

Entanglement Generation Protocol: Notes on Link+Physical Layer

Stephanie, Axel, Matthew, Leon, Erwin, Ronald

May 29, 2018

The objective of this document is to define the link layer in quantum networks connecting quantum processing nodes, and to propose two concrete link layer protocols based on an existing implementation of the physical layer with certain properties. In analogy to classical networks, the objective of the link layer will be to enable communication between two nodes A and B connected by a *link* on the same network. Here, enabling communication corresponds to producing entanglement between A and B , and we will hence refer to such protocols as Entanglement Generation Protocols (EGP). We propose the desired service, interface to the higher layer, as well as two possible EGPs that are closely related design alternatives. We first discuss an EGP between two nodes A and B , and discuss extensions to a proposed architecture connecting many nodes at the end.

To fit the link layer EGP into the future envisioned network stack we briefly sketch the stack framework here, going from higher to lower layer:

QTP - Qubit transport protocol (Transport Layer) Responsible for the end to end transmission of qubits.

EMP - Entanglement Management Protocol (Network Layer) Responsible for the generation of entanglement between two nodes that are not directly connected by a link, i.e. not on the same local network.

EGP - Entanglement Generation Protocol (Link Layer) Responsible for the generation of entanglement between to nodes connect by a direct link.

1 Entanglement Generation Protocols

Let us first describe the interface, service, as well as performance criteria of entanglement generation protocols.

1.1 Higher layer to EGP

An EGP supports a single command from the higher layer, namely a request to produce entanglement, which we call a CREATE command. This command includes some desired properties of the entanglement such as for example a minimum fidelity, and a maximum waiting time. In an actual physical implementation, there is a tradeoff between these parameters. More time, for example, may allow the underlying implementation to use entanglement distillation to produce higher quality pairs.

CREATE Produce entanglement with a node on the same network (i.e. connected by a link). Arguments supplied are:

Partner ID	ID of the node to generate entanglement with.
Number k	Number of pairs we want to create.
F_{\min}	Minimum acceptable fidelity (with high confidence).
t_{\max}	Maximum acceptable waiting time before request is completed.
Purpose ID	Identifying the purpose or application at this node (optional, default 0).
Priority	Manual setting of a priority for entanglement production (optional).
create ID	Sequence number identifying this CREATE command.

1.2 EGP to higher layer

Following the reception of the CREATE command, several actions of the EGP are possible. Let us start with the positive outcome, and then consider possible errors.

OK Entangled pair has successfully been produced deterministically (heralded). One message per pair created, delivered immediately (best effort) following pair creation. With high confidence, the minimum acceptable fidelity F_{\min} has been met, and the entanglement has been generated within the specified time frame t_{\max} . Information about the entanglement generation is provided, including an entanglement identifier. This identifier is required to be globally unique, and agreed upon by A and B . That is, A and B can locally use this entanglement identifier to determine which of their qubits is entangled with the remote node, and also which qubit belongs to which entangled pair. Entanglement identifiers are meant to be shared in the network by higher layer protocols and carry meaning beyond the nodes A and B . An entanglement identifier (Ent_{ID}) consists of:

(Node A ID, Node B ID)	IDs of the two nodes between which this entanglement is shared.
seqID	Sequence number. Unique (up to wrap around) between A and B , and globally unique when combined with the node IDs.
Goodness	Heuristic estimate for the fidelity of the generated pair.
t_{Goodness}	Time when this goodness was established (in EGP, usually the same as generation time).
t_{Create}	Time the pair was produced.

In addition the OK message also includes the following local information. We remark that Qubit IDs are exclusively local information (akin to the memory address in a computer) and not in general shared between network nodes.

Qubit ID Logical Qubit ID of the entangled pair can locally be found.

Entanglement generation may fail for wide number of reasons, some of which form an immediate error. It may also be that the entanglement later expires, or is discarded of which the EGP will inform the higher layer. Let us start by listing the immediate failure modes, where in all such cases the create ID will be included allowing the higher layer to identify which request has failed.

ERR_UNSUPP Operation not supported. For example, creation of entanglement with the specified minimum fidelity is unattainable, or unattainable within the given time frame, even if the node is not loaded.

