



Insights from Developing Hybrid Quantum-Classical Algorithms for Biomarker Discovery in Multimodal Cancer Data

Samantha Riesenfeld



THE UNIVERSITY OF CHICAGO
PRITZKER SCHOOL OF MOLECULAR ENGINEERING

Additional Affiliations

Department of Medicine
Committee on Data Science
Committee on Immunology
Comprehensive Cancer Center

Institute for Biophysical Dynamics
Biohub, Chicago
NSF-Simons National Institute for
Theory and Mathematics in Biology

How my background and academic interests intersect with workshop topics

- BA in Mathematics (with minor in Computer Science)
- PhD in Theoretical Computer Science
- Postdocs* in Computational Biology & Systems Immunology
- Asst. Prof. of Molecular Engineering and of Medicine, UChicago



- Discrete optimization, graphs algorithms, metric embeddings
- Analysis of single-cell 'omics data
- Tissue-specific immune responses
- Multi-modal AI and image-based ML in cancer and cancer immunology
- Riemannian manifold learning
- Transcriptional dynamics (e.g., improving RNA velocity)
- Leveraging spatial transcriptomics
- ...
- Today: quantum computing for biomarker discovery

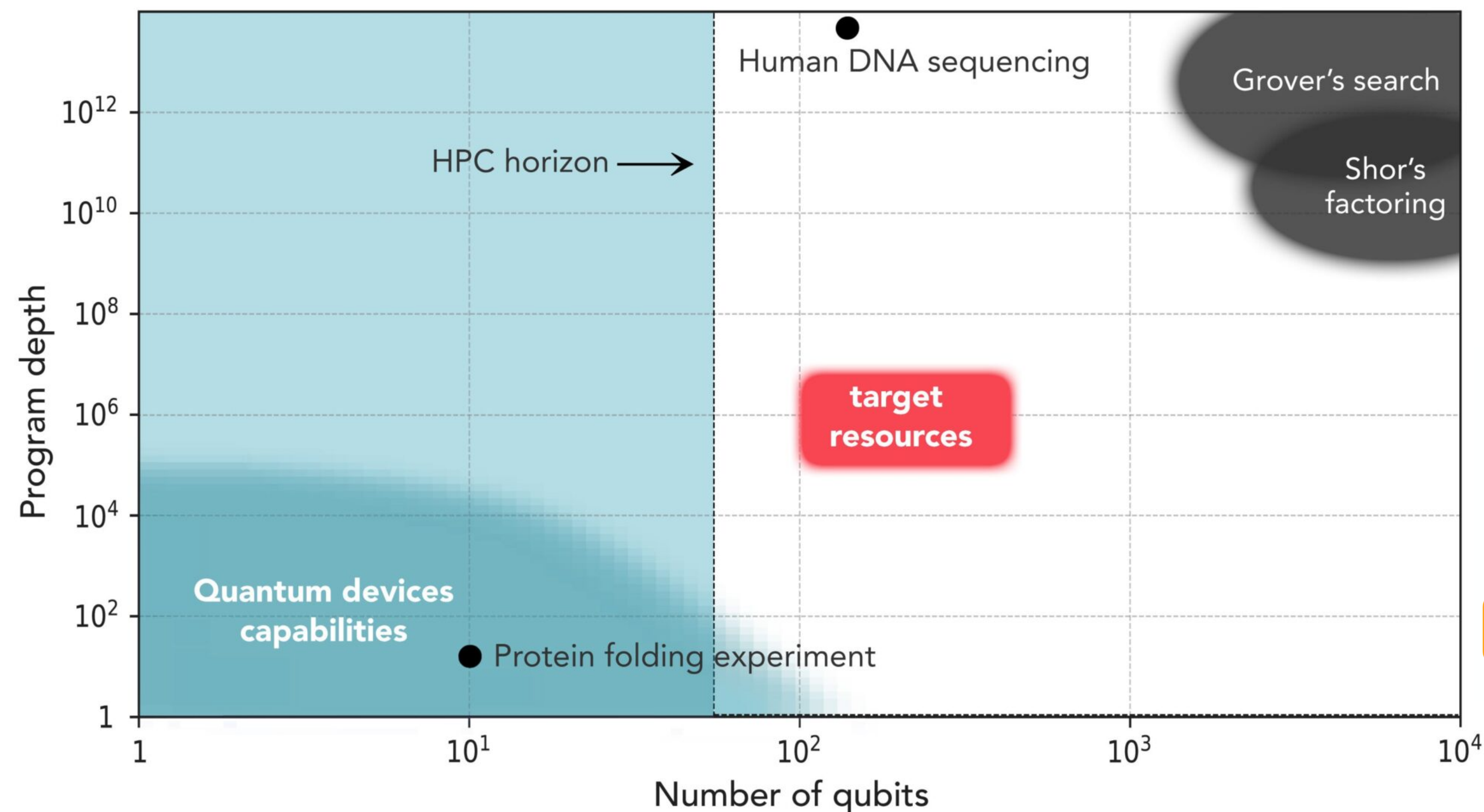
* First experience analyzing real-world data, working in biology, and collaborating with experimental scientists



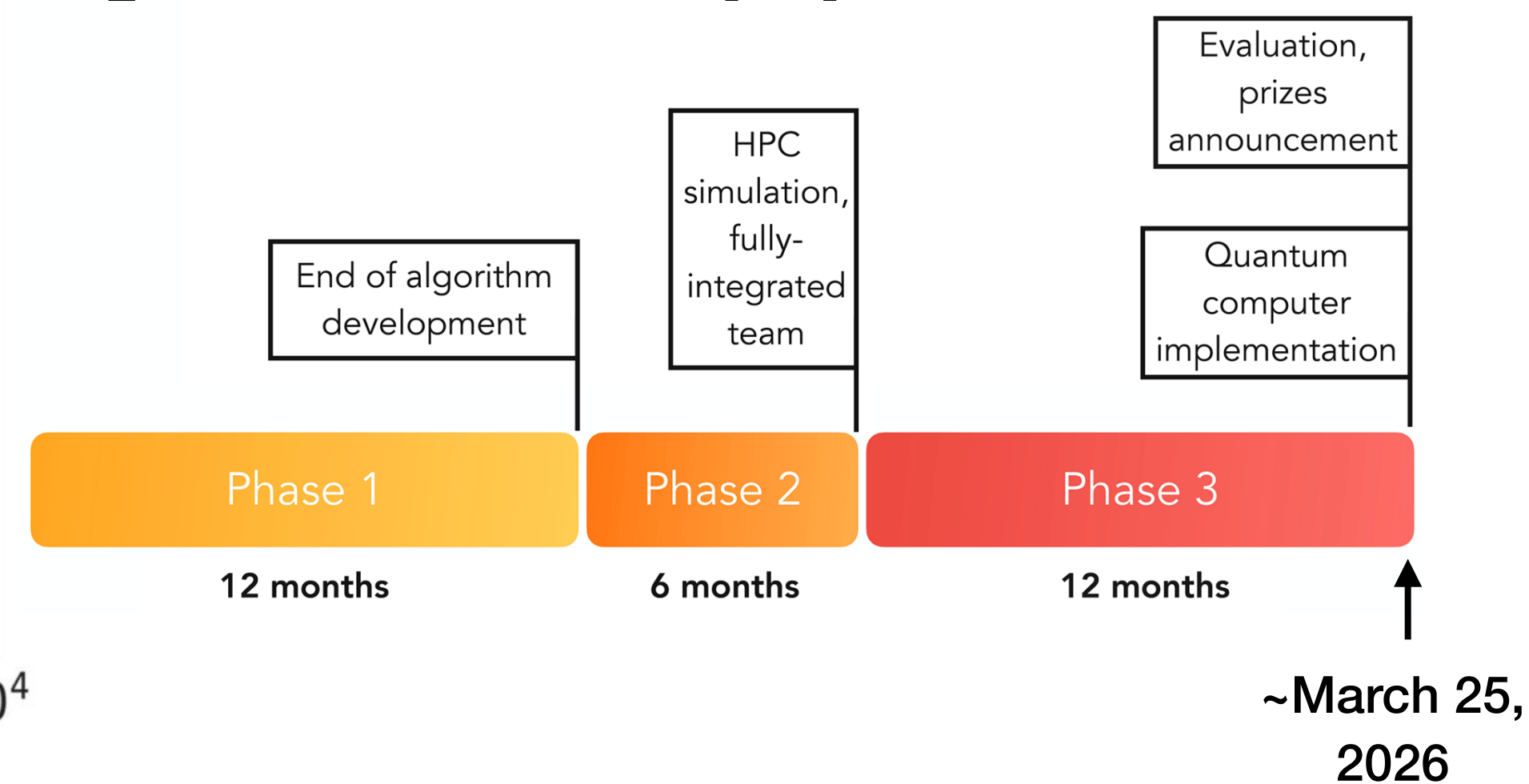
Origin of the research I'll be talking about today...



“Wellcome Leap’s Supported Challenge Program in Quantum for Bio is focused on identifying, developing, and demonstrating biology and healthcare applications that will benefit from the quantum computers expected to emerge in the next 3-5 years.”



“...The early days of any new computational method benefit from the co-development of application, software, and hardware – allowing early optimizations with not-yet-generalizable, early systems.”



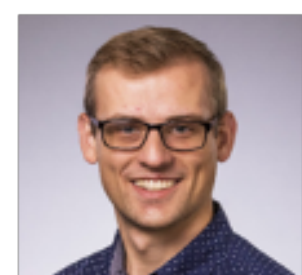
<https://wellcomeleap.org/q4bio/program/>

Multidisciplinary, cross-institutional, highly collaborative team



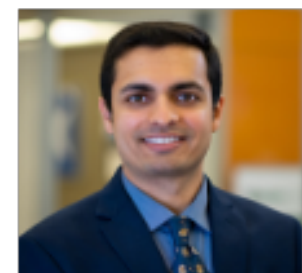
Frederic T. Chong, PhD

PI; Chief Scientist for Quantum Software; UChicago Seymour Goodman Professor



Teague Tomesh, PhD

Co-PI; Manager of Quantum Software Engineering



Pranav Gokhale, PhD

Chief Technology Officer



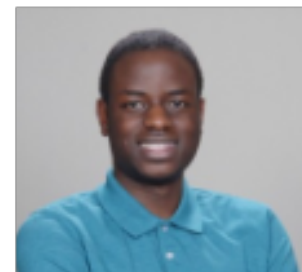
Peter Noell

Senior Quantum Solutions Manager



Colin Campbell

Quantum Applications Engineer



Victory Omole

Quantum Software Engineer



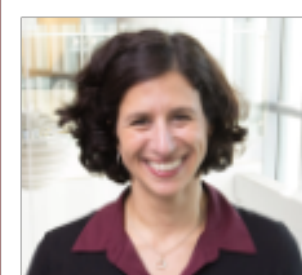
Bharath Thotakura

Quantum Software Engineer



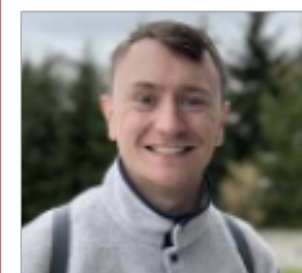
Alexander T. Pearson, MD, PhD

Co-PI; Practicing Medical Oncologist & Statistician



Samantha J. Riesenfeld, PhD

Co-PI; Asst. Prof. of Molecular Engineering & Medicine; Genomics-based ML Expert



Ryan Robinett, PhD

Postdoc, Information Theory & Manifold Learning Expert



Sid Ramesh, MD

UChicago Pritzker School of Medicine



Sophia Madejski

Computational Experimentalist, Pritzker School of Molecular Engineering



Zachary Morrell

PhD Candidate, Combinatorial Optimization Expert



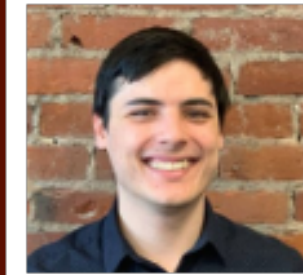
Willers Muye Yang

PhD Candidate, Quantum Computing Researcher



Aram Harrow, PhD

Co-PI; Prof. of Physics; Quantum Information and Computing Researcher



Eric Anshuetz, PhD

Quantum Algorithms Researcher



Jin-Peng Liu, PhD

Postdoc, Quantum Algorithms Researcher



Hanrui Wang

PhD Candidate, Quantum Computing Architecture and Machine Learning



Mariesa Teo

PhD Candidate, Quantum Computing Researcher



Dhirpal Shah

PhD Student, Quantum Computing Researcher



Tina Oberol

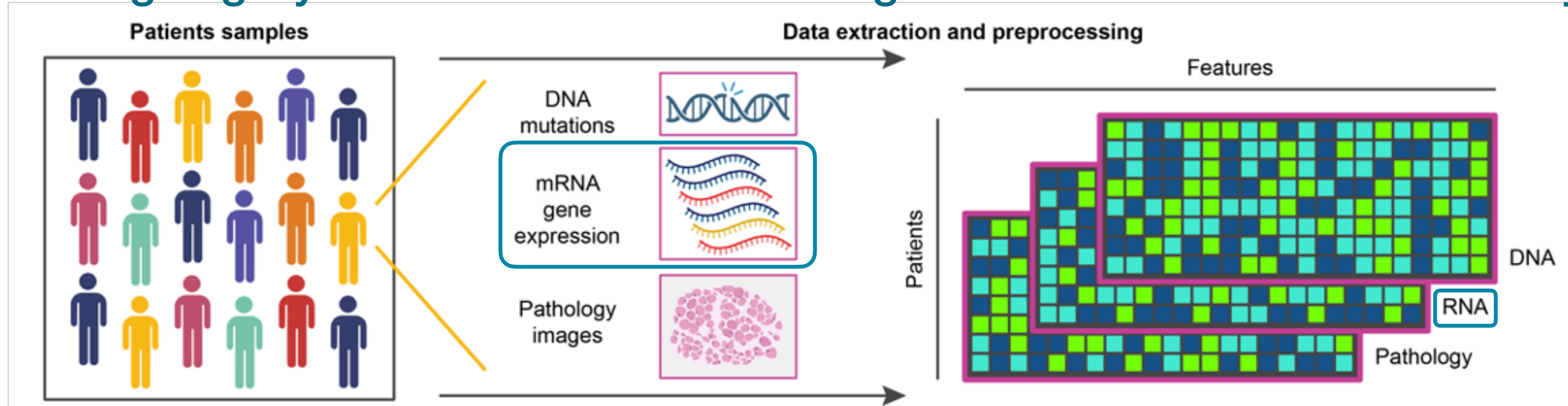
PhD Candidate, Quantum Computing Researcher



Ali-Javadi Abhari
Nate Earnest
Simon Martiel
Kevin Sung



Designing Hybrid Quantum-Classical Algorithms for Biomarker Discovery



Goal: Use these data to identify good biomarkers, i.e., measurable features indicating a clinically relevant state or outcome

Why biomarkers?

