

Efficient cavity control with SNAP gates

Thomas Fösel,^{1,2,3} Stefan Krastanov,^{2,3,4,5} Florian Marquardt,^{1,6} and Liang Jiang^{2,3,7}

¹*Max Planck Institute for the Science of Light, Staudtstr. 2, 91058 Erlangen, Germany*

²*Department of Applied Physics and Physics, Yale University, New Haven, CT, USA*

³*Yale Quantum Institute, Yale University, New Haven, CT, USA*

⁴*Department of Electrical Engineering and Computer Science,
Massachusetts Institute of Technology, Cambridge, MA 02139, USA*

⁵*John A. Paulson School of Engineering and Applied Sciences,
Harvard University, Cambridge, MA 02138, USA*

⁶*Physics Department, University of Erlangen-Nuremberg, Staudtstr. 5, 91058 Erlangen, Germany*

⁷*Pritzker School of Molecular Engineering, The University of Chicago, Chicago, Illinois 60637, USA*

Microwave cavities coupled to superconducting qubits have been demonstrated to be a promising platform for quantum information processing. A major challenge in this setup is to realize universal control over the cavity. A promising approach are selective number-dependent arbitrary phase (SNAP) gates combined with cavity displacements. It has been proven that this is a universal gate set, but a central question remained open so far: how can a given target operation be realized efficiently with a sequence of these operations. In this work, we present a practical scheme to address this problem. It involves a hierarchical strategy to insert new gates into a sequence, followed by a co-optimization of the control parameters, which generates short high-fidelity sequences. For a broad range of experimentally relevant applications, we find that they can be implemented with 3 to 4 SNAP gates, compared to up to 50 with previously known techniques.

I. INTRODUCTION

In the past decade, remarkable progress has happened in the field of quantum information processing [1, 2]. One of the leading platforms is circuit QED with microwave cavities coupled to superconducting qubits. While the (three-dimensional) cavities offer relatively long coherence times of around 1...10 ms, they allow only for very restricted direct control via cavity displacements (driven with a microwave pulse). In contrast, the non-linearity of superconducting qubits like transmons [3–5] allows for universal control over them, but their coherence time is much shorter (up to 100 μ s). To combine the respective strengths of these systems, i. e. the long lifetime and the good controllability, the quantum information can be stored in (a single mode of) the cavity, and the transmon is used as an ancilla to manipulate the cavity state. Impressive recent demonstrations using this platform are the realization of quantum error correction beyond the breakeven point [6], entanglement between remote systems [7] and deterministic teleportation of quantum gates [8].

Currently, there are two competing approaches towards achieving universal control over the cavity with the help of a transmon. The first one makes use of the full freedom of the control space in order to find suitable control pulses [9], which is typically done using numerical optimal-control techniques like gradient ascent pulse engineering (GRAPE [10]). The second approach is to build combinations from a restricted set of well-controllable gates. A promising candidate for

such a gate set is the selective number-dependent arbitrary phase (SNAP) gate [11] combined with cavity displacements, for which universality has been proven [12]. The main advantage of this second approach is its modularity, which results in better interpretability and simplifies the optimization of its constituents. In particular, the SNAP gate can be made robust against transmon decay via error transparency [13, 14].

One central challenge in this gate-based approach is to find suitable combinations for the (continuous) control parameters of these gates, which represents a complex task due to the huge size of the search space. Furthermore, the generated gate sequences have to be efficient, i. e. the number of gates needs to be small, as the following consideration shows.

The previously best known technique from [12] generates sequences with (at least) N^2 SNAP gates interleaved with $N^2 + 1$ displacements to implement an operation on N different Fock states. E. g., for $N = 10$, this means that 100 SNAP gates would be required. With 3...5 μ s per SNAP gate and less than 0.1 μ s per displacement [7], the total execution time would be 300...500 μ s, resulting in significant corruption of the quantum information given the lifetime of the cavity (1...10 ms).

A simple parameter-counting argument suggests that the actually required number of gates should be significantly lower. We compare the number of conditions to match an operation on N Fock states, which is $\dim(\text{SU}(N)) = N^2 - 1 = \mathcal{O}(N^2)$, to the effective number of free parameters for each pair of displacement and SNAP gate, which is $\mathcal{O}(N)$. From the quotient of these

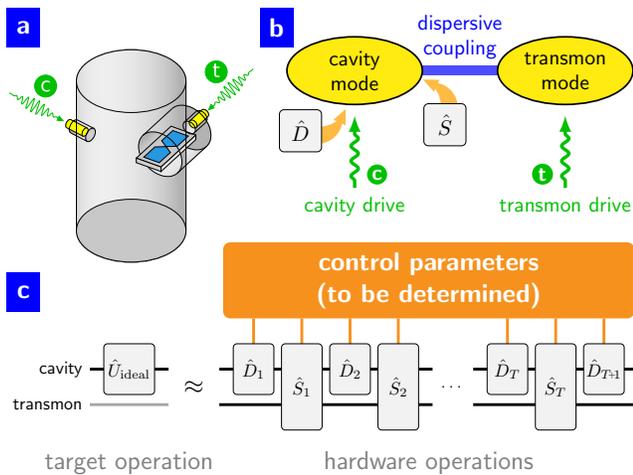


Figure 1. Setup and goal. a) Circuit QED platform: a microwave cavity (gray) coupled to a transmon qubit (blue). b) Modes and drives. The cavity and the transmon each contribute one mode to the system. The cavity mode is used to store the quantum information. The cavity drive allows to perform displacements \hat{D} on the cavity mode. Via the dispersive coupling, also the transmon drive indirectly affects the cavity mode; in this application, this is used to implement SNAP gates [11, 12] with the action $\hat{S}(\theta) = \sum_{n=0}^{\infty} e^{i\theta(n)} |n\rangle\langle n|$ on the cavity. c) Universal control schemes. An arbitrary target operation \hat{U}_{ideal} can be approximately decomposed into an alternate series of SNAP gates \hat{S} and displacements \hat{D} , as guaranteed by an existence theorem [12]. We present a practical construction scheme for such gate sequences, including their relevant control parameters. We try to achieve high fidelity while minimizing the length of the gate sequence (quantified by T , the number of SNAP gates).

values, we see that it should be possible to realize such an operation already within $\mathcal{O}(N)$ gates. This would be a quadratic speedup over the technique from [12]. For details, see appendix A 3.

In this work, we present a practical scheme to determine the control parameters for sequences of SNAP gates and displacements to match a given target operation. The results indicate that our scheme indeed follows this linear scaling. Further, we apply our technique to optimize several experimentally relevant applications. We find that, for example, correcting errors in a cavity code can be performed with 4 SNAP gates (compared to up to 50 with the old technique). Therefore, our work promises to increase the fidelity for operations in circuit QED and thereby paves the way for novel experiments.

II. PROBLEM SPECIFICATION

Our goal is to (approximately) decompose a given target operation into a series of hardware operations. In our circuit QED setup, these hardware operations are SNAP gates and displacements. As a practical way to find gate sequences with high fidelity under realistic experimental conditions, we aim to achieve a good trade-off between the following three objectives: (i) a large overlap between actual and desired output states disregarding noise and gate imperfections, (ii) a small number of gates, i.e. an efficient sequence, and (iii) low photon number at intermediate times. In this section, we will motivate this approach and make these criteria precise.

We start by introducing the two native hardware operations: SNAP gates $\hat{S}(\theta)$ and displacements $\hat{D}(\alpha)$. In this work, we will not address their hardware implementation (see [11–13] for further reading), but work directly with their resulting action on the cavity:

$$\hat{S}(\theta) = \sum_{n=0}^{\infty} e^{i\theta(n)} |n\rangle\langle n| \quad \hat{D}(\alpha) = \exp(\alpha\hat{a}^\dagger + \alpha^*\hat{a}). \quad (1)$$

These operators are idealized in the sense that noise and gate imperfections are not taken into account directly. However, we indirectly tackle these effects by searching for gate sequences with small number of gates and low photon number.

The interesting sequences are those which interleave the two gate types, as equal pairs could always be contracted to a single gate. This leads to the ansatz

$$\hat{U} = \hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1). \quad (2)$$

Here, the integer T counts the number of SNAP gates and thereby characterizes the length of the sequence. While Eq. (2) determines the sequence structure, we still need to find explicit choices for the control parameters, α_t and $\vec{\theta}_t$, of the gates in the sequence.

The action of the target operation does not need to be defined on the entire Hilbert space \mathcal{H} , but can be restricted to a logical subspace of dimension L . In the simplest case, we might only care about a single state $|\psi_{\text{in}}\rangle$ which should be transformed into another state $|\psi_{\text{out}}\rangle$. In this case, we have $L = 1$. As we can also see from this example, the input and output space do not have to coincide; here, it changes from $\text{span}(|\psi_{\text{in}}\rangle)$ to $\text{span}(|\psi_{\text{out}}\rangle)$.

In general, we have a (pre-determined) target operation defining an isometric mapping from an input space $\mathcal{X} \subseteq \mathcal{H}$ to an output space $\mathcal{Y} \subseteq \mathcal{H}$. We fix an orthonormal basis $\{|x_1\rangle, \dots, |x_L\rangle\} \subset \mathcal{X}$ of the input space, and

its image $\{|y_1\rangle, \dots, |y_L\rangle\} \subset \mathcal{Y}$ w.r.t. the target operation, which is a basis of the output space. With this, we can define the operator

$$\hat{V} := \sum_{l=1}^L |y_l\rangle\langle x_l| \quad (3)$$

to uniquely characterize the target operation (\hat{V} does not depend on the choice of the basis).

For \hat{U} , we are restricted to unitary transformations, but \hat{V} itself is not unitary unless the target operation is defined on the entire Hilbert space. However, \hat{V} can always be extended to unitary transformations which satisfy

$$\hat{U}_{\text{ideal}}|x\rangle = \hat{V}|x\rangle \in \mathcal{Y} \quad (4)$$

for all $|x\rangle \in \mathcal{X}$. The transformation $e^{i\varphi}\hat{U}_{\text{ideal}}$ would perform the target operation perfectly, where φ represents a possible global phase.

In practice, it might not be possible to implement $e^{i\varphi}\hat{U}_{\text{ideal}}$ precisely. Even assuming ideal gates, as in Eq. (1), most target operations cannot be matched exactly with a (finite) sequence of hardware operations. This error is addressed by objective (i). In a real experiment, noise and gate imperfections will further decrease the fidelity. Because these effects accumulate with each new gate, their error decreases for shorter sequences, which leads us to demand objective (ii). In addition, also the mean photon number matters, as this determines the rate of photon losses, the dominant noise process for a cavity. The photon number can change under displacements. So, even though the initial and final states are fixed, higher Fock states can get populated at intermediate points of the sequence. This cannot completely be avoided, but we can strongly suppress it. To this end, we also aim for objective (iii) to penalize high mean photon numbers.

Finally, we state explicit criteria for our three objectives. To this end, we introduce the coefficient vector $\vec{c} \in \mathbb{C}^L$ for decomposing the input (output) states into the $\{|x_l\rangle\}$ ($\{|y_l\rangle\}$) basis defined above. As a shorthand notation, we write $|x(\vec{c})\rangle = \sum_{l=1}^L c_l|x_l\rangle \in \mathcal{X}$ and $|y(\vec{c})\rangle = \sum_{l=1}^L c_l|y_l\rangle \in \mathcal{Y}$.

