to, the EGP makes use of a *distributed queue* which shares request information between peers. Management of the distributed queue is performed by the Distributed Queue Protocol (DQP). In addition to storing the parameters supplied with a CREATE request generated by the layers above the link layer<sup>7</sup>, DQP will keep additional information about each entanglement request including its *create\_time*, *min\_time* at which the request may be executed and *MHP timeout cycle* by which the entanglement request will time out. We proceed with the introduction of the DQP by describing the structure of priority queues, followed by the queue establishment process, DQP message sequence diagram and DQP associated messages.

*E.1.1 Priority Queues.* Priorities are necessary to fulfill the use case requirements outlined in section 3. This is accomplished by adding requests to different types of queues  $Q = \{Q_1, \dots, Q_L\}$ , where  $L$  is the total number of queues in the distributed

<sup>5</sup>Fibers measured for QL2020 have been found to have this loss level.

<sup>6</sup>Which is consistent with the measurements, e.g. in [30, Section 4].

<sup>7</sup>Refer to Section 4 of the main paper for the details of the CREATE request.

queue. Each queue can contain a maximum of  $x$  items simultaneously (in other words  $x$  is the maximum size of each individual queue), where an item is an individual entanglement request with its associated metadata, e.g. create\_time, min\_time, MHP timeout cycle. Each CREATE request is assigned a queue number by the scheduler (see Section E.3.1 below), and receives an absolute queue ID which is a tuple  $(j, i_j)$  where  $j$  indicates the designated queue  $Q_j$  (or, more abstractly, the *queue ID* of the entanglement request) and  $i_j$  is a unique ID within  $Q_j$ . Equivalently, for a finite number of queues we will denote  $(j, i_j)$  as the *absolute queue ID* or  $a_{ID}$ , and use  $(j, i_j) \in Q$  to indicate the ID of the request.

The queue ID must obey the following properties:

- *Total order*: Items on each queue follow a total order of items waiting in the queue determined by  $i_j$ .
- *Arrival time*: ID of an entanglement (CREATE) request is a function of its arrival time. Let  $t_1$  and  $t_2$  denote the create\_time of entanglement requests 1 and 2, respectively. Then, let  $i_1$  and  $i_2$  denote their respective queue ID's. If both requests are added to the same queue  $Q_j$ , and  $t_1 < t_2$ , then  $(i_2 - i_1) \bmod x > 0$ . That is, if a CREATE request arrives earlier, it will also receive a lower queue ID.

We will now outline the distributed queue establishment within DQP. For simplicity of exposition, we now assume there is only one queue, i.e.,  $L = 1$ .

*E.1.2 DQP Queue Establishment.* The core objective of DQP is to obtain shared queues at both nodes, i.e. the items and the order of the elements in the queues are agreed upon. That is, both controllable end nodes  $A$  and  $B$  hold local queues  $Q^A = \{Q_1^A, \dots, Q_L^A\}$  and  $Q^B = \{Q_1^B, \dots, Q_L^B\}$  respectively, which are synchronized using the DQP. CREATE request additions to the queue  $Q_j$  can be made by either  $A$  or  $B$  invoking the DQP by the function  $\text{ADD}(j, c_r)$ , where  $c_r$  is the entanglement request by CREATE message. ADD returns a tuple  $(i_j, R)$  where  $R$  indicated success or failure. Failure can occur if

- no acknowledgments are received within a certain time frame, i.e. a timeout occurs,
- the remote node rejects addition to the queue, or
- the queue is full.

Success means that the request to create entanglement is placed into  $Q^A$  and  $Q^B$  such that the following properties are satisfied:

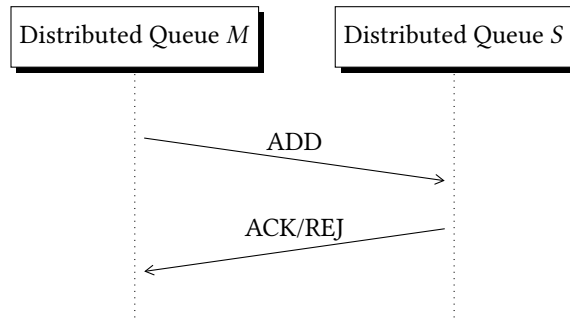
- *Equal queue number*: If a request is added by  $A$  as  $(j, i_j) \in Q^A$ , then it will (eventually) be added at  $B$  with the same absolute queue ID  $(j, i_j) \in Q^B$  (and vice versa);
- *Uniqueness of queue ID*: If a request is placed into the queue by either  $A$  or  $B$ , then it is assigned a unique queue number. That is, if  $(j, a) \in Q^A$  and  $(j, a') \in Q^A$  reference two distinct CREATE requests, then  $a \neq a'$ ;
- *Consistency*: If  $(j, i_j) \in Q^A$  and  $(j, i_j) \in Q^B$  then both absolute queue IDs refer to the same request<sup>8</sup>;
- *Fairness*: If  $A$  (or  $B$ ) is issuing requests continuously, then also the other node  $B$  (or  $A$ ) will get to add items to the queue after  $A$  (or  $B$ ) in a “fair manner“ as determined by the window size, denoted as  $W_A$  ( $W_B$ ). More precisely, if  $a_1, \dots, a_N$  are CREATE requests submitted at  $A$ , and  $b_1, \dots, b_M$  are CREATE requests submitted at  $B$  with  $N > W_A$  and  $M > W_B$ —all assigned to the same queue  $Q_j$  but not yet added—then the final ordering of the requests on the queue obeys  $a_1, \dots, a_m, b_1, \dots, b_k, a_{m+1}, \dots$  with  $m \leq W_A$  and  $k \leq W_B$ .

