

And now for something completely different:

Quantum Support Vector Machines

*but this time without quantum databases, or
quantum-linear-systems-HHL quantum linear algebra tricks*

this one could work on near-term quantum computers

BASED ON :

Supervised learning with quantum enhanced feature spaces

Havlicek, Córcoles, Temme, Harrow, Kandala, Chow, Gambetta

Nature. vol. 567, pp. 209-212 (2019)

& a bunch of other literature...

1) *Background*

- *limitations of QCs*
- *Support Vector Machines, QSVM (1) and the kernel trick*
- *Variational (parametrized) quantum circuits*

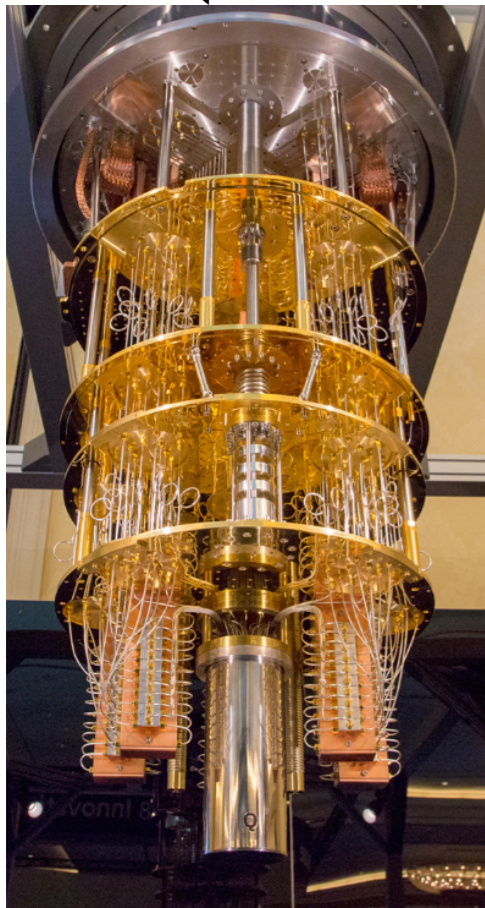
2) *Support Vector Machine with quantum kernels*

- 1) *version 1: “quantum-assisted”*
- 2) *version 2: “full quantum”*

3) *Some results*

1) Limitations of real-world QCs

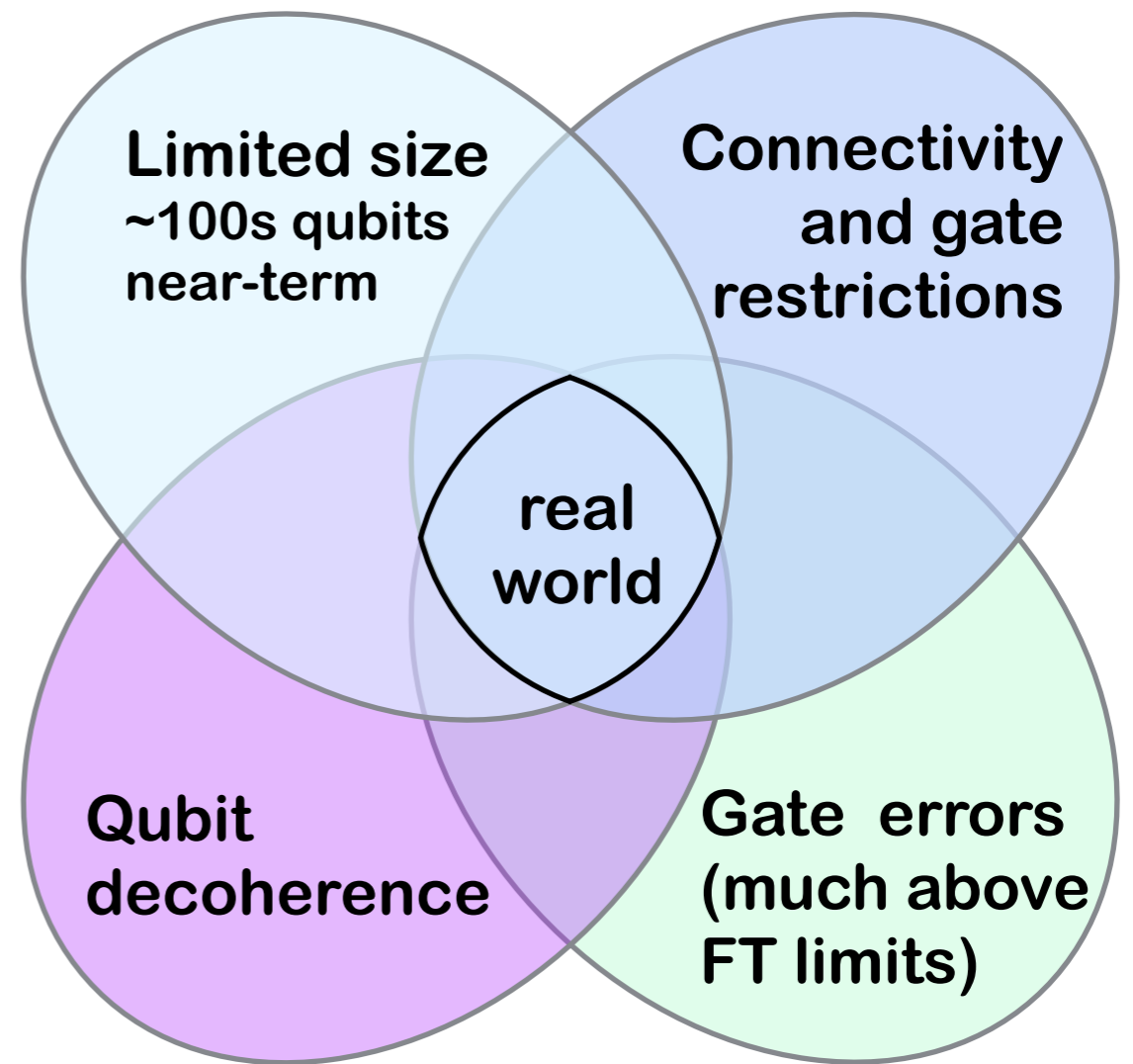
20 qubit "QC"



