

Is Circuit Cutting Scalable for Practical Quantum Applications?

Peter Yang¹ & Prakash Murali²

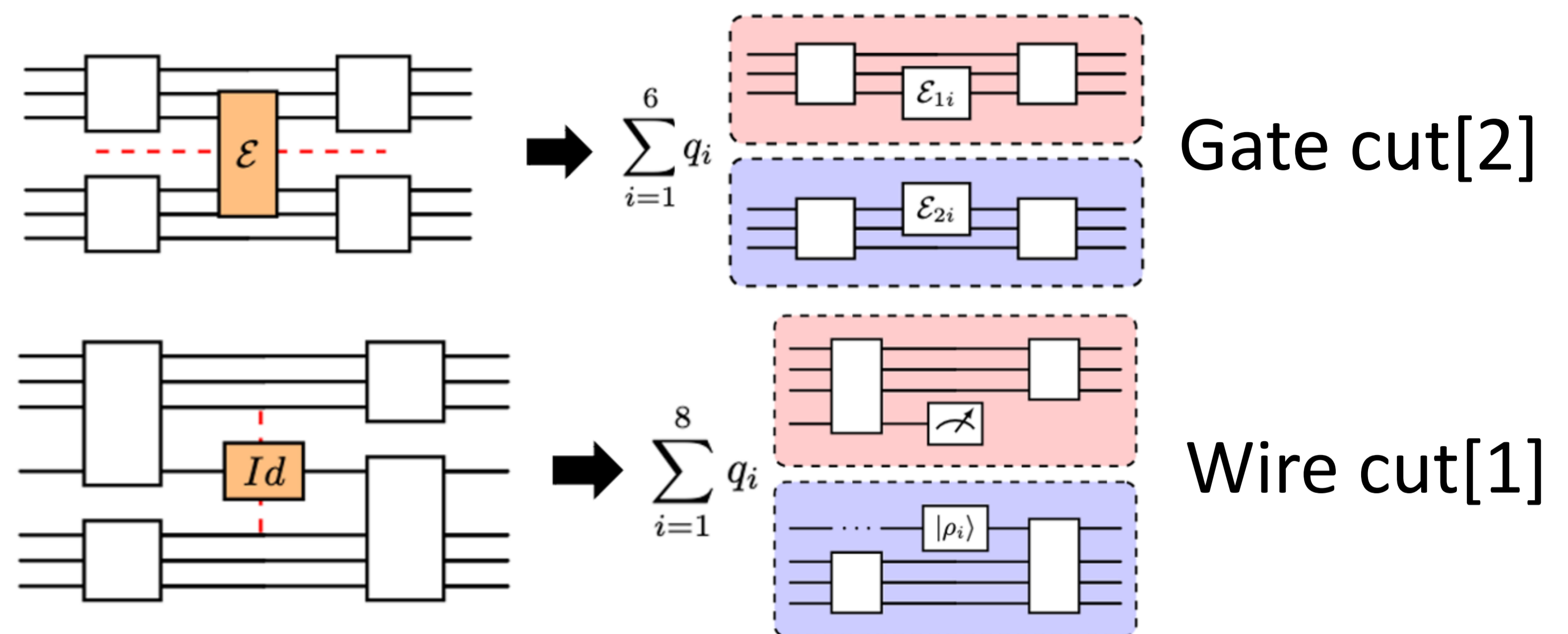