Recall that each request receives a minimum to be executed time—a time buffer before the request may begin processing which takes into account the processing time to add it into the queue (denoted by the min\_time)—which we will choose to be the expected propagation delay between  $A$  and  $B$ . The purpose of this minimum time is to decrease the likelihood  $A$  or  $B$  wants to produce an entanglement before the other node is ready. If either  $A$  or  $B$  begins processing early, no penalty other than reduced performance due to increased decoherence of the quantum memory results. Refer to Section E.1.4 on how this minimum time is passed between nodes.

We recall that in the current implementation of quantum network we have two nodes only. This implies that the queue establishment can be realized by one node being the master controller of the queue marshaling access to the queue, and the other the slave controller. Extensions to multiple nodes are more complex, and a motivation to consider heralding station-centric protocols in the future versions of the protocol. Also, as we have two nodes only, there is no need for the introduction of leader election or a network discovery mechanism. We leave this as future work.

*E.1.3 DQP Sequence Diagrams.* Figure 23 shows a DQP sequence diagram of adding an item to the queue containing a request to the distributed queue. Specifically, an item is an entanglement create request with its associated properties that is passed inside an ADD message within its REQ field, refer to Figure 24 for details.

<sup>8</sup>This is implied by the previous two conditions, but added for clarity.



**Figure 23: DQP operation timeline.** User *M* (master) adds an item to the distributed queue by sending an ADD message to peer node *S* (slave). *S* either acknowledges or rejects the request using ACK and REJ messages, respectively. Note that this process is symmetric when *S* attempts to add an item to the queue. For a definition of all messages refer to Figure 24.

Upon receiving an ADD message from master *M*, a slave *S* may choose to acknowledge the item with an ACK message, should validation pass, or reject it with the REJ message for any of the previously mentioned reasons. In the case that master *M* never receives an acknowledgment (ACK) or rejection (REJ) message after a timeout, the item will be removed from the queue and no processing will occur on the request. Loss of ADD, REJ, and ACK messages in the distributed queue protocol result in retransmissions of the original ADD to guarantee the receipt of rejection and acknowledgement messages.

|   |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
|---|--|--|--|-----|--|------|--|----|--|-----|--|----|--|------|--|------------|--|--|--|------------|--|--|--|--|--|--|--|--|--|--|--|
| 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| OPT (reserved)  |  |  |  | FT  |  | CSEQ |  |    |  | QID |  |    |  | QSEQ |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Schedule Cycle  |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Timeout   |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Minimum Fidelity  |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Purpose ID  |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  | Create ID  |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Number of Pairs   |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  | Priority   |  |  |  | (reserved) |  |  |  |  |  |  |  |  |  |  |  |
| Initial Virtual Finish  |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| Estimated Cycles/Pair   |  |  |  |     |  |      |  |    |  |     |  |    |  |      |  |            |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |
| STR   |  |  |  | ATM |  |      |  | MD |  |     |  | MR |  |      |  | (reserved) |  |  |  |            |  |  |  |  |  |  |  |  |  |  |  |

**Figure 24: Packet format for ADD, ACK, and REJ.** Explanation of the message fields—OPT: field reserved for future options, FT: frame type (00: ADD, 01: ACK, 10: REJ), CSEQ: the communication sequence number of the transmitted message; QID: the ID of the queue to add the request to; QSEQ: the sequence number within the specified queue to assign the request; Schedule Cycle: the first MHP cycle when the request may begin (equivalent to min\_time); Timeout: The MHP cycle when the request will time out; Initial Virtual Finish: Scheduling information for weighted fair queuing; STR: store flag, ATM: atomic flag, MD: measure directly flag, MR: master request flag.

When the slave *S* wishes to add an item to the queue, a message containing the request information and desired queue is included within the messages. Because the master controller has the final say on the state of the queue, a sequence number within the specified queue will be transmitted in return to the slave such that absolute queue IDs are consistent between the nodes.

*E.1.4 DQP Packet Formats.* Figure 24 presents the packet format for messages exchanged in the DQP. Schedule Cycle and Timeout Cycle of 64 bits is governed by the maximum number of MHP cycles in the scheduler. Purpose ID of 16 bits enables pointing to  $2^{16}$  different applications and the total number of uniquely addressed applications and follows from the number chosen for IPv4. Create ID defines the identifier of locally created request. Number of pairs enables to request up to  $2^{16}$  pairs. Priority field of 4 bits is used as we enable 16 local queues composing the distributed queue and each one represents a priority lane. Initial Virtual Finish is used for weighted fair queuing.

## E.2 Midpoint Heralding Protocol

The purpose of MHP is to create entanglement using a midpoint heralding protocol. The operation of the MHP is defined by Protocol E.2.

