

Quantum Link Prediction in Complex Networks

João P. Moutinho^{1, 2, *}, André Melo³, Bruno Coutinho¹, István A. Kovács^{4, 5, 6}, and Yasser Omar^{1, 2, 7}

¹Physics of Information and Quantum Technologies Group, Instituto de Telecomunicações, Lisbon, Portugal

²Instituto Superior Técnico, Universidade de Lisboa, Lisbon, Portugal

³Kavli Institute of Nanoscience, Delft University of Technology, Delft, The Netherlands

⁴Department of Physics and Astronomy, Northwestern University, Evanston, IL, USA

⁵Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL, USA

⁶Central European University, Budapest, Hungary

⁷Portuguese Quantum Institute, Portugal

Abstract

Predicting new links in physical, biological, social, or technological networks has a significant scientific and societal impact. Network-based link prediction methods utilize topological patterns in a network to infer new or unobserved links. Here, we propose a quantum algorithm for link prediction, QLP, which uses quantum walks to infer unknown links based on even and odd length paths. By sampling new links from quantum measurements, QLP avoids the need to explicitly calculate all pairwise scores in the network. We study the complexity of QLP and discuss in which cases one may achieve a polynomial speedup over classical link prediction methods. Furthermore, tests with real-world datasets show that QLP is at least as precise as state-of-the-art classical link prediction methods, both in cross-validation tests and in the prediction of experimentally verified protein-protein interactions.

1 Introduction

From genes and proteins that govern our cellular function, to our everyday use of the Internet, Nature and our lives are surrounded by interconnected systems [1]. Network science aims to study these complex networks, and provides a powerful framework to understand their structure, function, dynamics, and growth. Studies in network science typically have a substantial computational component, borrowing tools from graph theory to extract relevant information about the underlying system. With the advent of quantum computation, a natural question to ask is which problems in network science can be explored with this new computing paradigm, and what benefits it can yield. This question can be interpreted in at least two different ways. First, there is a large body of work in quantum algorithms for graph theoretical problems, some examples being Refs. [2, 3, 4, 5], which may have their own applications in network science problems. However, network science algorithms often look for specific patterns or organizing principles based on empirical observations from the real underlying systems, which may warrant the development of problem-specific quantum algorithms. This constitutes a novel research direction, different from the development of more general graph-theoretical algorithms. Previous connections have been made between quantum phenomena and complex networks, both by using quantum tools to study complex networks [6, 7, 8, 9] and by using complex network tools to

*Electronic Address: joao.p.moutinho@tecnico.ulisboa.pt

study quantum systems [10]. Nevertheless, to our knowledge, potential quantum speedups for network science problems have not been addressed.

In this work, we propose a quantum algorithm for the problem of link prediction in complex networks based on Continuous-Time Quantum Walks (CTQW) [11, 12], and discuss potential quantum speedups over classical algorithms. The objective in link prediction is to identify unknown connections in a network [13, 14, 15, 16, 17]. For example, in social networks, we aim to predict which individuals will develop shared friendships, professional relations, exchange of goods and services or others [13, 14]. In biological networks, the main focus is the issue of data incompleteness, which hinders our understanding of complex biological function. For example, in protein-protein interaction (PPI) networks link prediction methods have already proven to be a valuable tool in mapping out the large amount of missing data [18, 19]. While there are many approaches to the problem of link prediction, such as using machine learning techniques [20] or studying global perturbations [21], other methods focus on simple topological features like paths of different length between nodes, which we describe next.

1.1 Classical Link Prediction

Network-based link prediction methods take as input a graph $G(\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the set of nodes with size $N = |\mathcal{V}|$ and \mathcal{E} is the set of undirected links, and output a matrix of predictions $P \in \mathbf{R}^{N \times N}$ where each entry p_{ij} is a score value quantifying the likelihood of a link existing between nodes i and j (see Figure 1). Each method computes P differently, depending on the assumptions made about the network and its emergent topological features. Most methods are based on the Triadic Closure Principle (TCP), assuming that two nodes are more likely to connect the more similar they are [17, 18]. Similarity is often quantified based on the number of shared connections, i.e., paths of length two between two nodes, or in general of even length. It has been shown that, despite its dominant use in biological networks, the TCP approach is not valid for most protein pairs [18]. Instead, in [18], a link prediction method (L3) is proposed without the assumption that node similarity correlates with connectivity. L3 is based on the assumption that a candidate partner is similar to the existing partners of a node, $P = AS$, as illustrated in Figure 1 c). These results [18], and follow-up studies [22, 23, 24], show that the L3 method significantly outperforms other link prediction methods, for example, in some complementarity driven networks.

Our quantum approach takes inspiration from both these paradigms, utilizing even-length (TCP) and odd-length (L3-like) paths. One of the main reasons why link prediction may prove suitable to be tackled with a quantum computer is the realisation that in practice we are not interested in knowing the scores of all pairs of nodes, but we simply wish to know which ones have the highest score up to a certain cut-off threshold, as illustrated in Figure 1 d) and e). By encoding the prediction scores in the amplitudes of a quantum superposition and performing quantum measurements on the system, the predictions with the highest score will be naturally sampled with higher probability, which can potentially be advantageous compared to the classical case of explicitly computing all scores. We proceed now in Section 2 with the description of the quantum method and discuss in Section 3 the comparison with classical path-based methods both in terms of prediction precision and resource complexity.

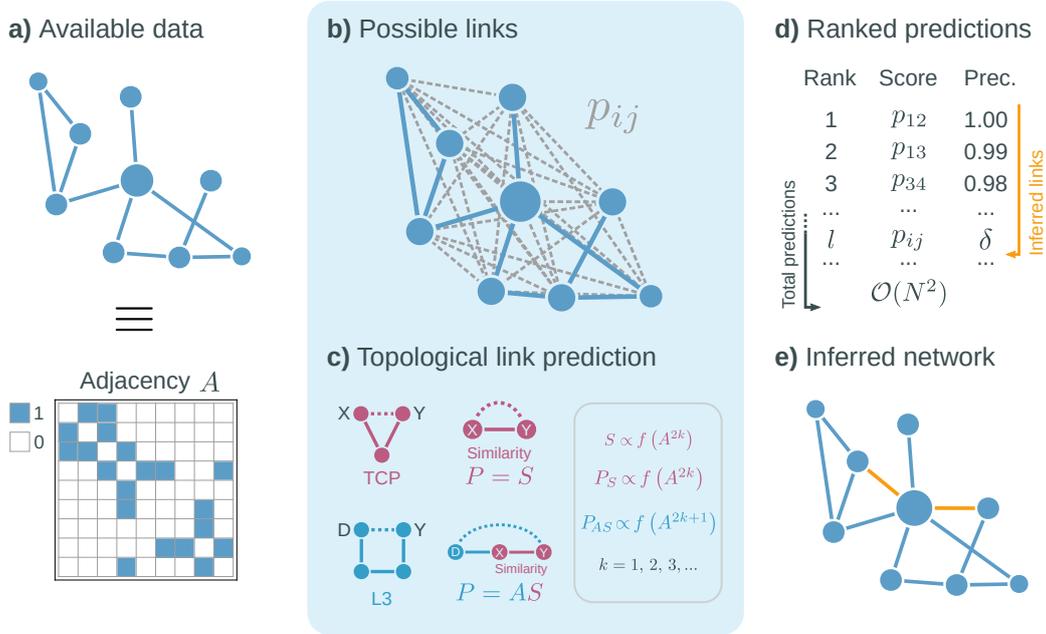


Figure 1: **Classical network-based link prediction.** **a)** Given a complex network described by a graph with a corresponding adjacency matrix A , **b)** one can predict new links by associating a prediction value p_{ij} , or *score*, to every pair of nodes $\{i, j\}$, such that a higher value p_{ij} correlates to a higher probability of the link $\{i, j\}$ appearing. **c)** Predictions based on TCP rely on similarity (matrix S) between nodes, quantified in the simplest case as $P \sim S \sim A^2$, counting paths of length 2 between pairs of nodes. As an alternative, proteins often connect to others that are similar to their neighbours, but not necessarily similar to themselves [18]. In the simplest case, the authors in [18] quantify this principle by taking $P \sim AS \sim A^3$, counting paths of length 3 between nodes. A possible extension of these principles is to quantify direct similarity with even powers of A and neighbour similarity with odd powers of A . **d)** Most classical link prediction methods output the full matrix P , often dense, organized in a ranked list of scores from highest to lowest, where the relevant top l predictions are those where the precision is above a user-determined threshold δ , contributing to the final inferred network (**e**). The need to calculate all pairwise scores constitutes a computational burden which can become intractable as we scale the methods to larger networks.