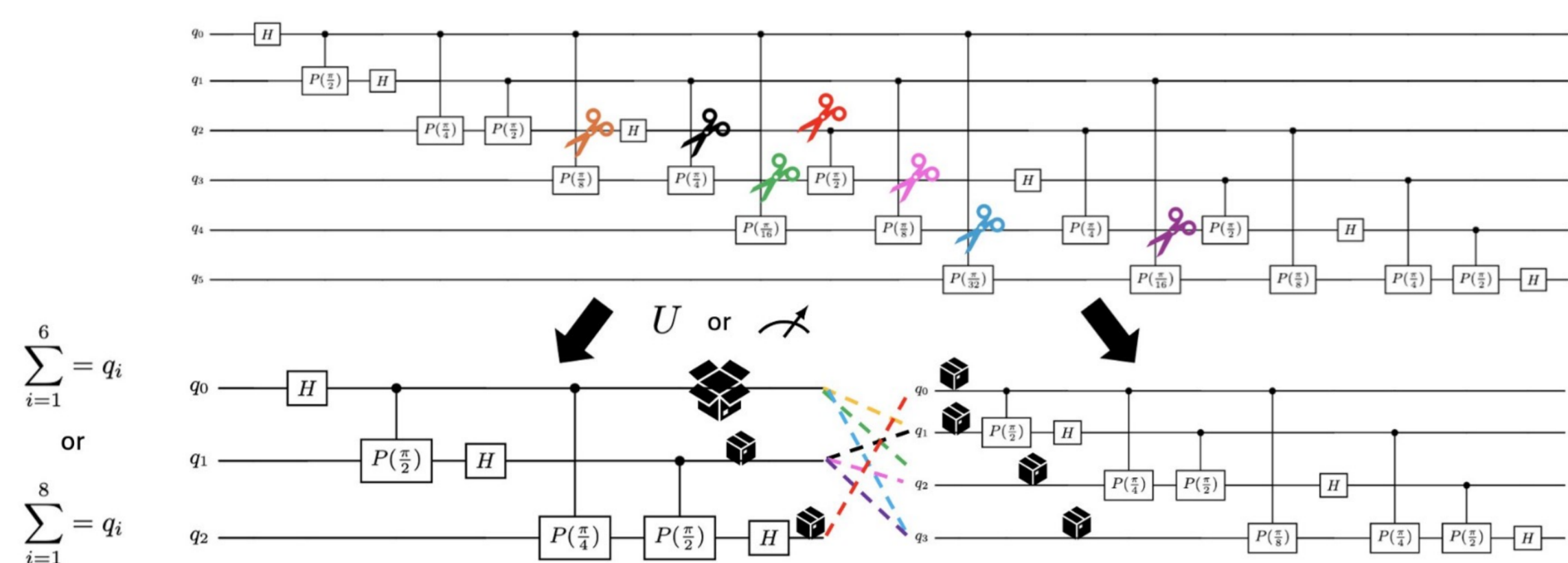
[1]Cavendish, Department of Physics, University of Cambridge