---

### Protocol 1 MHP for use with the Node-Centric EGP

---

*Definition of functions and variables.*

- POLLEGP: process to poll for entanglement parameters from EGP; it returns:
  - flag: true/false indicating whether entanglement should be attempted or not;
  - PSEQ: The pulse sequence identifier that should be issued to the hardware to initialize communication qubit and produce spin-photon entanglement. May also instruct the hardware to store the spin state within a storage qubit.
  - $a_{ID}$ : Absolute queue ID, i.e.  $(j, i_j)$ , of the request entanglement is being attempted for ( $a_{ID,A}$  and  $a_{ID,B}$  for nodes  $A$  and  $B$ , respectively);
  - params: parameters to use for the entanglement attempt such as bright state population  $\alpha$ ;
- $mhp_{err}$ : error in MHP reported to EGP through REPLY message ( $REPLY_A$  and  $REPLY_B$  sent to nodes  $A$  and  $B$ , respectively), which can take the following values:
  - GEN\_FAIL: general failure that occurs locally at the MHP (failed qubit initialization; other errors). Note: this error message is passed to EGP locally and not included in the REPLY message (see Figure 28);
  - QUEUE\_MISMATCH: an error sent by the midpoint when  $a_{ID}$  included in frame from  $A$  does not match  $a_{ID}$  included in frame from  $B$ ;
  - TIME\_MISMATCH: when messages from  $A$  and  $B$  does not arrive at midpoint within the same time interval;
  - NO\_MESSAGE\_OTHER: when the midpoint receives a message from only one of  $A$  or  $B$ ;
- GEN: the frame sent by  $A$  and  $B$  to the midpoint requesting entanglement. The contents include:
  - $a_{ID}$ : same as  $a_{ID}$  above;
- $REPLY_A$  and  $REPLY_B$ : the REPLY frames sent to  $A$  and  $B$  respectively. The contents include:
  - outcome: the outcome of the attempted entanglement at the midpoint, also encodes the error that occurred for the attempt at entanglement (see errors listed above); stored locally at mid-point.
  - $seq_{MHP}$ : the sequence number from the MHP;
  - $a_{ID,receiver}$ : The absolute queue ID that was submitted by the node receiving the REPLY.
  - $a_{ID,peer}$ : The absolute queue ID that was submitted by peer node.

*Initialization.* Initialize sequence numbers (set initial  $seq_{MHP} = 0$  at  $H$ ). Start timer using a global synchronized clock.

*The protocol, executed at each time step:*

(1) **Executed at Node A or Node B:**

- (a) Poll EGP, i.e.  $POLLEGP = (flag, PSEQ, a_{ID}, params)$
- (b) If  $flag = true$ , i.e., we want to make entanglement:
  - (i) Issue  $PSEQ$  to hardware to initialize communication qubit and produce spin-photon entanglement.  $PSEQ$  may also instruct the hardware to store spin state in a storage qubit. If any failures occur, send  $mhp_{err} = GEN\_FAIL$  back to the EGP and skip to next time step.
  - (ii) Use  $GEN = (a_{ID})$  and transmit to midpoint upon photon emission.

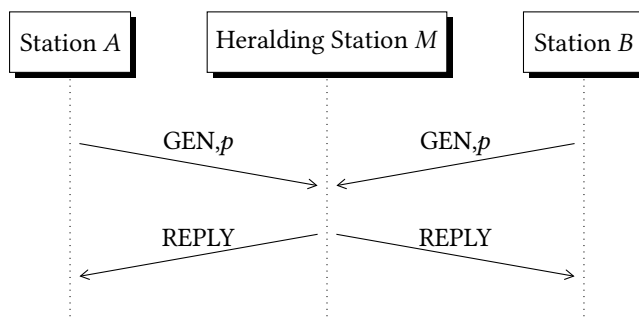
**Protocol 1** (cont.) MHP for use with the Node-Centric EGP**(2) Heralding station  $H$ :**

- (a) Perform the following upon receipt of GEN messages:
- (i) If messages from  $A$  and  $B$  do not arrive within the same time interval, let  $mhp_{err} = \text{TIME\_MISMATCH}$  and send  $\text{REPLY}_A = (mhp_{err}, seq_{MHP}, a_{ID,A}, ID_B)$  to  $A$  and  $\text{REPLY}_B = (mhp_{err}, seq_{MHP}, a_{ID,B}, a_{ID,A})$  to  $B$ .
  - (ii) If  $a_{ID,A} \neq a_{ID,B}$ , then set  $mhp_{err} = \text{QUEUE\_MISMATCH}$  and send  $\text{REPLY}_A = (mhp_{err}, seq_{MHP}, a_{ID,A}, a_{ID,B})$  to  $A$  and  $\text{REPLY}_B = (mhp_{err}, seq_{MHP}, a_{ID,B}, a_{ID,A})$  to  $B$ .
  - (iii) If GEN arrives only from  $A$ , set  $mhp_{err} = \text{NO\_MESSAGE\_OTHER}$  and send  $\text{REPLY} = (mhp_{err}, seq_{MHP}, a_{ID,A}, a_{ID,B}=\text{null})$  to  $A$ , where  $a_{ID,B}=\text{null}$  indicates leaving the field as the zero string. Perform vice versa if GEN arrives only from  $B$ .
  - (iv) If no errors occurred then execute quantum swap. Inspect detection result within corresponding time window with  $r \in \{0, 1, 2\}$  where 0 denotes failure and 1 and 2 denote the creation of states one and two respectively. If  $r \in \{1, 2\}$ , a unique and increasing sequence number  $seq_{MHP}$  is chosen by the heralding station (incrementing a counter) to be sent to both  $A$  and  $B$ . Midpoint sends  $\text{REPLY} = (\text{outcome}, seq_{MHP}, a_{ID}, a_{ID})$  to  $A$  and  $B$ .