Banana for scale



Also one qubit is *much* more expensive than a banana...

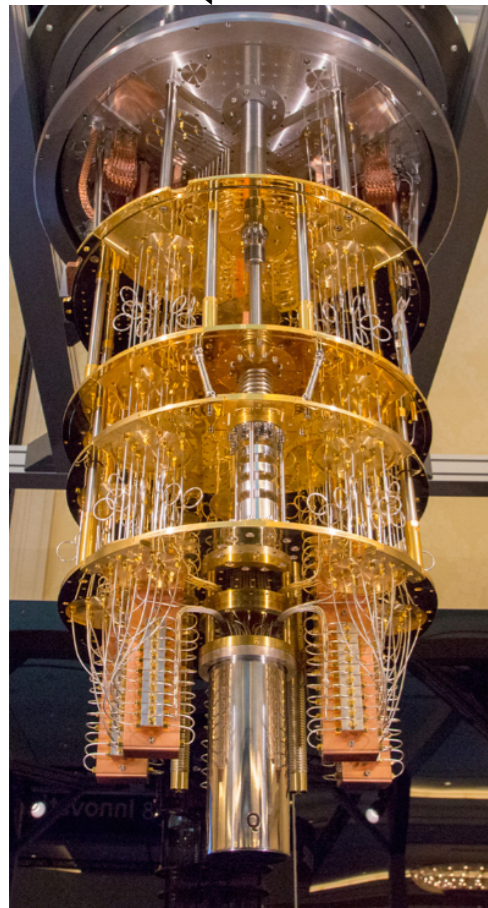


Decoherence: effects leading to degradation of qubit(s) state usually from “coupling to environment” and *relaxation*

- dephasing (environment “measures” qubit)
- de-polarization (gets noisy)
- relaxation (collapses to “ground state”)
- dissipation