2 Quantum Link Prediction

We now describe our method for quantum link prediction, denoted as QLP, which we summarize at the end. We base our approach on a Continuous-Time Quantum Walk (CTQW) [11, 12], where the Hilbert space of the quantum walker is defined by the orthonormal basis set $\{|j\rangle\}_{j \in \mathcal{V}}$, with each $|j\rangle$ corresponding to a localized state at a node j . For simplicity we consider the Hamiltonian of the evolution as the adjacency matrix of the graph, $H = A$, but later extend this selection to include a degree normalization. In Figure 2 we show the main structure of the QLP circuit using a qubit representation. In the simplest case, we require $n = \log_2 N$ qubits to add a binary label to each of the N nodes, hereafter marked by the subscript n , and we consider an extra ancilla qubit q_a that doubles the Hilbert space of the quantum walk, such that any node j has two associated basis states $|0\rangle_a |j\rangle_n$ and $|1\rangle_a |j\rangle_n$. For an initial state $|\psi_j(0)\rangle = |0\rangle_a |j\rangle_n$, the first step in the circuit of Figure 2 is to apply an Hadamard gate (Figure 2 c) to q_a , which creates the superposition $\frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)_a |j\rangle_n$. A conditional CTQW is then applied which evolves the $q_a = |0\rangle$ subspace with e^{-iAt} and the $q_a = |1\rangle$ subspace with e^{+iAt} . Finally, a second Hadamard gate is applied to q_a to interfere the two quantum walks in the computational basis, leading to the state

$$|\psi_j(t)\rangle = |0\rangle_a \left(\frac{e^{-iAt} + e^{iAt}}{2} \right) |j\rangle_n + |1\rangle_a \left(\frac{e^{-iAt} - e^{iAt}}{2} \right) |j\rangle_n. \quad (1)$$

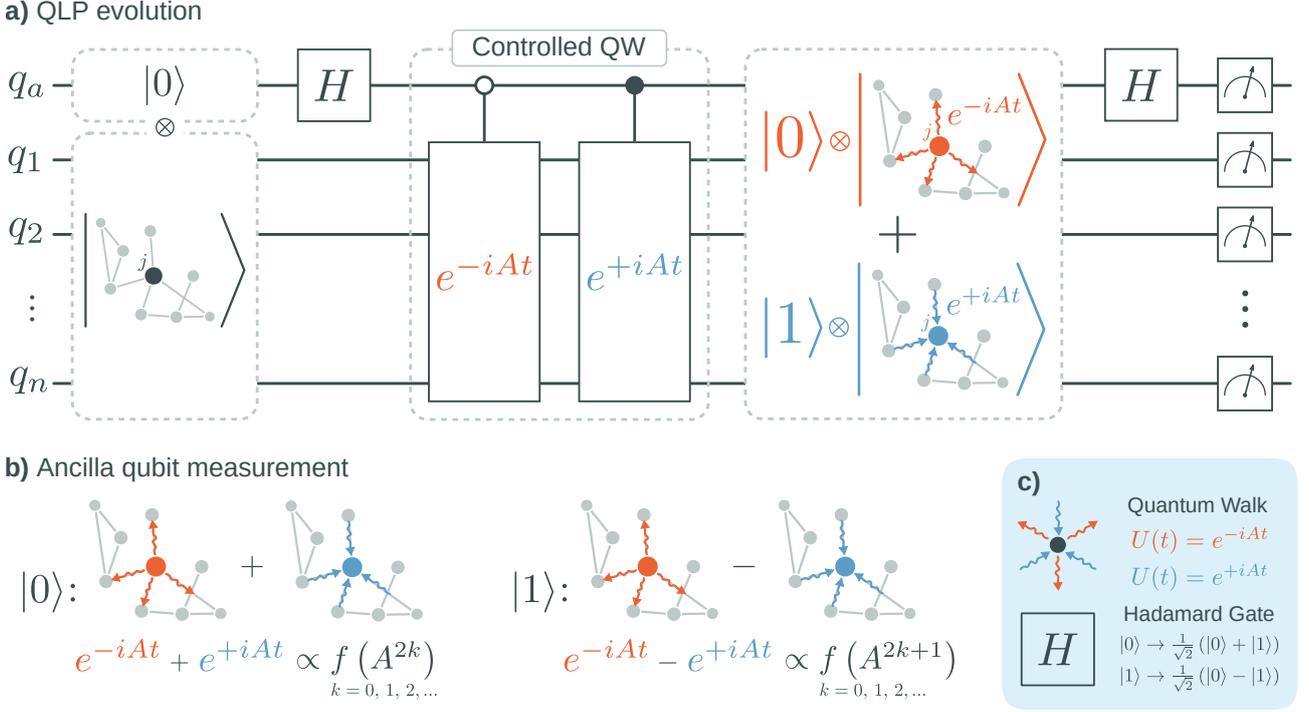


Figure 2: **Quantum link prediction (QLP)**. **a)** A set of $n = \log_2(N)$ qubits creates a Hilbert space that encodes each of the N nodes as a basis state. With an extra ancilla qubit q_a in a superposition of $|0\rangle$ and $|1\rangle$, created by an Hadamard gate (**c**), and an initial state marking a node j , the $q_a = |0\rangle$ subspace is evolved with a e^{-iAt} quantum walk, and the $q_a = |1\rangle$ subspace is evolved with the conjugate e^{+iAt} quantum walk (**c**). A second Hadamard gate applied to the ancilla qubit mixes the two subspaces together and creates an interference between the two quantum walks. **b)** Finally, by measuring q_a the state of the network collapses to one of two possible cases, imposing either a sum or subtraction of the two conjugate evolutions, which encodes even powers of A (even predictions) for $q_a = |0\rangle$ and odd powers of A (odd predictions) for $q_a = |1\rangle$. The measurement of the remaining n qubits returns a bit string marking a certain node i , which together with the initial node j forms a sample of a link (i, j) .

To make the connection with link prediction more evident, it is useful to rewrite the previous expression as

$$|\psi_j(t)\rangle = |0\rangle_a \left(\sum_{k=0}^{+\infty} c_{\text{even}}(k, t) A^{2k} \right) |j\rangle_n + i |1\rangle_a \left(\sum_{k=0}^{+\infty} c_{\text{odd}}(k, t) A^{2k+1} \right) |j\rangle_n, \quad (2)$$

where we have replaced the exponential terms with their respective power series, and defined the time-dependent coefficients as $c_{\text{even}}(k, t) = (-1)^k t^{2k} / (2k)!$ and $c_{\text{odd}}(k, t) = (-1)^{k+1} t^{2k+1} / (2k+1)!$. A detailed calculation leading to Eq. 10 can be found in Supplementary Note 1. Given some initial node j , Eq. 10 describes the state that is created following the QLP evolution. This state has two entangled components, one with a linear combination of even powers of A for $q_a = |0\rangle$, and another with odd powers of A for $q_a = |1\rangle$. The time t of the quantum walk defines the linear weights, and acts as a hyperparameter in the model. This describes the unitary part of the protocol. To obtain relevant predictions from this state we must perform repeated measurements on the system to draw multiple samples, as we now describe.