A reasonable figure of merit to quantify coherent errors is given by the mean overlap

$$\mathcal{F}[\hat{U}] := \max_{\varphi \in \mathbb{R}} \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \\ \|\vec{c}\|=1}} \text{Re}(\langle y(\vec{c}) | e^{i\varphi}\hat{U} | x(\vec{c}) \rangle). \quad (5)$$

We use this to quantify objective (i), with \hat{U} chosen according to Eqs. (1) and (2), i.e. neglecting noise and gate imperfections. Property (ii), to aim for efficient

sequences, simply means that T should be small. Because SNAP gates take much longer than displacements, objective (iii) can be achieved by minimizing $\sum_{t=1}^T \bar{n}_t$ where \bar{n}_t denotes the mean photon number while executing $\hat{S}(\vec{\theta}_t)$, averaged over all input states. Since SNAP gates conserve the photon number, \bar{n}_t can be computed as

$$\bar{n}_t = \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \\ \|\vec{c}\|=1}} \langle x_t(\vec{c}) | \hat{n} | x_t(\vec{c}) \rangle \quad (6)$$

where

$$|x_t(\vec{c})\rangle = \hat{D}(\alpha_t) \cdot \hat{S}(\vec{\theta}_{t-1}) \cdots \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1) |x(\vec{c})\rangle \quad (7)$$

represents the state immediately before applying $\hat{S}(\vec{\theta}_t)$.

III. RESULTS

Before we describe in Sec. IV how our technique works, we will present here results for a range of typical applications. Our intention behind this ordering is to first develop more intuition for the problem, and to show in advance some use cases which will help to understand certain design aspects for our technique.

A. Operations on a 10-dimensional Fock space

As a first benchmark for our technique, we consider transformations of the entire subspace spanned by the 10 lowest Fock states. To demonstrate generality, we choose eight target operations from different problem classes (cmp. appendix E1). For each of them, we generate gate sequences of varying length. We expect that adding more gates, i.e. more control parameters, improves the ability to approximate the target operation. The achieved mean overlap for these gate sequences is plotted in Fig. 2a. Indeed, we find that the overlap steeply increases with the sequence length. In addition, we see that the same level of fidelity can be achieved with a comparable number of gates across the different examples.

These results constitute a significant improvement compared to the technique in [12], which here could only generate sequences with at least 100 SNAP gates. For comparable examples, this yielded a mean overlap in the order of 0.999, which can be reached with ca. 10 SNAP gates by our technique, i.e. with ca. 10 times less gates!

In an experiment, one would choose a sequence length that minimizes the overall infidelity, including noise and gate imperfections. Because these effects grow with each new gate, the optimum is shifted towards shorter

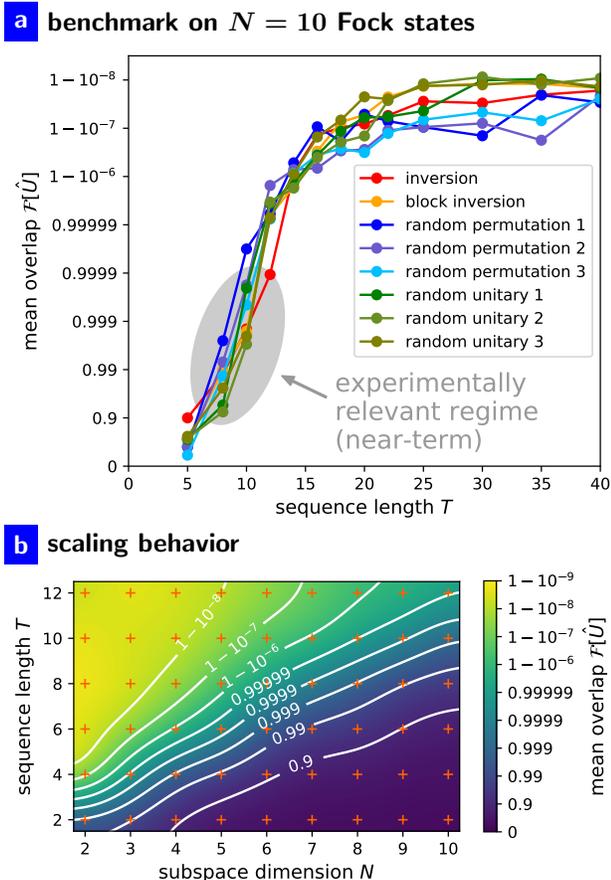


Figure 2. Results for operations on Fock subspaces. a) Benchmark for our technique, showing significantly reduced sequence lengths. The target operations are unitary transformations on the $N = 10$ lowest Fock states. We apply our scheme to these target operations at different sequence lengths T , and plot the achieved mean overlap \mathcal{F} as a function of T . We see that longer sequences yield higher values for \mathcal{F} , which is computed disregarding noise and gate imperfections. We highlight the regime with a good overall trade-off considering these effects. To improve visibility, the y axis is scaled such that the error $1 - \mathcal{F}$ decreases logarithmically. b) Scaling behavior. The colors indicate how well sequences with T SNAP gates can approximate target operations on the N lowest Fock states. The isocurves roughly follow lines through the origin, which indicates that the sequence length for our method scales linearly in N . The values at the points marked with a cross are based on simulations (see Sec. III B), and we use bicubic interpolation to extend this discrete grid of points to the entire plane.

sequences. For an error probability per SNAP gate in a broad range between 0.01% and 1%, the best trade-off is achieved around $T = 8$ to $T = 10$. This regime is highlighted as “experimentally relevant” in Fig. 2a.

We observe that the mean overlap saturates at a level between $1 - 10^{-7}$ and $1 - 10^{-9}$. Similar behavior will also occur in the following examples. Because we are not aware of any reason for this fidelity regime to be the ultimate limit, we assume that the saturation is a property of our numerical scheme. However, this question is of minor practical relevance: for a gate sequence with this precision, the other error sources like noise and gate imperfections will become the limiting factors for the overall fidelity, so further optimization of this contribution would effectively not help. The saturation would only become a bottleneck for an error probability per gate of less than 10^{-7} .

Besides the sequence length and the fidelity, the third important property of a gate sequence is the photon number, which determines the rate of photon losses. To improve this aspect, we have equipped our technique with a mechanism to prefer low photon numbers during optimization (cmp. Sec. IV B). In order to quantify its success, we consider $\sum_{t=1}^T \bar{n}_t$ where \bar{n}_t is the mean photon number while executing the t -th SNAP gate (see Eq. (6)). For our examples, we find that $\frac{1}{T} \sum_{t=1}^T \bar{n}_t$ ranges from 4.51 to 5.76, and is less than 5.0 in 80% of the cases. For comparison, the minimum value for \bar{n}_t to store 10 states is 4.5.

B. Scaling behavior

In the introduction, we have argued that for fixed fidelity, the minimum sequence length, characterized by T , should scale linearly with the subspace dimension N . We now analyse the scaling behavior for the gate sequences generated by our technique.

To this end, we choose three random unitary transformations of the N lowest Fock states as the target operations, for each value of N between 2 and 10. The target operations for $N = 10$ coincide with the random unitaries from Fig. 2a. We apply our technique to these target operations for all even sequence lengths T between 2 and 12. Then, we compute their mean overlap \mathcal{F} , and determine for all tuples (N, T) an average value according to¹ $-\ln(1 - \bar{\mathcal{F}}) = \langle -\ln(1 - \mathcal{F}_j) \rangle_j$. The result is plotted in Fig. 2b.

Fig. 2b indicates that, at least for this class of operations, the sequence length to achieve similar overlap indeed grows linearly with N , since the isocurves roughly follow lines through the origin. This linear scaling behavior coincides with the expectation for the minimally

¹ The linear average gives a similar picture, but is more prone to random deviations.

required sequence length according to our parameter-counting argument. However, this does not yet mean that our sequences are optimum (or close-to-optimum), as there could be a difference in the constant of proportionality.

Based on our results, we can compute a rough estimate for the mean overlap as $\mathcal{F} \approx 1 - e^{-6.49 \cdot (T/N)^{1.91}}$ (cmp. Fig. S2). This relation is an empirical observation from the data points, but we do not have a model to predict this behavior.

C. Applied problems

We will now consider several experimentally relevant problems for quantum information processing in cavities, and demonstrate how our technique can be beneficial for them.

In circuit QED, the quantum information is typically stored in the cavity. In the simplest case, we have a single logical qubit which is encoded in the two lowest Fock states:

$$|\psi_{\alpha\beta}^{(\text{triv})}\rangle = \alpha|0\rangle + \beta|1\rangle. \quad (8)$$

We refer to this as the trivial code. This code has the issue that it is very vulnerable to noise. To overcome this problem, several error-correction codes have been developed in the recent years. Among the most important codes for cavities are the binomial (or *kitten*) codes[15, 16]. In this work, we will specifically consider the binomial code $|\psi_{\alpha\beta}^{(\text{bin})}\rangle = \alpha|b_0\rangle + \beta|b_1\rangle$ with

$$|b_0\rangle = \frac{|0\rangle + \sqrt{3}|6\rangle}{2} \quad |b_1\rangle = \frac{\sqrt{3}|3\rangle + |9\rangle}{2} \quad (9)$$

which protects against single- and double-photon loss. In principle, this code could also tolerate a single dephasing error, but this noise channel plays only a negligible role if the cavity evolves freely.

The first application we consider is how to prepare the cavity in one of these binomial code states $|\psi_{\alpha\beta}^{(\text{bin})}\rangle$, starting from the vacuum state $|0\rangle$. This is an example where the transformation of only one state matters, so we have a one-dimensional logical subspace ($L = 1$). The achieved mean overlap for preparing several superpositions of $|0\rangle$ and $|1\rangle$ is shown in Fig. 3a, as a function of the sequence length. We see that high fidelity can be reached within only 3 or 4 SNAP gates!

As a second application, we consider recovery operations. The main benefit of an error-correction code is to protect the quantum information against certain noise events, in the case of the binomial code from Eq. (9) against up to two photon losses. To avoid that over

time more of these events accumulate than which can be corrected for, the error syndrome (here: photon number modulo 3) has to be monitored frequently. If such a syndrome detection reveals a single- or double-photon loss, the corresponding number of photons needs to be re-pumped to return into the original, robust code. Even if no photon loss is detected, a small correction needs to be applied to counter the effect of amplitude decay, i. e. that the amplitudes for higher Fock states decay faster (cmp. appendix E2). Hence, an active recovery operation is required in all three cases. This operation must not depend on the logical state, which is unknown and not allowed to be measured, so a two-dimensional subspace needs to be transformed ($L = 2$). Again, our technique can be used to generate efficient gate sequences, as shown in Fig. 3b: 4 SNAP gates are sufficient to perform these operations with high fidelity.

Another related task is to perform logical operations, i. e. to effect a unitary transformation of the coefficients α and β , here on the binomial code $\alpha|b_0\rangle + \beta|b_1\rangle$. As for the previous problem, the operation must not depend on the incoming logical state, so again the logical subspace is two-dimensional ($L = 2$). A special case are rotations around the logical z axis which can be realized with a single SNAP gate. All other operations, however, require more complex control schemes. Fig. 3c shows the results for five selected examples. We reach good fidelity with 4, in one case 3, SNAP gates.

We observe that all these operations can be realized with between 3 and 4 SNAP gates using our technique. This value is exceeded considerably by the *minimum* sequence length for the method from [12] in almost every case, as shown in appendix E3. The largest improvement can be obtained for the error correction examples (Fig. 3b), which are also the most time-critical operations. Here, the minimum number of SNAP gates for the technique in [12] would be 16 for the syndrome $\mathbb{1}$, and 49 for the syndroms \hat{a} and \hat{a}^2 .

In addition, the sequences here are significantly shorter than what would be expected from Sec. III B, taking into account their respective number of involved Fock states. This is possible because here only one or two states within this subspace are actually relevant. In these cases, we benefit from our scheme to focus solely on the action on the input states (cmp. Sec. IV C).