Classify

Cancer subtypes, e.g., with different prognoses

Predict

If an individual will respond well to a specific treatment

Hypothesize

What mechanism leads to differential responses?

Manipulate

Alter response by targeting the mechanism

Why interpretable biomarkers?

- Easier to deploy clinically
- Less likely to overfit

- More directly yield biological understanding
- Generate experimentally testable hypotheses



Drawbacks of Current Biomarker Approaches

1. Complex, dense, highly parameterized models

- ✓ • Learn from large-scale data
- ✓ • Useful indicators in research studies
- ✗ • Very limited interpretability, clinical deployment
- ✗ • Can overfit to training data

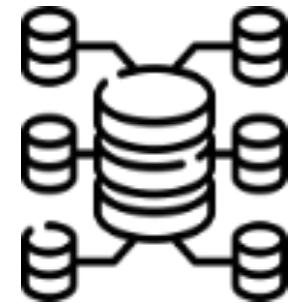
2. Very simple models (e.g., 1–3 features)

- ✓ • Interpretable, relatively broad clinical deployment
- ✓ • Biological associations can lead to therapies
- ✗ • Poor performance (false positive/negative predictions)
- ✗ • Too simple to capture complex biology

Feature selection could address many of these challenges and complement current approaches, but optimizing feature sets is hard with classical computers

Challenges in Core Problem of Biomarker Discovery

Why is identifying clinically viable, multimodal cancer biomarkers by finding compact, informative feature sets in vast search spaces hard?



Sample Scarcity

Limited patient samples (due to cost, privacy, IP, etc.) relative to feature space → **overfitting**, poor generalization



Complex Synergies & Interactions

Higher-order interactions within and across data modalities (e.g., DNA, RNA, pathology features)



Computational Goals & Constraints

High premium on **interpretability** and **accuracy/precision** (rather than speed), yet classical solvers struggle

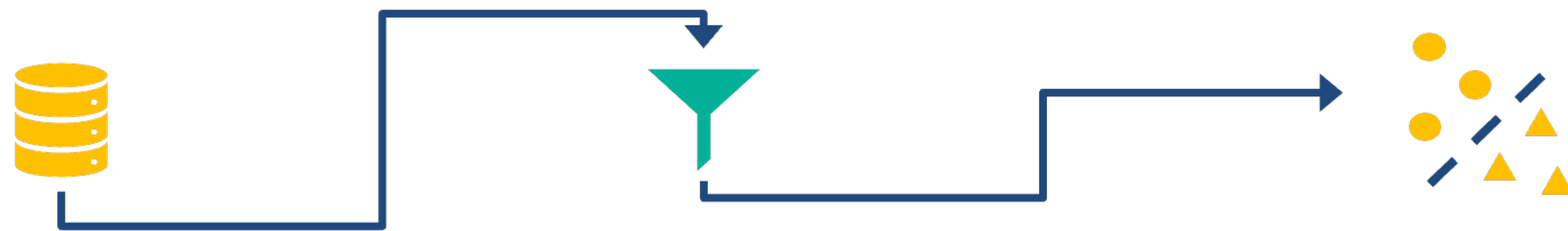


Data Explosion

New modalities generate 1000s of new features → Exponential search space

Designing Hybrid Quantum-Classical Algorithms for Biomarker Discovery

Goal: Identify a small feature set that maximizes performance by a downstream learning model on a particular task



1. Data (e.g., from public databases) is collected, cleaned, and preprocessed

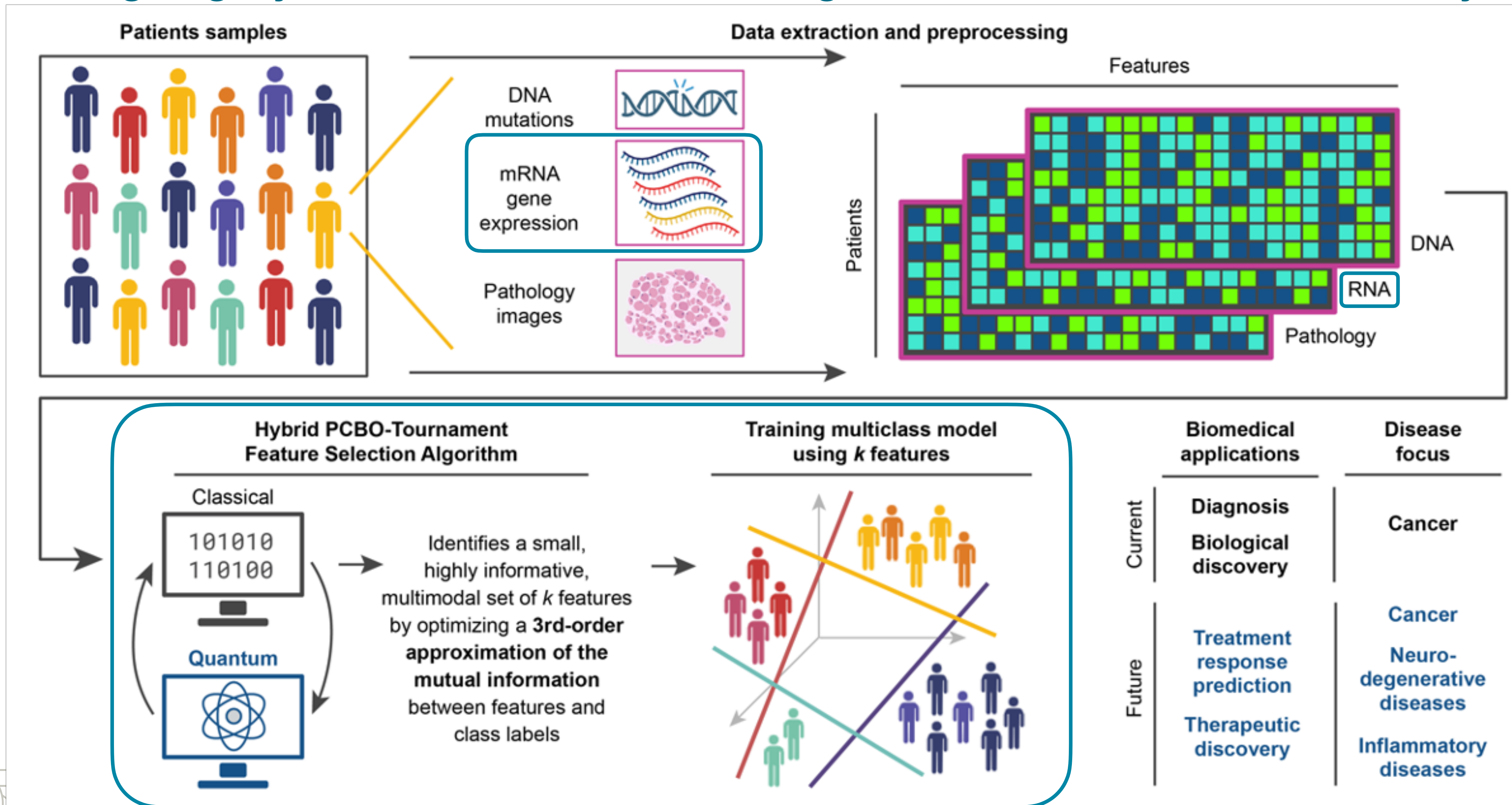
2. Feature selection performed, with the aid of quantum computers

3. Simple (currently, classical) machine learning models are trained

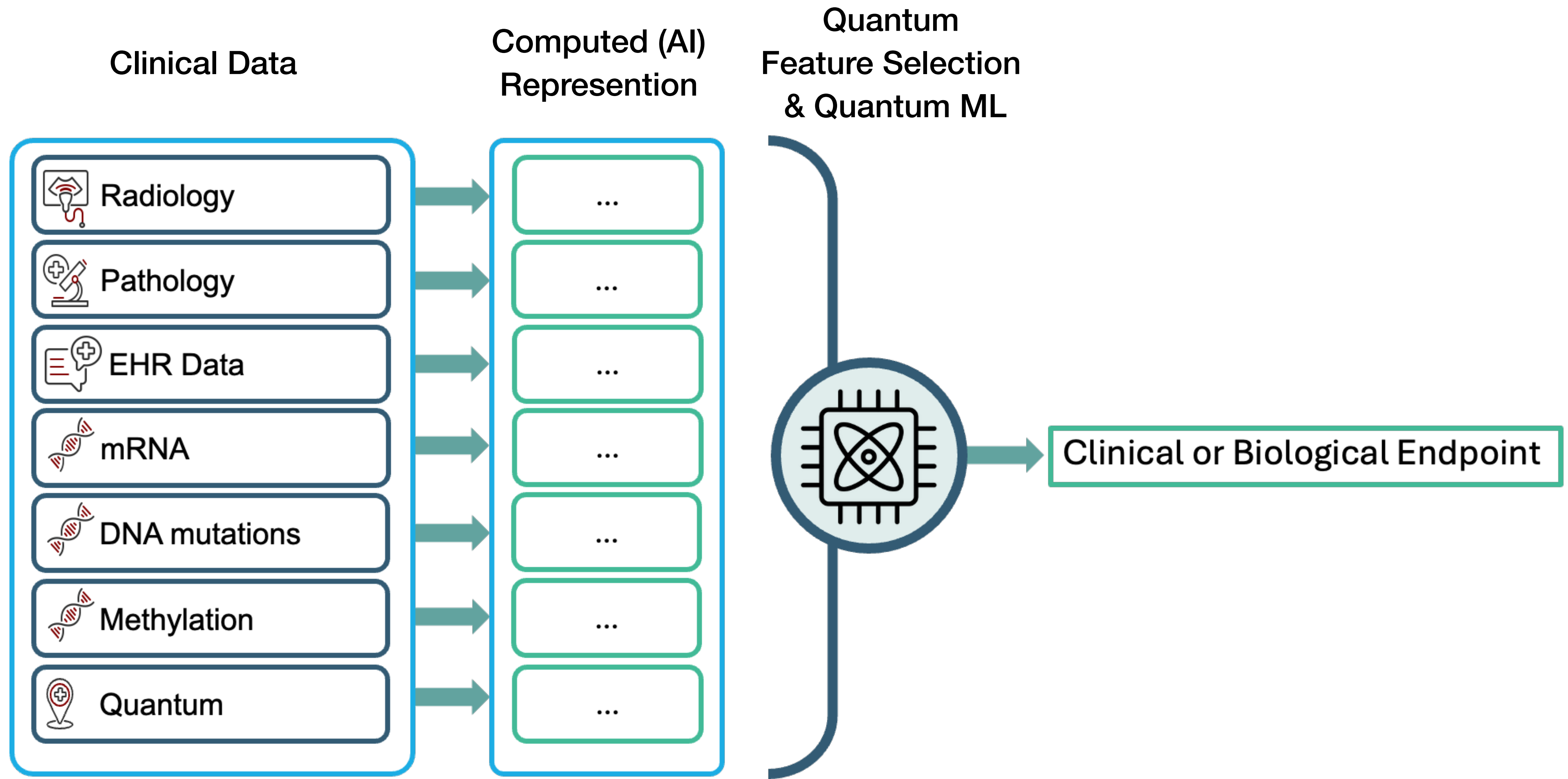
The pipeline is composed of both **classical** and **quantum** steps

A quantum computing framework supports feature selection through effective combinatorial optimization

Designing Hybrid Quantum-Classical Algorithms for Biomarker Discovery



Future Goals...



Formulate Feature Selection as Information-Theoretic Optimization Problem

Each dataset consists of:

- m patients, n initial features: f_1, \dots, f_n (we assume to be discrete, random, variables)
 - Transcript counts discretized into quintiles, i.e., $f_i \in \{0,1,2,3,4\}$, computed per gene across samples
- Target y (label/class/prediction)
 - Initially we focus on $y \in \{0, \dots, 9\}$ as an indicator of cancer tissue of origin among 9 cancer types
- Each patient is considered to be a sample of the **joint variable** (f_1, \dots, f_n, y)
- Our high-level problem:
$$\arg \max_{\{i_1, \dots, i_k\} \subset [1, n]} I\left(y; f_{i_1}, f_{i_2}, \dots, f_{i_k}\right)$$
 Mutual information (MI) of target and feature set

MI well established in prior biological work for measuring relevancy of and reducing redundancy among features

How might we make this optimization more feasible?