ERR_NOTIME Cannot meet the desired fidelity demand within the given time frame due to high load.

ERR_NORES No resources (such as qubits to store the entanglement) available.

ERR_TIMEOUT Failure to produce entanglement within the specified time frame.

ERR_OTHER Failure for unspecified reasons, such as hardware failures.

In addition, the following failure mode can occur later when an entangled pair is expired. The primary use case of this will be to deal with extremely improbable failures in which recognition of the failure only becomes available after the higher layer has already received an OK message. This allows for a tradeoff between speed and certainty in recognizing failure modes. Since entanglement is very short lived, increased certainty can if desired be sacrificed for speed.

EXPIRE Expire Qubit ID. Any entanglement associated with Qubit ID has become unavailable.

1.2.1 Questions

- The term "High confidence" is not defined and we need to decide what we mean by that, and also if this is some parameter where/by whom it is determined.

1.3 Performance metrics

Apart from correctly fulfilling requests, a variety of performance metrics can be considered for EGPs. Not all of these can be simultaneously optimized, but occasionally impose tradeoffs. We hereby also draw a distinction between performance metrics of interest to a specific “user” requesting entanglement from the EGP, and the overall performance of the network. Evidently, for all metrics below adverage, variance, and worst case behaviour is of interest. Once more data is available on how quantum networks are used in practise, one may also consider “typical” values for these metrics.

Let us first consider “user” centric metrics, measuring the experience of one invidual user rather than a behaviour of the network as a whole. We remark that nevertheless these metrics are a consequence of the total usage of the network.

Fidelity Quality of the entanglement produces. By design the fidelity has to exceed the minium requested fidelity F_{\min} .

Latency Time between submission of a CREATE request, and an OK response when successful. By design this time may not exceed t_{\max} .

In addition, we can consider measures defined by the behaviour of the network when dealing with a large number of requests.

Throughput Number of pairs/s. Refined variants of throughput to be measured include: instantaneous throughput and sustained throughput.

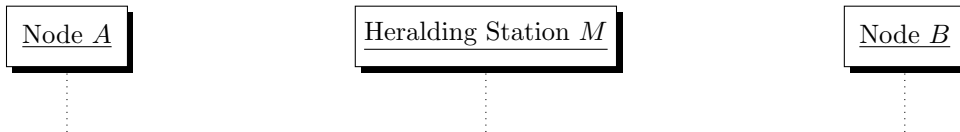
Fairness Difference in performance metrics between requests originating at A and B .

Availability Availability is a concern here if a network node requires resetting and two nodes require resynchronization at certain time intervals.

We remark that measured values like throughput evidently depend on the request behaviour, including what we will call the *request ratio*, i.e. the number of pairs requested/number of requests total.

2 EGPs based on midpoint heralding protocols (MHPs)

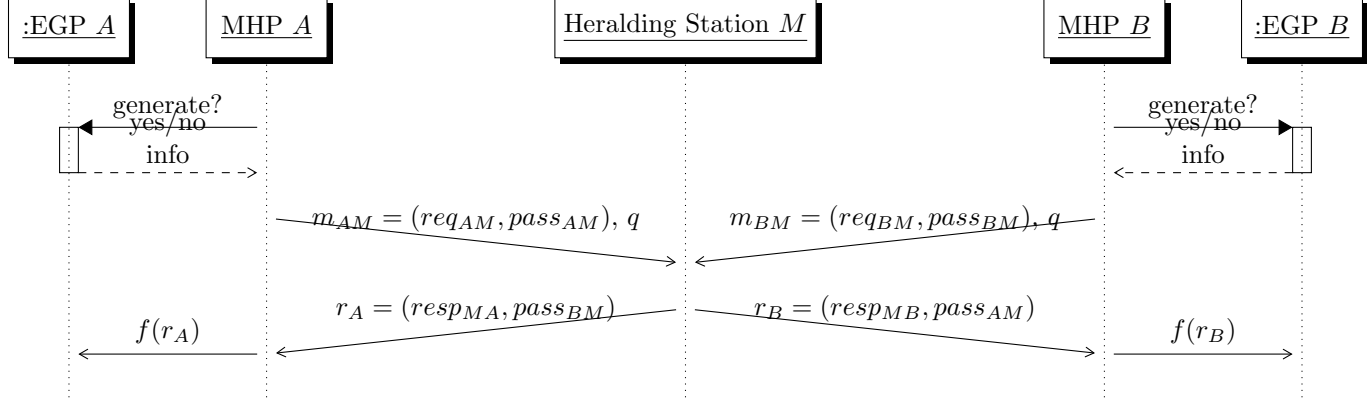
Before proposing a specific EGP, let us first consider a general class of EGPs that are build on top of a system supporting heralded entanglement generation at the physical layer. More precisely, we will consider physical layer protocols that produce entanglement between two nodes A and B , by means of an heralding station M between them, in the following configuration:



Several variations of such schemes exist, such as the single-click, the Barrket-Kok (BK) protocol or even a memory assisted protocol in which entanglement distillation is performed. For simplicity, we will assume a single-click/BK type scheme below, but the following also applies to memory assisted schemes with minor modifications. Evidently, the choice of the physical layer entanglement generation scheme effects the overall performance metrics stated above, as well as the possibility to service requests for entanglement above a specific fidelity F_{\min} . For example, certain fidelities may only be attainable by performing entanglement distillation.

Nevertheless, we may cast such physical layer generation schemes in the same following abstract form, upon which our EGPs will be built. We remark that MHPs in our network stack make no high level decisions on when to produce entanglement, scheduling, etc. Also the decision on which MHP to use - eg whether to perform single click, or distill - or even what parameters to set in the single click protocol (such as the bright state population α) are not part of the MHP, but such decisions are left to the link layer EGP, that will use the appropriate MHP (parameters) to obtain the desired service from the physical layer.

Let us thus first give a very abstract description of such protocols, before specifying some assumptions and design alternatives, and stating desired requirements. Here, we have subdivided node into the different units, EGP and MHP.



As a simple example, consider the single click protocol with some fixed parameters. Here, m_{AM}, m_{BM} are empty, $resp_{MA}, resp_{MB} \in \{OK, FAIL\}$. Some assumptions and choices were made in the above description:

- Assumptions**
1. An (essentially) instantaneous association between the classical control messages m below, and the entanglement generation. For this reason, we will write a classical message transmission as simply m , and q for arbitrary quantum signal q , and assume simultaneous arrival. This could be realized by forming a timing association between the classical and quantum signals. To make it clear, how we will use this abstract description later, we will always take $m = (req, pass)$ where req is request information only for the mid point and later protocol specific, and $pass$ is something that will by default always be passed onto the other side (also protocol specific).
 2. Midpoint will provide a response $resp$, which includes at least the following information to both A and B : (1) Success or failure of entanglement generation. (2) In case different types of states can be produced, the identity of the state made. (3) A sequence number or equivalent information that allows A and B to determine the ordering of generated pairs.

2.1 Design considerations

Before considering possible protocols, let us abstract some design considerations resulting from the implementation for the non quantum reader:

Basic facts Nodes A and B in the considered implementation are few qubit processors capable of storing and manipulating qubits. Entanglement can be produced probabilistically between one of these qubits, and a travelling photon (q in the above cartoon). When photons from both side arrive simultaneously at the heralding station, a measurement is performed. This measurement will produce 2 possible entangled states with a certain probability p_{succ} , which we will call states 1 and 2 below. This constitutes successful generation. The two types of states can be converted into each other by applying local gates at either A (or B) and thus both states can be considered equally valuable. The measurement at the heralding station can also fail, in which case another attempt must be made. Typical values are $p_{succ} = 1/2$. The information about success - incl. whether we have state 1 or 2 - or failure is only available at the heralding station, until it is transmitted to A and B .

Not all qubits are created equal, in the sense that not all of them can be directly entangled with a traveling photon. We will refer to those that can as communication qubits, and the other as storage qubits. One can transfer the state of a communication qubit to a storage qubit at any time (where we remark that this operation evidently also costs time and introduces noise). In NV in diamond, there is one communication qubits (electron spin) and several storage qubits (called nuclear spins or carbon spins).