The examples so far show how our technique can be used to efficiently implement the relevant operations for one binomial code. However, similar problems occur also under other circumstances. Recovery operations as in Fig. 3b are needed for every cavity code that requires active error correction. Logical operations as in Fig. 3c occur for every encoding, whether error-corrected or not. State preparation as in Fig. 3a has many applications even outside of quantum information processing. Our

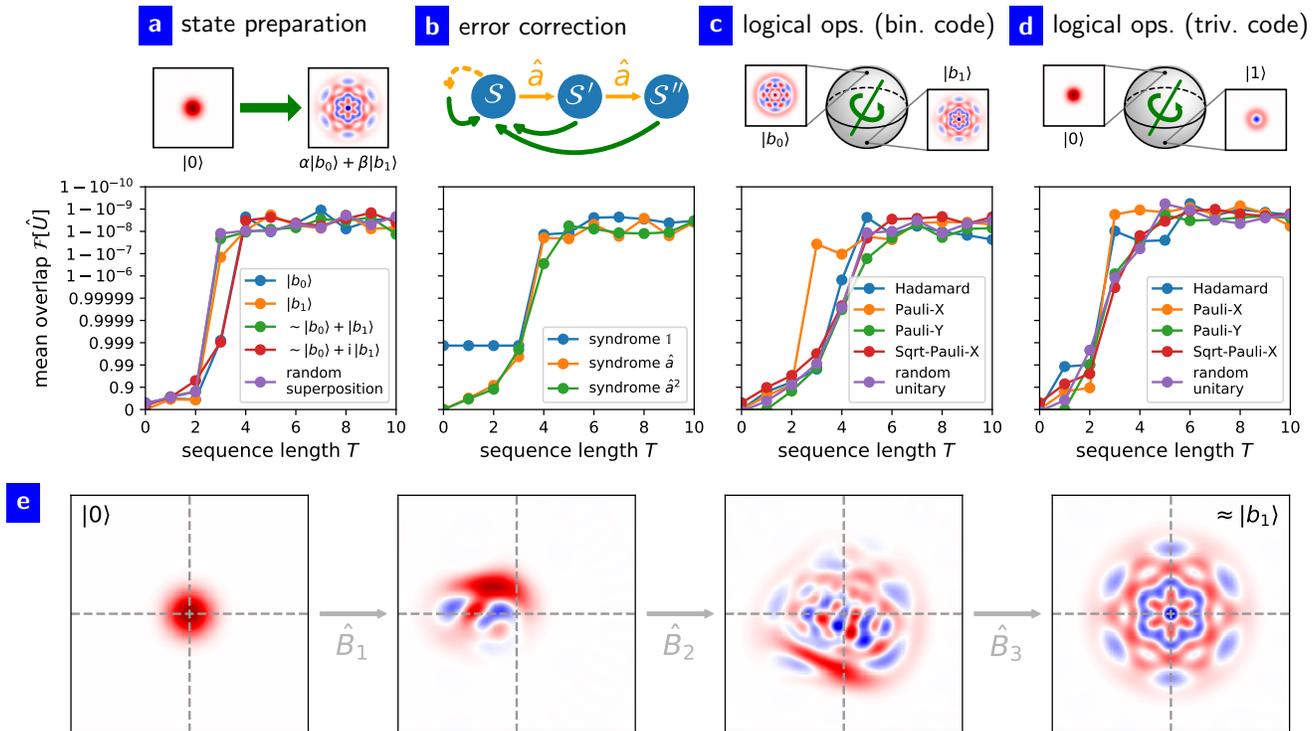


Figure 3. Results for experimentally relevant applications, showing that high fidelity for these operations can be achieved with a remarkably small number of gates. a) State preparation. A binomial code state $\alpha|b_0\rangle + \beta|b_1\rangle$ is generated from the vacuum state $|0\rangle$. b) Error correction for the binomial code $\alpha|b_0\rangle + \beta|b_1\rangle$. After a syndrome detection, the system needs to be transformed back into the code space $\mathcal{S} = \text{span}(|b_0\rangle, |b_1\rangle)$; also the $\mathbb{1}$ syndrome requires a small correction, see main text. c) Logical operations on the binomial code $\alpha|b_0\rangle + \beta|b_1\rangle$. The logical state, encoded in the coefficients α and β , is modified via a unitary transformation of (α, β) . This operation must be performed without knowledge of the logical state. d) Logical operations on the trivial code $\alpha|0\rangle + \beta|1\rangle$. As in (c), but for a different encoding. The upper row in (a)-(d) illustrates the target operation, and the lower row plots the achieved fidelities (disregarding noise and gate imperfections) as a function of the sequence length. e) Evolution of the Wigner density for preparation of the binomial code state $|b_1\rangle$ (cmp. a). The pictures display snapshots between different building blocks (cmp. Eq. (10)). With every step, the structure of the state gains complexity, and also the volume in phase space is growing.

technique is also useful for those cases, and by no means restricted to the binomial code.

As an example, we consider one more application, namely logical operations on the trivial code as in Eq. (8). We consider the same logical operations as in Fig. 3c, but on the code $\alpha|0\rangle + \beta|1\rangle$ instead of $\alpha|b_0\rangle + \beta|b_1\rangle$. The results, as plotted in Fig. 3d, indicate that sequences with 3 SNAP gates can perform these operations with sufficiently high fidelity.

IV. TECHNIQUE

After presenting the results, we now turn to a description of our technique. It consists of two parts: an initialization step, followed by a finetuning step. Whereas

this basic concept has already been used in [12], we have fundamentally redesigned the approaches towards both of these steps. These changes enable us to arrive at significantly better results.

In the initialization step, we use a heuristic strategy to generate a gate sequence which roughly approximates the target operation. The goal is not yet to achieve high fidelity, as eventually required in the experiment, and usually we obtain just a coarse approximation (cmp. Fig. 4b). Instead, we delegate the responsibility for optimizing the fidelity mainly into the finetuning step, and are here content with providing a reasonable starting point for this. In contrast to the fidelity, the length of the sequence will not be changed during finetuning. Our new initialization scheme can construct suitable sequences of arbitrary length, in particular also at rela-

tively small numbers of gates. This is a central advantage compared to its counterpart in [12] (which is based on $SO(2)$ rotations).

In the finetuning step, we use gradient descent to co-optimize the control parameters of the gate sequence. This is necessary to reach the high-fidelity regime (cmp. Fig. 4c) which is needed for application in an experiment. We require from this process that it can be launched

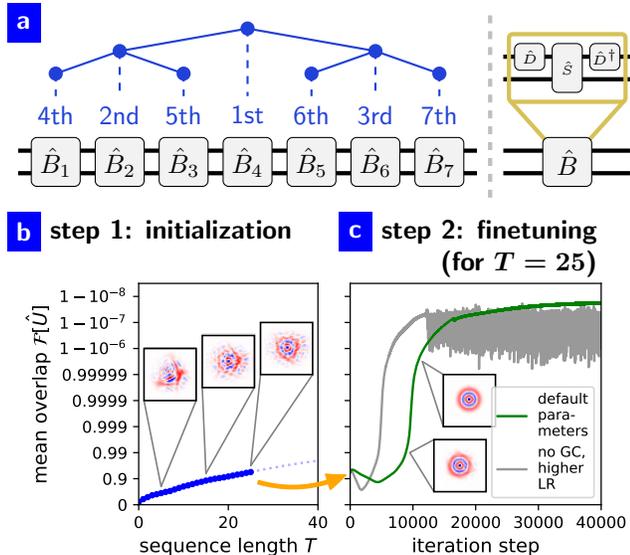


Figure 4. Construction process. a) Hierarchical insertion strategy for the initialization, which lets the fidelity improve progressively during this whole process. The order in which the building blocks $\hat{B}_t = \hat{D}^\dagger(\alpha_t)\hat{S}(\vec{\theta}_t)\hat{D}(\alpha_t)$ are inserted corresponds to the breadth-first traversal of a binary tree (here until sequence length $T = 7$). b) Progress during initialization, from an empty sequence towards a coarse approximation of the target operation. The curve shows how the mean overlap \mathcal{F} improves by inserting new gates into the sequence. c) Progress during finetuning, showing the attainment of high fidelity. After stopping initialization, here at sequence length $T = 25$, we further improve the fidelity by co-optimizing the control parameters. With our default hyperparameters (green curve), the fidelity grows steeply within the first 15000 iterations (after a first phase where mainly the photon number is optimized), surpassing a value of $1 - 10^{-7}$. Afterwards, the improvement gets much slower. For comparison, the gray curve shows the progress without gradient clipping and for a higher learning rate. While this parameter configuration learns faster in the beginning, it becomes unstable in the high-fidelity regime. However, post-selecting good intermediate values leads to comparable overall results as with the default parameters. The exact behavior can vary from case to case. The insets in (b) and (c) show the final state that should ideally be mapped to the state $|3\rangle$. The target operation is “random unitary 3” from Fig. 2a.

successfully from the starting point provided by the initialization step, for which the mean overlap is typically between 0.75 and 0.995.

During both steps, we will work with building blocks of the form

$$\hat{B}(\alpha, \vec{\theta}) = \hat{D}^\dagger(\alpha)\hat{S}(\vec{\theta})\hat{D}(\alpha) \quad (10)$$

where $\hat{D}(\alpha)$ denotes a cavity displacement and $\hat{S}(\vec{\theta})$ a SNAP gate. The direct output of our technique are sequences of these building blocks. To implement them on hardware, we can easily reparameterize them in terms of SNAP gates and displacements (cmp. appendix A 2). This equivalence also guarantees these building blocks to be computationally universal.

A. Initialization

As motivated above, the goal for the initialization step is rather to approximate the target operation than to aim directly for very high fidelity. The more important aspect is the flexibility regarding the length of the sequence. In the following, we describe our initialization scheme which starts from an empty sequence and expands it successively with the building blocks from Eq. (10).

We first consider, for simplicity, the special case of a gate sequence with only one building block $\hat{B}(\alpha, \vec{\theta})$. In this situation, we can directly compute the parameters α and $\vec{\theta}$ which maximize the mean overlap $\mathcal{F}[\hat{B}(\alpha, \vec{\theta})]$ as

$$\alpha = \operatorname{argmax}_{\tilde{\alpha}} \left[\sum_n |g_n(\tilde{\alpha})| \right] \quad \theta^{(n)} = \operatorname{arg}(g_n(\alpha)) \quad (11)$$

where $g_n(\alpha) = \langle n | \hat{D}(\alpha) \hat{V} \hat{D}^\dagger(\alpha) | n \rangle$ and \hat{V} encodes the target operation (cmp. Eq. (3)). For the derivation, see appendix C.

With a slight modification, we can use this scheme also to successively extend an existing gate sequence. In each step, one building block \hat{B}_t is inserted between the already fixed segments $\hat{B}_T, \dots, \hat{B}_{t+1}$ and $\hat{B}_{t-1}, \dots, \hat{B}_1$. We choose \hat{B}_t such that it greedily maximizes the mean overlap $\mathcal{F}[\hat{B}_T \dots \hat{B}_1]$ of the extended sequence. Afterwards, also this \hat{B}_t remains fixed for the rest of this procedure, and we continue with inserting the next building block (for a new subdivision into the two segments) until the desired sequence length is reached.

The greedy optimization of $\hat{B}_t = \hat{B}(\alpha_t, \vec{\theta}_t)$ is an equivalent problem to matching the input states $\hat{B}_{t-1} \dots \hat{B}_1 |x(\vec{c})\rangle$ to the output states $\hat{B}_{t+1}^\dagger \dots \hat{B}_T^\dagger |y(\vec{c})\rangle$ with a single building block. Hence, by replacing \hat{V} with $\hat{B}_{t+1}^\dagger \dots \hat{B}_T^\dagger \cdot \hat{V} \cdot \hat{B}_1 \dots \hat{B}_{t-1}$, we can reuse our solution for this situation: the parameters α_t and $\vec{\theta}_t$ can be

determined according to Eq. (11), where

$$g_n(\alpha) = \langle n | \hat{D}(\alpha) \cdot \hat{B}_{t+1}^\dagger \cdots \hat{B}_T^\dagger \cdot \hat{V} \cdot \hat{B}_1^\dagger \cdots \hat{B}_{t-1}^\dagger \cdot \hat{D}^\dagger(\alpha) | n \rangle. \quad (12)$$

This choice guarantees to never reduce the fidelity because there is always the trivial option $\hat{B}_t = \mathbb{1}$, so the old value for the fidelity is a lower bound for the new one. In practice, we find that typically the fidelity gradually increases with every new building block (cmp. Fig. 4b).