The first step is to measure q_a , yielding $|0\rangle$ or $|1\rangle$ and collapsing the state of the remaining qubits to $|\psi_j(t)\rangle_n^{\text{even}} \propto (\sum_k c_{\text{even}}(k, t) A^{2k}) |j\rangle$ or $|\psi_j(t)\rangle_n^{\text{odd}} \propto (\sum_k c_{\text{odd}}(k, t) A^{2k+1}) |j\rangle$, respectively, where we omitted the normalization. This effectively selects whether the link sampled will be drawn from a distribution encoding even or odd powers of A . The last step is then to measure the remaining qubits,

yielding a bit string corresponding to a sample of some node i with probability

$$p_{ij}^{\text{even}} \propto \left| \langle i | \left(\sum_{k=0}^{+\infty} c_{\text{even}}(k, t) A^{2k} \right) | j \rangle \right|^2 \quad \text{or} \quad p_{ij}^{\text{odd}} \propto \left| \langle i | \left(\sum_{k=0}^{+\infty} c_{\text{odd}}(k, t) A^{2k+1} \right) | j \rangle \right|^2, \quad (3)$$

which together with the initial node j forms a sample of a link (i, j) . The values p_{ij}^{even} and p_{ij}^{odd} encode the prediction scores of the link (i, j) , but these can not be directly extracted from the algorithm. Instead, what this algorithm allows is the repeated sampling of these distributions, yielding pairs of nodes (i, j) with probability proportional to p_{ij}^{even} or p_{ij}^{odd} . This is analogous to sampling entries (i, j) from the matrix of prediction scores P with probability proportional to $|P_{ij}|^2$. As discussed in Section 1.1, predictions coming from even or odd powers of A are typically useful in different types of networks. For a given network application of QLP, whether each sample obtained corresponds to an even or odd prediction depends on the value measured in the ancilla qubit, and this postselection can only be done probabilistically [25]. This is a potential sampling overhead, as unwanted predictions need to be discarded. Another overhead to consider is the possibility of sampling the initial node, or to sample already existing links, given the contribution of the identity I in p_{ij}^{even} and A in p_{ij}^{odd} , which must also be discarded. As stated, QLP uses a linear combination of powers of A weighted by the time t . A classical prediction method with a linear combination of odd powers of A was presented in [23], which was shown to sometimes improve the prediction precision compared to the original L3 method from [18] by also fitting an additional model parameter. Another popular link prediction method is the Katz index [26], which uses a linear combination of all powers of A .

We can now summarize the QLP algorithm. Firstly, an initial state $|\psi_j(0)\rangle = |0\rangle_a |j\rangle_n$ is prepared for a node j in the network. Secondly, the QLP evolution leading to Eq. 10 is performed for a specific time t . Finally, the ancilla and node qubits are measured to obtain a sample of a link (i, j) corresponding to an even or odd prediction, and the procedure is repeated. The number of samples that output a certain link (i, j) will follow the distributions described by Eq. 3, and thus represent a score for link (i, j) . Once predictions associated with node j are sufficiently characterized, the procedure can be repeated for other relevant nodes in the network.

3 Results and Discussion

3.1 Complexity analysis

To identify a potential quantum advantage, we briefly discuss how link prediction scales on a classical computer. Complex networks are typically sparse [1] with the average degree $k_{\text{av}} \ll N$, and thus there are $\mathcal{O}(N^2)$ potentially missing links. Thus, the general case of computing all possible scores leads to a classical complexity of at least $\mathcal{O}(N^2)$. Different methods scale differently depending on the assumptions made about the solution. For example, the scaling of simple length-2 based methods is $\mathcal{O}(N\langle k^2 \rangle)$ and the scaling of L3 [18] is upper bounded by $\mathcal{O}(N\langle k^3 \rangle)$, where $\langle k^n \rangle$ is the average of the n -th power of the degrees. For certain realistic values of γ , the exponent in the power-law degree distribution of a scale-free network, the moments $\langle k^2 \rangle$ and $\langle k^3 \rangle$ diverge with growing N , as we estimate in Supplementary Note 2. These methods do not calculate a score for every possible missing link, only for those corresponding to nodes at distance 2 or 3. However, other methods also surpass the $\mathcal{O}(N^2)$ scaling, as is the case of LO [23] that uses a matrix inversion to represent a linear combination of odd powers of A , and is one of the best performing classical methods tested. Complex networks can easily

reach sizes of up to millions or billions of nodes, consider for example online social and e-commerce networks, or the neuronal network in the human brain [27]. Improving these scalings may thus be decisive in the application of link prediction methods to larger networks in the future.

In order to estimate the complexity of QLP, we must estimate the number of samples required for a given network. It is often reasonable to assume the number of missing links at any given node j will be proportional to its observed degree k_j , which happens for instance when the missing links are removed randomly from the network. Then, each initial node j will require a number of repetitions of QLP proportional to k_j to sufficiently characterize the predictions associated with node j . This leads to $\mathcal{O}(Nk_{\text{av}})$ total samples for a network of N nodes and average degree k_{av} . To provide a more detailed estimate we must analyse the cost of implementing the unitary e^{-iHt} on a quantum computer, representing the CTQW used to obtain each sample. For this we can look at results from quantum simulation using a d -sparse Hamiltonian model, meaning that H has at most d entries in any given row. A state of the art result [28] shows that implementing e^{-iHt} scales as $\mathcal{O}\left(dt\|H\|_{\text{max}} + \frac{\log(1/\epsilon)}{\log\log(1/\epsilon)}\right)$, where t is the time interval of the evolution, ϵ is the allowed error, and $\|H\|_{\text{max}}$ is the maximum entry in absolute value. In our case, $d = k_{\text{max}}$, the maximum degree of the network, and $\|H\|_{\text{max}} = 1$ for $H = A$, which allows us to write the complexity of implementing e^{-iHt} as $\tilde{\mathcal{O}}(k_{\text{max}}t)$, omitting logarithmic factors. We can thus write a direct complexity estimate for QLP as $\tilde{\mathcal{O}}(Nk_{\text{av}}k_{\text{max}}t)$.

The most meaningful complexity comparison we can make is between methods that make similar assumptions. In that sense, both QLP and LO assume the solution is a linear combination of powers of the adjacency matrix, and as we will see in the next section, these methods are often the best performing. Here, we can see that QLP has a potential quantum speedup given the polynomially lower dependence on N but with an extra $k_{\text{av}}k_{\text{max}}$ factor. Comparing QLP to simple length-2 and length-3 based methods is less straightforward, as the difference is solely based on the degree factors. As mentioned earlier, while the moments $\langle k^2 \rangle$ and $\langle k^3 \rangle$ diverge with growing N , the average degree remains finite. This hints at a potential quantum speedup, but it is not clear if the dependence on k_{max} will spoil the difference. We note that the dependence on k_{max} comes from assuming a circuit-based simulation of the quantum walk. However, our method is general and can also admit an analog quantum walk implementation, which would require a different complexity analysis.

