# Cost-Based Query Optimization for Quantum Computation
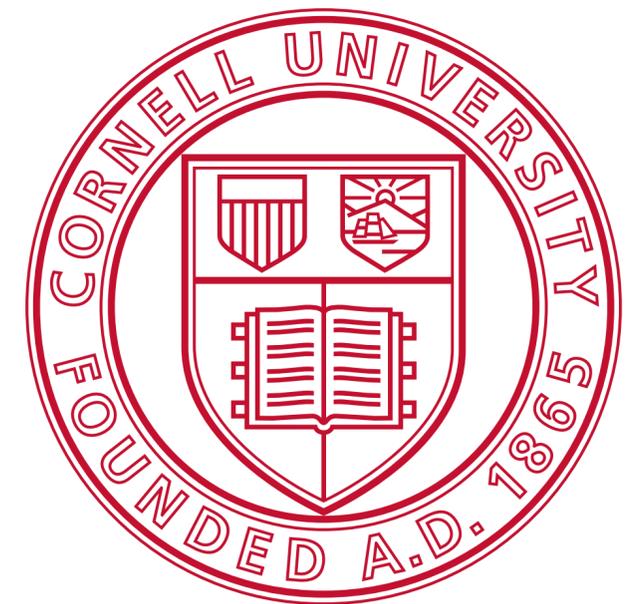
**Immanuel Trummer**

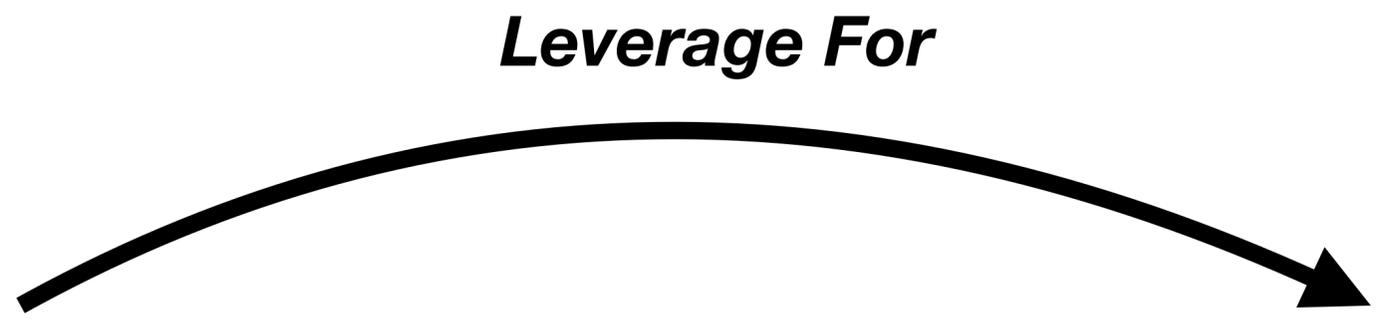**Quantum Computing** *Leverage For* → **Query Optimization**

**QUBO Problem** → Decomposition → Embedding → Solver → Post-Proc. → **Solution**

**QUBO Problem** → **Decomposition** → **Embedding** → **Solver** → **Post-Proc.** → *Solution*

**Choice of Method**

**QUBO Problem** → **Decomposition** → **Embedding** → **Solver** → **Post-Proc.** → **Solution**

**Choice of Method**

Decomposition Algorithms for Scalable Quantum Annealing

Decomposition Algorithms for Solving NP-hard Problems on a Quantum Annealer

...

QUBO Problem → Decomposition → Embedding → Solver → Post-Proc. → *Solution*

**Choice of Method**

**Decomposition Algorithms for Scalable Quantum Annealing**

Elijah Pelofske — Los Alamos National Laboratory, Los Alamos NM 87545, USA — epelofske@lanl.gov

Georg Hahn — Lancaster University, Lancaster LA1 4YW, U.K.

Hristo Djidjev — Los Alamos National Laboratory, Los Alamos, USA — djidjev@lanl.gov

**Integer programming techniques for minor-embedding in quantum annealers**

**Decomposition Algorithms for Solving NP-hard Problems on a Quantum Annealer**

**CHARME: A chain-based reinforcement learning approach for the minor embedding problem**

**QUBO Problem** → Decomposition → Embedding → Solver → Post-Proc. → **Solution**

Choice of Method

...                    ...                    ...

**QUBO Problem** → Decomposition → Embedding → Solver → Post-Proc. → **Solution**

Choice of Method



· · ·          · · ·          · · ·          · · ·

**QUBO Problem** → **Decomposition** → **Embedding** → **Solver** → **Post-Proc.** → **Solution**

Choice of Method

Method Parameters

...

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

↓

*Quantum Optimizer*

↓

**Quantum Computation Plan**
**(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance
User Preferences**

↓

**Quantum Optimizer**

**Optimization Algorithm**

**Cost & Quality Model**

↓

**Quantum Computation Plan
(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

**Quantum Optimizer**

**Optimization Algorithm**

*Multi-Objective Optimization*

**Cost & Quality Model**

**Quantum Computation Plan**
**(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

**Quantum Optimizer**

*Multi-Objective Optimization*

**Optimization Algorithm**

**Cost & Quality Model**

**Quantum Computation Plan**
**(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

**Quantum Optimizer**

**Optimization Algorithm**

**Cost & Quality Model**

**Multi-Objective Optimization**

**Approximation Algorithms**

**Quantum Computation Plan**
**(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

**Quantum Optimizer**

**Optimization Algorithm**

**Multi-Objective Optimization**

**Cost & Quality Model**

**No Analytical Formulas**

**Approximation Algorithms**

**Quantum Computation Plan (Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**

Multiobjective Query Optimization

*Quantum Optimizer*

**Optimization Algorithm**

*Multi-Objective Optimization*

**Cost & Quality Model**

*No Analytical Formulas*

*Approximation Algorithms*

Neo: A Learned Query Optimizer

**Quantum Computation Plan**
**(Methods & Parameters)**

# "Query Optimization" for Quantum Computing

**Problem Instance**
**User Preferences**



*Quantum Optimizer*

**Optimization Algorithm**

**Cost & Quality Model**

*Multi-Objective Optimization*

*No Analytical Formulas*

*Approximation Algorithms*

*Learned Optimization*

**Quantum Computation Plan**
**(Methods & Parameters)**

# Summary

- Use **query optimization ideas** for multi-step quantum computation

  - Approximate multi-objective optimization

  - Learned cost and quality models