One degree of freedom in this procedure is the order in which the building blocks are inserted. We find that for a linear arrangement, i.e. building the sequence either from start to end or from end to start, the improvement of the fidelity would stagnate after a few steps because only operations close to identity will be inserted. Besides its lower fidelity, such a sequence has the issue that it does not at all make efficient use of its expressive power, which impedes subsequent finetuning. This problem can be circumvented by inserting the building blocks hierarchically between the previously added ones (cmp. Fig. 4a). We follow this strategy during the initialization for all results shown in Sec. III.

For one example, the recovery of the $\mathbb{1}$ syndrome for a binomial code (cmp. Fig. 3b), we have to slightly deviate from the scheme as described so far. Because all $g_n(\alpha)$ take real values here, we would run into the issue that only identity operations would be inserted. We break this blockade by starting the initialization with one randomly chosen building block and then continue as usual.

For this strategy of expanding sequences stepwise, it is very advantageous to have the building blocks parameterized as $\hat{D}^\dagger(\alpha) \hat{S}(\vec{\theta}) \hat{D}(\alpha)$. In particular, such building blocks do not generate a net displacement, in contrast to other parameterizations like $\hat{S}(\vec{\theta}) \hat{D}(\alpha)$, and therefore already isolated building blocks can be useful for the important case of displacement-free target operations. The importance of choosing a good parameterization has been observed in other platforms before [17].

B. Finetuning

After the initialization, we perform an additional finetuning step where all the control parameters in the gate sequence (displacements α_t and SNAP angles $\vec{\theta}_t$) are co-optimized. We accomplish this using gradient descent, using the cost function

$$C(\alpha_1, \vec{\theta}_1, \dots, \alpha_T, \vec{\theta}_T) = \ln(1 - \mathcal{F}[\hat{B}_T \cdots \hat{B}_1]) + \lambda \sum_{t=1}^T \frac{\bar{n}_t + \bar{n}'_t}{2} \quad (13)$$

where \mathcal{F} is the mean overlap from Eq. (5), \bar{n}_t is the mean photon number while executing $\hat{S}(\vec{\theta}_t)$ as defined

in Eq. (6), and \bar{n}'_t is the corresponding quantity for the inverse sequence starting from the desired output states $|y(\vec{c})\rangle$ (see appendix B 1).

The term $\ln(1 - \mathcal{F}[\hat{B}_T \cdots \hat{B}_1])$ is responsible for matching the target operation. We use this logarithmic expression instead of $-\mathcal{F}[\hat{B}_T \cdots \hat{B}_1]$ to avoid that the gradient diminishes in the high-fidelity regime, i.e. when $\mathcal{F}[\hat{B}_T \cdots \hat{B}_1]$ gets close to 1. In practice, this significantly accelerates the training process.

The term $\sum_{t=1}^T \frac{\bar{n}_t + \bar{n}'_t}{2}$ penalizes larger photon numbers in the cavity, which would amplify the rate of photon losses. Without this contribution to the cost function, or if its coefficient λ was chosen too small, higher Fock states would get populated stronger than necessary during the gate sequence. Increasing λ helps to reduce this effect, and often even generates sequences with a higher mean overlap \mathcal{F} . However, for λ too large, reducing the photon number would become at some point more favorable than matching the target operation. We have always found good intermediate values for λ with high fidelity and low photon number, but this range can vary from example to example (for us, between $\lambda = 0.16$ and $\lambda = 2.4$, cmp. appendix F 1).

We compute the parameter updates from the gradient according to the Adam method [18]. In addition, we use gradient clipping to stabilize the training process (cmp. Fig. 4c). These are two well-established concepts from the field of machine learning, and they help us to reach convergence with a significantly reduced number of iterations. To keep the runtime of the simulations moderate, we need in addition an efficient numerical scheme to compute the gradient (cmp. appendix D 1).

C. Discussion

Our scheme is carefully designed such that it focuses only on the relevant aspect of the operation, namely on the mapping between the specified input and output states, and does on purpose not consider the action on any other state. This goes further than merely ignoring the action on Fock states which have no overlap with any input or output state: The previous scheme from [12] requires to fix explicitly a unitary on the full subspace of involved Fock states, even if the aim was only to transform a single input state into one desired output state. In contrast, our approach keeps the number of conditions minimal, which allows to improve the behavior on the relevant states. On the technical side, this means that we must not use an explicit extension of \hat{V} (from Eq. (3)) to a unitary operator \hat{V}_{ext} , e.g. for a figure of merit like $|\text{tr}[\hat{V}_{\text{ext}}^\dagger \hat{U}]|$. Instead, both during initialization and finetuning, we maximize $\mathcal{F}[\hat{U}] \sim |\text{tr}[\hat{V}^\dagger \hat{U}]|$ (cmp.

appendix B 1) which is directly based on \hat{V} . A similar strategy is followed by [9] in their optimal-control pulse engineering.

Our technique involves only a few hyperparameters. The only one which is changed between the different examples is λ , the coefficient for the photon number cost (see Sec. IV B). Depending on the application, the range of good values can span a factor between 2 and more than 100. For all other hyperparameters, the same values are used for all results in Sec. III. All hyperparameter values are listed in appendix F 1.

Even though we do not employ neural networks, we can still benefit a lot from machine-learning techniques: our finetuning approach makes use of the Adam optimizer, gradient clipping, and a backpropagation-like scheme to compute the gradient. In this sense, our work is related to previous applications of machine learning for problems related to quantum computing [19–21], and especially to quantum control [22–25]. However, there is still a big difference in the concrete methodology, as these works use neural networks and/or reinforcement learning. Another aspect that sets us apart (from the mentioned works on quantum control) is the considered control problem: their goal is to engineer pulses for single gates, whereas here we focus on building sequences by combining multiple gates.

V. CONCLUSION

In this work, we have presented a practical scheme to construct efficient sequences of SNAP gates and displacements for given target operations. We have demonstrated that our scheme works robustly for a broad range of applications, including tasks of direct experimental relevance. Moreover, our scheme is directly applicable for experimentalists: all necessary gates have already been demonstrated, and our technique outputs explicitly their relevant control parameters.

Like other gradient-descent techniques, our technique likely does not find the global optimum, and in principle we cannot even be sure that it finds close-to-optimal solutions. Nevertheless, in practice, our scheme can generate much shorter gate sequences compared to previously known techniques, with higher fidelity. The consequence

of this improvement is that realizing universal control over cavities with SNAP gates becomes experimentally feasible: As we have demonstrated above, we can obtain accelerations by one order of magnitude, which turns out to take us into the regime where gate times are sufficiently short given the coherence of current devices. In this way, our work also reveals the power and flexibility of the SNAP gate when used properly in combination with displacements.

In our simulations, we have not considered directly the effect of gate imperfections and noise. However, our technique is able to considerably reduce the sequence length; this already mitigates their influence very effectively. A more rigorous approach is to include these effects into the simulation, which is subject to future work.

Our work is compatible with the error-transparency scheme in [13, 14], from which we can gain robustness against transmon decay to first order. Accelerating the SNAP gates while correcting for coherent errors, e. g. using a Magnus-based approach [26], could further improve the overall fidelity.

We have applied our technique explicitly for sequences of SNAP gates and displacements, due to their direct relevance for circuit QED experiments. However, variations of our technique may be used to optimize control schemes for other platforms, too. Moreover, our combination of initialization and finetuning is a more general concept which has possible applications even outside of quantum control, namely for any continuous control problem on discrete time steps without feedback.

Data and code availability The data and the code that support the plots within this paper will be made publicly available. Other findings of this study are available from the corresponding author on request.

Acknowledgements We thank Hugo Ribeiro, Thales Figueiredo Roque, Wenlong Ma and Christopher Wang for fruitful discussions. T.F. thanks Yale Quantum Institute for hosting during his research stay. L.J. acknowledges supports from the ARL-CDQI (W911NF-15-2-0067), ARO (W911NF-18-1-0020, W911NF-18-1-0212), ARO MURI (W911NF-16-1-0349), AFOSR MURI (FA9550-15-1-0015, FA9550-19-1-0399), DOE (DE-SC0019406), NSF (EFMA-1640959, OMA-1936118), and the Packard Foundation (2013-39273).

-
- [1] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, Cambridge ; New York, anniversary edition edition, January 2011.
- [2] Thaddeus D Ladd, Fedor Jelezko, Raymond Laflamme, Yasunobu Nakamura, Christopher Monroe, and

- Jeremy Lloyd O’Brien. Quantum computers. *Nature*, 464(7285):45–53, 2010.
- [3] RJ Schoelkopf and SM Girvin. Wiring up quantum systems. *Nature*, 451(7179):664–669, 2008.
- [4] Michel H Devoret and Robert J Schoelkopf. Superconducting circuits for quantum information: an outlook.