Formulate Feature Selection as Information-Theoretic Optimization Problem

$$I\left(y; f_{i_1}, f_{i_2}, \dots, f_{i_k}\right) = H\left(f_{i_1}, f_{i_2}, \dots, f_{i_k}\right) - H\left(f_{i_1}, f_{i_2}, \dots, f_{i_k} | y\right)$$

For a partition \mathcal{P} of the k features into subsets $S_i, 1 \leq i \leq l$,

$$\begin{aligned} H\left(f_{i_1}, f_{i_2}, \dots, f_{i_k}\right) &= H\left(S_1 \cup \dots \cup S_l\right) \\ &= H\left(S_1\right) + H\left(S_2 | S_1\right) + \dots + H\left(S_l | S_1 \cup \dots \cup S_{l-1}\right) \end{aligned} \quad \text{Chain rule of entropy}$$

We target a trade-off in expressivity and efficiency by assuming that feature interactions have order ≤ 3 , i.e., there is some partition \mathcal{P} with $|S_i| \leq 3$ for all i , such that:

$$\begin{aligned} H\left(f_{i_1}, f_{i_2}, \dots, f_{i_k}\right) &= H\left(S_1\right) + H\left(S_2\right) + \dots + H\left(S_l\right) \\ &\stackrel{?}{\approx} \frac{c}{\binom{k}{3}} \sum_{\substack{S \subset \{i_1, \dots, i_k\}, \\ |S| = 3}} H(S) \end{aligned} \quad \begin{array}{l} \text{Average over all partitions of features into triplets,} \\ \text{since } \mathcal{P} \text{ unknown} \end{array}$$

Formulate Feature Selection as a Polynomial Unconstrained Boolean Optimization (PUBO) Problem

$$\arg \max_{\{i_1, \dots, i_k\} \subset \{1, \dots, n\}} I(y; f_{i_1}, f_{i_2}, \dots, f_{i_k}) \stackrel{?}{\approx} \arg \max_{\substack{\{i_1, \dots, i_k\} \subset [1, n] \\ |S| = 3}} \frac{c}{\binom{k}{3}} \sum_{S \subset \{i_1, \dots, i_k\}} (H(S) - H(S|y))$$

$$\propto \arg \max_{\vec{x} \in \mathbb{Z}_2^n, \|\vec{x}\|_1 = k} \sum_{\substack{\{i_1, i_2, i_3\} \subset [1, n] \\ i_1 < i_2 < i_3}} \left(H(f_{i_1}, f_{i_2}, f_{i_3}) - H(f_{i_1}, f_{i_2}, f_{i_3} | y) \right) x_{i_1} x_{i_2} x_{i_3}$$

PCBO problem

$$= \arg \max_{\vec{x} \in \mathbb{Z}_2^n, \|\vec{x}\|_1 = k} \sum_{\substack{\{i_1, i_2, i_3\} \subset [1, n] \\ i_1 < i_2 < i_3}} C_{i_1 i_2 i_3} x_{i_1} x_{i_2} x_{i_3} \quad (\text{generalization of a Sherrington-Kirkpatrick spin glass})$$

“entropy-CUBO”,
aka “H-CUBO”

$$Q_{\text{H-CUBO}}(\vec{x}, k) = - \sum_{\substack{\{i_1, i_2, i_3\} \subset [1, n] \\ i_1 < i_2 < i_3}} C_{i_1 i_2 i_3} x_{i_1} x_{i_2} x_{i_3}$$

PUBO problem, as a Hamiltonian
(Hard constraint of PCBO can be absorbed as a soft quadratic penalty with a user-set Lagrange multiplier)

H-CUBO makes features selection amenable to quantum algorithms designed for PCBO-type problems

We focus on a novel variant of the Quantum Approximation Optimization Algorithm (QAOA),
a heuristic approach for combinatorial optimization (Farhi, Goldstone. 2014)

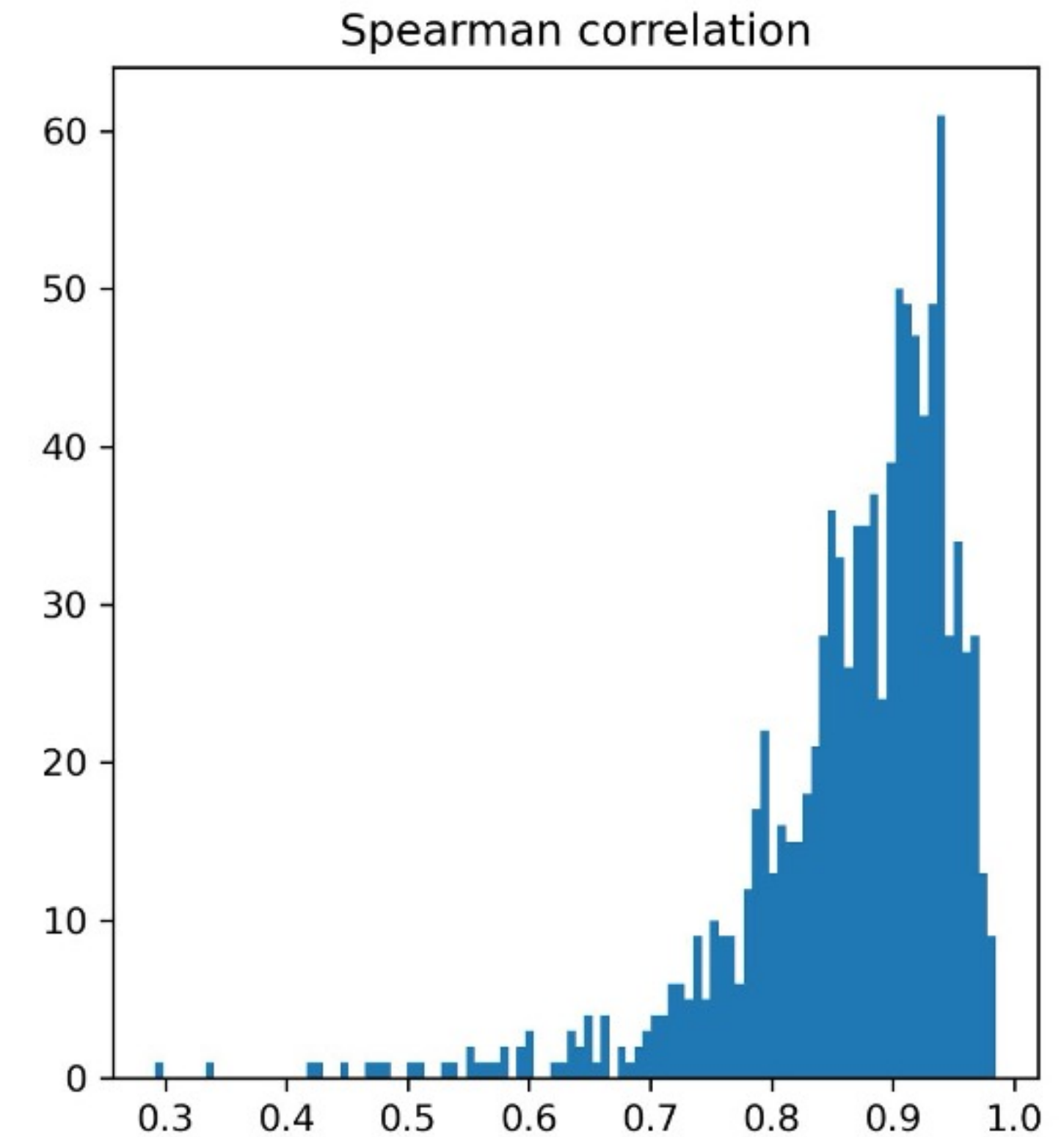
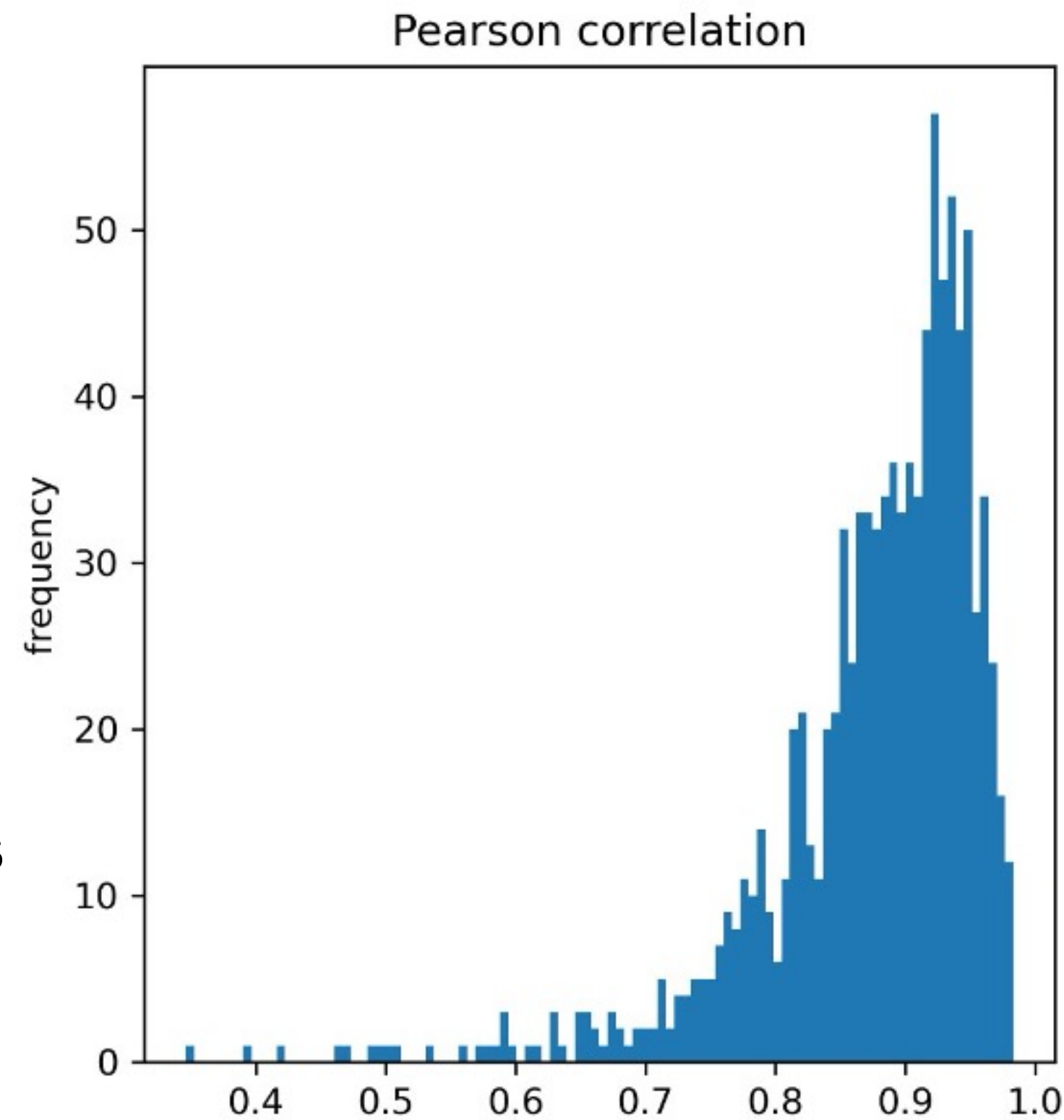
Empirical sanity checks suggest H-CUBO is a reasonable substitute for MI

Example:

For 1,000 problems of size $\binom{12}{6}$
(i.e., $n=12$, $k=6$), we computed the correlation of empirical MI and energy of H-CUBO across candidate solutions

- To enable MI computations, we use small problem sizes, and expression is discretized into tertiles (not quintiles)

Results suggest reasonable correlation





High-level view of quantum computing, with an eye toward QAOA

A *qubit* refers to a quantum system whose classical state set is $\{0,1\}$.

“Pure” quantum states are *superpositions*, essentially linear combinations, of the classic (z -basis) states:

$$|\psi\rangle = \alpha |0\rangle + \beta |1\rangle \text{ for } \alpha, \beta \in \mathbb{C}$$

($|0\rangle$ and $|1\rangle$ could be any orthonormal basis for the 2-dimensional qubit Hilbert space)

If a quantum state is measured, each classical state of the system appears with probability equal to the *absolute value squared* of the entry in the quantum state vector corresponding to that classical state.