1) Limitations of real-world QCs

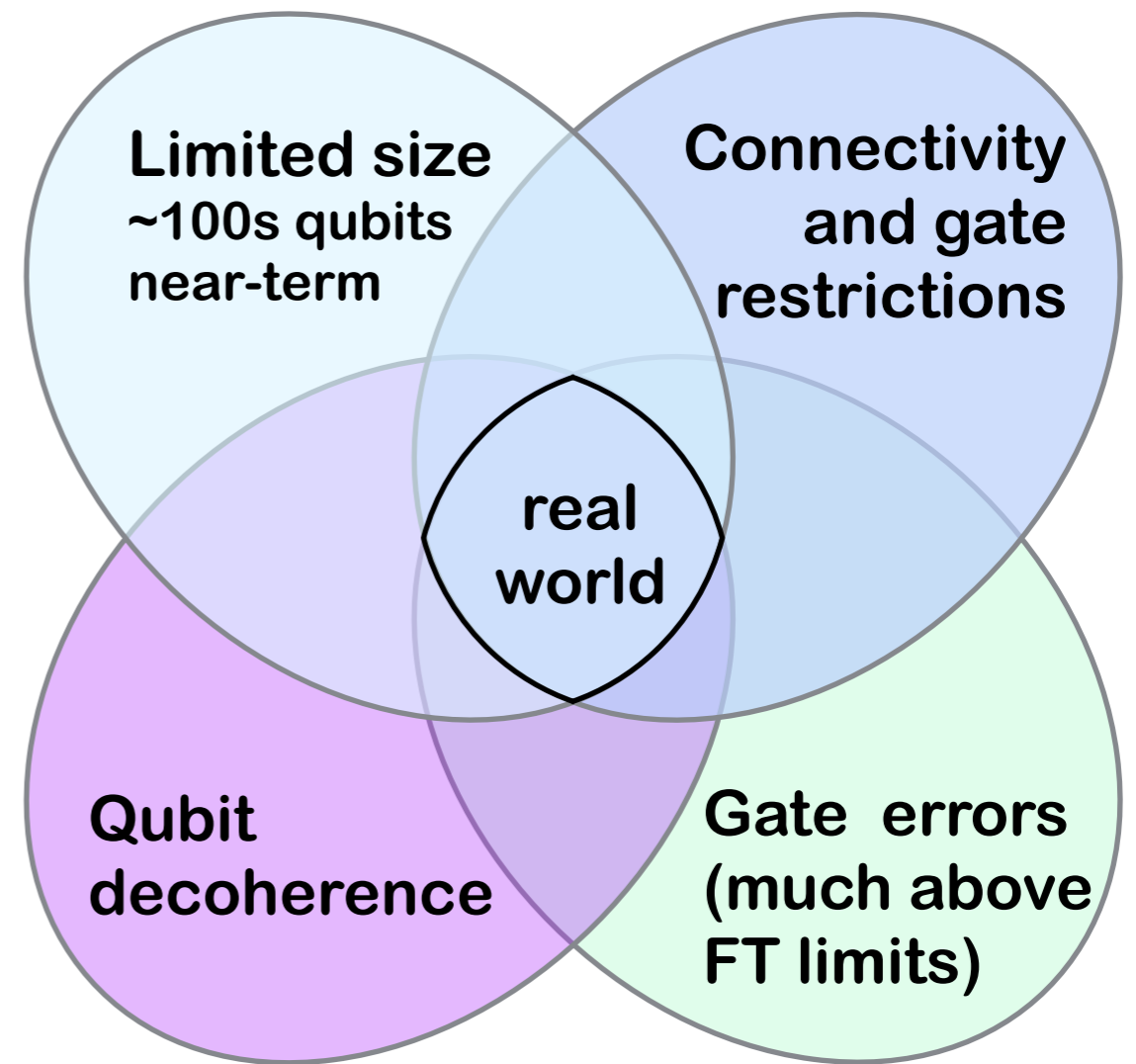
20 qubit "QC"



Banana for scale



Also one qubit is *much* more expensive than a banana...