- Science*, 339(6124):1169–1174, 2013.
- [5] G Wendin. Quantum information processing with superconducting circuits: a review. *Reports on Progress in Physics*, 80(10):106001, 2017.
- [6] Nissim Ofek, Andrei Petrenko, Reinier Heeres, Philip Reinhold, Zaki Leghtas, Brian Vlastakis, Yehan Liu, Luigi Frunzio, SM Girvin, L Jiang, et al. Extending the lifetime of a quantum bit with error correction in superconducting circuits. *Nature*, 536(7617):441, 2016.
- [7] Christopher J Axline, Luke D Burkhardt, Wolfgang Pfaff, Mengzhen Zhang, Kevin Chou, Philippe Campagne-Ibarcq, Philip Reinhold, Luigi Frunzio, SM Girvin, Liang Jiang, et al. On-demand quantum state transfer and entanglement between remote microwave cavity memories. *Nature Physics*, 14(7):705, 2018.
- [8] Kevin S Chou, Jacob Z Blumoff, Christopher S Wang, Philip C Reinhold, Christopher J Axline, Yvonne Y Gao, Luigi Frunzio, MH Devoret, Liang Jiang, and RJ Schoelkopf. Deterministic teleportation of a quantum gate between two logical qubits. *Nature*, 561(7723):368, 2018.
- [9] Reinier W Heeres, Philip Reinhold, Nissim Ofek, Luigi Frunzio, Liang Jiang, Michel H Devoret, and Robert J Schoelkopf. Implementing a universal gate set on a logical qubit encoded in an oscillator. *Nature Communications*, 8(1):94, 2017.
- [10] Navin Khaneja, Timo Reiss, Cindie Kehlet, Thomas Schulte-Herbrüggen, and Steffen J Glaser. Optimal control of coupled spin dynamics: design of nmr pulse sequences by gradient ascent algorithms. *Journal of Magnetic Resonance*, 172(2):296–305, 2005.
- [11] Reinier W Heeres, Brian Vlastakis, Eric Holland, Stefan Krastanov, Victor V Albert, Luigi Frunzio, Liang Jiang, and Robert J Schoelkopf. Cavity state manipulation using photon-number selective phase gates. *Physical Review Letters*, 115(13):137002, 2015.
- [12] Stefan Krastanov, Victor V Albert, Chao Shen, Chang-Ling Zou, Reinier W Heeres, Brian Vlastakis, Robert J Schoelkopf, and Liang Jiang. Universal control of an oscillator with dispersive coupling to a qubit. *Physical Review A*, 92(4):040303, 2015.
- [13] Philip Reinhold, Serge Rosenblum, Wen-Long Ma, Luigi Frunzio, Liang Jiang, and Robert J Schoelkopf. Error-corrected gates on an encoded qubit. *arXiv preprint arXiv:1907.12327*, 2019.
- [14] Wen-Long Ma, Mengzhen Zhang, Yat Wong, Kyungjoo Noh, Serge Rosenblum, Philip Reinhold, Robert J Schoelkopf, and Liang Jiang. Path-independent quantum gates with noisy ancilla. *arXiv preprint arXiv:1911.12240*, 2019.
- [15] Marios H Michael, Matti Silveri, RT Brierley, Victor V Albert, Juha Salmilehto, Liang Jiang, and Steven M Girvin. New class of quantum error-correcting codes for a bosonic mode. *Physical Review X*, 6(3):031006, 2016.
- [16] L Hu, Y Ma, W Cai, X Mu, Y Xu, W Wang, Y Wu, H Wang, YP Song, C-L Zou, et al. Quantum error correction and universal gate set operation on a binomial bosonic logical qubit. *Nature Physics*, 15(5):503, 2019.
- [17] Daniel Zeuch and NE Bonesteel. Efficient two-qubit pulse sequences beyond cnot. *arXiv preprint arXiv:2001.09341*, 2020.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Giacomo Torlai and Roger G Melko. Neural decoder for topological codes. *Physical review letters*, 119(3):030501, 2017.
- [20] Alexey A Melnikov, Hendrik Poulsen Nautrup, Mario Krenn, Vedran Dunjko, Markus Tiersch, Anton Zeilinger, and Hans J Briegel. Active learning machine learns to create new quantum experiments. *Proceedings of the National Academy of Sciences*, 115(6):1221–1226, 2018.
- [21] Thomas Fösel, Petru Tighineanu, Talitha Weiss, and Florian Marquardt. Reinforcement learning with neural networks for quantum feedback. *Physical Review X*, 8(3):031084, 2018.
- [22] Marin Bukov, Alexandre GR Day, Dries Sels, Phillip Weinberg, Anatoli Polkovnikov, and Pankaj Mehta. Reinforcement learning in different phases of quantum control. *Physical Review X*, 8(3):031086, 2018.
- [23] Moritz August and José Miguel Hernández-Lobato. Taking gradients through experiments: Lstms and memory proximal policy optimization for black-box quantum control. In *International Conference on High Performance Computing*, pages 591–613. Springer, 2018.
- [24] Murphy Yuezheng Niu, Sergio Boixo, Vadim N Smelyanskiy, and Hartmut Neven. Universal quantum control through deep reinforcement learning. *npj Quantum Information*, 5(1):1–8, 2019.
- [25] Riccardo Porotti, Dario Tamascelli, Marcello Restelli, and Enrico Prati. Coherent transport of quantum states by deep reinforcement learning. *Communications Physics*, 2(1):1–9, 2019.
- [26] Thales Figueiredo Roque, Aashish A Clerk, and Hugo Ribeiro. Engineering fast high-fidelity quantum operations with constrained interactions. *arXiv preprint arXiv:2003.12096*, 2020.

Appendix A: Gate sequences

1. Ansatz

In order to motivate the ansatz in Eq. (2), we start with the observation that two (or multiple) subsequent SNAP gates can be contracted into a single one, and the same statement holds also for displacements:

$$\hat{S}(\vec{\theta}_2) \cdot \hat{S}(\vec{\theta}_1) = \hat{S}(\vec{\theta}_1 + \vec{\theta}_2) \qquad \hat{D}(\alpha_2) \cdot \hat{D}(\alpha_1) = \hat{D}(\alpha_1 + \alpha_2).$$

Hence, the most efficient sequences are those which always alternate between these two gate types:

$$\hat{U} = \dots \cdot \hat{S}(\vec{\theta}_{t+1}) \cdot \hat{D}(\alpha_{t+1}) \cdot \hat{S}(\vec{\theta}_t) \cdot \hat{D}(\alpha_t) \cdot \dots$$

The sequence structure is not yet characterized uniquely, as we can still decide whether to start, and also to end, such a sequence either with a SNAP gate or with a displacement. Because both the execution time and the gate errors for displacements are negligible compared to the corresponding values for SNAP gates, it makes sense to both start and end with a displacement, as these operations add a new degree of freedom to the sequence at very low costs:

$$\hat{U} = \hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1)$$

The control parameters α_t can be restricted to real values because any complex phase can be absorbed into the neighboring SNAP gates, as implied by the relation $\hat{D}(e^{i\varphi}\alpha) = \hat{S}(\vec{\theta}_{\text{rot}}(\varphi)) \cdot \hat{D}(\alpha) \cdot \hat{S}(-\vec{\theta}_{\text{rot}}(\varphi))$ with $\theta_{\text{rot}}^{(n)}(\varphi) = n\varphi$. This argument applies to all displacement gates except for the first and the last one. In total, the full gate sequence loses only a single degree of freedom by the restriction to real α_t .

2. Reparameterization using building blocks

Here, we justify the usage of the building blocks $\hat{B}(\alpha, \vec{\theta}) = \hat{D}^\dagger(\alpha)\hat{S}(\vec{\theta})\hat{D}(\alpha)$ as defined in Eq. (10). These building blocks have several practical advantages over working directly with the native hardware gates $\hat{S}(\vec{\theta})$ and $\hat{D}(\alpha)$ during the construction of the gate sequence.

We can rewrite the gate sequence from Eq. (2) as

$$\begin{aligned} & \hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1) = \\ & = \hat{D}(\alpha'_{T+1}) \cdot \hat{B}(\alpha'_T, \vec{\theta}_T) \cdot \dots \cdot \hat{B}(\alpha'_2, \vec{\theta}_2) \cdot \hat{B}(\alpha'_1, \vec{\theta}_1) \end{aligned}$$

with $\alpha'_t = \sum_{s=1}^t \alpha_s$.

A displacement operation $\hat{D}(\alpha)$ can be composed using two building blocks as

$$\hat{D}(\alpha) = \hat{B}\left(-\frac{\alpha - \delta\alpha}{4}, \vec{\theta}_{\text{rot}}(\pi)\right) \cdot \hat{B}\left(\frac{\alpha + \delta\alpha}{4}, \vec{\theta}_{\text{rot}}(\pi)\right)$$

for $\vec{\theta}_{\text{rot}}^{(n)}(\pi) = n\pi$ and any $\delta\alpha \in \mathbb{C}$ (set to 0 in the next equation). This relation can be verified by expanding the right-hand side and using $\hat{S}(\vec{\theta}_{\text{rot}}(\pi))\hat{D}(\alpha)\hat{S}(\vec{\theta}_{\text{rot}}(\pi)) = \hat{D}(-\alpha)$. This consideration guarantees that any sequence of SNAP gates and displacements can be exactly reparameterized with these building blocks:

$$\begin{aligned} & \hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1) = \\ & = \hat{B}\left(-\alpha'_{T+1}/4, \vec{\theta}_{\text{rot}}(\pi)\right) \cdot \hat{B}\left(\alpha'_{T+1}/4, \vec{\theta}_{\text{rot}}(\pi)\right) \cdot \hat{B}(\alpha'_T, \vec{\theta}_T) \cdot \dots \cdot \hat{B}(\alpha'_2, \vec{\theta}_2) \cdot \hat{B}(\alpha'_1, \vec{\theta}_1). \end{aligned}$$

For our purposes, however, it is more practical to put the mild constraint that the sum of all displacements vanishes, i. e. $\alpha'_{T+1} = \sum_{s=1}^{T+1} \alpha_s \stackrel{!}{=} 0$. This condition fixes only one degree of freedom of the gate sequence, and hence does not take away much of its expressive power.

3. Length estimate

Here, we provide more details on the statement in the introduction that the optimum sequence length should scale linear with the number of involved Fock states.

This statement refers primarily to the situation where an entire Fock subspace is transformed, i. e. if the input space coincides with the output space, and this subspace is the span of Fock states (and not just embedded in the span of all Fock states with non-vanishing overlap), like for the examples in Fig. 2. If this is not the case, like for the examples in Fig. 3, a length estimate becomes more difficult. For instance, also the dimension of the logical subspace will play a role then. In this case, the estimate for the transformation of the entire Fock subspace needs to be understood as an upper bound.

Focusing on the simple situation, our parameter-counting argument, here in more detail, works as follows: Searching for a gate sequence to realize a target operation can be interpreted as solving a high-dimensional (non-linear) system of equations. We can find a solution if the number of constraints does not exceed the number of free variables. There are as many constraints as angles to parameterize an operation on N Fock states, which is $\dim(\text{SU}(N)) = N^2 - 1 = \mathcal{O}(N^2)$. The free variables are the control parameters of the gate sequence. The control parameters for a SNAP gate are the phase shifts for all the Fock levels, so, in principle, there are infinitely many of them. In practice, however, only the $\mathcal{O}(N)$ Fock levels close to the input/output space will get populated significantly, so effectively each SNAP gate contributes only $\mathcal{O}(N)$ free variables. In addition, each displacement contributes one more free parameter (the absolute of α , as the complex phase can be absorbed into the neighboring SNAP gates, see appendix A 1). So, in total, we obtain $\mathcal{O}(N)$ free variables per pair of SNAP gate and displacement. Hence, we need $\mathcal{O}(N^2)/\mathcal{O}(N) = \mathcal{O}(N)$ pairs of these gates to arrive at a sequence with sufficiently many control parameters.

This optimistic estimate, promising quadratic speedup over the previous technique from [12], was actually the starting point for this project. It motivated us to look deeper into optimizing this control problem of combining SNAP gates and displacements, which has eventually led to the technique described in this paper.

Appendix B: Objective functions

1. Overview

To quantify the mismatch between the target operation and the actual transformation of the cavity state (disregarding noise and gate imperfections), we consider the mean overlap

$$\begin{aligned} \mathcal{F}[\hat{U}] &:= \max_{\varphi \in \mathbb{R}} \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s. t.} \\ \sum_i |c_i|^2 = 1}} \text{Re}(\langle y(\vec{c}) | e^{i\varphi} \hat{U} | x(\vec{c}) \rangle) = \\ &= \frac{1}{L} |\langle \hat{V}, \hat{U} \rangle_{\text{HS}}| \end{aligned}$$

as defined in Eq. (5). For \hat{U} , insert $\hat{U} = \hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1)$ as in Eq. (2), or its reparameterization $\hat{U} = \hat{B}_T \cdot \dots \cdot \hat{B}_1$ using the building blocks defined in Eq. (10). The interpretation as the mean overlap gets obvious from the first line. The second line, using the Hilbert-Schmidt product $\langle \hat{A}, \hat{B} \rangle_{\text{HS}} = \text{tr}[\hat{A}^\dagger \hat{B}]$, is more useful to perform calculations. Their equivalence is shown in appendix B 2.

Another important property of the gate sequences is the mean photon number, which determines the rate of photon losses. For this purpose, we introduce \bar{n}_t as the mean photon number during $\hat{S}(\vec{\theta}_t)$, i. e. after executing $\hat{D}(\alpha_{t-1}) \cdot \hat{S}(\vec{\theta}_{t-2}) \cdot \dots \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1)$, averaged over all input states $|x(\vec{c})\rangle$.