Operations on quantum state vectors are represented by *unitary* matrices

A quantum algorithm often involves preparation of a state and then measuring it, which can be akin to sampling from a potentially complex distribution

QAOA generates a quantum circuit whose output distributionally approaches the ground state of the given Hamiltonian

“Cost” Hamiltonian (i.e., for H-CUBO):

$$U(C, \gamma) = e^{-i\gamma C}$$

“Mixer” Hamiltonian:

$$U(B, \beta) = e^{-i\beta B} \text{ for } B = \sum_{j=1}^n \sigma_j^x$$

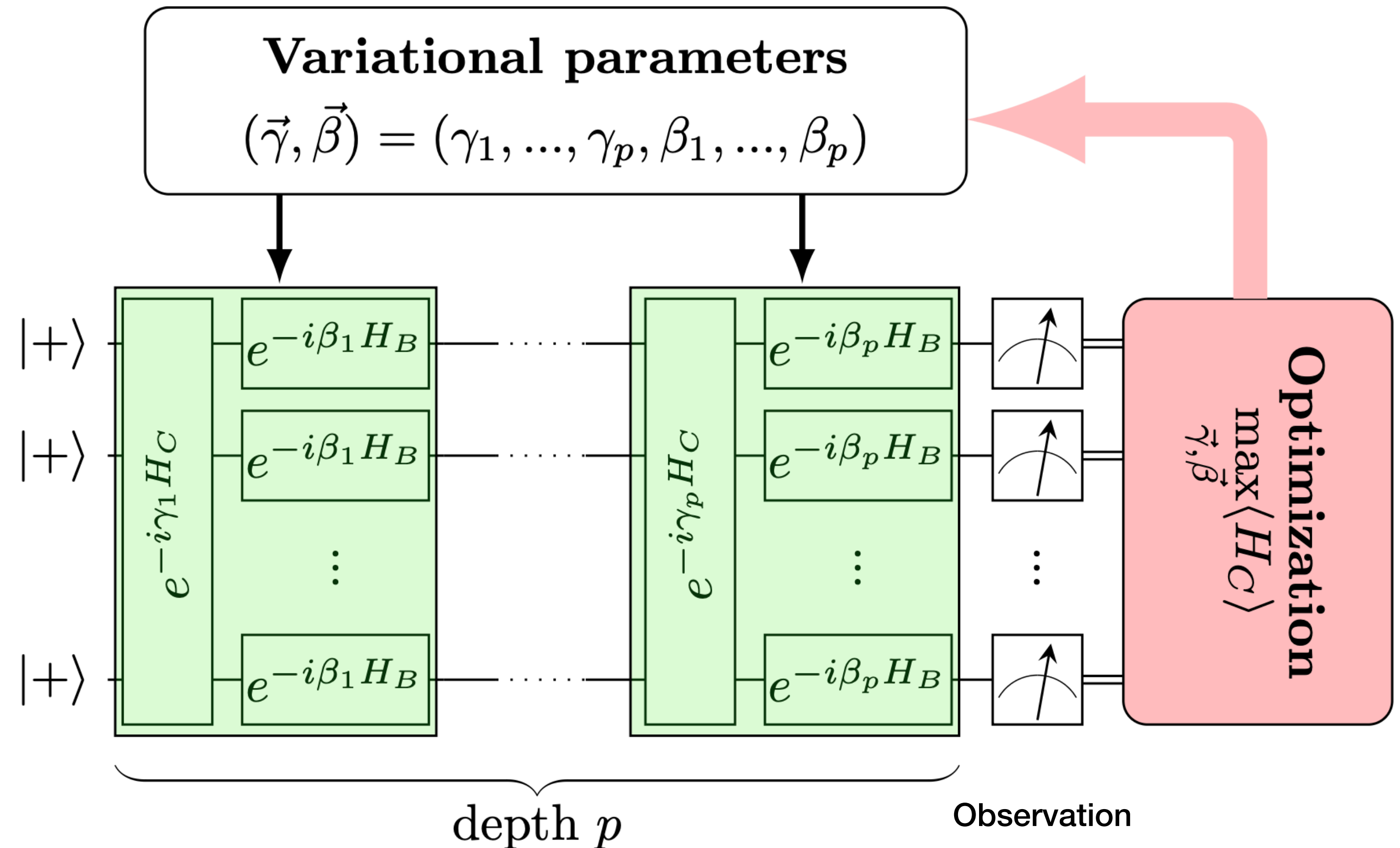


Image from Lee, et al. arXiv 2108.05288v1

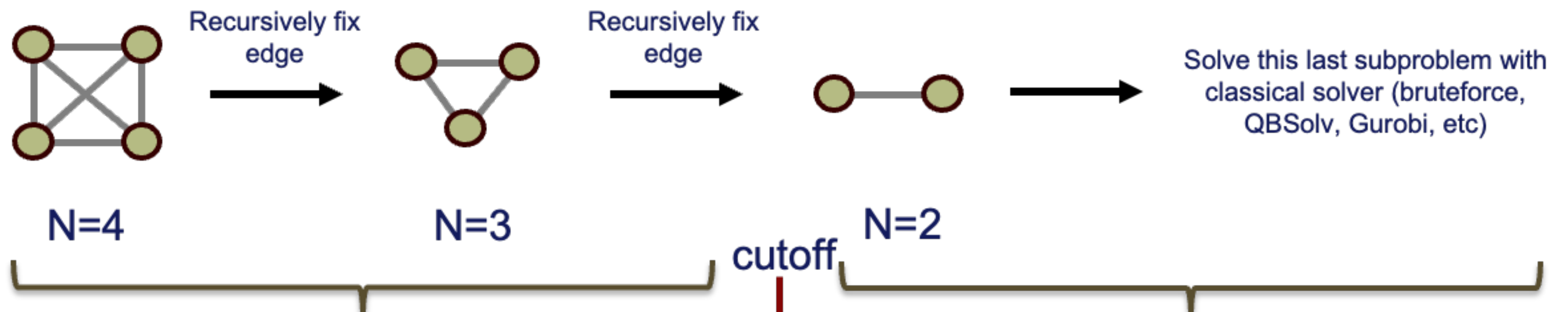
Problem: Optimization of variational parameters is hard/expensive

Potential solution: Recursive QAOA (RQAOA) (Bravyi, et al. Phys Rev Letters 2020)

How Recursive QAOA (RQAOA) helps leverage quantum resources

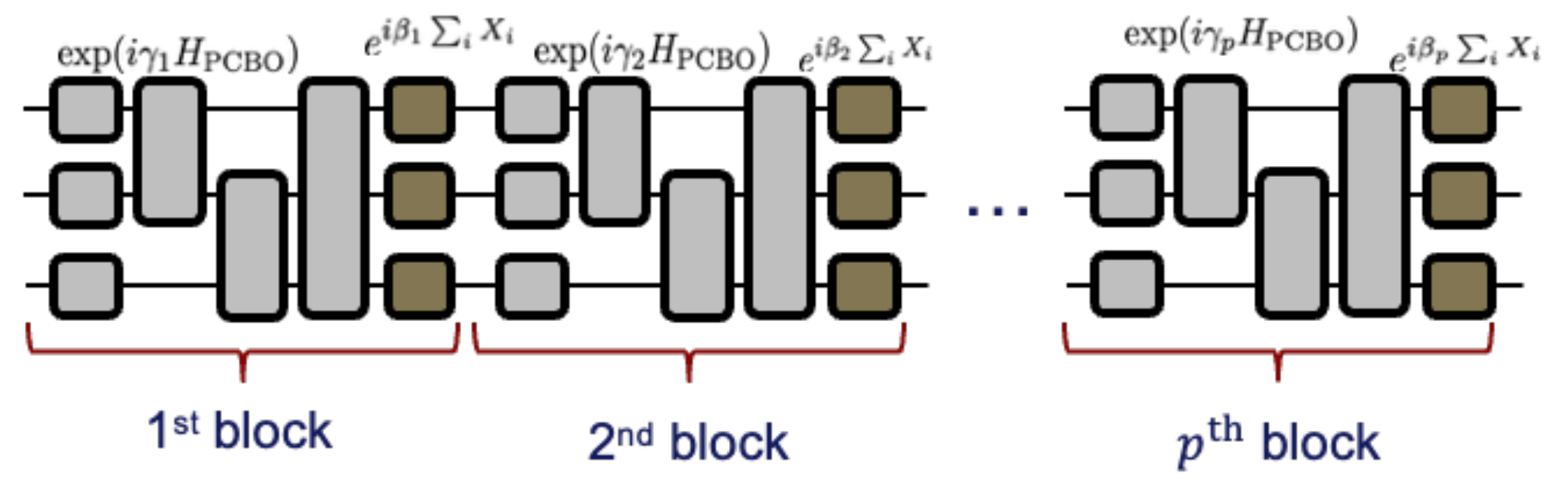
$$H = \sum_{i < j < k} C_{ijk} Z_i Z_j Z_k$$

Third-order "edge"