**(3) Executed at Node  $A$  or Node  $B$**  (here  $a_{ID,local}=a_{ID,A}$ ,  $a_{ID,peer}=a_{ID,B}$  in  $A$  and vice versa in  $B$ ):

- (a) If  $\text{REPLY} = (\text{outcome}, seq_{MHP}, a_{ID,local}, a_{ID,peer})$  returns from midpoint:
  - (i) Set  $\text{RESULT} = (\text{outcome}, seq_{MHP}, a_{ID,local}, err=000, a_{ID,peer})$  and pass to EGP.
- (b) Else if  $\text{REPLY} = (mhp_{err}, seq_{MHP}, a_{ID,local}, a_{ID,peer})$  returns from the midpoint:
  - (i) Set  $\text{RESULT} = (\text{outcome}=0, seq_{MHP}, a_{ID,local}, mhp_{err}, a_{ID,peer})$  and pass to EGP.

*E.2.1 MHP Sequence Diagrams.* The MHP sequence diagram is defined by two cases: the successful—see Figure 25, and unsuccessful one—see Figure 26. Specifically, there are two failure scenarios that may occur in the MHP protocol: queue mismatch error (Figure 26a)—where the message consistency check fails at the midpoint—and single-sided transmission error (Figure 26b).



**Figure 25: Timeline of the MHP message exchange with a successful reply by the heralding station;  $p$  is a photon associated with the GEN message. For a definition of GEN and REPLY message refer to Figure 27 and Figure 28, respectively.**

*E.2.2 MHP Packet Formats.* MHP relies on the exchange of the packets listed in the MHP sequence diagrams, see Figure 25 and Figure 26: GEN and REPLY.

GEN packet (Figure 27) is used by the midpoint to determine whether the nodes are consistent in their local information regarding their knowledge of the attempt at entanglement.

REPLY packet (Figure 28) is sent by the midpoint in the case of no error. It will include the senders' submitted absolute queue ID (i.e. QID and QSEQ) and additionally pass on the submitted queue ID of the peer node (i.e. QIDP and QSEQP). The sequence number, SEQ, denotes the number of successful heralded entanglement generations that have occurred at the

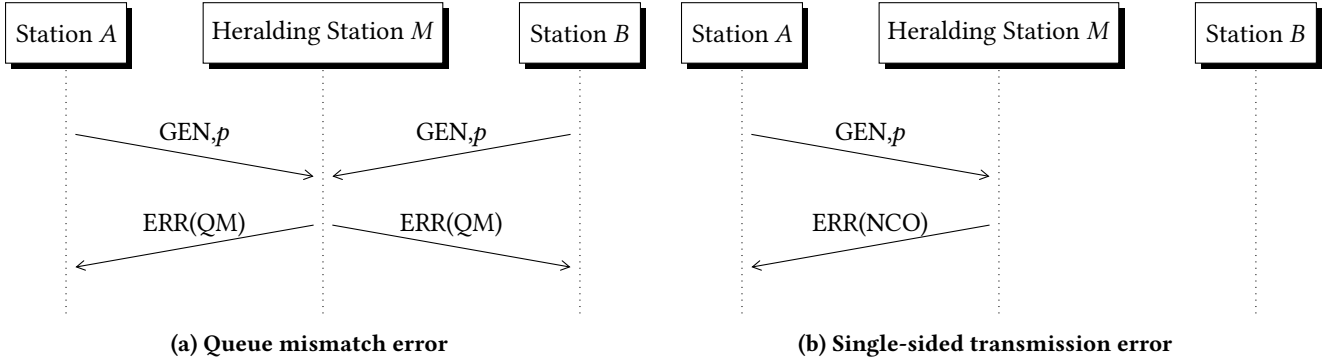


Figure 26: Timeline of two types of errors within MHP. For a definition of GEN and REPLY message refer to Figure 27 and Figure 28, respectively. QM and NCO refer to specific fields of the REPLY message (i.e. OT field), i.e. QUEUE\_MISMATCH and NO\_MESSAGE\_OTHER, respectively; both error types are explained in Protocol E.2.

midpoint heralding station and allows the end nodes to keep track of the number of entangled pairs that have been generated. OT encodes the heralding signal from the midpoint upon successful operation and encodes errors in case of failures.

|        |   |   |   |   |   |   |   |   |   |    |    |     |    |    |    |      |    |    |    |            |    |    |    |    |    |    |    |    |    |    |    |
|--------|---|---|---|---|---|---|---|---|---|----|----|-----|----|----|----|------|----|----|----|------------|----|----|----|----|----|----|----|----|----|----|----|
| 0      | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12  | 13 | 14 | 15 | 16   | 17 | 18 | 19 | 20         | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| Header |   |   |   |   |   |   |   |   |   |    |    | QID |    |    |    | QSEQ |    |    |    | (reserved) |    |    |    |    |    |    |    |    |    |    |    |

Figure 27: GEN packet format (used in MHP) sent by end stations to heralding station (midpoint). The pair (QID, QSEQ) represents the absolute queue ID, in other words they map to  $(j, i_j)$ —see Section E.1.1.

|        |   |   |   |      |   |   |   |      |   |       |    |            |    |     |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
|--------|---|---|---|------|---|---|---|------|---|-------|----|------------|----|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0      | 1 | 2 | 3 | 4    | 5 | 6 | 7 | 8    | 9 | 10    | 11 | 12         | 13 | 14  | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| Header |   |   |   |      |   |   |   |      |   |       |    | OT         |    | SEQ |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| QID    |   |   |   | QSEQ |   |   |   | QIDP |   | QSEQP |    | (reserved) |    |     |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |

Figure 28: REPLY/ERR packet format (used in MHP) for replies by midpoint with no error. OT: outcome/error reported by midpoint; SEQ: sequence number; QIDP: QID Peer; QSEQP: QSEQ Peer; QID and QSEQ are defined the same as in for GEN message—see Figure 27. Error codes include (001: QUEUE\_MISMATCH; 010: TIME\_MISMATCH, 100: NO\_MESSAGE\_OTHER, (refer to Protocol E.2 for the above error description).