To obtain an expression for the photon number cost that is more symmetric than $\sum_{t=1}^T \bar{n}_t$, we follow the idea that mapping the input states $|x(\vec{c})\rangle$ to the output states $|y(\vec{c})\rangle$ with the gate sequence $\hat{D}(\alpha_{T+1}) \cdot \hat{S}(\vec{\theta}_T) \cdot \hat{D}(\alpha_T) \cdot \dots \cdot \hat{S}(\vec{\theta}_2) \cdot \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1)$ is equivalent to mapping the output states $|y(\vec{c})\rangle$ to the input states $|x(\vec{c})\rangle$ with the inverse sequence $\hat{D}^\dagger(\alpha_1) \cdot \hat{S}^\dagger(\vec{\theta}_1) \cdot \hat{D}^\dagger(\alpha_2) \cdot \hat{S}^\dagger(\vec{\theta}_2) \cdot \dots \cdot \hat{D}^\dagger(\alpha_T) \cdot \hat{S}^\dagger(\vec{\theta}_T) \cdot \hat{D}^\dagger(\alpha_{T+1})$. Therefore, we minimize $\sum_{t=1}^T \frac{\bar{n}_t + \bar{n}'_t}{2}$ rather than $\sum_{t=1}^T \bar{n}_t$, where \bar{n}'_t is the mean photon number after executing $\hat{D}^\dagger(\alpha_t) \cdot \hat{S}^\dagger(\vec{\theta}_{t+2}) \cdot \dots \cdot \hat{D}^\dagger(\alpha_T) \cdot \hat{S}^\dagger(\vec{\theta}_T) \cdot \hat{D}^\dagger(\alpha_{T+1})$, starting from the output states $|y(\vec{c})\rangle$.

The quantities \bar{n}_t and \bar{n}'_t can be computed as

$$\begin{aligned}\bar{n}_t &:= \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \langle x_t(\vec{c}) | \hat{n} | x_t(\vec{c}) \rangle = \\ &= \frac{1}{L} \langle \hat{V}^\dagger \hat{V}, \hat{B}_1^\dagger \cdots \hat{B}_{t-1}^\dagger \cdot \hat{D}(\alpha_t) \cdot \hat{n} \cdot \hat{D}^\dagger(\alpha_t) \cdot \hat{B}_{t-1} \cdots \hat{B}_1 \rangle_{\text{HS}} \\ \bar{n}'_t &:= \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \langle y_t(\vec{c}) | \hat{n} | y_t(\vec{c}) \rangle = \\ &= \frac{1}{L} \langle \hat{V}^\dagger \hat{V}, \hat{B}_T \cdots \hat{B}_{t+1} \cdot \hat{D}(\alpha_t) \cdot \hat{n} \cdot \hat{D}^\dagger(\alpha_t) \cdot \hat{B}_{t+1}^\dagger \cdots \hat{B}_T^\dagger \rangle_{\text{HS}}\end{aligned}$$

with $|x_t(\vec{c})\rangle = \hat{D}(\alpha_t) \cdot \hat{S}(\vec{\theta}_{t-1}) \cdots \hat{D}(\alpha_2) \cdot \hat{S}(\vec{\theta}_1) \cdot \hat{D}(\alpha_1) |x(\vec{c})\rangle = \hat{D}^\dagger(\alpha_t) \cdot \hat{B}_{t-1} \cdots \hat{B}_1 |x(\vec{c})\rangle$ and $|y_t(\vec{c})\rangle = \hat{D}^\dagger(\alpha_t) \cdot \hat{B}_{t+1}^\dagger \cdots \hat{B}_T^\dagger |y(\vec{c})\rangle$.

2. Calculation for the mean overlap

Here, we demonstrate how the definition for $\mathcal{F}[\hat{U}]$ can be rewritten into a more practical form. We start by evaluating the maximization over φ :

$$\begin{aligned}\mathcal{F}[\hat{U}] &= \max_{\varphi \in \mathbb{R}} \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \text{Re}(\langle y(\vec{c}) | e^{i\varphi} \hat{U} | x(\vec{c}) \rangle) = \\ &= \max_{\varphi \in \mathbb{R}} \text{Re} \left(e^{i\varphi} \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \langle y(\vec{c}) | \hat{U} | x(\vec{c}) \rangle \right) = \\ &= \left| \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \langle y(\vec{c}) | \hat{U} | x(\vec{c}) \rangle \right|\end{aligned}$$

By inserting $|x(\vec{c})\rangle = \sum_{l=1}^L c_l |x(\vec{e}_l)\rangle$ and $|y(\vec{c})\rangle = \sum_{l'=1}^L c_{l'} |y(\vec{e}_{l'})\rangle$, we can further simplify this expression:

$$\begin{aligned}\mathcal{F}[\hat{U}] &= \left| \text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} \langle y(\vec{c}) | \hat{U} | x(\vec{c}) \rangle \right| = \\ &= \left| \sum_{l, l'=1}^L \overbrace{\text{avg}_{\substack{\vec{c} \in \mathbb{C}^L \text{ s.t.} \\ \sum_i |c_i|^2 = 1}} [c_l^* c_{l'}]}^{=\delta_{ll'}/L} \langle y(\vec{e}_{l'}) | \hat{U} | x(\vec{e}_l) \rangle \right| = \\ &= \frac{1}{L} \left| \sum_{l=1}^L \langle y(\vec{e}_l) | \hat{U} | x(\vec{e}_l) \rangle \right| = \\ &= \frac{1}{L} \left| \text{tr} \left[\underbrace{\left(\sum_{l=1}^L |x(\vec{e}_l)\rangle \langle y(\vec{e}_l)| \right)}_{=\hat{V}^\dagger} \hat{U} \right] \right| = \\ &= \frac{1}{L} |\langle \hat{V}, \hat{U} \rangle_{\text{HS}}|\end{aligned}$$

Appendix C: Initialization

In Sec. IV A, we describe a greedy strategy to initialize a gate sequence. Each step in this approach can be reduced to determining the optimum combination of α and $\vec{\theta}$ that maximizes mean overlap $\mathcal{F}[\hat{B}(\alpha, \vec{\theta})]$ for a single building

block $\hat{B}(\alpha, \vec{\theta}) = \hat{D}^\dagger(\alpha) \hat{S}(\vec{\theta}) \hat{D}(\alpha)$. Here, we derive the result as stated in Sec. IV A.

First, we compute

$$\mathcal{F}[\hat{B}(\alpha, \vec{\theta})] = \frac{1}{L} |\langle \hat{V}, \hat{B}(\alpha, \vec{\theta}) \rangle_{\text{HS}}| = \dots = \frac{1}{L} \left| \sum_{n=0}^{\infty} e^{i\theta^{(n)}} g_n^*(\alpha) \right|$$

with $g_n(\alpha) = \langle n | \hat{D}(\alpha) \hat{V} \hat{D}^\dagger(\alpha) | n \rangle$ as defined in Sec. IV A. Now, we can optimize $\vec{\theta}$ for an arbitrary, but fixed value for α :

$$\begin{aligned} \operatorname{argmax}_{\theta^{(n)}} \mathcal{F}[\hat{B}(\alpha, \vec{\theta})] &= \bar{\theta} + \arg(g_n(\alpha)) + 2\pi\mathbb{Z} \\ \max_{\vec{\theta}} \mathcal{F}[\hat{B}(\alpha, \vec{\theta})] &= \frac{1}{L} \sum_{n=0}^{\infty} |g_n(\alpha)|. \end{aligned}$$

Hence, the optimum value for α is

$$\alpha_{\text{opt}} = \operatorname{argmax}_{\alpha} \left[\sum_{n=0}^{\infty} |g_n(\alpha)| \right],$$

and the optimum value for $\vec{\theta}$ is $\vec{\theta}_{\text{opt}} = \bar{\theta} + \arg(g_n(\alpha_{\text{opt}})) + 2\pi\mathbb{Z}$. By choosing $\bar{\theta} = 0$, we arrive at Eq. (11) in the main text.

This recipe can be used to extend a sequence of building blocks $\hat{B}(\alpha, \vec{\theta})$ to any desired length in a meaningful way, as explained in Sec. IV A.

Appendix D: Finetuning

During finetuning, we take the gradient of the cost function in Eq. (13) w. r. t. the control parameters α_t and $\vec{\theta}_t$ of the building blocks $\hat{B}_t = \hat{B}(\alpha_t, \vec{\theta}_t)$. Here, we provide explicit expressions for the gradients. We also demonstrate how these gradients can be computed efficiently using a recursive scheme which is analogous to backpropagation in machine learning.

1. Gradient for the overlap cost

The gradient for the overlap cost is given by

$$\begin{aligned} \frac{\partial}{\partial \alpha_t} \ln \left(1 - \mathcal{F}[\hat{B}_T \dots \hat{B}_1] \right) &= \operatorname{Re} \left(\left\langle \hat{G}_t, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \alpha_t} \right\rangle_{\text{HS}} \right) \\ \frac{\partial}{\partial \theta_t^{(n)}} \ln \left(1 - \mathcal{F}[\hat{B}_T \dots \hat{B}_1] \right) &= \operatorname{Re} \left(\left\langle \hat{G}_t, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \theta_t^{(n)}} \right\rangle_{\text{HS}} \right) \end{aligned}$$

with

$$\hat{G}_t := - \frac{\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}}{|\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}|} \cdot \frac{\hat{B}_{t+1}^\dagger \dots \hat{B}_T^\dagger \cdot \hat{V} \cdot \hat{B}_1^\dagger \dots \hat{B}_{t-1}^\dagger}{L - |\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}|}.$$

\hat{G}_t can be efficiently computed using the recursive scheme

$$\hat{G}_1 = - \frac{\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}}{|\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}|} \cdot \frac{\hat{B}_2^\dagger \dots \hat{B}_T^\dagger \cdot \hat{V}}{L - |\langle \hat{V}, \hat{B}_T \dots \hat{B}_1 \rangle_{\text{HS}}|} \quad \hat{G}_{t+1} = \hat{B}_{t+1} \hat{G}_t \hat{B}_t^\dagger.$$

2. Gradient for the photon-number cost

The gradient for the photon-number cost can be computed from

$$\begin{aligned} \frac{\partial}{\partial \alpha_t} \left[\sum_{s=1}^T \bar{n}_s \right] &= \text{Re} \left(\left\langle 2\hat{X}_t \hat{B}_t \hat{\rho}_t^{(x)}, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \alpha_t} \right\rangle_{\text{HS}} + \left\langle 2\hat{\rho}_t^{(x)}, \hat{D}^\dagger(\alpha_t) \hat{n} \frac{\partial \hat{D}(\alpha_t)}{\partial \alpha_t} \right\rangle_{\text{HS}} \right) \\ \frac{\partial}{\partial \theta_t^{(n)}} \left[\sum_{s=1}^T \bar{n}_s \right] &= \text{Re} \left(\left\langle 2\hat{X}_t \hat{B}_t \hat{\rho}_t^{(x)}, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \theta_t^{(n)}} \right\rangle_{\text{HS}} \right) \\ \frac{\partial}{\partial \alpha_t} \left[\sum_{s=1}^T \bar{n}'_s \right] &= \text{Re} \left(\left\langle 2\hat{\rho}_t^{(y)} \hat{B}_t \hat{Y}_t, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \alpha_t} \right\rangle_{\text{HS}} + \left\langle 2\hat{\rho}_t^{(y)}, \hat{D}^\dagger(\alpha_t) \hat{n} \frac{\partial \hat{D}(\alpha_t)}{\partial \alpha_t} \right\rangle_{\text{HS}} \right) \\ \frac{\partial}{\partial \theta_t^{(n)}} \left[\sum_{s=1}^T \bar{n}'_s \right] &= \text{Re} \left(\left\langle 2\hat{\rho}_t^{(y)} \hat{B}_t \hat{Y}_t, \frac{\partial \hat{B}(\alpha_t, \vec{\theta}_t)}{\partial \theta_t^{(n)}} \right\rangle_{\text{HS}} \right). \end{aligned}$$

with

$$\begin{aligned} \hat{\rho}_t^{(x)} &= \frac{1}{L} \cdot \hat{B}_{t-1} \cdots \hat{B}_1 \cdot \hat{V}^\dagger \cdot \hat{V} \cdot \hat{B}_1^\dagger \cdots \hat{B}_{t-1}^\dagger \\ \hat{\rho}_t^{(y)} &= \frac{1}{L} \cdot \hat{B}_{t+1}^\dagger \cdots \hat{B}_T^\dagger \cdot \hat{V} \cdot \hat{V}^\dagger \cdot \hat{B}_T \cdots \hat{B}_{t+1} \\ \hat{X}_t &= \sum_{s=t+1}^T \hat{B}_{t+1}^\dagger \cdots \hat{B}_{s-1}^\dagger \cdot \hat{D}^\dagger(\alpha_s) \cdot \hat{n} \cdot \hat{D}(\alpha_s) \cdot \hat{B}_{s-1} \cdots \hat{B}_{t+1} \\ \hat{Y}_t &= \sum_{s=1}^{t-1} \hat{B}_{t-1} \cdots \hat{B}_{s+1} \cdot \hat{D}^\dagger(\alpha_s) \cdot \hat{n} \cdot \hat{D}(\alpha_s) \cdot \hat{B}_{s+1}^\dagger \cdots \hat{B}_{t-1}^\dagger. \end{aligned}$$