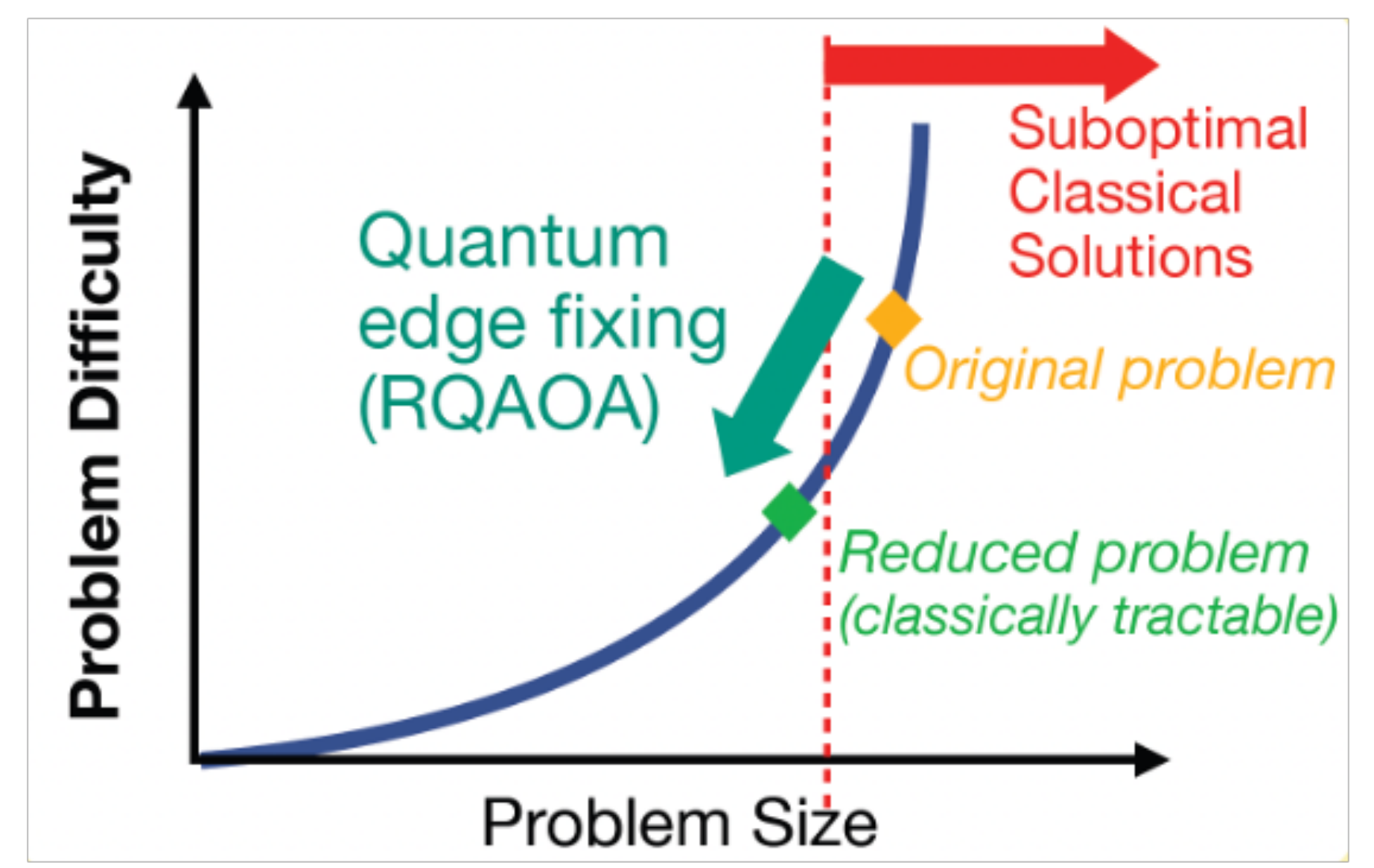


Quantum

Classical

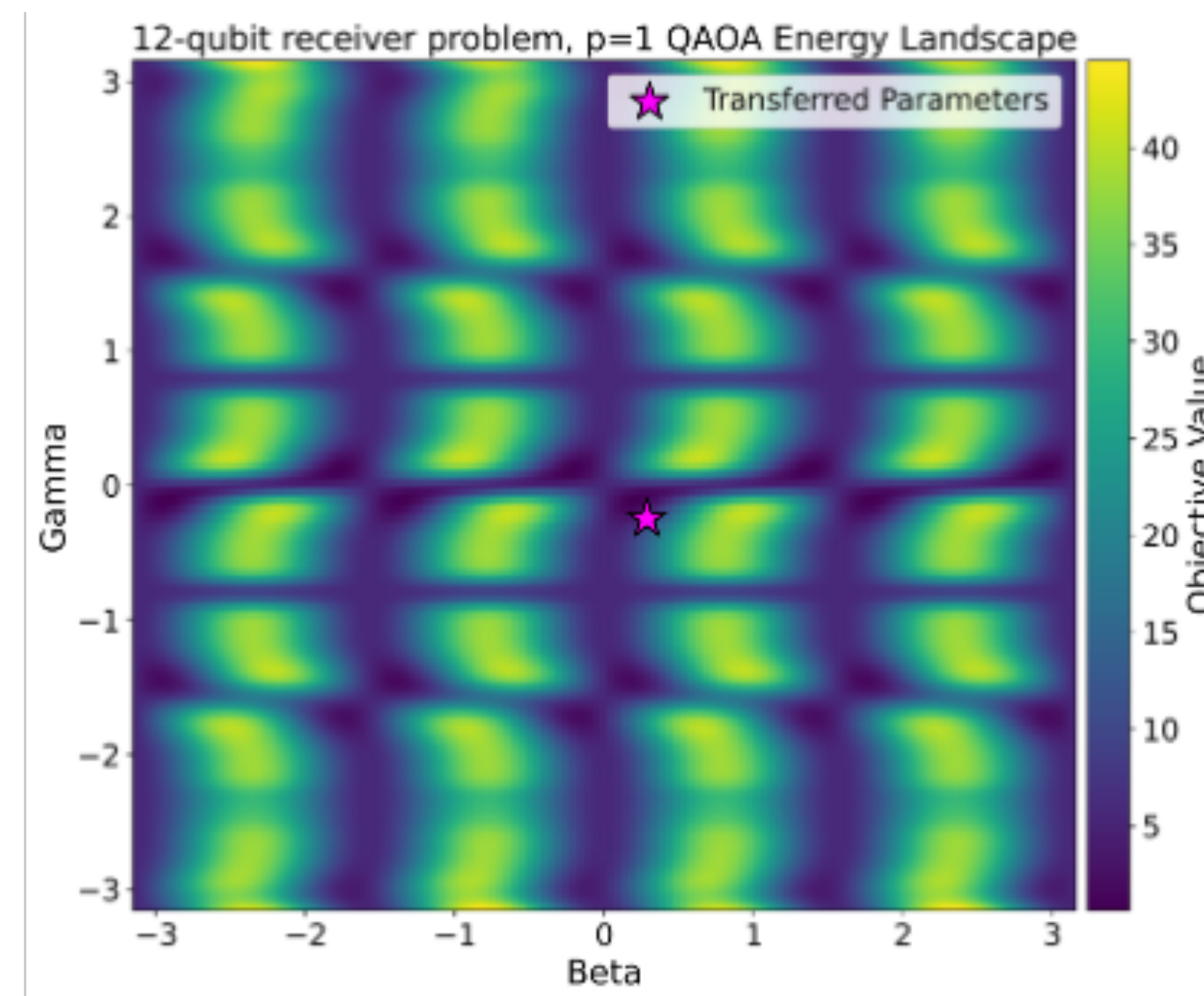
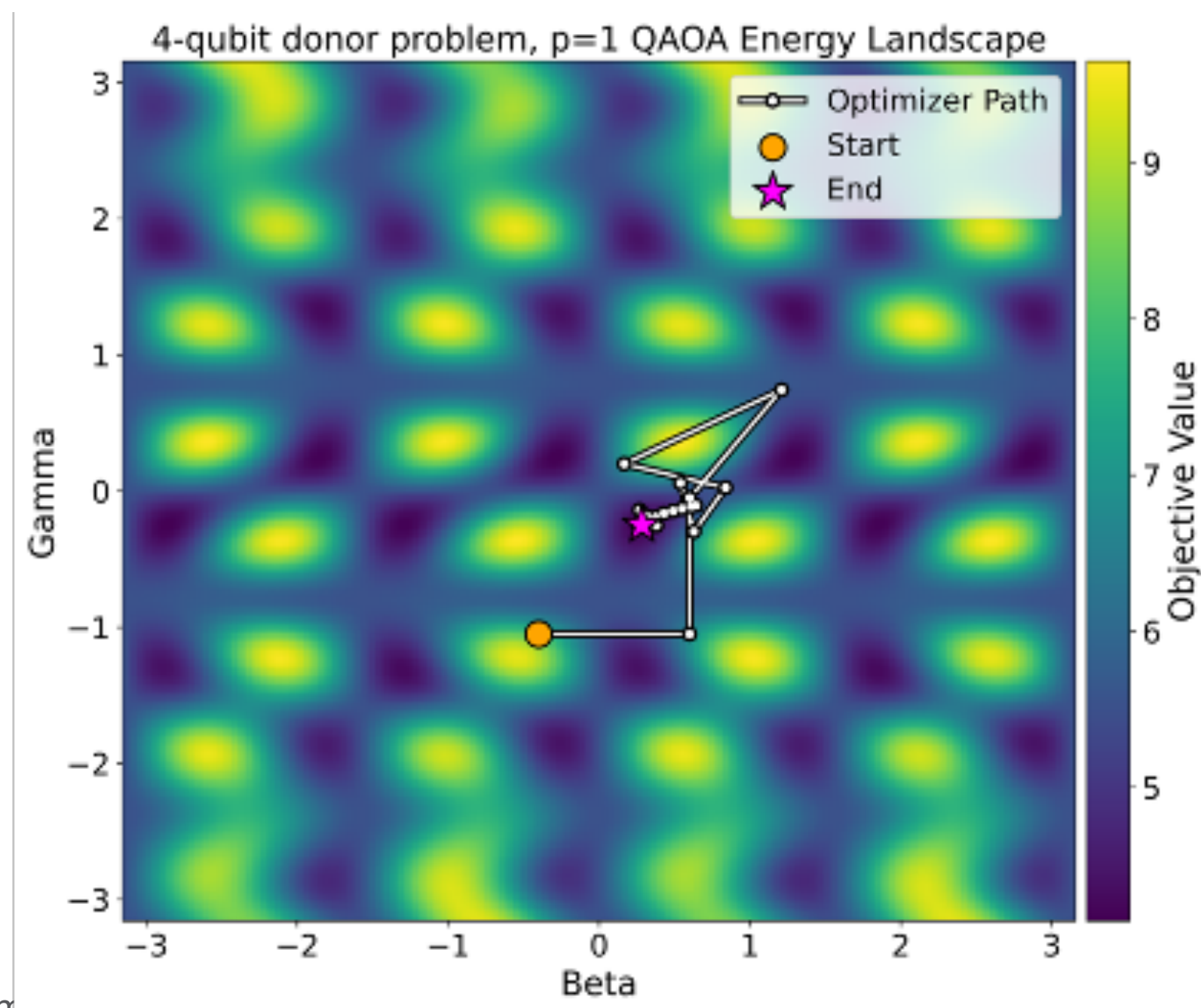
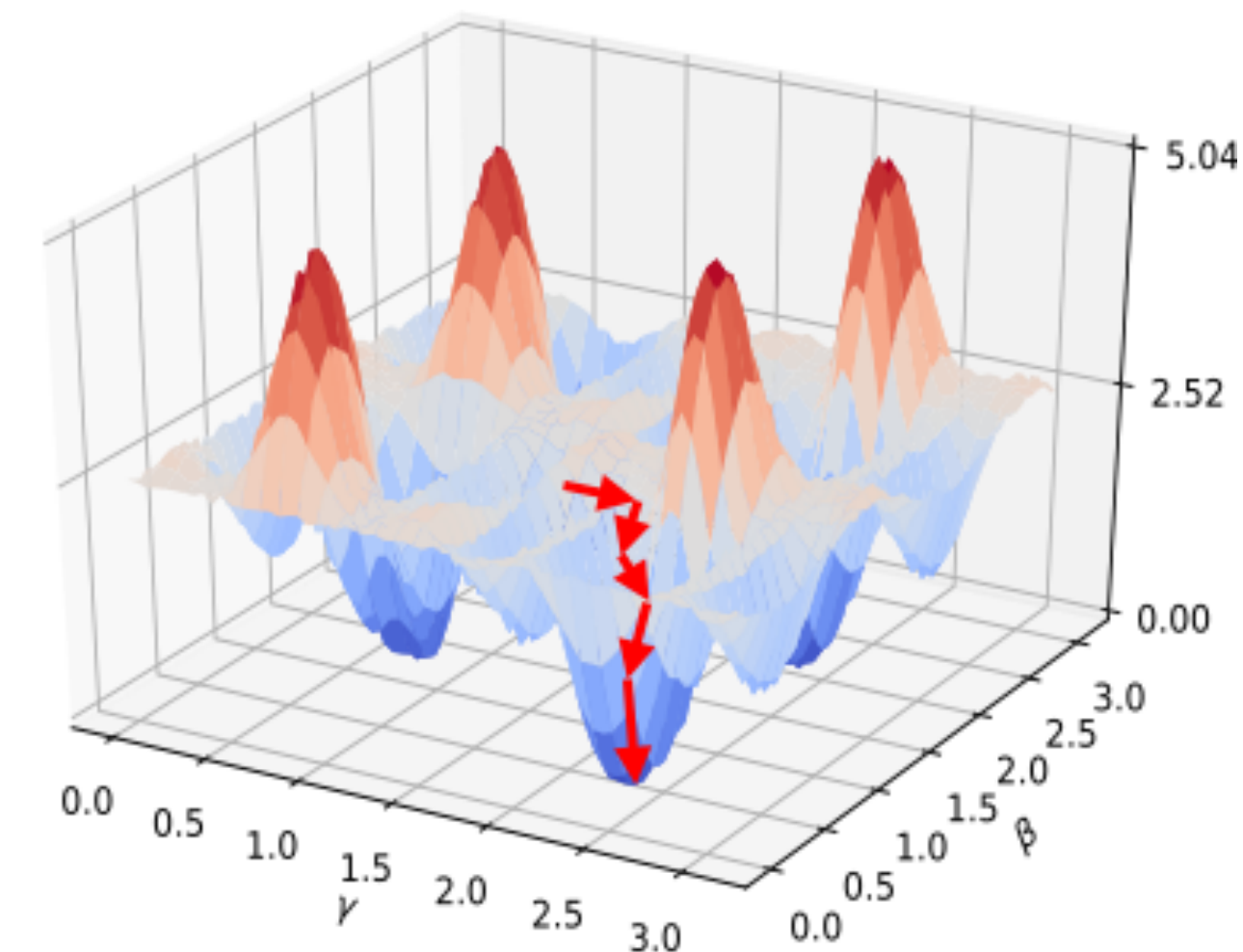


- We use the probability distribution, obtained by sampling from the QC many times, to determine which edge to fix



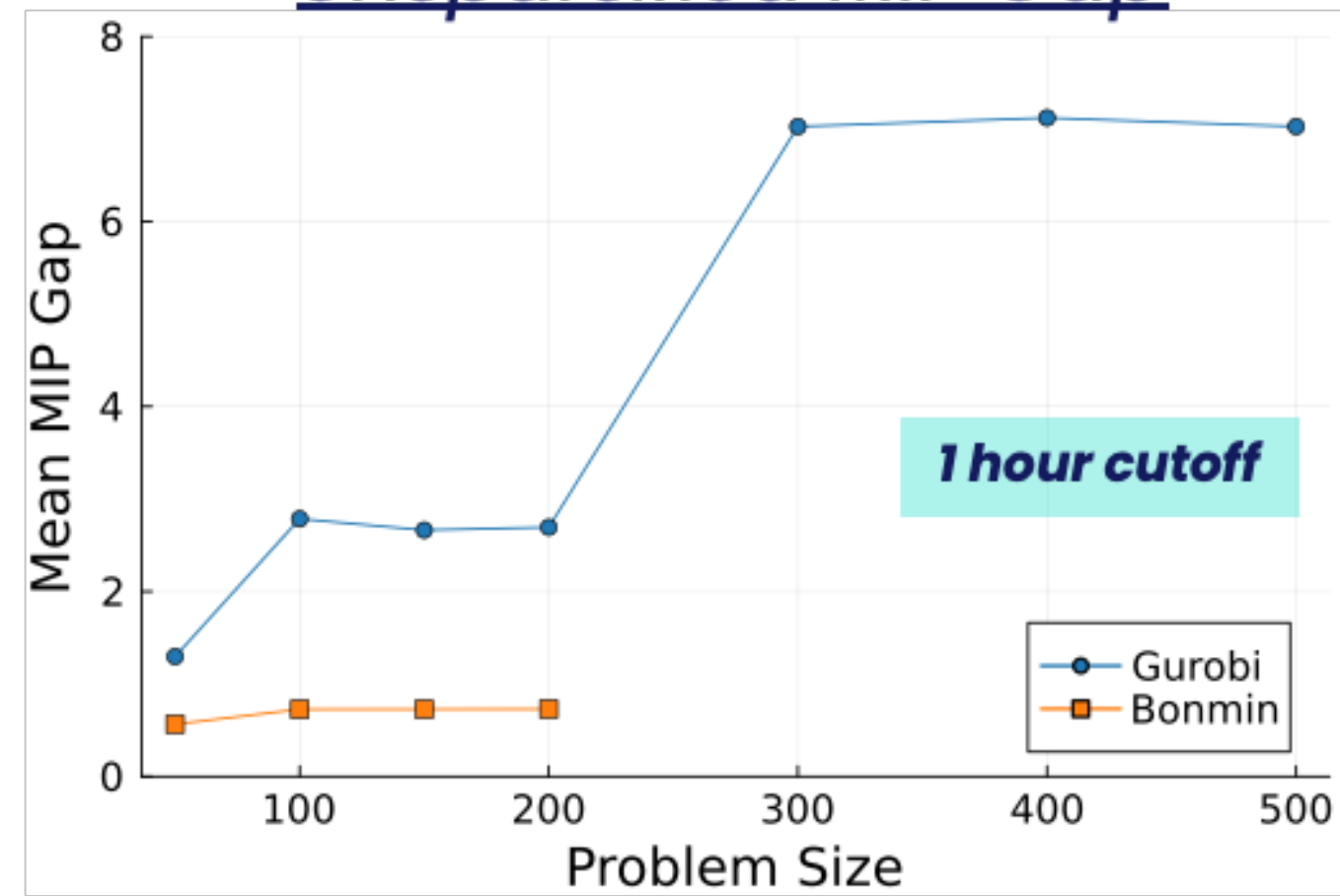
We use parameter transfer to address trainability challenges

- Parameter transfer is a powerful technique for **warm starting** the ansatz to circumvent barren plateau results
- Allows us to initialize in the vicinity of a good local minimum
- Empirically, this is exactly what we observe
- Reduces # iterations by 1000X