Triggering generation Generation of entanglement requires a time synchronized trigger at both A and B that will result in the probabilistic creation of entanglement between the electron spin, and the photon (q in the cartoon above) traveling to the midpoint. If the trigger only occurs at one node, no entanglement can be made. Agreement to produce entanglement at specific time steps thus already has to be realized, requiring some communication ahead of time.

Here, we will assume that all low level timing and other synchronization is implemented in the MHP, which is then able to produce entanglement at certain “pre-agreed” i.e. synchronized time instances without additional communication between A and B . As such, the EGP only performs higher level processing. This motivates the choice above that the MHP will poll the EGP, e.g. by reading a specific state variable, on whether entanglement generation is required at a specific synchronized time step. This is in contrast to the EGP sending a signal to the MHP to produce a pair. Since the EGP does not deal with timing synchronization, it cannot know when the actual physical trigger should be produced and hence the MHP could then only save the request until pair production is timely. We remark that this would mean that the MHP would need to keep state of how many outstanding triggers there are, which is not desirable from a design point of view where if t_{\max} has elapsed the EGP may no longer want generation by the MHP. Consequently, we here choose for the MHP to poll the EGP, which does have state on desired generation.

Noise due to generation One may wonder, why entanglement generation is not enabled continuously. That is, attempts to produce entanglement are made all the time, and then the entanglement is either discarded or directly available to the EGP. Two reasons motivated by the physical implementation considered (NV in diamond) make this undesirable: First, triggering an attempt to produce entanglement causes additional noise on the storage qubits. This means that the storage is significantly degraded by trying to produce further entanglement. As such, it is desirable that triggers to attempt generation are only made whenever entanglement is indeed wanted. Second, there are only a small number of storage qubits (presently, 4). If we produce entanglement quickly, by for example triggering an attempt and then immediately transferring the state of the communication qubit to the storage qubit and then proceeding with the next attempt before having heard back from the heralding station, then several storage qubits are needed to support this, making the memory unavailable to other purposes.

For these reasons, MHP will inquire with EGP whether more entanglement is desired, and only then commence production.

Scheduling and flow control The EGP will be responsible for all higher level logic, including scheduling requests. A form of scheduling is flow control which controls the speed of entanglement generation, which hence also falls under higher level logic (see questions in EGP section below).

Memory allocation Decisions on which qubits to use for what purpose lies in the domain of higher level logic, where more information is available on - for example - the number of outstanding requests allowing scheduling decisions including the scheduling of how memory is best used. MHP will hence also not perform memory allocation, i.e., determine which communication qubits or storage qubits to use.

This impacts the types of information included in “info” in the protocol above, which we later take to be of the form (Physical ID Communication Qubit, Physical ID Storage Qubit, Parameters)

2.2 Sending classical messages

Above we assumed that there exists a means to transmit classical data between A , B and M . How this is realized is not the objective of this document, and it could be achieved both by a dedicated fiber (possibly using two wavelength for bidirectional communication), or interspersed with quantum signals on the same fiber.

Here, of interest are merely standard numbers - and the system will need to be implemented to ensure a quality that yields good performance in our link layer protocol. We hence for now consider merely standard numbers:

- Classical channels are bidirectional, meaning data could in principle be sent in both direction at the same time (and, as a relevant consequence, messages can cross and are not ordered in both directions)
- Likelihood of losses: p_{loss} probability of loss (e.g. following standard fiber loss plus electronics if applicable).
- Likelihood of errors: p_{err} probability of error - where we remark that as in other classical communication burst errors are probably dominant.
- Standard delays of interest: propagation delay (over the fiber), transmission delay (incl. delays of the electronics in putting the packets on the fiber), and processing delay, if known. We will assume that given the highly sophisticated electronics and the fact that the rate of classical communication is low due the relatively low repetition rate of entanglement generation attempts , the transmission and processing delay are essentially negligible.

2.2.1 Enhanced situation

Two standard methods exist to enhance this situation to the following, whose exact form and choice depends on the parameters above. We may also consider running an existing link layer protocol, such as Ethernet over fiber - we do however note that authentication is highly desirable due to control messages leading to local spin photon entanglement generation, and hence otherwise allow unauthenticated messages to manipulate the state of the quantum processing node.