[2]Department of Computer Science and Technology, University of Cambridge

(E-mail: sqhy2@cam.ac.uk)

What is circuit cutting?

- A technique to run large circuit on small quantum computers
- Three phases: cut, run on quantum device, reconstruct on classical device



Challenges in circuit cutting

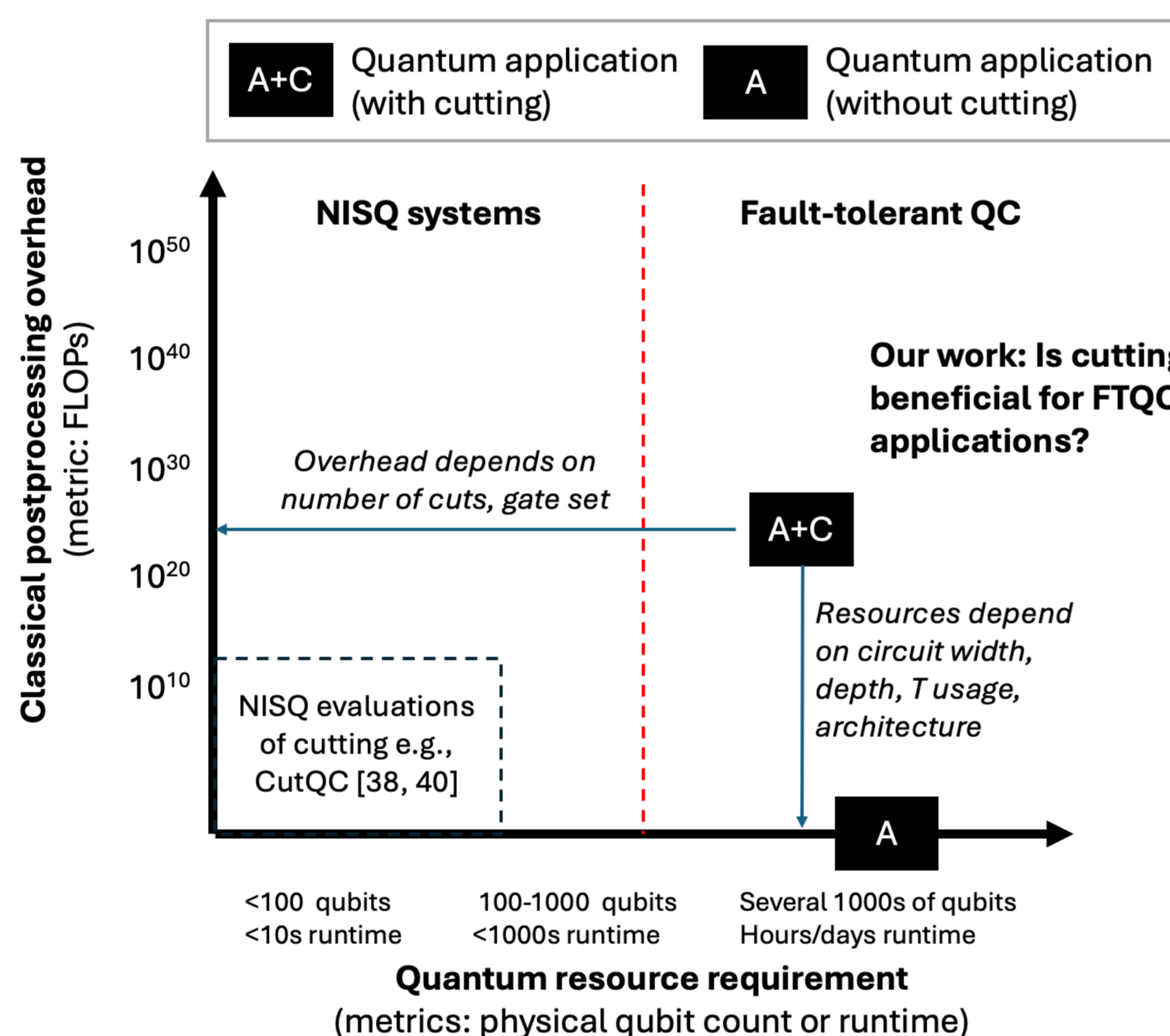
Sampling overhead is exponential in number of cuts, is this scalable for real applications?