3.2 Cross-validation tests

In Figure 3 we compare the prediction precision of QLP with classical link prediction methods using the standard link prediction benchmark of cross-validation on a selection of networks from different applications. Here we used a degree-normalized adjacency matrix as the Hamiltonian for QLP, $H = D^{-\frac{1}{4}}AD^{-\frac{1}{4}}$, with the final predictions mapped back to A as $P = D^{\frac{1}{4}}\tilde{P}D^{\frac{1}{4}}$, where D is a diagonal matrix with each entry k_i being the degree of node i . This penalizes the counting of paths that go through hubs in the network [29], which, for the purpose of link prediction, can introduce superfluous shortcuts in the network [18]. The scores used for QLP were an exact calculation of the distributions in Eq. 3 by computing the full evolution of the quantum walkers. For each network, we selected the time t that maximizes the prediction precision in the first 10% of the plotted ranks. As shown in Figure 3, we can conclude that QLP matches the best performing classical link prediction methods tested in terms of prediction precision for a range of real life complex networks [19, 30, 31, 32, 33, 34]. In most cases, we observe that both QLP-Odd and LO stand out as the best performing methods, a result which further affirms the case that there can be advantages in including higher order powers of

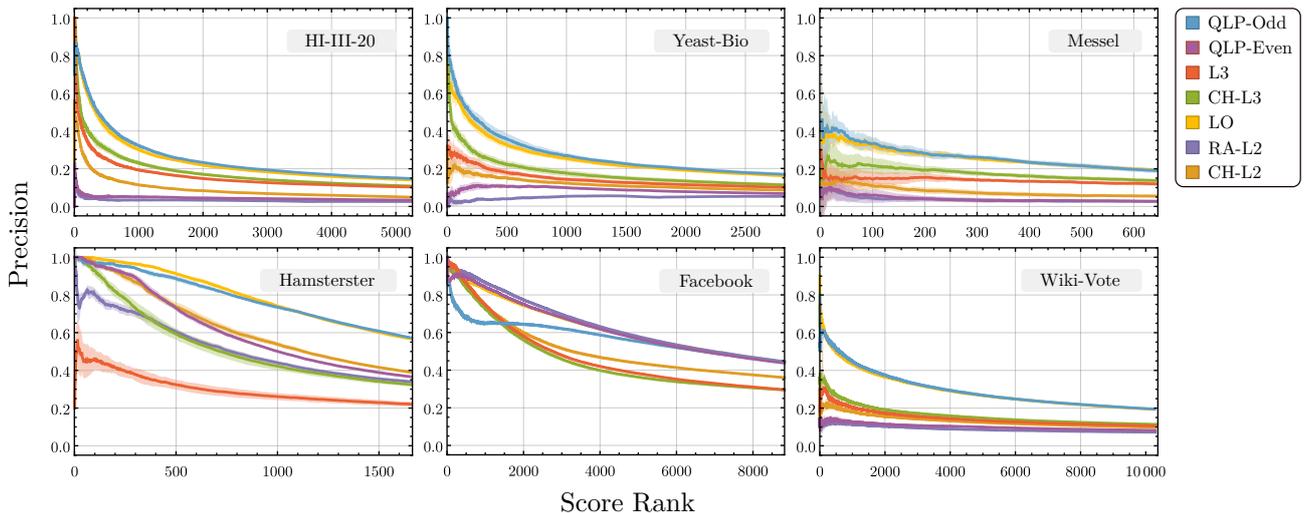


Figure 3: **Computational cross-validation.** Cumulative precision over the top 5% ranked predictions out of Nk_{av} scores for each network, averaged over a 10-fold cross validation procedure. The shaded regions correspond to the standard deviation. In each trial 10% of the links were randomly removed and the remainder used as input to the link prediction methods. The networks used correspond to the PPI networks HI-III-20, the most recent PPI mapping of the human interactome [19], Yeast-Bio, a PPI network of a yeast organism [30], Messel, a food web [31], Hamsterster [32] and Facebook [33], two online social networks, and Wiki-Vote, a vote network between users for adminship of Wikipedia [34]. For comparison, we implemented five classical link prediction methods: the L3 method [18], the LO method [23], the CH-L3 method [22], and two even power methods, RA-L2 (resource allocation) [35], and CH-L2 [36, 22]. The dataset parameters characterizing each network are shown in Supplementary Table 1, and the values selected for the optimal parameters t in the QLP method and α in the LO method are shown in Supplementary Table 2.

the adjacency matrix in the predictions [23]. In some cases, QLP-Odd has a slight advantage over LO. We highlight the results for the Facebook dataset as the only case where even power methods stand out (RA-L2 and QLP-Even), although matched by LO, and also the only case where there is a clear difference between QLP-Odd and LO. While QLP-Odd and LO, when expressed as a power series, have a similar form, the predictions they produce are in fact different, as shown in the Supplementary Note 3 and Table 3. Further results for the cross-validation benchmark are shown in Supplementary Figure 1, as well as detailed results for each of the experimental screens that contribute to the full HI-III-20 network in Supplementary Figure 2. Here, we predict interactions that have been obtained by independent, full experimental screens, simulating the case of real life performance against future experiments.

4 Conclusions

To the best of our knowledge, QLP is the first quantum algorithm for link prediction in complex networks, and the first to potentially offer a quantum speedup for a practical network science problem. Furthermore, the inclusion of even and odd paths allows QLP to make both TCP-like and L3-like predictions, thus encompassing all types of networks where these topological patterns play a role. Our results serve as a proof of principle for potential future applications of QLP in large complex networks using quantum hardware. Recently, a 62-node network CTQW was demonstrated experimentally [37], an important first step towards this goal. Besides the potential improvement in complexity when sampling from the quantum solution, we should also note that a classical simulation of QLP relies on the diagonalization of the adjacency matrix, and thus it has a comparable classical complexity to

other classical link prediction methods. This makes QLP easier to be further developed with a focus on immediately relevant real world applications, while at the same time exploring other ways in which quantum features of QLP can be advantageous when quantum hardware becomes more widely available. These findings open the way to explore this novel frontier of quantum-computational advantage for complex network applications.

Acknowledgements

The authors thank Albert-László Barabási for the useful discussion, and acknowledge the support from the JTF project *The Nature of Quantum Networks* (ID 60478). JPM, BC and YO thank the support from Fundação para a Ciência e a Tecnologia (FCT, Portugal), namely through project UIDB/50008/2020, as well as from projects TheBlinQC and QuantHEP supported by the EU H2020 QuantERA ERA-NET Cofund in Quantum Technologies and by FCT (QuantERA/0001/2017 and QuantERA/0001/2019, respectively), and from the EU H2020 Quantum Flagship project QMiCS (820505). JPM acknowledges the support of FCT through scholarships SFRH/BD/144151/2019, and BC acknowledges the support of FCT through project CEECINST/00117/2018/CP1495.

References

- [1] Barabási, A.-L. *et al. Network science* (Cambridge University Press, 2016).
- [2] Dürr, C., Heiligman, M., Hoyer, P. & Mhalla, M. Quantum query complexity of some graph problems. *SIAM Journal on Computing* **35**, 1310–1328 (2006).
- [3] Ambainis, A. & Špalek, R. Quantum algorithms for matching and network flows. In *Annual Symposium on Theoretical Aspects of Computer Science*, 172–183 (Springer, 2006).
- [4] Chakraborty, S., Novo, L., Ambainis, A. & Omar, Y. Spatial search by quantum walk is optimal for almost all graphs. *Physical Review Letters* **116**, 100501 (2016).
- [5] Chakraborty, S., Novo, L., Di Giorgio, S. & Omar, Y. Optimal quantum spatial search on random temporal networks. *Physical Review Letters* **119**, 220503 (2017).
- [6] Tsomokos, D. I. Quantum walks on complex networks with connection instabilities and community structure. *Physical Review A* **83**, 052315 (2011).
- [7] Sánchez-Burillo, E., Duch, J., Gómez-Gardenes, J. & Zueco, D. Quantum navigation and ranking in complex networks. *Scientific reports* **2**, 1–8 (2012).
- [8] Faccin, M., Johnson, T., Biamonte, J., Kais, S. & Migdał, P. Degree distribution in quantum walks on complex networks. *Physical Review X* **3**, 041007 (2013).
- [9] Mukai, K. & Hatano, N. Discrete-time quantum walk on complex networks for community detection. *Physical Review Research* **2**, 023378 (2020).
- [10] Faccin, M., Migdał, P., Johnson, T. H., Bergholm, V. & Biamonte, J. D. Community detection in quantum complex networks. *Physical Review X* **4**, 041012 (2014).