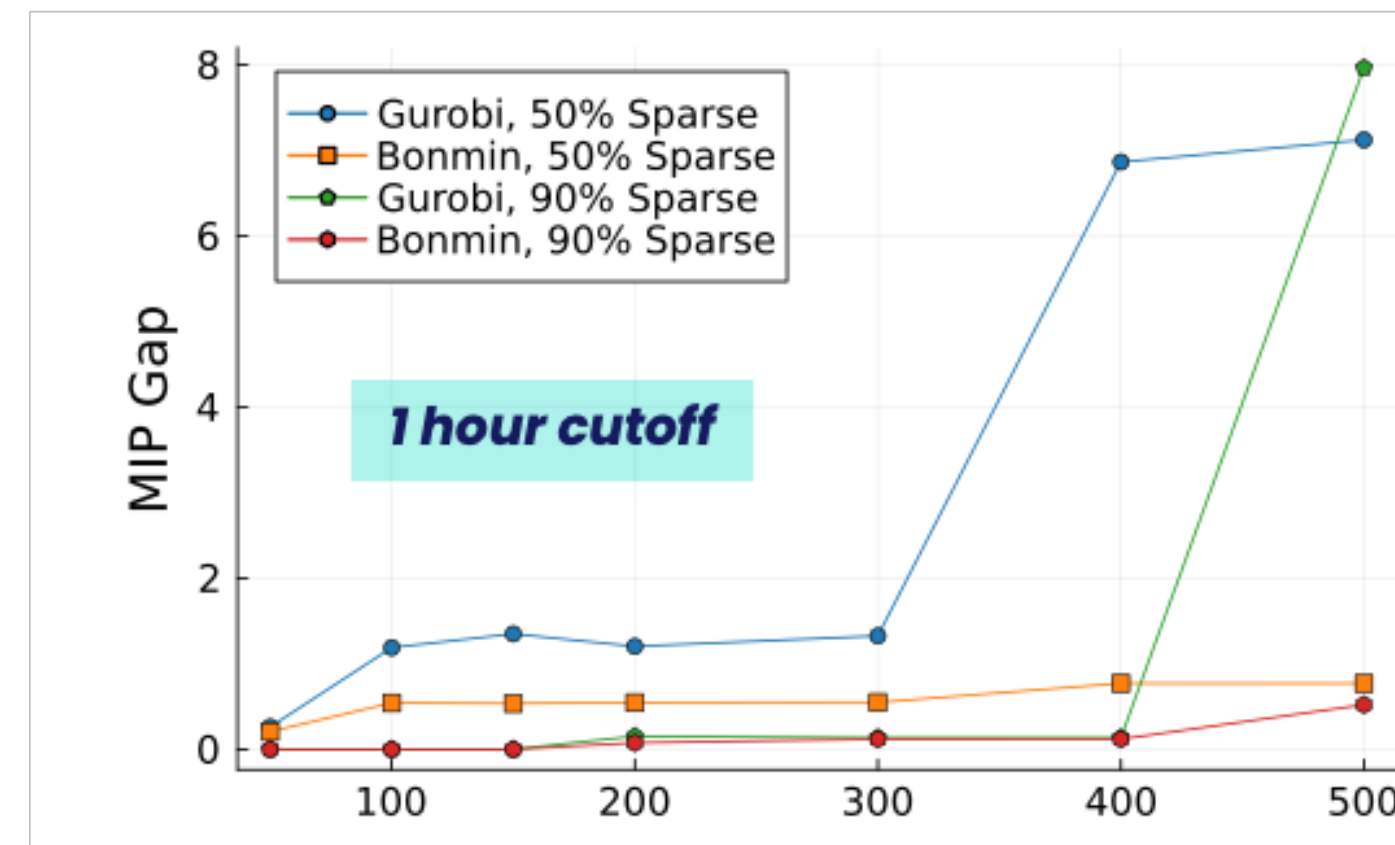


“Architecture-topology-aware” problem sparsification improves efficiency, but approximation gap saturates with increasing problem size