Again, there are recursive schemes to compute the quantities $\hat{\rho}_t^{(x)}$, $\hat{\rho}_t^{(y)}$, \hat{X}_t and \hat{Y}_t efficiently:

$$\begin{aligned} \hat{\rho}_1^{(x)} &= \frac{1}{L} \hat{V}^\dagger \hat{V} & \hat{\rho}_{t+1}^{(x)} &= \hat{B}_t \hat{\rho}_t^{(x)} \hat{B}_t^\dagger \\ \hat{\rho}_T^{(y)} &= \frac{1}{L} \hat{V} \hat{V}^\dagger & \hat{\rho}_{t-1}^{(y)} &= \hat{B}_t^\dagger \hat{\rho}_t^{(y)} \hat{B}_t \\ \hat{X}_T &= 0 & \hat{X}_{t-1} &= \hat{B}_t^\dagger \hat{X}_t \hat{B}_t + \hat{D}^\dagger(\alpha_t) \hat{n} \hat{D}(\alpha_t) \\ \hat{Y}_1 &= 0 & \hat{Y}_{t+1} &= \hat{B}_t \hat{Y}_t \hat{B}_t^\dagger + \hat{D}^\dagger(\alpha_t) \hat{n} \hat{D}(\alpha_t) \end{aligned}$$

Appendix E: Benchmarks

1. Transformations of the 10 lowest Fock states

Here, we provide explicit formulas for the operators \hat{V} representing the target operation (cmp. Eq. (2)), for all examples in Fig. 2a.

Inversion: $\hat{V} = \sum_{n=0}^9 |9-n\rangle\langle n|$

Block inversion: $\hat{V} = \sum_{n=0}^9 |\text{mod}(n+5, 10)\rangle\langle n|$

Random permutations: $\hat{V} = \sum_{n=0}^9 |p(n)\rangle\langle n|$ where p are random permutations of the numbers $0, \dots, 9$.

Random unitaries: $\hat{V} = \sum_{m,n=0}^9 v_{mn} |m\rangle\langle n|$ where v is a unitary 10×10 matrix generated by diagonalizing random Hermitian matrices.

2. Applied problems

Here, we provide explicit formulas for the operators \hat{V} representing the target operation (cmp. Eq. (2)), for all examples in Fig. 3.

State preparation (Fig. 3a):

$$\hat{V} = (\alpha|b_0\rangle + \beta|b_1\rangle)\langle 0|$$

for the stated values of the coefficients α and β . The numbers for “odd superposition” are $\alpha = \sqrt{(1 + \sin 0.72104)}/2$ and $\beta = \sqrt{(1 - \sin 0.72104)}/2 \cdot e^{-1.27275i}$.

Error correction (Fig. 3b):

$$\hat{V}_s = |\psi_{1,0}^{(1)}(0)\rangle\langle\psi_{1,0}^{(s)}(t)| + |\psi_{0,1}^{(1)}(0)\rangle\langle\psi_{0,1}^{(s)}(t)|$$

with $\Gamma t = 0.02$ and

$$\begin{aligned} |\psi_{\alpha\beta}^{(1)}(t)\rangle &= \frac{\alpha \cdot (|0\rangle + \sqrt{3}e^{-6\Gamma t}|6\rangle) + \beta \cdot (\sqrt{3}e^{-3\Gamma t}|3\rangle + e^{-9\Gamma t}|9\rangle)}{\sqrt{|\alpha|^2 \cdot (1 + 3e^{-12\Gamma t}) + |\beta|^2 \cdot (3e^{-6\Gamma t} + e^{-18\Gamma t})}} \\ |\psi_{\alpha\beta}^{(\hat{a})}(t)\rangle &= \frac{\alpha \cdot (\sqrt{2}e^{-3\Gamma t}|5\rangle) + \beta \cdot (|2\rangle + e^{-6\Gamma t}|8\rangle)}{\sqrt{|\alpha|^2 \cdot (2e^{-6\Gamma t}) + |\beta|^2 \cdot (1 + e^{-12\Gamma t})}} \\ |\psi_{\alpha\beta}^{(\hat{a}^2)}(t)\rangle &= \frac{\alpha \cdot (\sqrt{5}e^{-3\Gamma t}|4\rangle) + \beta \cdot (|1\rangle + 2e^{-6\Gamma t}|7\rangle)}{\sqrt{|\alpha|^2 \cdot (5e^{-6\Gamma t}) + |\beta|^2 \cdot (1 + 4e^{-12\Gamma t})}} \end{aligned}$$

These states can be computed as follows: To solve the Lindblad master equation for a cavity under photon loss,

$$\frac{d\hat{\rho}_{\alpha\beta}(t)}{dt} = \Gamma \cdot (2\hat{a}\hat{\rho}_{\alpha\beta}(t)\hat{a}^\dagger - \{\hat{a}^\dagger\hat{a}, \hat{\rho}_{\alpha\beta}(t)\}),$$

with initial state $\hat{\rho}_{\alpha\beta}(0) = (\alpha|b_0\rangle + \beta|b_1\rangle)(\alpha|b_0\rangle + \beta|b_1\rangle)^\dagger$, we decompose $\hat{\rho}_{\alpha\beta}(t)$ as $\hat{\rho}_{\alpha\beta}(t) = \sum_{j=0}^9 \hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t)$ with

$$\begin{aligned} \frac{1}{\Gamma} \frac{d\hat{\rho}_{\alpha\beta}^{(1)}(t)}{dt} &= -\{\hat{n}, \hat{\rho}_{\alpha\beta}^{(1)}(t)\} & \hat{\rho}_{\alpha\beta}^{(1)}(0) &= \mathbb{1} \\ j \geq 1 : \quad \frac{1}{\Gamma} \frac{d\hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t)}{dt} &= -\{\hat{n}, \hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t)\} + 2\hat{a}\hat{\rho}_{\alpha\beta}^{(\hat{a}^{j-1})}(t)\hat{a}^\dagger & \hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(0) &= 0. \end{aligned}$$

The solutions $\hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t)$ all have rank 1 (unless $t = 0$). $|\psi_{\alpha\beta}^{(1)}(t)\rangle$, $|\psi_{\alpha\beta}^{(\hat{a})}(t)\rangle$ and $|\psi_{\alpha\beta}^{(\hat{a}^2)}(t)\rangle$ are chosen such that they satisfy $|\psi_{\alpha\beta}^{(\hat{a}^j)}(t)\rangle\langle\psi_{\alpha\beta}^{(\hat{a}^j)}(t)| = \hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t) / \text{tr}[\hat{\rho}_{\alpha\beta}^{(\hat{a}^j)}(t)]$.

Note that $|\psi_{1,0}^{(1)}(0)\rangle = |b_0\rangle$ and $|\psi_{0,1}^{(1)}(0)\rangle = |b_1\rangle$.

Logical operations on binomial code (Fig. 3c):

$$\hat{V} = \sum_{j,k=0}^1 v_{jk}|b_j\rangle\langle b_k|$$

with

$$\begin{aligned} v &= \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} & v &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} & v &= \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} & v &= \frac{1}{2} \begin{pmatrix} 1-i & 1+i \\ 1+i & 1-i \end{pmatrix} \\ & \text{(Hadamard)} & \text{(Pauli-X)} & \text{(Pauli-Y)} & \text{(Sqrt-Pauli-X)} \end{aligned}$$

v for “odd operation” is a 2×2 matrix that was generated by diagonalizing a random Hermitian matrix.

Logical operations on trivial code (Fig. 3d):

$$\hat{V} = \sum_{m,n=0}^1 v_{mn}|m\rangle\langle n|$$

with the same v matrices as for Fig. 3c.

3. Comparison with previous technique

The following table compares the sequence length between our technique and the technique in [12], for the target operations in Fig. 3. For [12], we give the minimum sequence lengths, which is the squared number of involved Fock states. For our technique, we state the smallest value for T such that \mathcal{F} exceeds a value of 0.999.

target operation		involved Fock levels	T_{\min} for [12]	$T_{\mathcal{F} \geq 0.999}$ here
state preparation (Fig. 3a)	$ b_0\rangle$	0, 6	4	4
	$ b_1\rangle$	0, 3, 9	9	3
	superpositions	0, 3, 6, 9	16	3...4
error correction (Fig. 3b)	syndrome $\mathbb{1}$	0, 3, 6, 9	16	4
	syndrome \hat{a}	0, 2, 3, 5, 6, 8, 9	49	
	syndrome \hat{a}^2	0, 1, 3, 4, 6, 7, 9	49	
logical operations on binomial code (Fig. 3c)		0, 3, 9	16	3...4
logical operations on trivial code (Fig. 3d)		0, 1	4	3

Appendix F: Simulations

1. Implementation details and hyperparameters

Both during initialization and finetuning:

- The Hilbert space of a cavity is infinite-dimensional, so in order to be able to simulate it, we have to truncate it at some point. In all our simulations, we neglect those Fock states $|n\rangle$ with $n \geq 100$. We choose this limit deliberately high to rule out any non-negligible leakage of the relevant input states out of this space, guaranteeing for correctness of our calculations. With 100 simulated Fock states, leakage to higher states is much smaller than the numerical floating-point error, so our choice is very generous. This goes at the expense of the simulation runtime (which scales quadratically with the truncation limit), but this we can afford.

During initialization only:

- To find the optimal values for the control parameters α and $\vec{\theta}$ of an inserted building block according to Eq. (11), we have to take into account in principle every real number for α . In our implementation, we restrict our search to a fixed set of candidates for α . Explicitly, we choose this set as $\{-2, -1.8, \dots, 1.8, 2\}$. Intentionally, this set excludes values with $|\alpha|$ above a certain threshold (here 2), as large displacements lead to high values for the average mean photon number \bar{n}_t (which we want to avoid).
- All SNAP phase angles for the Fock levels $n \geq 15$ are set to 0 during initialization. Even though this is not absolutely necessary, it helps to launch finetuning with a lower value for the average mean photon number \bar{n}_t . However, also during initialization, we simulate the cavity up to Fock state $n = 100$ (see above) since also the Fock states ≥ 15 are still populated. Moreover, during finetuning, we allow non-vanishing SNAP phase angles also for the Fock states between 15 and 99 (even though in practice these phase angles do not change significantly during finetuning).

During finetuning only:

- We run our finetuning procedure for 100 000 iterations. The results in Fig. 2 and Fig. 3 always refer to the parameter configuration after the last iteration step. Due to fluctuations in the cost function, its value might have been slightly better at previous steps, but the difference is not large enough such that a post-selection over these intermediate configurations would make sense.