$$N_s = O\left(\gamma^{2n} \times \frac{1}{\epsilon^2}\right)$$

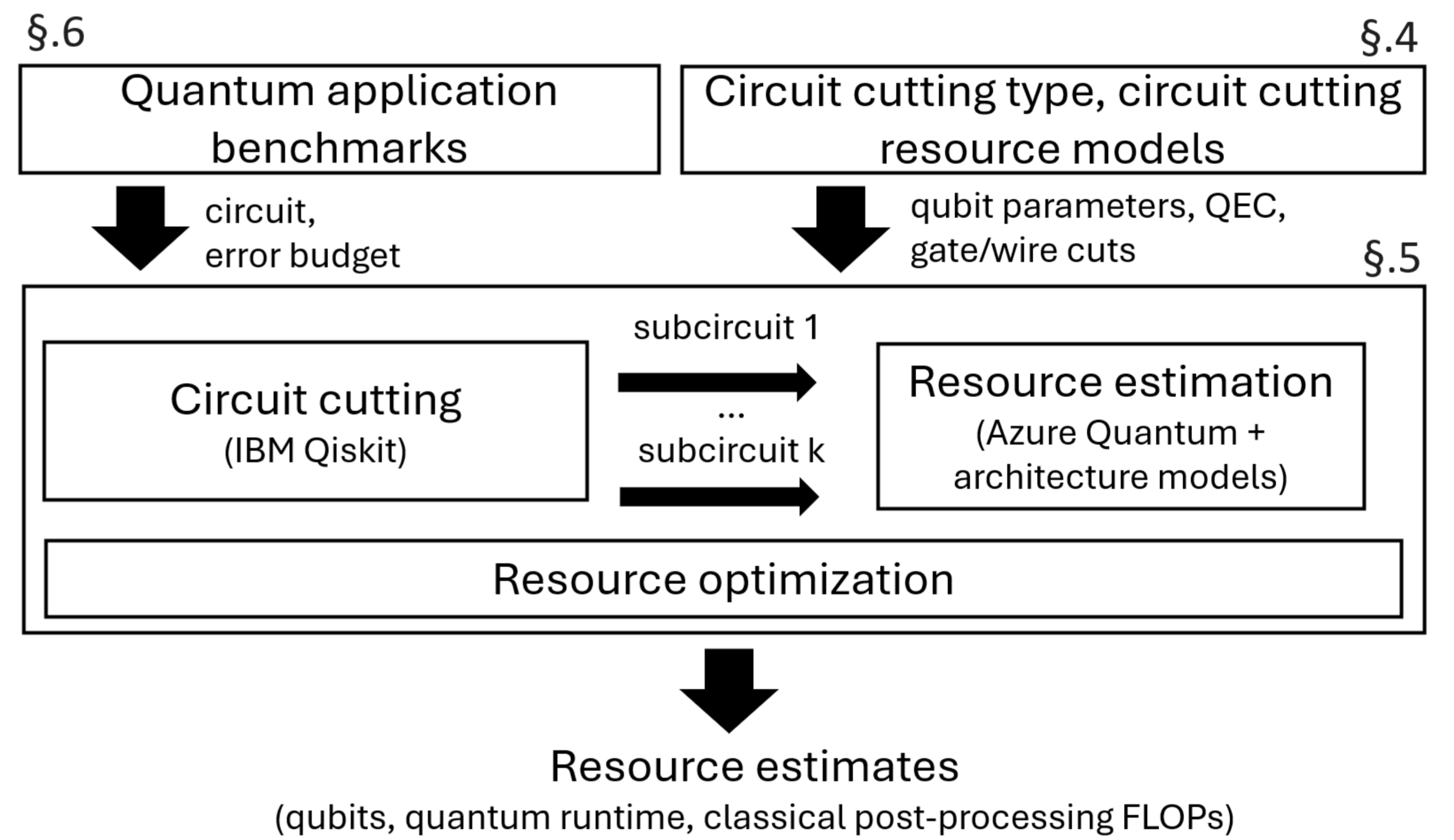
1. Dependence on application structure
 - Clustered applications – easy to cut
 - Applications with more uniform gate structure – hard to cut
2. Different types of cuts:
 - Gate and wire cuts
 - Architecture-dependent cuts
3. In a fault-tolerant setting, resource depend on qubits, QEC, distillation etc.

Instruction	Sampling overhead factor
CSGate, CSXGate	$3 + 2\sqrt{2} \approx 5.828$
CXGate, CHGate	$3^2 = 9$
iSwapGate, SwapGate	$7^2 = 49$
RXXGate, RZXGate	$[1 + 2 \sin(\theta)]^2$
CRXGate, CPhaseGate	$[1 + 2 \sin(\theta/2)]^2$