- Error detection: This can be achieved using standard methods, where probably a simple CRC depending on the length of the headers is most appropriate. This will add a number of bits to the messages headers below if employed. For example, for a standard CRC-32 as used in Ethernet, the CRC is computed over the message and stored in a 32 bit header field.
- Message authentication: In this case, this refers to the fact that A knows the messages originate with B (and vice versa). Similarly, M can authenticate messages if desired. Such authentication can be realized using a message authentication code (MAC) (see eg ?). These can be realize with varying levels of security. If A and B share consumable key (such as for example generated by QKD), they can afford to use one-time MAC which - similar to a one-time pad - offers information-theoretic security. Such MACS, for example based on two-universal hashing, can in principle be fast (see e.g. ? needing however various levels of key), although it is a question whether they are fast enough to be useful in this context. We remark that also weaker forms of authentication are acceptable as they merely form a means of protection and would need to be broken in a very short time frame to have a practical impact.

2.2.2 Enhanced situation

Based on the general shape of such protocols above, one can now consider a slight “enhancement” of a protocol of this form - like single-click - that makes explicit some (probably obvious) failure modes, and produces a total ordering of pairs that A and B agree upon, even if some messages may go missing.

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
Rest of header																															
Error detection CRC																															
Message authentication (MAC)																															

} To be filled later

3 Candidate EGPs

We will now consider three classes of link layer EGPs build on top of MHPs. Strictly speaking, we will consider two non trivial implementations, and one rather adhoc one with rather trivial features for comparison.

The present objective is to implement these in simulation to:

1. Assess their relative performance - also with respect to each other. Performance metrics have been specified above.
2. Assess the effect of specific design choices made in each of these protocols.
3. Study the relevance of parameter demands, specifically also the required quality of the classical communication and its timing so it can be implemented.

The main difference between the more sophisticated EGPs is where queuing and scheduling is done, demanding more or less power at the heralding station.

3.1 Node Centric: Distributed Queue

In the first scenario, the heralding station is essentially passive and performs no higher level functions. Before sketching the protocol, let us first outline its ingredients next to the MHP.

3.1.1 Priority Queues (PQ)

In principle, we may want to give priorities to certain requests. This will be accomplished by adding them to different types of queues. For simplicity, we will for now assume there is only one.

3.1.2 Distributed Queue Protocol (DQP)

The objective of the DQP is to obtain a shared - i.e. agreed upon queue - at both nodes. That is, both A and B hold a local queue, which is synchronized using the DQP. Specifically, any request to the EGP at node A or B to perform entanglement is placed into the local queues such that the following are satisfied:

- (Unique queue numbers) If a request is placed into the queue by either A or B, then it is assigned a unique queue number (modulo the maximum queue number).
- (Total agreed ordering) A and B agree on the ordering of their queue: If request with queue number j exists in the queue of A, then eventually (whp within the usual propagation delay between A and B) the item also exists in the queue of B with number j (and vice versa).
- (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) has added a maximum of `queue_window_size` many items.
- (Scheduling information) When items are added to the queue, they are marked ready if acks have been received (ie we can reasonably assume they are in both queues). In this case, the item receives a minimum to be executed time, corresponding to the expected propagation delay between A and B, decreasing the likelihood A or B wants to produce before the other node is ready as well. No penalty other than reduced performance due to increased decoherence results if one of them goes too early.

Given we have only two nodes, the above can easily be realized by one node being the master controller of the queue marshalling access to the queue, and the other the slave controller. (see `DistributedQueue` for implementation)

3.1.3 Scheduler

The scheduler fulfills a number of important arbitrating functions:

- Determines deterministically which queue to serve next. In the case there is only one queue, the top most item is returned if it is ready to be scheduled.

- Determines which parameters to use in the MHP depending on the number of type of outstanding requests.
- In cooperation with QMM determine storage qubits.
- Includes the flow controller below.

3.1.4 Quantum Memory Management (QMM)

While not part of the actual protocol, the QMM can be asked to (smartly) allocate a fresh qubit and to convert logical qubit IDs to physical Qubit IDs and vice versa.