- We use the Adam optimizer [18] to compute the updates for α and $\vec{\theta}$ from the corresponding gradients of the cost function. We set the Adam hyperparameters to $\eta = 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$.
- Gradient clipping is implemented by truncating the values for $|\partial C/\partial \alpha_t|$ at 100.0 and for $|\partial C/\partial \theta_t^{(n)}|$ at 50.0 before passing them to the Adam optimizer.
- The only hyperparameter which is changed between the different examples is λ , the coefficient for the photon number cost (see Sec. IV B). Explicitly, we use the following values:

Fig. 2a:

target operation	$T \leq 8$	$T \geq 10$
inversion	$\lambda = 0.16$	$\lambda = 0.4$
block inversion	$\lambda = 0.8$	
random permutations	$\lambda = 1.6$	
random unitaries		

Fig. 2b:

N	λ
2, 3	$\lambda = 2.4$
4, 5	$\lambda = 1.8$
≥ 6	$\lambda = 1.6$

Fig. 3:

application	λ	conditions
state preparation (Fig. 3a)	$\lambda = 0.6$	*
error correction (Fig. 3b)	$\lambda = 0.6$ $\lambda = 0.4$	syndrome \hat{a} and $T \leq 4$ otherwise
logical operation on binomial code (Fig. 3c)	$\lambda = 1.0$ $\lambda = 0.32$	operation Pauli-X and $T \leq 3$ otherwise
logical operation on trivial code (Fig. 3d)	$\lambda = 2.4$	*

- For the “no GC, higher LR” curve in Fig. 4c, we did not apply gradient clipping, and in addition changed two Adam hyperparameters: $\eta = 2.5 \cdot 10^{-4}$ and $\beta_2 = 0.99$. For this parameter configuration, post-selection would be required to reliably find well-performing sequences. Therefore, we find our default hyperparameters preferable which avoid this additional step of complexity.

2. Computational costs

Because we need to construct a large number of gate sequences (for different target operations and sequence lengths), we can obtain the strongest acceleration from embarrassing parallelization, and therefore run each instance single-threaded on one CPU core. For each job, initialization is completed in less than 1 min, and finetuning takes between 2 h and 40 h (for 100 000 iterations). If needed, this runtime could potentially be reduced significantly with a more efficient implementation.

The RAM requirements are very modest. For the example with the highest memory usage, it takes 32 KiB to store the control parameters for one sequence, and between 10 MB and 20 MB to compute the gradient during finetuning. The computations were distributed over 8 nodes à 32 cores.

Supplementary Material

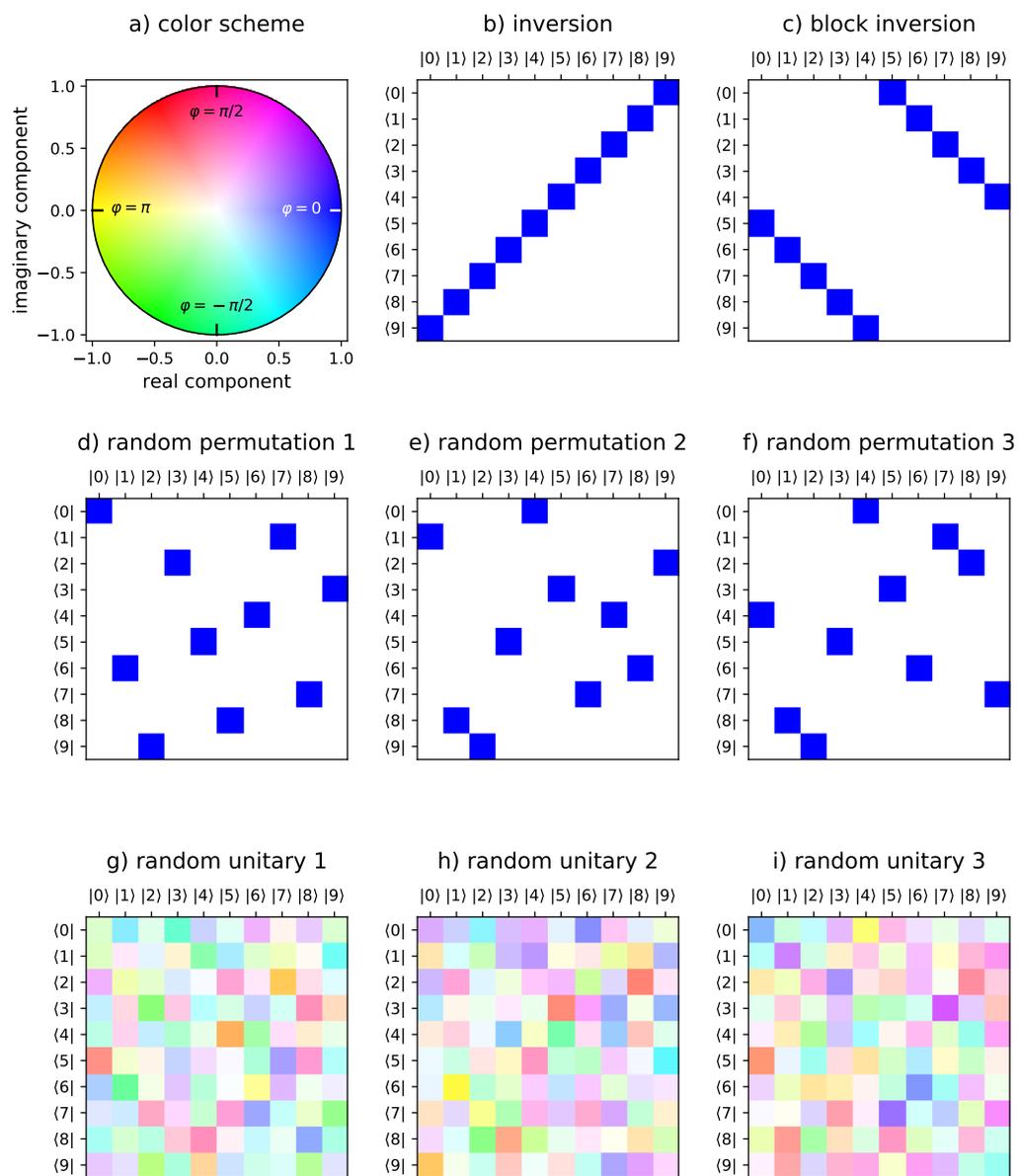


Figure S1. Supplementary information to Fig. 2a. a) Color scheme to represent complex numbers in the following subfigures. b-i) Graphical representation of the target operations used for Fig. 2a.

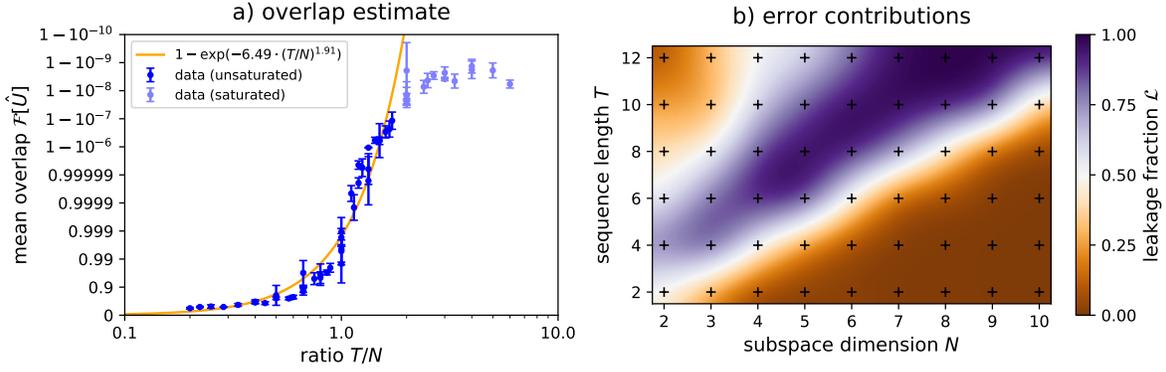


Figure S2. Supplementary information to Fig. 2b. a) Estimate for the mean overlap $\mathcal{F}[\hat{U}]$ obtained for target operations on the N lowest Fock states. Our result from Fig. 2b, that T , the number of building blocks in a sequence, should scale linear in N to achieve the same mean overlap $\mathcal{F}[\hat{U}]$, means that $\mathcal{F}[\hat{U}]$ is primarily a function of the ratio T/N . Therefore, we plot here the averaged values for $\mathcal{F}[\hat{U}]$ from Fig. 2b against T/N . We find that the data points indeed lie on a curve, up to small deviations. To get more insight into the dependency between $\mathcal{F}[\hat{U}]$ and T/N , we make a fit to these data points; the orange curve shows the result. For this fit, we have left out the data points in the saturation region, where $\mathcal{F}[\hat{U}]$ visibly follows a different behavior than for smaller values of T/N . b) Error contributions. The total error can be divided into two parts: leakage of the actual output states out of the desired output space, and a mismatch within this desired output space. To learn how strong these contributions are, we plot the value of $\mathcal{L} = (1 - \mathcal{F}'[\hat{U}]) / (1 - \mathcal{F}[\hat{U}])$ into the same coordinate system as in Fig. 2b. Here, $\mathcal{F}'[\hat{U}] = \frac{1}{L} \|\hat{V}\hat{U}^\dagger\hat{V}\|_1$ quantifies the probability that the actual output states have not leaked out of the desired output space, averaged over all relevant input states. Hence, \mathcal{L} gives the fraction to which leakage contributes to the total error. As the plot shows, leakage has a minor effect both for small and large values of T/N , whereas it is the dominating factor for intermediate values of T/N . This behavior could be a consequence of our particular technique to construct the gate sequences.

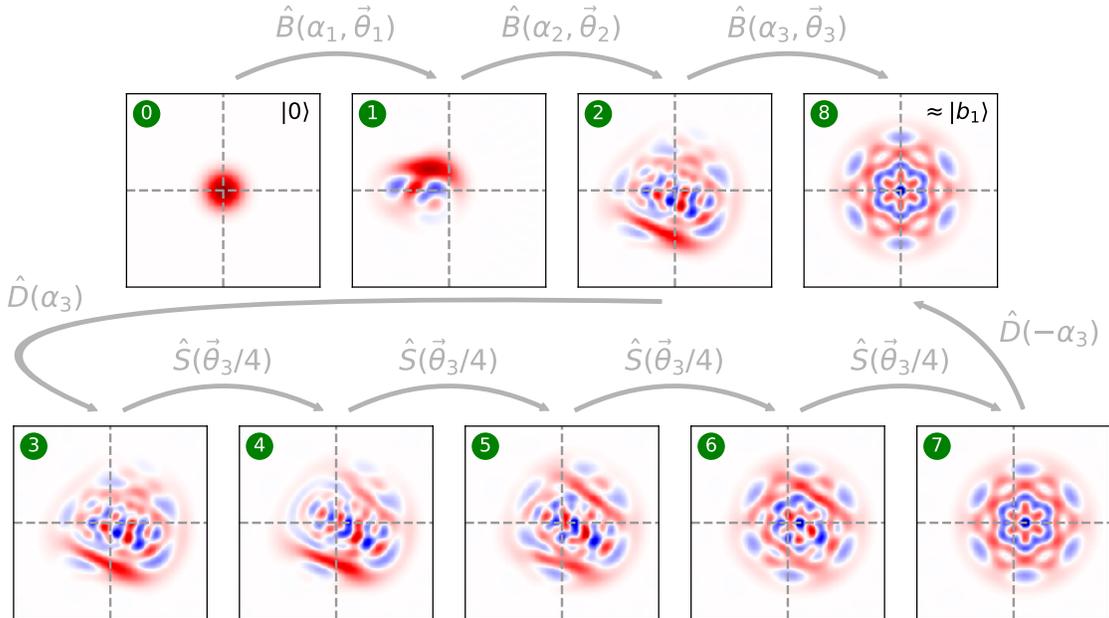


Figure S3. Supplementary information to Fig. 3e. The first row is equal to Fig. 3e, the second row shows additional Wigner density plots for the substeps of the final building block. The displacements ($2 \rightarrow 3$, $7 \rightarrow 8$) just shift the Wigner density along x , whereas the SNAP gate ($3 \rightarrow \dots \rightarrow 7$) transforms the Wigner density in a complex manner.