### E.3 Entanglement Generation Protocol

The role of the Entanglement Generation Protocol (EGP) is to produce the required entanglement between two end nodes or otherwise declare failure.

*E.3.1 Entanglement Generation Scheduler.* We now proceed with the description of the scheduler—refer to Protocol E.3.1 for details. The EGP scheduler fulfills the following arbitrage functions, where we remark that for CREATE requests that demand multiple EPR pairs, only one request is added to the queue, and hence NEXT (function to select the next request from the local set of queues, see below) will return multiple pairs to be produced for the same request when called successively.

- GET\_QUEUE(creq): Once a request has been submitted, GET\_QUEUE deterministically chooses which queue  $Q_j$  to assign the CREATE request creq to. This may depend on the details of the request, such as for example  $t_{\max}$ , or  $F_{\min}$  as well as the purpose ID and priority.
- NEXT: Selects the next request from the local set of queues  $Q$  to serve, if any. Specifically, NEXT will determine:
  - Flag, set to true when a request is ready to be served;
  - Absolute queue ID (and corresponding request details) of request to be served;
  - Parameters to use in the MHP depending on the number of type of outstanding requests;
  - Communication and storage qubits, determined in cooperation with QMM.

---

**Protocol 2 EGP - Node A (B analogous exchanging A and B)**


---

*Definition of functions and variables.*

- Node ID: the ID of the peer we want to create entanglement with;
- $n$ : the number of entangled pairs we wish to create;
- $F_{\min}$ : the minimum acceptable fidelity required for the generated pairs;
- $t_{\max}$ : the maximum amount of time the higher layer is willing to wait for the entanglement to be created
- Purpose ID: port/application ID the requested this entanglement—used for forwarding OK messages to the appropriate application;
- priority: The priority of the request;
- $\text{seq}_{\text{expected}}$ : The expected sequence number from the midpoint. Initially set to 1;
- $j = \text{GET\_QUEUE}(c_r)$ : The call to the scheduler to obtain the queue ID of  $Q_j$  where the request should be placed in the distributed queue.
- $(i_j, ok) = \text{ADD}(j, c_{req})$ : The call to the distributed queue to add the create request  $c_r$  to  $Q_j$ .  $i_j$  is the unique ID of the request within  $Q_j$  and  $ok$  is a status code of performing the ADD to the distributed queue. Can take the values success (item added), timeout (communication timeout with peer while adding), or reject (the peer rejected adding the item to the queue);
- $\text{ERR\_NOTIME}$ : Error issued to higher layers by the EGP upon receiving  $\text{OK} = \text{timeout}$  from adding item to queue;
- $\text{ERR\_REJECT}$ : Error issued to higher layers by the EGP upon receiving  $\text{OK} = \text{reject}$  from adding item to queue;
- Trigger pair: equivalent to the POLLEGP call within the MHP protocol outline;
- $(\text{flag}, (j, i_j), \text{PSEQ}, \text{params}) = \text{NEXT}()$ : The call to the scheduler to obtain information for the next entanglement generation attempt where  $\text{flag} = \text{True/False}$  indicates whether entanglement should be attempted (same in MHP outline),  $(j, i_j)$  is the absolute queue ID of the create request being served, PSEQ is a pulse sequence identifier encoding the communication and storage qubit information to use for entanglement attempts, and  $\text{params}$  encodes the parameters to use for entanglement attempts (same as  $\text{params}$  in MHP outline);
- $\text{proto\_err}$ : Status of the attempt at entanglement, encodes the  $\text{mhp}_{\text{err}}$  from the MHP outline if an error occurred, 0 if no errors happened;
- $\text{create\_time}$ : timestamp of when the entanglement was generated;
- $F_{\text{est}}$ : the goodness passed in the frame, the estimate of the fidelity of the entangled qubits;
- $\text{logical\_id}$ : The storage qubit ID where the entangled qubit is stored, for use by higher layers;
- $t_{\text{Goodness}}$  - A timestamp of when  $F_{\text{est}}$  was record.
- $t_{\text{Create}}$  - A timestamp of when the entanglement was created.
- $k$ : the number of pairs left to generate for the request.

**(1) Adding to Queue:**

- (a) Ask scheduler which queue this request should be added to:  $j = \text{GET\_QUEUE}(c_r)$ .
- (b) Try to add request to the queue using the DQP:  $(i_j, ok) = \text{ADD}(j, c_r)$ .
- (c) If  $ok = \text{timeout}$  error, issue  $\text{ERR\_NOTIME}$  and stop.
- (d) If  $ok = \text{reject}$  error, issue  $\text{ERR\_REJECTED}$  and stop.
- (e) Otherwise the request has been added to the Distributed Queue.