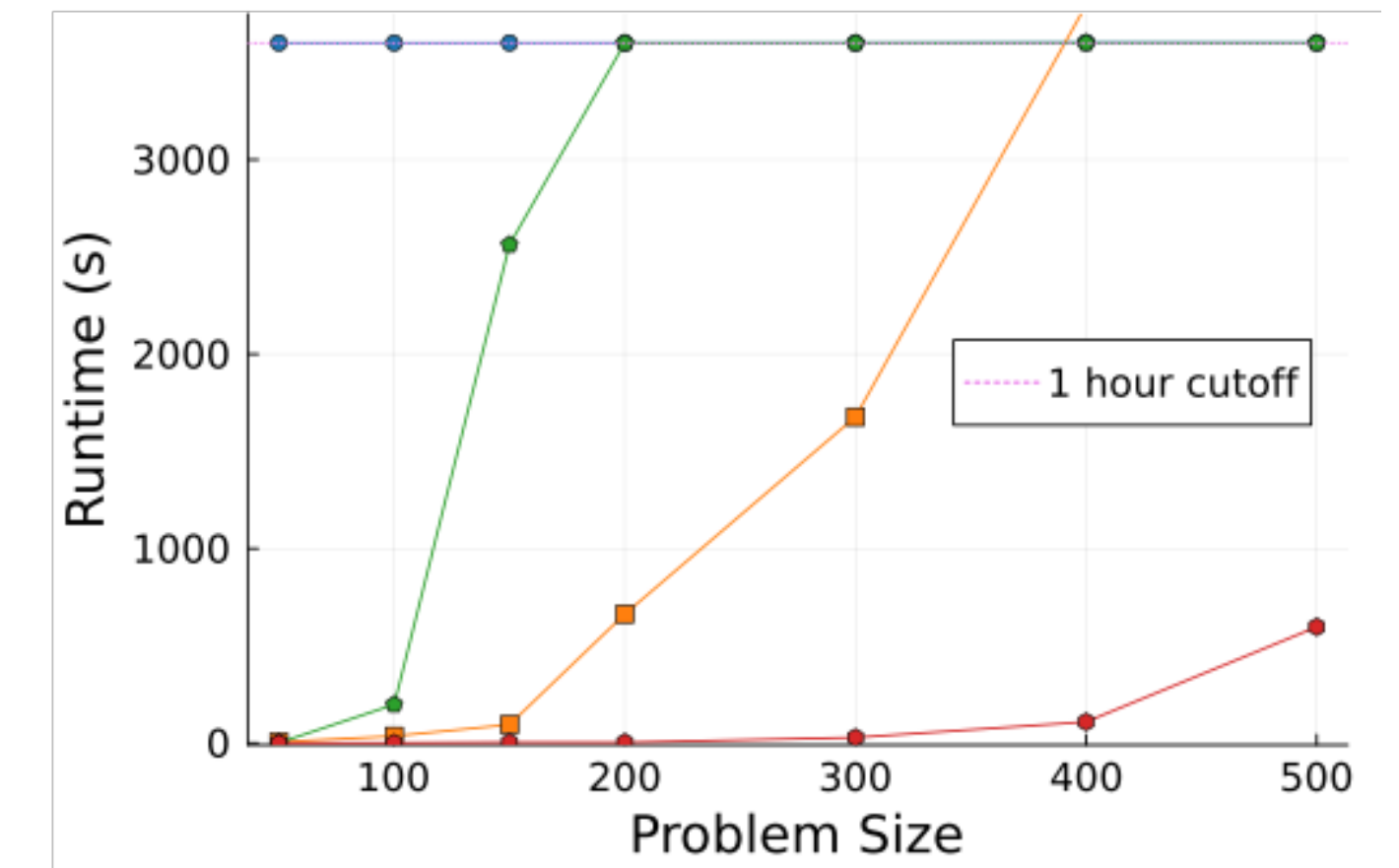
Unsparsified MIP Gap



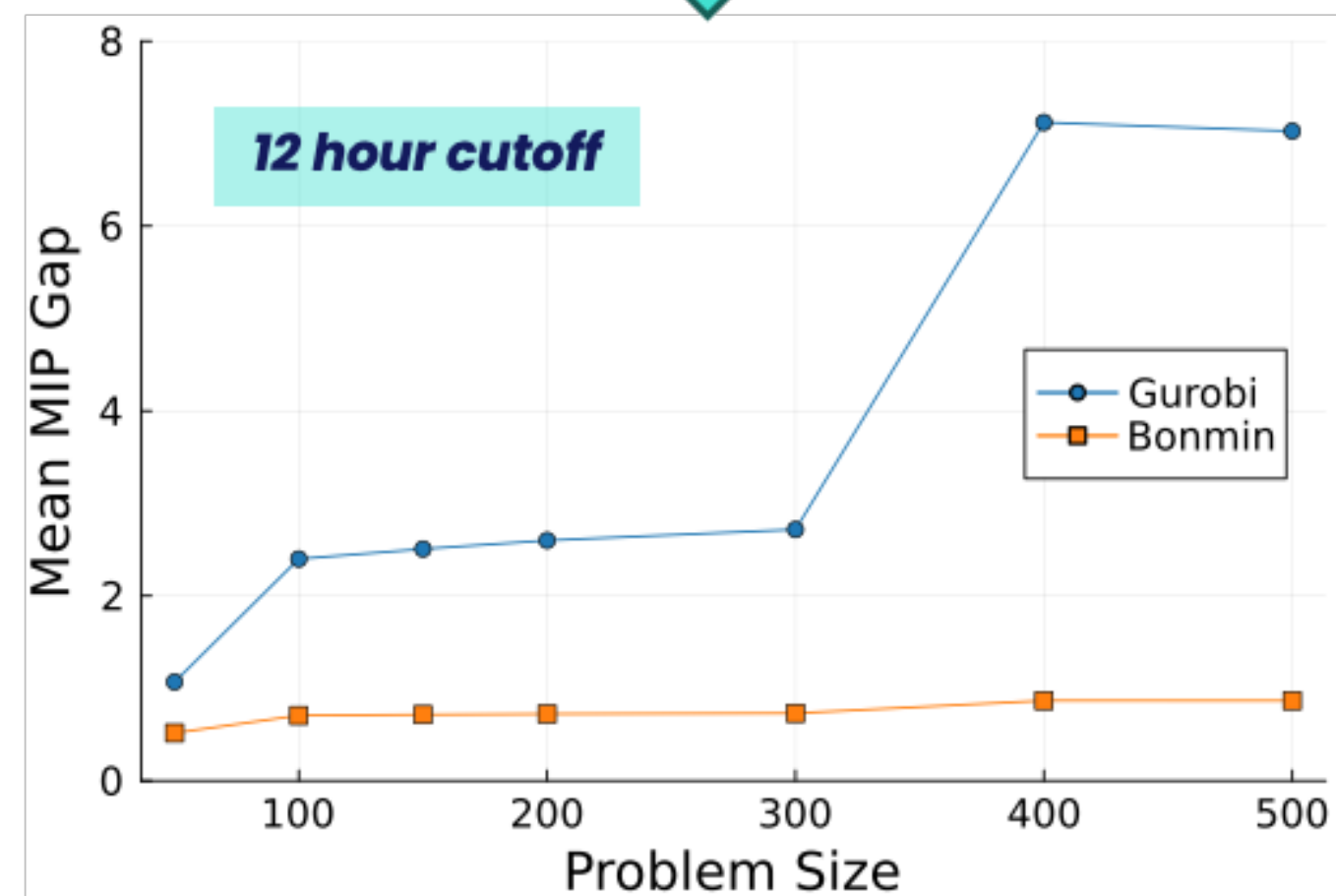
Sparsified MIP Gap



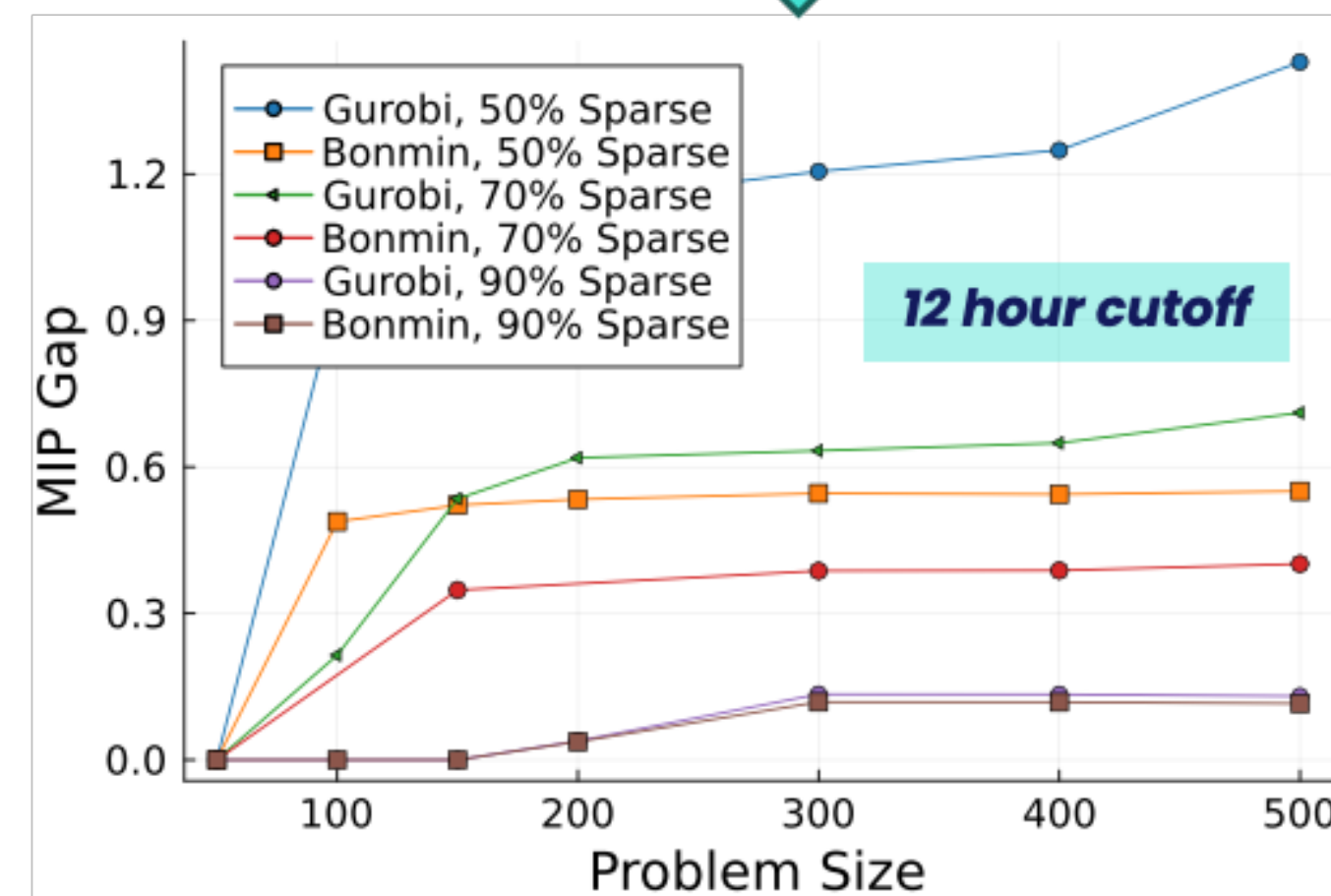
Sparsified Runtime



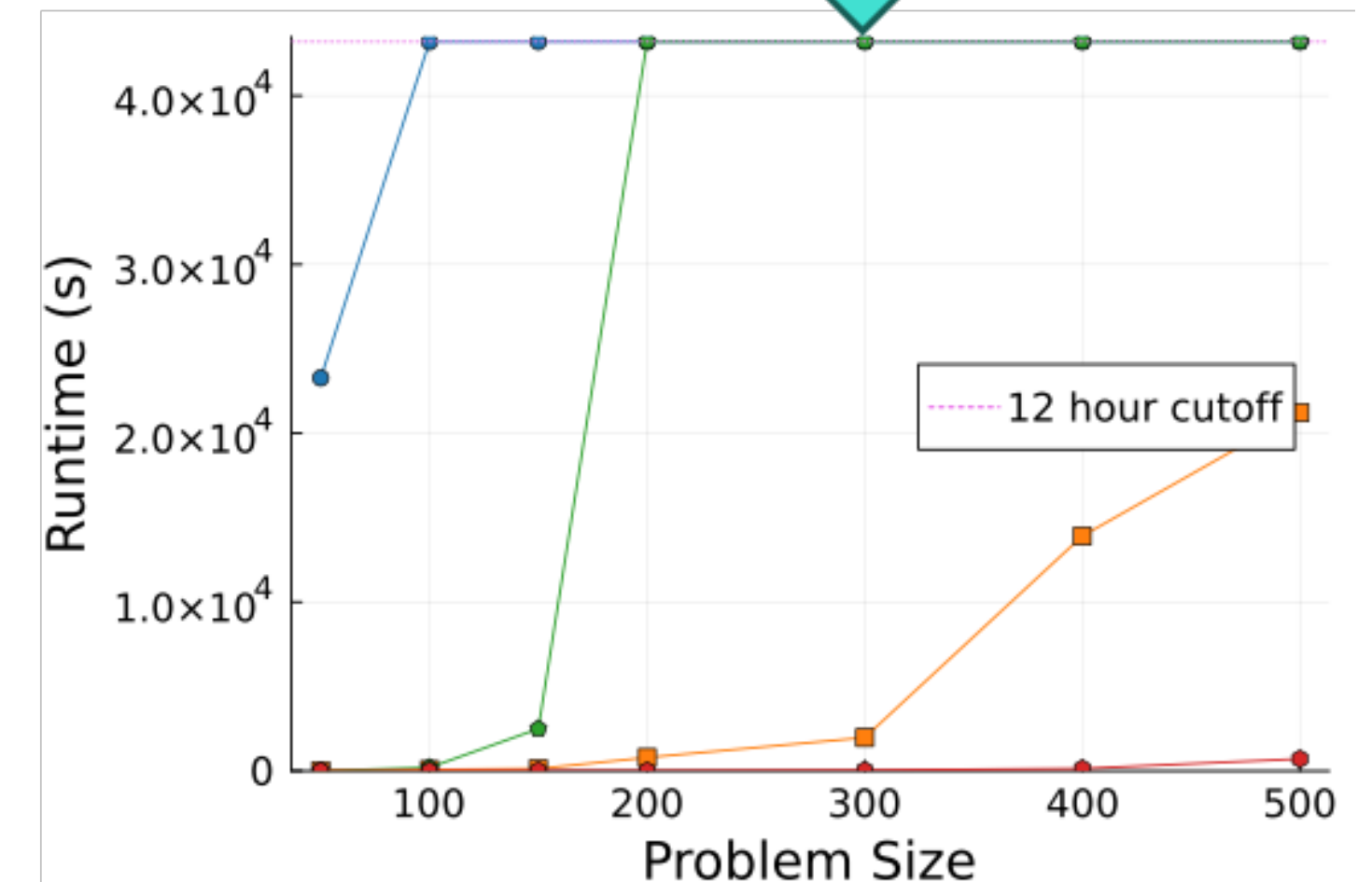
12 hour cutoff



12 hour cutoff



12 hour cutoff

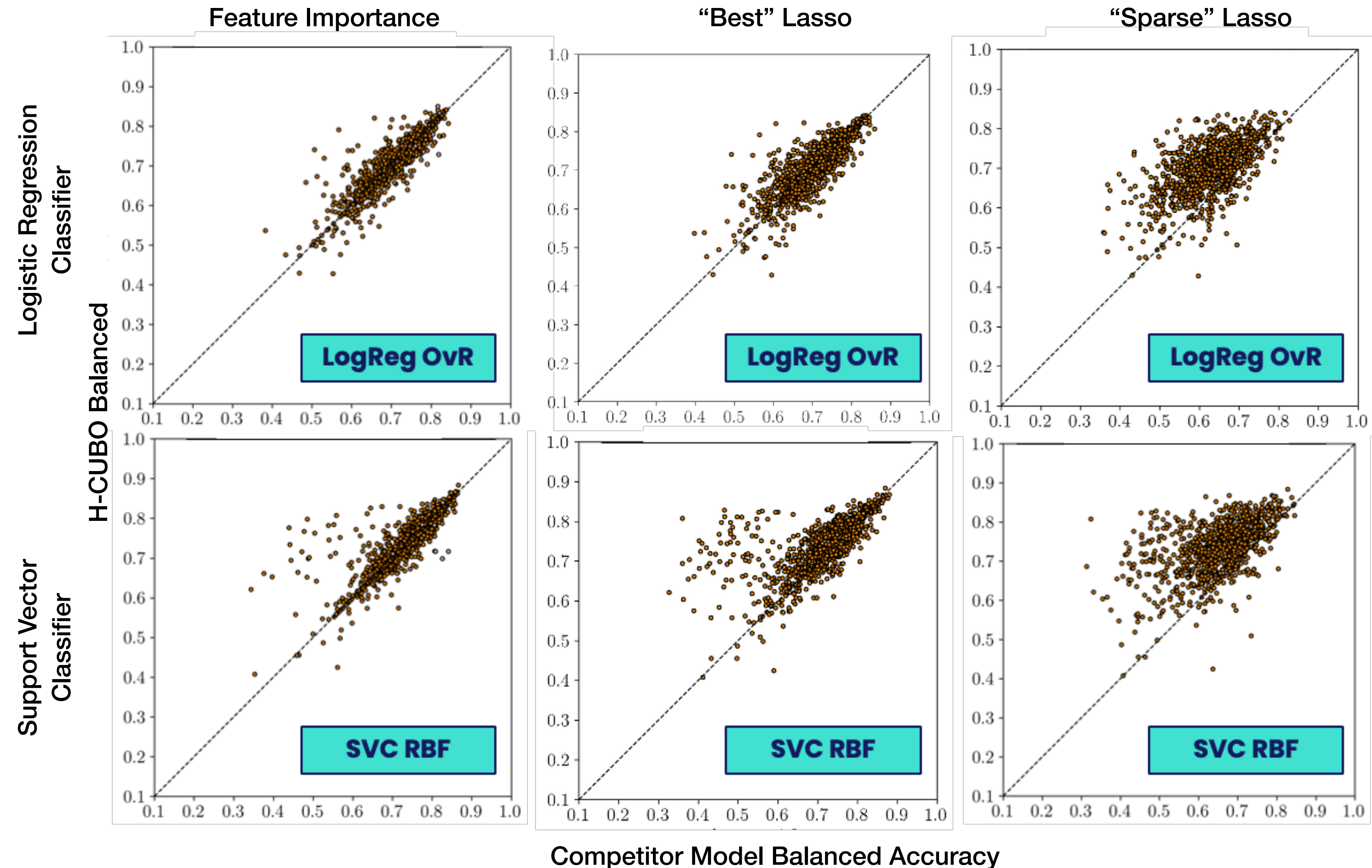


H-CUBO Performs Well Against Classical Feature Selection Methods

Classifier models were trained on 10-feature solutions selected from 20-feature problems (randomly selected from genes in data), by 4 feature selection methods:

- Entropy-CUBO (y axes)
- Feature Importance
- “Best” Lasso: find optimal fit, take top k features
- “Sparse” Lasso: constrain model until only k features have non-zero weight

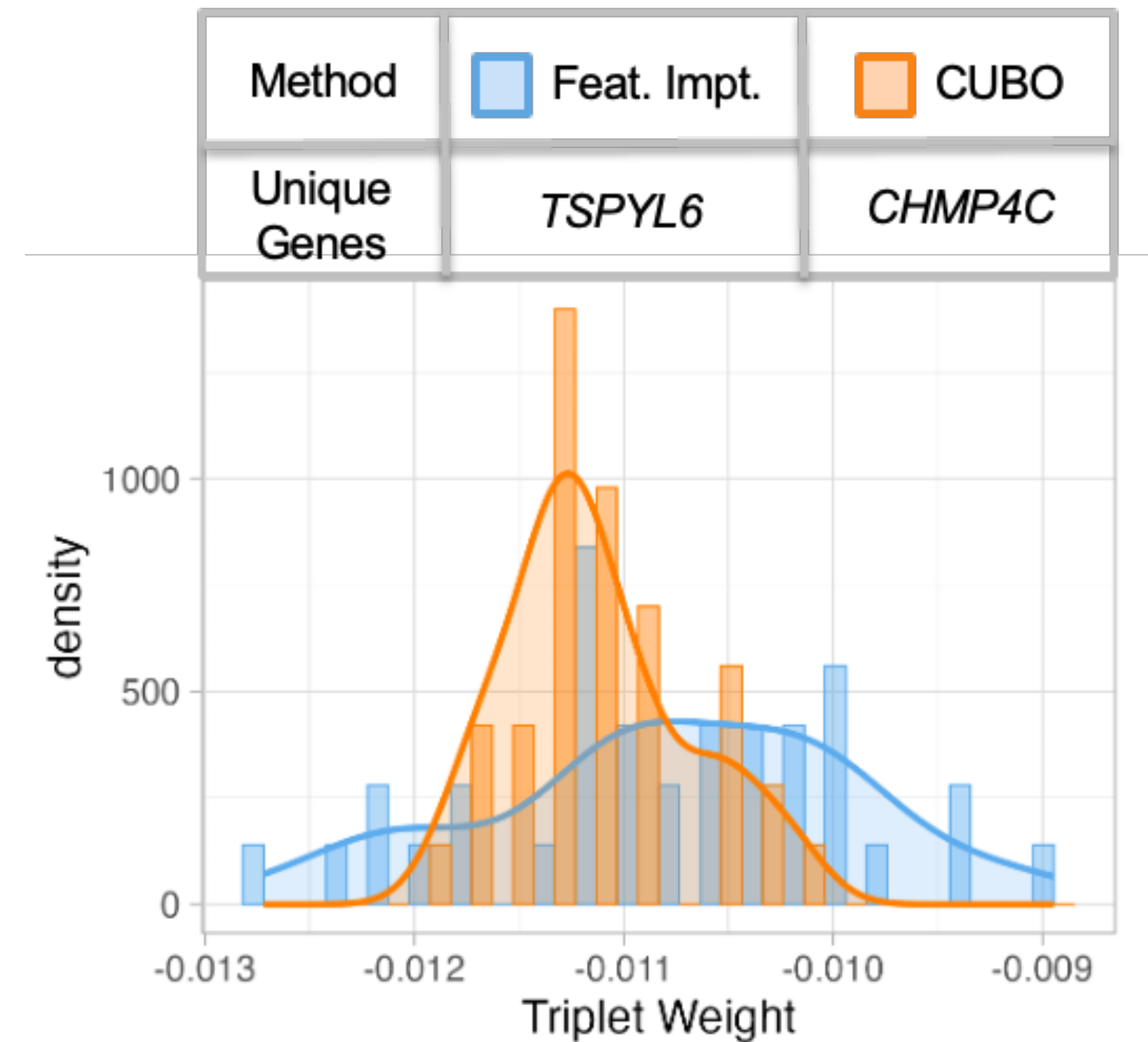
For 1000 problems, training data are from TCGA (large public database); all models tested in CPTAC (different public database); patient samples represent 9 cancer types (target labels)



Small taste of what we hope to dig more into in the results...