- [11] Farhi, E. & Gutmann, S. Quantum computation and decision trees. *Physical Review A* **58**, 915 (1998).
- [12] Kempe, J. Quantum random walks: an introductory overview. *Contemporary Physics* **44**, 307–327 (2003).
- [13] Liben-Nowell, D. & Kleinberg, J. The link-prediction problem for social networks. *Journal of the American society for information science and technology* **58**, 1019–1031 (2007).
- [14] Wang, P., Xu, B., Wu, Y. & Zhou, X. Link prediction in social networks: the state-of-the-art. *Science China Information Sciences* **58**, 1–38 (2015).
- [15] Albert, I. & Albert, R. Conserved network motifs allow protein–protein interaction prediction. *Bioinformatics* **20**, 3346–3352 (2004).
- [16] Getoor, L. & Diehl, C. P. Link mining: a survey. *ACM Sigkdd Explorations Newsletter* **7**, 3–12 (2005).
- [17] Lü, L. & Zhou, T. Link prediction in complex networks: A survey. *Physica A: statistical mechanics and its applications* **390**, 1150–1170 (2011).
- [18] Kovács, I. A. *et al.* Network-based prediction of protein interactions. *Nature Communications* **10**, 1–8 (2019).
- [19] Luck, K. *et al.* A reference map of the human binary protein interactome. *Nature* **580**, 402–408 (2020).
- [20] Al Hasan, M., Chaoji, V., Salem, S. & Zaki, M. Link prediction using supervised learning. In *SDM06: workshop on link analysis, counter-terrorism and security*, vol. 30, 798–805 (2006).
- [21] Lü, L., Pan, L., Zhou, T., Zhang, Y.-C. & Stanley, H. E. Toward link predictability of complex networks. *Proceedings of the National Academy of Sciences* **112**, 2325–2330 (2015).
- [22] Muscoloni, A., Abdelhamid, I. & Cannistraci, C. V. Local-community network automata modelling based on length-three-paths for prediction of complex network structures in protein interactomes, food webs and more. *bioRxiv preprint* (2018).
- [23] Pech, R., Hao, D., Lee, Y.-L., Yuan, Y. & Zhou, T. Link prediction via linear optimization. *Physica A: Statistical Mechanics and its Applications* **528**, 121319 (2019).
- [24] Kitsak, M. Latent geometry for complementarity-driven networks. *arXiv preprint arXiv:2003.06665* (2020).
- [25] Kothari, R. Efficient algorithms in quantum query complexity (2014).
- [26] Katz, L. A new status index derived from sociometric analysis. *Psychometrika* **18**, 39–43 (1953).
- [27] Azevedo, F. A. *et al.* Equal numbers of neuronal and nonneuronal cells make the human brain an isometrically scaled-up primate brain. *Journal of Comparative Neurology* **513**, 532–541 (2009).
- [28] Low, G. H. & Chuang, I. L. Optimal hamiltonian simulation by quantum signal processing. *Physical review letters* **118**, 010501 (2017).

- [29] Barabási, A.-L. & Albert, R. Emergence of scaling in random networks. *Science* **286**, 509–512 (1999).
- [30] Stark, C. *et al.* Biogrid: a general repository for interaction datasets. *Nucleic acids research* **34**, D535–D539 (2006).
- [31] Dunne, J. A., Labandeira, C. C. & Williams, R. J. Highly resolved early eocene food webs show development of modern trophic structure after the end-cretaceous extinction. *Proceedings of the Royal Society B: Biological Sciences* **281**, 20133280 (2014).
- [32] Kunegis, J. Konect: the koblenz network collection. In *Proceedings of the 22nd international conference on world wide web*, 1343–1350 (2013).
- [33] McAuley, J. J. & Leskovec, J. Learning to discover social circles in ego networks. In *NIPS*, vol. 2012, 548–56 (Citeseer, 2012).
- [34] Leskovec, J., Huttenlocher, D. & Kleinberg, J. Signed networks in social media. In *Proceedings of the SIGCHI conference on human factors in computing systems*, 1361–1370 (2010).
- [35] Zhou, T., Lü, L. & Zhang, Y.-C. Predicting missing links via local information. *The European Physical Journal B* **71**, 623–630 (2009).
- [36] Cannistraci, C. V., Alanis-Lobato, G. & Ravasi, T. From link-prediction in brain connectomes and protein interactomes to the local-community-paradigm in complex networks. *Scientific Reports* **3**, 1613 (2013).
- [37] Gong, M. *et al.* Quantum walks on a programmable two-dimensional 62-qubit superconducting processor. *Science* **372**, 948–952 (2021).
- [38] Fiol, M. A. & Garriga, E. Number of walks and degree powers in a graph. *Discrete Mathematics* **309**, 2613–2614 (2009).
- [39] Interactome, A. I. M. C. A. Mapping consortium evidence for network evolution in an arabidopsis interactome map. *Science* **333**, 601–607 (2011).
- [40] Leskovec, J., Kleinberg, J. & Faloutsos, C. Graph evolution: Densification and shrinking diameters. *ACM transactions on Knowledge Discovery from Data (TKDD)* **1**, 2–es (2007).
- [41] Sen, P. *et al.* Collective classification in network data. *AI magazine* **29**, 93–93 (2008).

Supplementary Note 1 - QLP method

We consider the usual continuous-time quantum walk (CTQW) model, where the Hilbert space of the quantum walker is defined by the orthonormal basis set $\{|j\rangle\}_{j \in \mathcal{V}}$, each basis state $|j\rangle$ corresponding to a localized state at a node j in the network, and the hamiltonian of the evolution given by the adjacency matrix of the graph,

$$|\psi(t)\rangle = e^{-iAt} |\psi(0)\rangle, \quad (4)$$

By taking the power series of the time evolution operator we can immediately make the connection to link prediction,

$$e^{-iAt} = \sum_{k=0}^{+\infty} \frac{1}{k!} (-it)^k A^k, \quad (5)$$

as each power A^k encodes the number of paths of length k between any two nodes in the graph. Furthermore, we note that the imaginary term i^k adds a phase to the quantum evolution that separates the sum over even powers in the real part of the evolution and the sum over odd powers in the imaginary part. To proceed we wish to separate the evolution over even powers from the evolution over odd powers, and for that it is useful to consider a qubit representation of the graph, as seen in Fig. 2 of the main text. We now define an operator CQW(t) corresponding to a controlled quantum walk which applies a normal or conjugate evolution operator on the node qubits depending on the value of q_a ,

$$\text{CQW}(t) = |0\rangle\langle 0|_a (e^{-iAt})_n + |1\rangle\langle 1|_a (e^{+iAt})_n. \quad (6)$$