**(2) Trigger pair (polled by MHP):**

- (a) Ask the scheduler whether entanglement should be made and for which request:  $\text{NEXT} = (\text{flag}, (j, i_j) \equiv \text{req}, \text{param}, \text{PSEQ})$ . For this end, the scheduler will employ its priority policy, as well as perform flow control depending on whether  $B$  is likely to produce entanglement.
  - (b) If there is a generation waiting to be satisfied:
    - (i) Construct response for MHP  $\text{POLLEGP}() = (\text{flag}, (j, i_j), \text{PSEQ}, \text{params})$ .
    - (ii) Provide the response to the MHP.
  - (c) Otherwise if there are no generations to perform provide  $\text{POLLEGP}() = (\text{flag} = \text{False}, \text{aID}_{\text{local}} = \text{null}, \text{PSEQ} = \text{null}, \text{params} = \text{null})$ .
-



**Protocol 2** (cont.) EGP - Node A (*B* analogous exchanging A and B)

---

**(3) Handle reply (message from MHP):**

- (a) Retrieve message from MHP including: result of generation  $r \in \{0, 1, 2\}$ ,  $\text{seq}_{\text{MHP}}$ , absolute queue id  $(j, i_j)$ , and protocol error flag  $\text{proto\_err}$ .
  - (b) If the specified absolute queue ID is not found locally then the request may have timed out or expired:
    - (i) Free the reserved communication/storage qubit in the Quantum Memory Manager.
    - (ii) Update  $\text{seq}_{\text{expected}}$  with  $\text{seq}_{\text{MHP}}+1$  and stop handling reply.
  - (c) Otherwise there is an absolute queue ID included in the response that is associated with an active request:
    - (i) If  $\text{proto\_err}$  is not OK, a non quantum error occurred and no entanglement was produced. Update expected sequence number with  $\text{seq}_{\text{MHP}}$  and stop handling reply.
    - (ii) If  $r = 0$  we failed to produce entanglement, stop handling reply.
    - (iii) Process  $\text{seq}_{\text{MHP}}$ :
      - (A) If the  $\text{seq}_{\text{MHP}}$  is larger than the expected sequence number, (partially) ERR\_EXPIRE the request to higher layer and send EXPIRE message to peer. Stop handling reply.
      - (B) Else, if  $\text{seq}_{\text{MHP}}$  is smaller than the expected sequence number, ignore reply. Stop handling reply.
      - (C) Else, update next expected  $\text{seq}_{\text{MHP}}$  sequence number (increment current one modulo  $2^{16}$ ).
    - (iv) A pair is established. If  $r = 2$  and we are the origin of this request, apply correction information to transform state  $|\Psi^-\rangle$  to state  $|\Psi^+\rangle$ , if we are the peer then we instruct the scheduler to suspend subsequent generation attempts until we believe the request originator has completed correction.
    - (v) Look up queue item  $(j, i_j)$ :
      - (A) If  $\text{create\_time} + t_{\text{max}} >$  current time or the request is not stored locally anymore, issue ERR\_TIMEOUT to higher layer and remove item from queue. Stop handling reply.
      - (B) Get fidelity estimate  $F_{\text{est}}$  from Fidelity Estimation Unit.
      - (C) Issue OK with Entanglement ID:  $(A, B, \text{seq}_{\text{MHP}})$ , logical\_id, Goodness  $F_{\text{est}}$ , and  $t_{\text{Goodness}} = t_{\text{Create}} = \text{now}$ .
      - (D) If  $k = 1$ , i.e. this was the last pair to be produced for this request, remove item from queue.
      - (E) If  $k > 2$ , decrement  $k$  on queue item.
- 

*E.3.2 EGP Sequence Diagrams.* We proceed with the introduction of all message passing sequence diagrams for the EGP.

Figure 29 presents a sequence diagram for the EGP operation when performing emission multiplexing. In some cases (such as the CM use case) it is not necessary to wait for the REPLY message from the midpoint before attempting entanglement generation again if one simply desires to generate correlated bit streams.

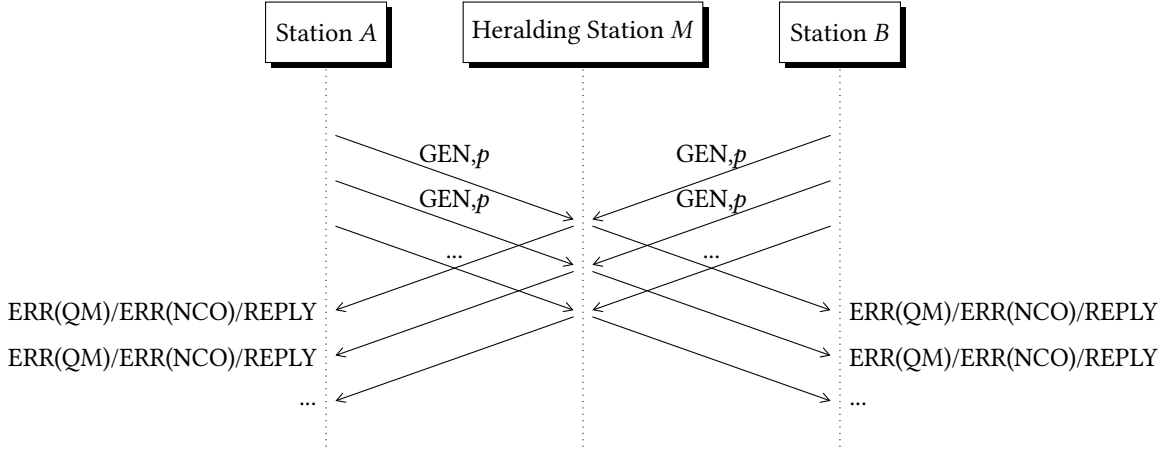
Figure 30 shows a sequence diagram detailing the message flow should station A obtain information that it is behind as well as a timeline of the message exchange when EGP processes at both nodes exchange available quantum memory information. Sharing this information allows both nodes to know whether there are available resources in order to proceed with satisfying an entanglement request. In the absence of resources of either peer there is no use in photon emission. Simply, both nodes must be able to emit photons for the protocol to operate properly.