Results Example

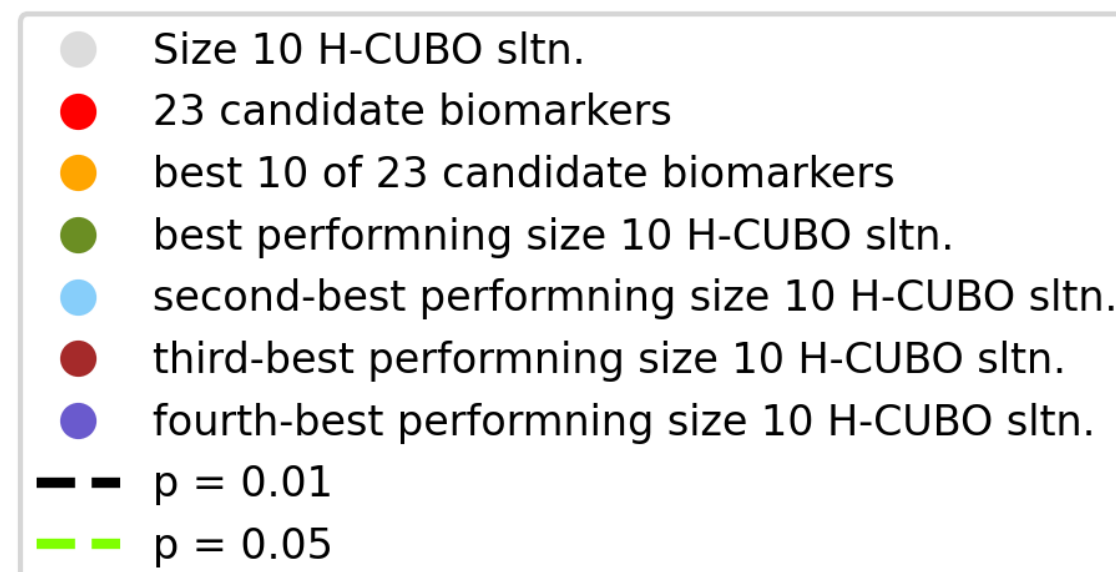
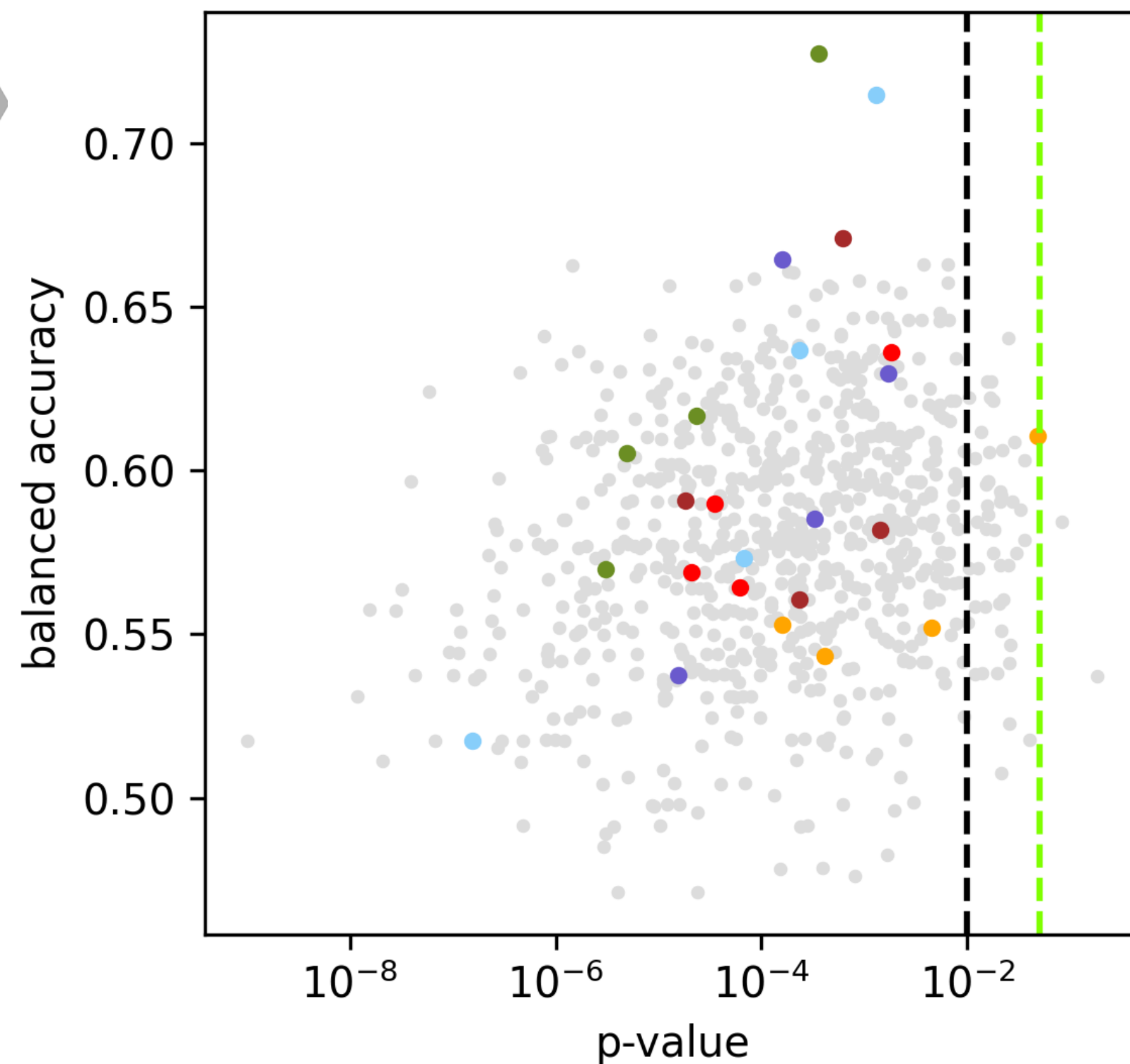
- CUBO and feature importance solution sets overlap
 - Genes expressed in specific cell types
 - Some known prognostic markers of cancers
 - **Differences?**
- **Classical selection: *TSPYL6***
 - Chromatin binder & histone binding
 - Individual high-effect gene
- **Quantum-hybrid selection: *CHMP4C***
 - Component in key protein transport machinery
 - **Larger joint mutual information** about target in triplets with other solution genes
 - **Associated with *unique processes, tissues, and cancers***, relative to other genes in solution
 - **Result: Covers *distinct facets of cancer biology***



Moving into more challenging applications...

Improving biomarkers for predicting tumor treatment response

logreg classifier performance on ISPY redux data



Data is from a seminal breast cancer study (I-SPY2); we focus on subset of patients highly enriched for triple-negative breast cancer (few targeted therapies)

Our prediction target: pathologic complete response (pCR) (across a few treatment arms)

Early results here are from 200 problems of size 30-choose-10, constructed via importance sampling of gene expression

Here, each 10-feature solution is evaluated using 4-fold cross-validation (with a learned model for each fold)

We also evaluated set of 23 candidate biomarkers (not individual genes), reportedly each associated with pCR

Results: Using only this simple, sparse strategy, we already identify 10-feature subsets that outperform existing biomarkers



Acknowledgments



Riesefeld Group

Ryan Robinett, PhD

Sophia Madejski

Ruxandra Tonea

Akansha Gupta

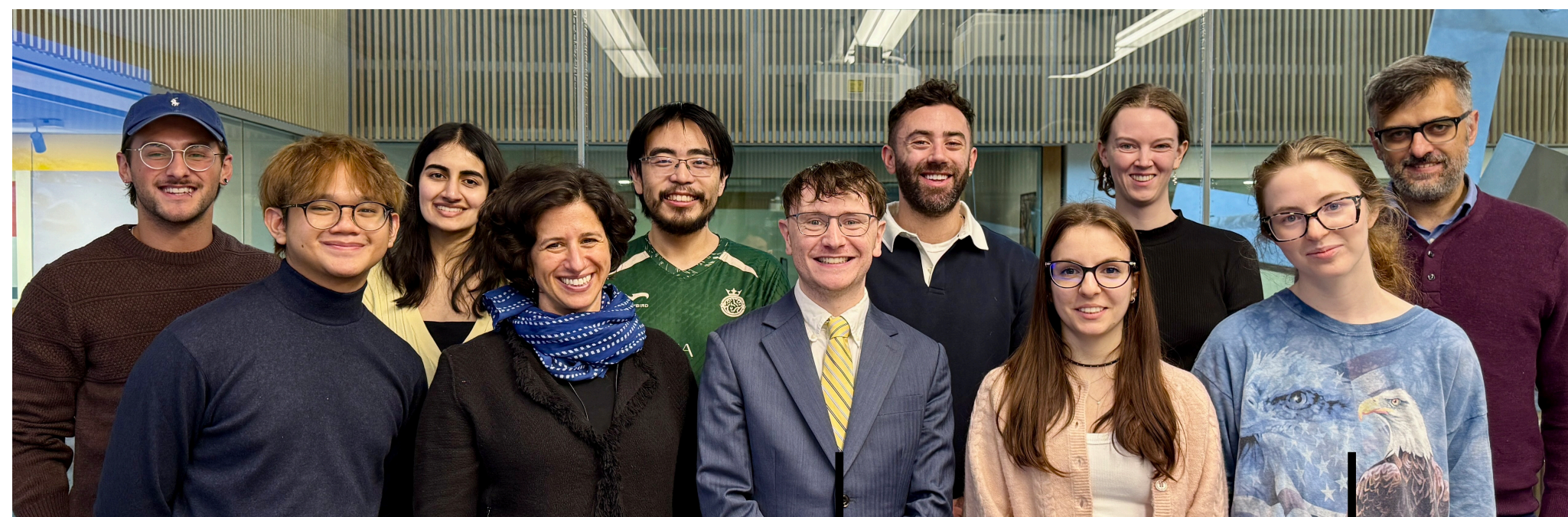
Hope Anderson

Joey Federico

Ryan Febryanto

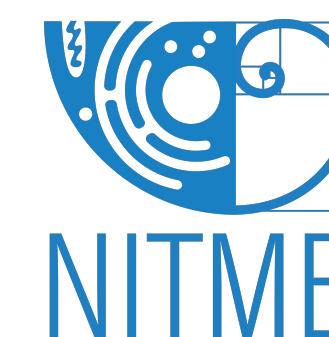
Joseph Sifakis, PhD

(Alum) Hanna Hieromnimon, PhD



Ryan

Sophia



Acknowledgments



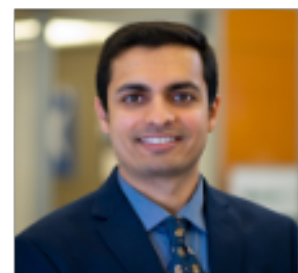
Frederic T. Chong, PhD

PI; Chief Scientist for Quantum Software;
UChicago Seymour Goodman Professor



Teague Tomesh, PhD

Co-PI; Manager of Quantum Software
Engineering



Pranav Gokhale, PhD

Chief Technology Officer



Peter Noell

Senior Quantum Solutions Manager



Colin Campbell

Quantum Applications Engineer



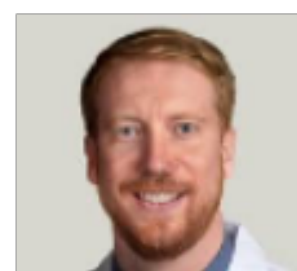
Victory Omole

Quantum Software Engineer



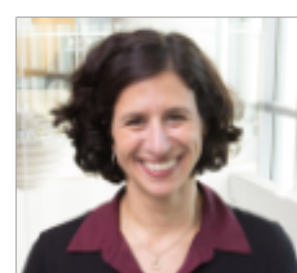
Bharath Thotakura

Quantum Software Engineer



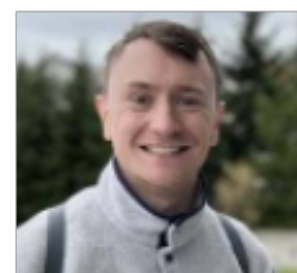
Alexander T. Pearson, MD, PhD

Co-PI; Practicing Medical Oncologist &
Statistician



Samantha J. Riesenfeld, PhD

Co-PI; Asst. Prof. of Molecular Engineering &
Medicine; Genomics-based ML Expert



Ryan Robinett, PhD

Postdoc, Information Theory & Manifold
Learning Expert



Sid Ramesh, MD

UChicago Pritzker School of Medicine



Sophia Madejski

Computational Experimentalist, Pritzker
School of Molecular Engineering



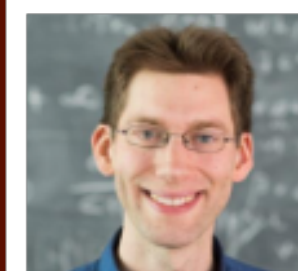
Zachary Morrell

PhD Candidate, Combinatorial
Optimization Expert



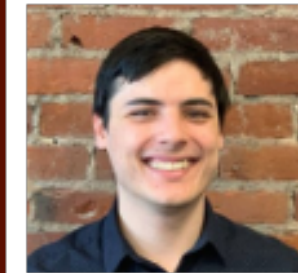
Willers Muye Yang

PhD Candidate, Quantum Computing
Researcher



Aram Harrow, PhD

Co-PI; Prof. of Physics; Quantum Information
and Computing Researcher



Eric Anshuetz, PhD

Quantum Algorithms Researcher



Jin-Peng Liu, PhD

Postdoc, Quantum Algorithms Researcher



Hanrui Wang

PhD Candidate, Quantum Computing
Architecture and Machine Learning



Mariesa Teo

PhD Candidate, Quantum Computing
Researcher



Dhirpal Shah

PhD Student, Quantum Computing
Researcher



Tina Oberol

PhD Candidate, Quantum Computing
Researcher

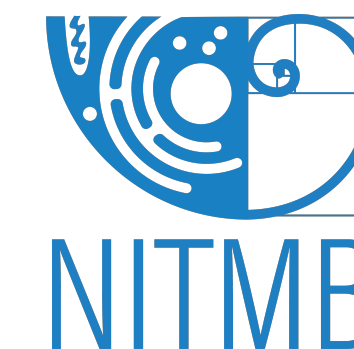
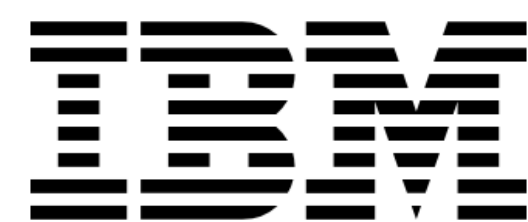


Ali-Javadi Abhari
Nate Earnest
Simon Martiel
Kevin Sung



Thank you

Questions?



THE UNIVERSITY OF CHICAGO



**PRITZKER SCHOOL OF
MOLECULAR ENGINEERING**

Additional Affiliations

Department of Medicine
Committee on Data Science
Committee on Immunology
Comprehensive Cancer Center

Institute for Biophysical Dynamics
Biohub, Chicago
NSF-Simons National Institute for
Theory and Mathematics in Biology