Considering now an initial state localized at node j in the form $|\psi_j(0)\rangle = |0\rangle_a |j\rangle_n$ we start by applying an Hadamard gate on the ancilla qubit,

$$H_a |0\rangle_a |j\rangle_n = \frac{1}{\sqrt{2}} (|0\rangle + |1\rangle)_a |j\rangle_n, \quad (7)$$

followed by the CQW(t) operator,

$$\text{CQW}(t) \left[\frac{1}{\sqrt{2}} (|0\rangle + |1\rangle)_a |j\rangle_n \right] = \frac{1}{\sqrt{2}} (|0\rangle_a e^{-iAt} |j\rangle_n + |1\rangle_a e^{+iAt} |j\rangle_n), \quad (8)$$

followed by a second Hadamard gate on the ancilla qubit, leading to the following expression after rearranging the terms:

$$|\psi_j(t)\rangle = \frac{1}{2} |0\rangle_a (e^{-iAt} + e^{iAt}) |j\rangle_n + |1\rangle_a (e^{-iAt} - e^{iAt}) |j\rangle_n. \quad (9)$$

Finally, taking the power series from Eq. 5 to replace the exponential terms we arrive at

$$|\psi_j(t)\rangle = |0\rangle_a \left(\sum_{k=0}^{+\infty} c_{\text{even}}(k, t) A^{2k} \right) |j\rangle_n + i |1\rangle_a \left(\sum_{k=0}^{+\infty} c_{\text{odd}}(k, t) A^{2k+1} \right) |j\rangle_n, \quad (10)$$

with $c_{\text{even}}(k, t) = (-1)^k t^{2k} / (2k)!$ and $c_{\text{odd}}(k, t) = (-1)^{k+1} t^{2k+1} / (2k+1)!$ being time-dependent coefficients.

Supplementary Note 2 - Classical complexity of link prediction

Consider a graph $G(\mathcal{V}, \mathcal{E})$ describing a complex network, where \mathcal{V} is the set of nodes with size $N = |\mathcal{V}|$ and \mathcal{E} is the set of undirected links. Link prediction on a classical computer requires $\frac{1}{2}N(N-1) - |\mathcal{E}|$ scores to be computed, one for each of the $\frac{1}{2}N(N-1)$ possible links, with the exception of those already present in the set of known links \mathcal{E} . Rewriting in terms of the average degree, $k_{\text{av}} = 2|\mathcal{E}|/N$, we have that the total number of scores scales as $\frac{1}{2}N^2 - \frac{1}{2}N(1 + k_{\text{av}})$. Real complex networks are typically sparse [1] with $k_{\text{av}} \ll N$, and thus $\mathcal{O}(N^2)$ scores are evaluated. Taking $\mathcal{O}(N^2)$ as an estimate for the complexity of a general classical link prediction method assumes two more things: that the method will indeed compute a score for every potential missing link, and that the cost of computing each score is $\mathcal{O}(1)$. In order to analyse these assumptions, let us pick a concrete method and study its complexity.

Common Neighbours (CN) is one of the simplest link prediction algorithms. It quantifies the likelihood of a link existing between two nodes i and j by the number of common neighbours they share, or in other words, by the number of paths of length 2 between i and j . While we don't use CN directly in the various simulations presented in this work, we used the method of Resource Allocation [35] (marked as RA-L2 in the plots), which is similar to CN with the addition of a degree normalization to each score. Adding the degree normalization does not affect the complexity significantly, and so we will analyse the simpler problem of counting paths of length 2. The objective of CN is to compute

$$p_{ij} = |\Gamma(i) \cap \Gamma(j)| \quad (11)$$

for every pair of nodes (i, j) where $|\Gamma(i) \cap \Gamma(j)| \neq 0$, $\Gamma(x)$ being the set of nodes neighbouring x . A simple algorithm to accomplish this iterates through all nodes z in the graph and adds a contribution to p_{ij} for each pair of nodes (i, j) neighbouring z . Such an algorithm will visit every path of length 2 in the graph and thus its complexity will be proportional to $\sum_{i,j=1}^N (A^2)_{ij}$. As detailed in [38] this sum can be simplified as

$$\sum_{i,j=1}^N (A^2)_{ij} = \sum_{i=1}^N k_i^2 = N \langle k^2 \rangle, \quad (12)$$

where $\langle k^2 \rangle$ is the average of the second power of the degrees in the graph. By assuming that the cost of accessing the graph data structure and adding the contributions to each p_{ij} is $\mathcal{O}(1)$ we can conclude that the CN method scales as $\mathcal{O}(N \langle k^2 \rangle)$.

Common Neighbours is a TCP based method, and as discussed in the main text, it is not able to match the precision of methods based on paths of length 3 in many complementarity driven networks. For that reason, let us see how the complexity changes when counting paths of length 3, which is the main computational cost behind the L3 method [18]. An algorithm to count paths of length 3 can be easily built with an extension of the CN algorithm, and using the same argument as before, its complexity will be proportional to $\sum_{i,j=1}^N (A^3)_{ij}$. This sum is not as easy to simplify, but the authors in [38] prove the following bound for a general power of A

$$\sum_{i,j=1}^N (A^n)_{ij} \leq \sum_{i=1}^N k_i^n = N \langle k^n \rangle. \quad (13)$$

With this information we can conclude that the complexity for counting paths of length 3 will be upper bounded by $\mathcal{O}(N \langle k^3 \rangle)$.

For scale-free networks, characterized by γ , the degree exponent in the degree power law distribution, we can analyse the moments $\langle k^n \rangle$ in terms of γ and N (see Section 4 of [1]). Typically, $\langle k \rangle$ (denoted

as k_{av} in the rest of the text) is much smaller than $\langle k^2 \rangle$ or $\langle k^3 \rangle$. For many scale-free networks γ is between 2 and 4. As N grows, $\langle k^2 \rangle$ diverges for $2 \leq \gamma \leq 3$ and $\langle k^3 \rangle$ diverges for $2 \leq \gamma \leq 4$, while $\langle k \rangle$ remains finite. These divergences can be seen in the expressions below from [1] which estimate the dependence of $\langle k^n \rangle$ with N :

$$\langle k^n \rangle \propto \frac{k_{\max}^{n-\gamma+1} - k_{\min}^{n-\gamma+1}}{n - \gamma + 1} \quad (14)$$

which together with the relation $k_{\max} = k_{\min} N^{\frac{1}{\gamma-1}}$ can be written as:

$$\langle k^n \rangle \propto \frac{k_{\min}^{n-\gamma+1}}{n - \gamma + 1} \left(N^{\frac{n-\gamma+1}{\gamma-1}} - 1 \right) \quad (15)$$

As stated before, out of the methods tested in this work, RA-L2 and L3 fall in the complexity categories of counting paths of length 2 and 3, respectively. CH-L2 and CH-L3 also have path counting as a base, but use a more complex structure of paths which has added complexity. LO, the best performing classical method tested, uses a matrix inversion for which the best algorithms scale between $\mathcal{O}(N^{2.3})$ and $\mathcal{O}(N^{2.4})$.