3.1.5 Flow Control

We will perform a form of flow control, that is we will speed up or slow down the production of entanglement. This will be implemented in the simulation for now as a stub, for first testing to be filled in later. Speeding up or slowing down is relevant if the local quantum memory is full, i.e., we can no longer produce entanglement since insufficiently many free qubits are available. To do this we remark that by a standard Hoeffding bounds, A (or B) can likely produce entanglement except with probability ε , if the number of qubits A (or B) has free m_A (or m_B) satisfies

$$n \leq \frac{m_A}{p_{\text{succ}} + \sqrt{\frac{1}{2n} \ln \left(\frac{1}{\varepsilon} \right)}}, \quad (1)$$

where p_{succ} is the probability of successful entanglement generation and n is the number of qubits currently in flight (i.e. we do not yet know the outcome of entanglement generation). Below we will say that A (or analogously B) can likely produce entanglement, if the above expression is satisfied.

3.1.6 Fidelity estimation unit

We are planning to perform error detection by including later an online fidelity estimation protocol, that together with already available information from the experiment allows us to make a guess for the fidelity communicated to higher layers via the Goodness in the Entanglement ID.

May intersperse CREATE requests of its own to update its online fidelity estimate of the EGP.

3.1.7 Protocol Sketch

Let us now describe the local protocol running at each node, once a CREATE request is received. For clarity we will describe it for node A, but it applies to both analogously:

1. Ask the scheduler which queue the request should be added to. Call that Q .
2. Add the request to the Distributed Queue Q (ie run the DQP to get it added to the consistent local queue)
3. Ask the scheduler whether there are ready items for which entanglement should be made ?
4. Ask the flow controller whether B is likely to receive?
5. If the answer to both of these questions is yes:
 - (a) Ask the scheduler to fetch a storage ID from the QMM
 - (b) Turn on the entanglement generation flag, with physical ID electron spin, storage ID the one just obtained from QMM, and parameters in accordance with the desired fidelity TBD by scheduler.

Let us now describe how we will use the MHP, recall that it will be automatically executed at specific time steps and poll the EGP for the entanglement generation flag (or more accurately list of outstanding “flags”). Here, we will assume that in addition to the flag y/n whether to produce entanglement, we will also supply the MHP with (1) the physical ID of the communication qubit (2) the physical ID of the storage qubit (3) any parameters such as the bright state population for entanglement generation (4) the current free memory size m_A .

We specify the MHP by filling in the default messages given above. Again we will for simplicity do this only for Alice but the same holds for Bob. It is unlikely we want the acking here, but I’ll include it anyways.

1. Initialize sequence numbers to 0.
2. Initialize list to be acked LA.
3. Use $req_{AM} = (Produce)$ and $pass_{AM} = (m_A, optionalackgenerationsequencenumberonLA)$. If no entanglement is desired to produce, set $req_{AM} = (Info)$ and use the same pass to update the remote node with the new memory size if any.
4. Use $resp_{MA} = (r, MHPseq, info)$ 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 MHP seq is chosen by the heralding station and send to both A and B. The heralding station may send additional information (eg the moment of detection) to allow a more accurate fidelity guess.
5. (Optional) Resend from previous round
6. Choose f to pass the following information up to the EGP: Current free memory of the remote node, r , and MHP seq.
7. (optional) Save seq onto LA to be acknowledged as part of the next generation round. For those outstanding too long, also include that information in the message back to the EGP.

Let us now specify what EGP does once an event of generation is communicated back by the MHP:

1. If the sequence number is in the expected order, and $r \in \{1, 2\}$, ask the scheduler what request to serve. Ask the fidelity estimation unit for an estimated fidelity based also on the info provided by the heralding station. If the request asks for many pairs, fulfill at least one request immediately. Send OK up and update the request and queue accordingly.
2. If the sequence number is not in expected order, we lost a heralding signal. EXPIRE the request that this would have belonged to.
3. If we do acks in the MHP, then check for too long outstanding acks - retroactively EXPIRE those requests.
4. Update our knowledge of the remote nodes memory size, to make an assessment whether it is likely to receive in the future.

3.2 Heralding Station Centric: Arbitrage by heralding station

3.3 Happy go lucky