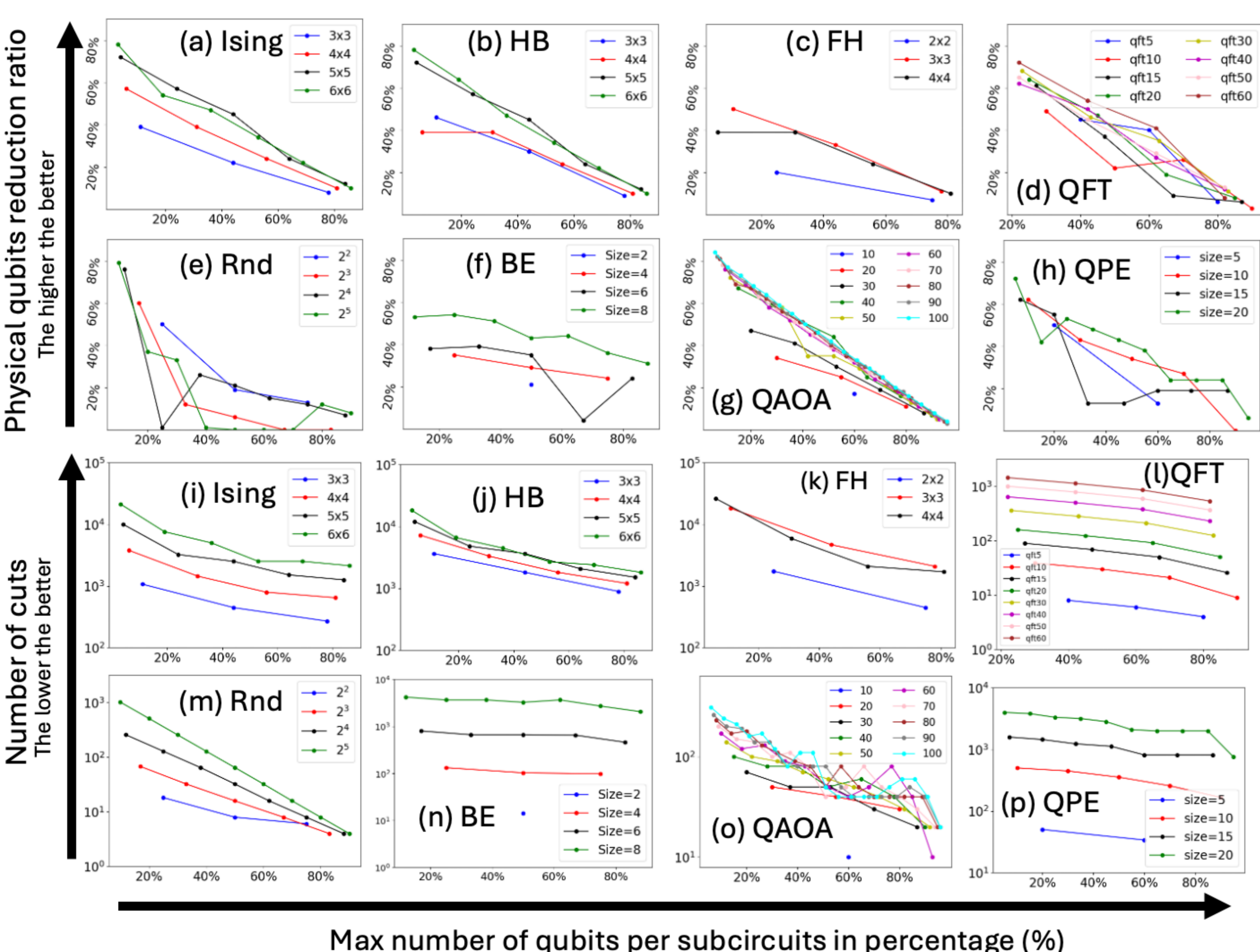
High-level conclusion & Contribution



Workflow



Results



How to model resource requirements?

- Space-efficient execution on a fault-tolerant QC system
- Run subcircuits one by one to reduce space, but increases time (best case for circuit cutting space reductions)

$$Q_{\max} = \max_{c \in \{\text{subexperiments}\}} Q_c$$

$$T_{\max} = \max_{c \in \{\text{subexperiments}\}} T_c \times N_s$$

$$\text{FLOPs required} \sim O(4^n)$$

References:

- [1] Peng, Tianyi, et al. "Simulating large quantum circuits on a small quantum computer." *Physical review letters* 125.15 (2020): 150504.
 [2] Piveteau, Christophe, and David Sutter. "Circuit knitting with classical communication." 2022. *arXiv preprint arXiv:2205.00016* (2022).