Imperfect message transmission may cause any of the GEN, REPLY, EXPIRE, REQ(E), and ACK messages to become lost or corrupted in transit between nodes. Depending on which messages are lost in the protocol, different actions are taken to prevent deadlock. A lost EXPIRE or corresponding ACK results in a retransmission of the EXPIRE to ensure that OK messages are properly revoked at the peer node. Loss of a REQ(E) message or its corresponding ACK results in a retransmission of the REQ(E) to make sure that both nodes have up-to-date information of available resources.

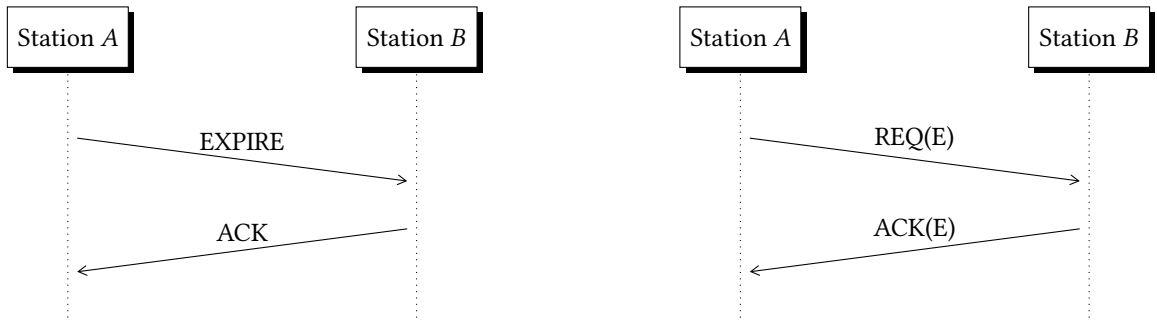
Losing a GEN message is handled by the midpoint heralding station when only one GEN message arrives from an end node. In this case the REPLY message containing NO\_CLASSICAL\_OTHER (see Protocol E.2) is issued to alert the nodes of the failure. In this case no attempt at entanglement is made and the sequence number at the midpoint remains the same.

Losing a REPLY message from the midpoint that contains an outcome of 0 has no impact on A or B as  $\text{seq}_{\text{expected}}$  is only updated when a successful attempt at entanglement occurs. When a REPLY message containing an outcome of 1 or 2 (for successful entanglement) is lost, the end node(s) that lost the message will continue attempting entanglement generation in subsequent polls by the MHP as there are outstanding pairs to be generated for the request. Upon successful receipt of any

message from the midpoint (REPLY), the included SEQ will be ahead of the receiving node(s)  $seq_{expected}$  and the loss will be detected. The detecting node will then transmit an EXPIRE message to its peer containing the old  $seq_{expected}$  that did not agree with the SEQ received from the midpoint along with its new  $seq_{expected}$ , so that any OK messages containing the missing set of sequence numbers are revoked at the peer.

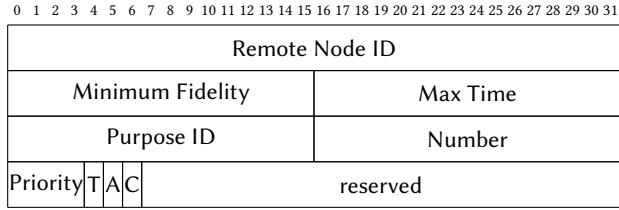


**Figure 29: Timeline of multiplexing photon emission in the MHP: multiple  $GEN,p$  messages are sent one by one for which any reply messages (ERR(QM), ERR(NCO) or REPLY) is possible. Compare this with the sequence diagram presented in Figure 26 and Figure 25 for the MHP.**

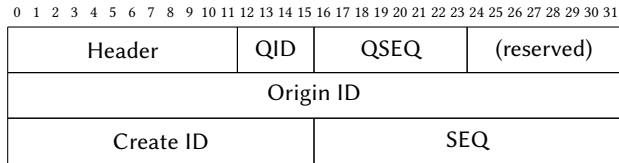


**Figure 30: (Left) Timeline of request expiration within EGP. Definitions of EXPIRE and ACK message are given in Figure 32 and Figure 33, respectively. (Right) Timeline of memory advertisement requests within EGP. Definition of REQ(E) and ACK(E) message is given in Figure 34.**

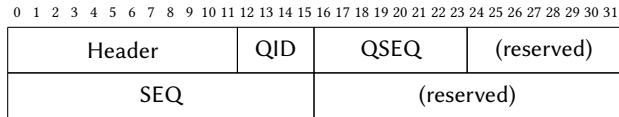
*E.3.3 EGP Packet Formats.* Finally, we present the definitions of all messages being used by the EGP. Figure 35 describes the information passed from the EGP to the MHP during each periodic cycle, while Figure 36 shows the packet format for replies from the physical layer MHP to the EGP. Figure 32 and Figure 33, define EXPIRE and ACK messages, respectively, exchanged in entanglement request expiration sequence diagram described in Figure 30. Figure 34 shows the REQ(E) and ACK(E) packet formats exchanged by memory advertisement requests made by the EGP sequence diagram described in Figure 30. Figure 37 and Figure 38 present the format of OK messages passed from the EGP to higher layers, in case of *create and keep request* and *measure directly request*, respectively.