Supplementary Table 1 - Dataset parameters

	Network	Ref.	$ V $	$ E $	k_{av}	ρ	d_{max}	d_{av}	C
Main Text - Fig. 3	HI-III-20	[19]	8275	52569	12.589	1.59×10^{-3}	12	3.844	5.92×10^{-2}
	Yeast-Bio	[30]	4885	28270	11.161	2.29×10^{-3}	10	3.603	1.20×10^{-1}
	Messel	[31]	700	6444	18.326	2.61×10^{-2}	6	2.632	1.04×10^{-1}
	Hamsterster	[32]	2426	16631	13.711	5.65×10^{-3}	10	3.589	5.38×10^{-1}
	Facebook	[33]	4039	88234	43.691	1.08×10^{-2}	8	3.693	6.06×10^{-1}
	Wiki-Vote	[34]	7115	103689	29.147	3.98×10^{-3}	7	3.248	1.41×10^{-1}
S.I. - Fig. 1	Arabidopsis	[39]	4865	11374	4.493	9.24×10^{-4}	14	5.180	9.82×10^{-2}
	Pombe	[30]	1929	3700	3.397	1.76×10^{-3}	14	4.671	6.37×10^{-2}
	AS Routes	[40]	6474	13.895	3.884	6.00×10^{-4}	9	3.705	2.52×10^{-1}
	Citeseer	[41]	3264	4536	2.779	8.518×10^{-4}	28	9.315	1.45×10^{-1}
	Cora	[41]	2708	5429	4.010	1.44×10^{-3}	19	6.311	2.41×10^{-1}
	P2P-Gnutella	[40]	10876	39994	7.355	6.78×10^{-4}	10	4.622	6.22×10^{-3}
S.I. - Fig. 2 - HuRI Screens	Screen 1	[19]	4643	16447	6.970	1.64×10^{-3}	12	4.094	5.13×10^{-2}
	Screen 2	[19]	4177	11644	5.467	1.31×10^{-3}	13	4.284	4.33×10^{-2}
	Screen 3	[19]	3807	10245	5.268	1.31×10^{-3}	15	4.456	4.16×10^{-2}
	Screen 4	[19]	3082	5685	3.655	1.19×10^{-3}	14	5.370	1.51×10^{-2}
	Screen 5	[19]	2712	4496	3.277	1.21×10^{-3}	15	5.560	1.25×10^{-2}
	Screen 6	[19]	3128	5981	3.774	1.21×10^{-3}	16	5.361	1.54×10^{-2}
	Screen 7	[19]	3508	7910	4.486	1.28×10^{-3}	15	5.465	9.73×10^{-3}
	Screen 8	[19]	3383	7533	4.436	1.31×10^{-3}	16	5.555	1.20×10^{-2}
	Screen 9	[19]	3404	7712	4.512	1.33×10^{-3}	15	5.520	9.40×10^{-3}

Table 1: Datasets and respective network parameters.

Supplementary Table 2 - Hyperparameter values

	Network	t_{Even} (QLP-Even)	t_{Odd} (QLP-Odd)	α (LO)
Main Text - Fig. 3	HI-III-20	3.0×10^{-6}	1.00	8.0×10^{-3}
	Yeast-Bio	7.0×10^{-1}	1.10	2.0×10^{-2}
	Messel	2.0×10^{-6}	1.20	2.0×10^{-2}
	Hamsterster	4.0×10^{-1}	1.60	3.0×10^{-2}
	Facebook	2.0×10^{-1}	1.00	7.0×10^{-3}
	Wiki-Vote	2.0×10^{-6}	1.00	4.0×10^{-3}
S.I. - Fig. 1	Arabidopsis	2.0×10^{-6}	1.00	2.0×10^{-2}
	Pombe	9.0×10^{-1}	1.00	4.0×10^{-2}
	AS Routes	5.0×10^{-4}	0.60	2.0×10^{-3}
	Citeseer	5.0×10^{-5}	1.40	8.0×10^{-2}
	Cora	1.0×10^{-6}	0.03	6.0×10^{-2}
	P2P-Gnutella	-	1.30	2.0×10^{-2}
S.I. - Fig. 2 - HuRI Screens	Screen 1	6.0×10^{-6}	1.10	1.0×10^{-2}
	Screen 2	4.0×10^{-6}	0.90	8.0×10^{-3}
	Screen 3	6.0×10^{-5}	0.90	1.0×10^{-2}
	Screen 4	9.0×10^{-5}	0.02	1.0×10^{-2}
	Screen 5	6.0×10^{-5}	0.03	3.0×10^{-3}
	Screen 6	1.0×10^{-6}	0.03	7.0×10^{-3}
	Screen 7	1.0×10^{-6}	0.70	7.0×10^{-4}
	Screen 8	2.0×10^{-6}	0.30	4.0×10^{-4}
	Screen 9	1.0×10^{-6}	0.80	6.0×10^{-4}

Table 2: Hyperparameter values used in each dataset. We omitted t_{Even} for P2P-Gnutella given that this is a bipartite network, and thus predictions based on even length paths have no meaning.

Supplementary Figure 1 - Extra cross-validation results

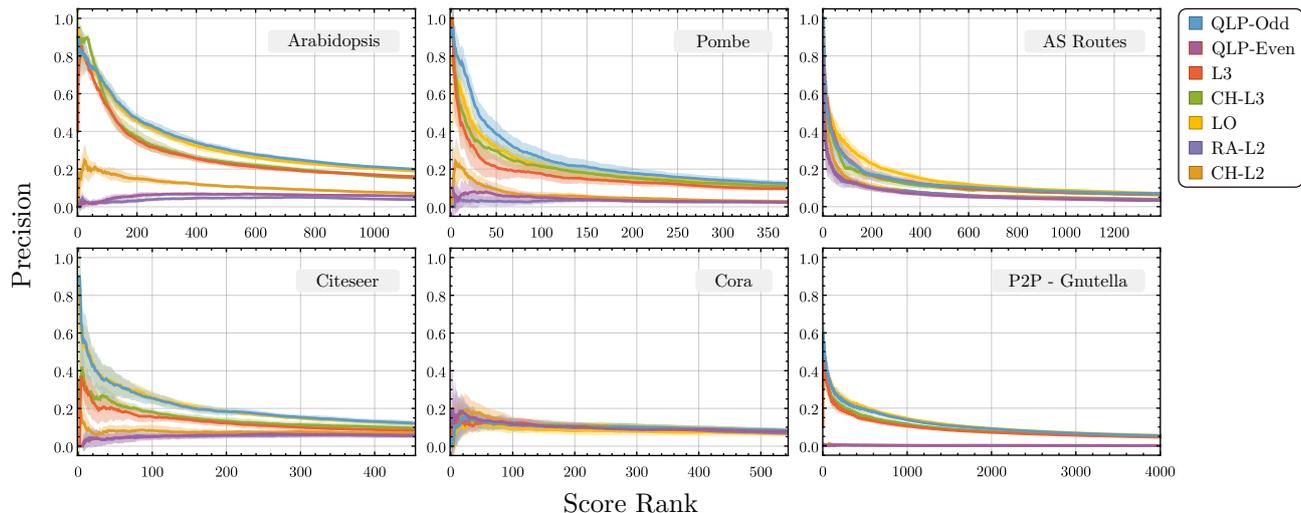


Figure 1: **Extra cross-validation results.** The results presented correspond to the cumulative precision over the top 1000 ranked predictions, averaged over ten different trials. The shaded regions correspond to the standard deviation. In each trial 50% of the links were randomly removed and the remainder used as input to the link prediction methods. The networks used correspond to the PPI network Arabidopsis [39], AS Routes [40], Facebook [33], Citeseer [41], Cora [41] and P2P-Gnutella [40]. We note that P2P-Gnutella is a bipartite network, and thus the TCP based methods QLP-Even, RA-L2 and CH-L2 have null precision.