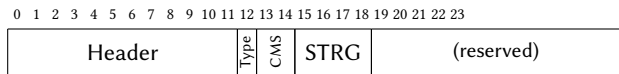
**Figure 31: Packet format for CREATE message to EGP. Explanation of the message fields—Remote Node ID: Used if the node is directly connected to multiple nodes. Indicates which node to generate entanglement with; Minimum Fidelity—The desired minimum fidelity, between 0 and 1, of the generated entangled pair; Max Time—The maximum time in seconds the higher layer is willing to wait for the request to be fulfilled; Purpose ID—Allows the higher layer to tag the request for a specific purpose; Number—The number of entangled pairs to generate; Priority—Can be used to indicate if this request is of high priority and should ideally be fulfilled early; T—the type of request, Either create and keep (K) or measure directly (M), where K stores the generated entanglement in memory and M measures the entanglement directly; A—atomic flag, indicates that the request should be satisfied as a whole, i.e. that all entangled pairs are available in memory at the same time; C—consecutive flag, indicates that the entangled pairs of the request should be generated close in time.**



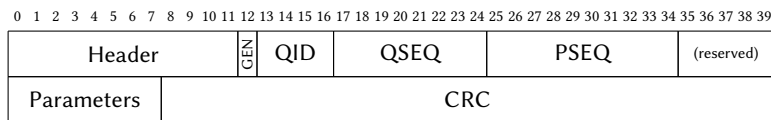
**Figure 32: Packet format for EXPIRE message. Explanation of the message fields—Origin ID: ID of the node the request originated from; Create ID: creation ID of the request. SEQ: up-to-date expected MHP sequence number at the node the EXPIRE originates from. Recall that (QID, QSEQ) represent the absolute queue ID.**



**Figure 33: Packet format for ACK message. Recall that (QID, QSEQ) is the absolute queue ID and SEQ is the acknowledger’s up-to-date expected MHP sequence number (the same as in the case of EXPIRE message, see Figure 32).**



**Figure 34: REQ(E)/ACK(E) Packet format for EGP memory requests. Explanation of the message fields—Type: message type (0: REQ(E), 1: ACK); CMS: the number of available communication qubits, STRG: the number of available storage qubits.**



**Figure 35: Packet format of POLLEGP messages sent from EGP to MHP. Explanation of the message fields—GEN: emit photon flag; PSEQ: pulse sequence identifier which instructs the underlying quantum communication device of the parameters to use when emitting the photon. See figure 5.**

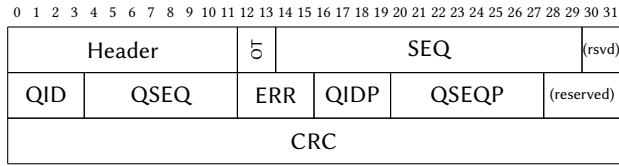


Figure 36: Packet format of messages from MHP to EGP. Explanation of the message fields—OT: measurement outcome (the same field as in REPLY message define in Figure 28). Error codes in ERR field encode the errors described in the MHP protocol (refer to Protocol E.2). See figure 5.

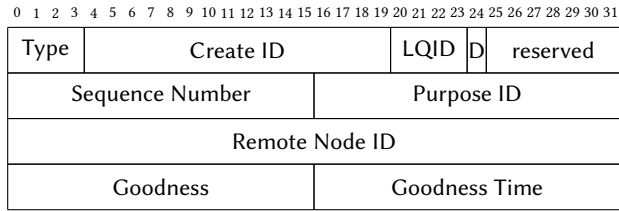


Figure 37: Packet format for OK message corresponding to a create and keep request. Explanation of the message fields—Type: Indicates that this is a create and keep OK; Create ID: The same as the Create ID returned to the requester; LQID: Logical Qubit ID where the entanglement is stored; D: Directionality flag indicating the source of the request; Sequence Number: A sequence number for identifying the entangled pair; Purpose ID: The purpose ID of the request; Remove Node ID: Used if connected to multiple nodes; Goodness: An estimate of the fidelity of the generated pair; Goodness Time: Time of the goodness estimate. See figure 5.

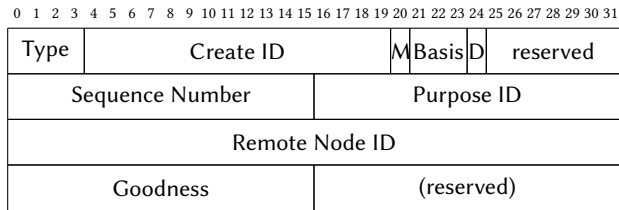


Figure 38: Packet format for OK message corresponding to a measure directly request. Explanation of the message fields—Type: Indicates this is an OK for a measure directly request; Create ID: The same Create ID given to the requester; M: Outcome of the measurement performed on the entangled pair; Basis: Which basis the entangled pair was measured in, used if the basis is random. The remainder of the fields are explained in Figure 37. See figure 5.

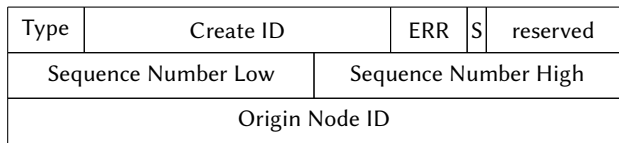


Figure 39: Packet format for ERR messages containing errors from EGP. Explanation of the message fields—Type: Indicates this is an ERR message; ERR: The error that occurred in the EGP, may be any of the errors described in Section E.3; S: Used by ERR\_EXPIRE, to specify whether a range of sequence numbers should be expired (S=1) or all sequence numbers associated with the given Create ID and Origin Node (S=0); Sequence Number Low/High: Use together to specify a range of sequence numbers to expire. The remaining fields are explained in figure 37.