Supplementary Figure 2 - HuRI screen vs screen results

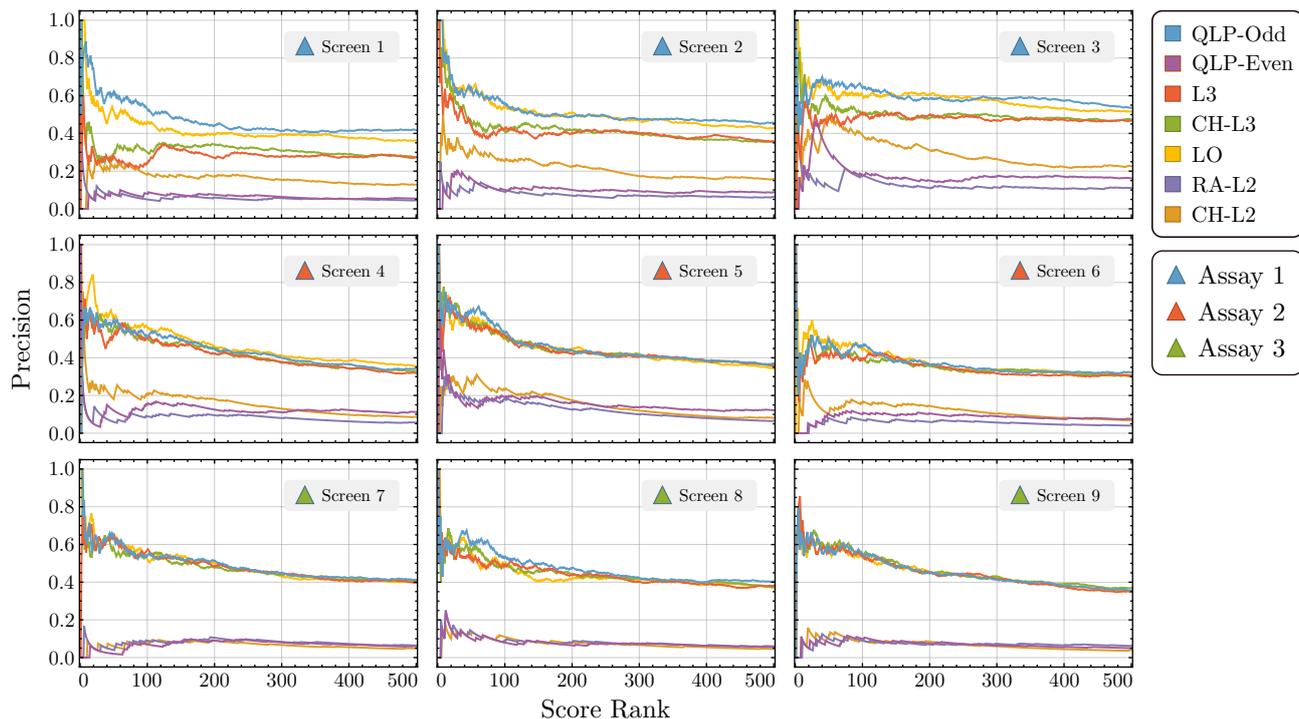


Figure 2: Prediction of missing PPIs in the human interactome validated against experimental screens. The results presented correspond to the cumulative precision over the top 500 ranked predictions, as detailed in the Methods section. The dataset used consists of nine different screens over a search space of human binary PPIs using a panel of three different assay versions [19]. For each of the nine plots we used the results of the respective screen as the input network to the link prediction methods and compared the predictions obtained to the PPIs identified in the remaining two screens from the same assay. For example, for the case of Screen 1, the predictions were compared with the PPIs identified in Screen 2 and Screen 3 combined. For the methods with a free parameter (QLP and LO) we randomly removed 50% of the input dataset and optimized the method to best predict the removed links by maximizing the area under the precision curve over the top 500 score ranks. The optimized parameter was then used for the results displayed.

In Fig. 2 we applied QLP to a dataset of the most recent experimental study of the human interactome network [19], named HI-III-20. In this study the authors presented a reference interactome map of human binary protein interactions with 52,569 protein-protein interactions involving 8,275 proteins. This map was generated by screening a search space of roughly 90% of the protein-coding genome a total of nine times with a panel of three different but complementary assay versions. We considered each screen as an input network, and validated the predictions obtained by comparison with the PPIs identified in the other two screens of the same assay combined. These results further complement our conclusion that QLP is competitive with the state of the art classical link prediction methods in terms of prediction precision.

Supplementary Note 3 - Comparison between QLP-Odd, LO, A^3 and L3

In this section we look to compare the predictions obtained from QLP-Odd, LO and L3. The motivation to do so is that QLP-Odd and LO admit a very similar power series, and produce very similar results in many of the tested datasets. Besides the clear difference of QLP-Odd being a quantum method, we wish to highlight other practical differences in the predictions obtained by each method. To do so, we first start by writing the power series of QLP-Odd and LO in matrix form,

$$\begin{aligned} P_{\text{QLP-Odd}} &\propto \left| -tA + \frac{t^3}{3!}A^3 - \frac{t^5}{5!}A^5 + \frac{t^7}{7!}A^7 - \dots \right|_{(ij)}^2 \\ P_{\text{LO}} &= \alpha A^3 - \alpha^2 A^5 + \alpha^3 A^7 - \dots \end{aligned} \quad (16)$$

where $P_{\text{QLP-Odd}}$ is written up to a normalization factor, and we denote $|\cdot|_{(ij)}^2$ as the entry-wise absolute value squared. P_{LO} is written according to [23]. Immediately we see that the differences are the presence of the linear term in QLP-Odd, which does not contribute to the predictions, the different weight definition, and the absolute value squared. To test their practical difference we present in the Supplementary Table 3 the number of matching predictions in the top-10000 scores between QLP-Odd, LO, A^3 and L3 for increasing values of t , and for three datasets: Yeast-Bio, Facebook and Wiki-Vote. Note that we simply compared the total number of common predictions in the top-10000 list, irrespective of the predictions being ranked the same. In fact, besides the case of very small t , we found that almost all predictions were ranked differently between the methods. We conclude that for the three tested datasets the predictions between all methods are only similar or equal for small values of t , when the predictions of QLP-Odd are indeed governed mostly by the third power of A .

Supplementary Table 3 - Matching predictions between QLP-Odd, LO, A^3 and L3

		Matching Predictions in the top-10000 scores			
	t_{Odd}	α	QLP-Odd vs LO	QLP-Odd vs A^3	QLP-Odd vs L3
Yeast-Bio	0.001	1.67×10^{-10}	10000	9972	8942
	0.01	1.67×10^{-7}	10000	9972	8941
	0.05	2.08×10^{-5}	9717	9664	8935
	0.1	1.67×10^{-4}	8911	8518	8911
	0.5	2.08×10^{-2}	4098	1095	8155
	0.8	8.53×10^{-2}	4460	567	6715
	1.0	1.67×10^{-1}	4753	336	5113
	1.2	2.88×10^{-1}	5291	211	3432
	1.5	5.63×10^{-1}	5796	186	2007
Facebook	0.001	1.67×10^{-10}	9998	9998	9758
	0.01	1.67×10^{-7}	9877	9874	9758
	0.05	2.08×10^{-5}	5975	5579	9764
	0.1	1.67×10^{-4}	5040	4187	9755
	0.5	2.08×10^{-2}	2827	503	5447
	0.8	8.53×10^{-2}	2410	317	3066
	1.0	1.67×10^{-1}	2311	272	3692
	1.2	2.88×10^{-1}	2249	267	3083
	1.5	5.63×10^{-1}	2141	217	3418
Wiki-Vote	0.001	1.67×10^{-10}	10000	9999	9435
	0.01	1.67×10^{-7}	9846	9841	9434
	0.05	2.08×10^{-5}	5363	4887	9409
	0.1	1.67×10^{-4}	3649	2066	9321
	0.5	2.08×10^{-2}	1765	530	5648
	0.8	8.53×10^{-2}	1926	411	2651
	1.0	1.67×10^{-1}	2531	414	2087
	1.2	2.88×10^{-1}	3408	326	1296
	1.5	5.63×10^{-1}	4597	384	723

Table 3: Matching predictions in the top-10000 scores between QLP-Odd, LO, A^3 and L3 for Yeast-Bio, Facebook and Wiki-Vote for varying values of t . The values of α in the LO method were chosen as $\alpha = t^3/3!$ in order to equal the leading prediction order to that of QLP-Odd in the power series. In the comparison with LO and A^3 we used $H = A$, and in the comparison with L3 we used $H = D^{-\frac{1}{4}}AD^{-\frac{1}{4}}$, thus matching the degree normalization between the methods.