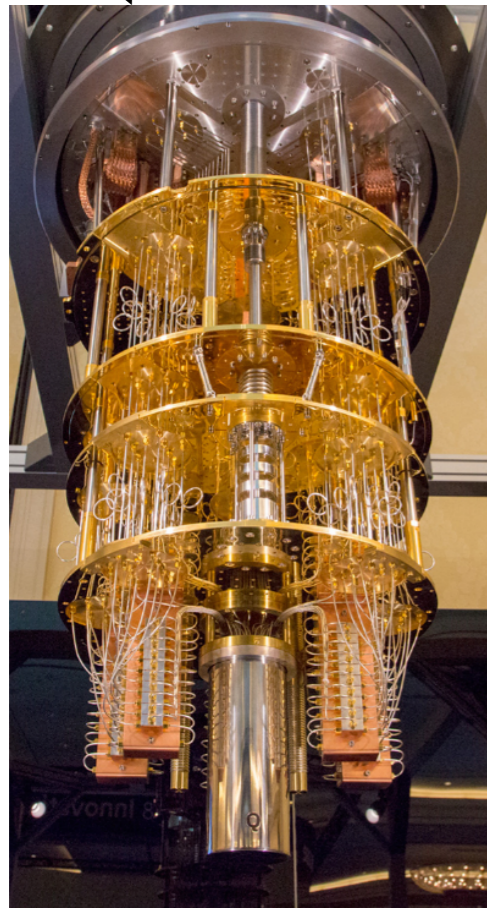


**In short: qubits have a life-time (half-life)....
up to milliseconds**

gates can take 10s-100s of nanoseconds

1) Limitations of real-world QCs

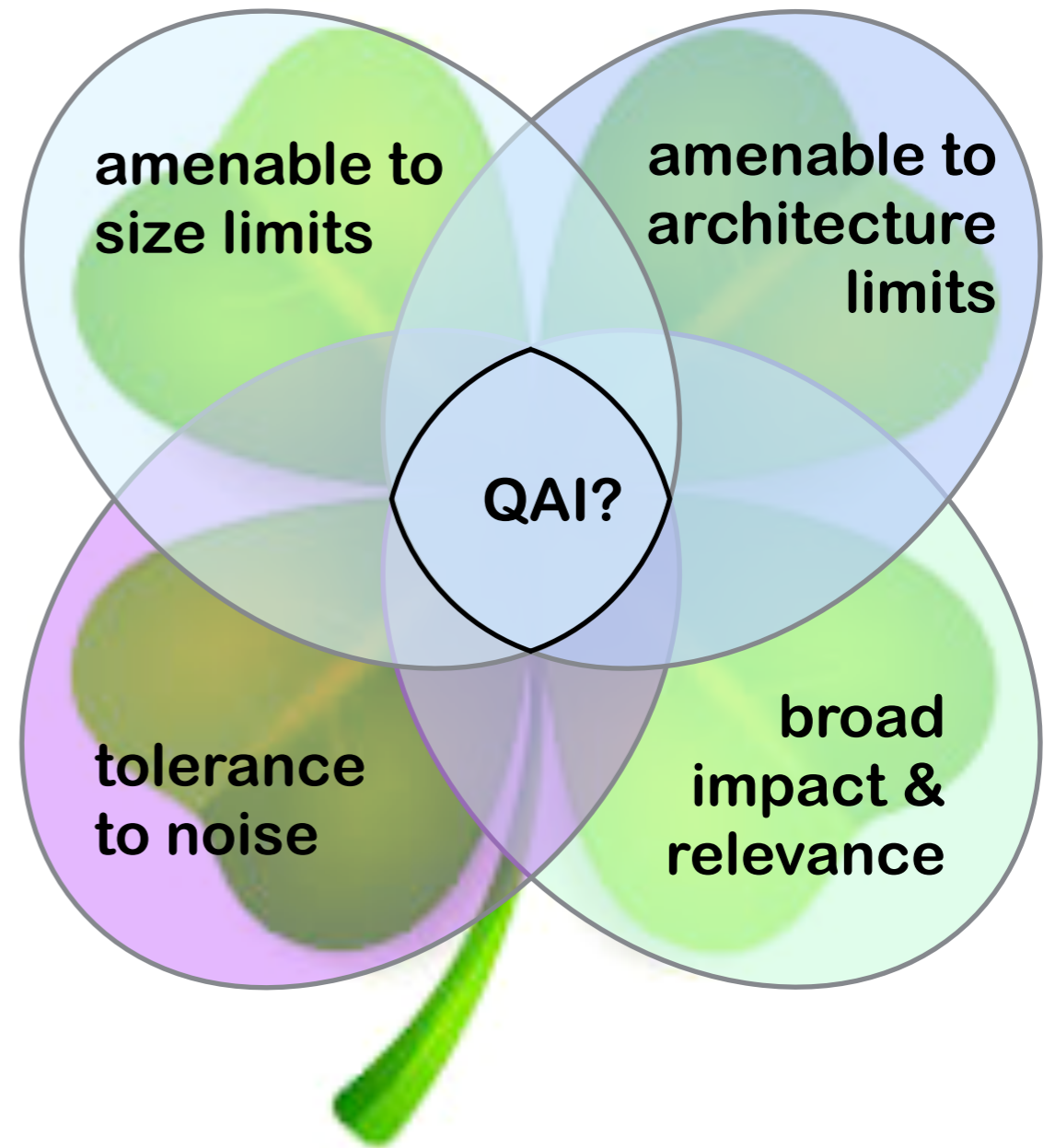
20 qubit "QC"



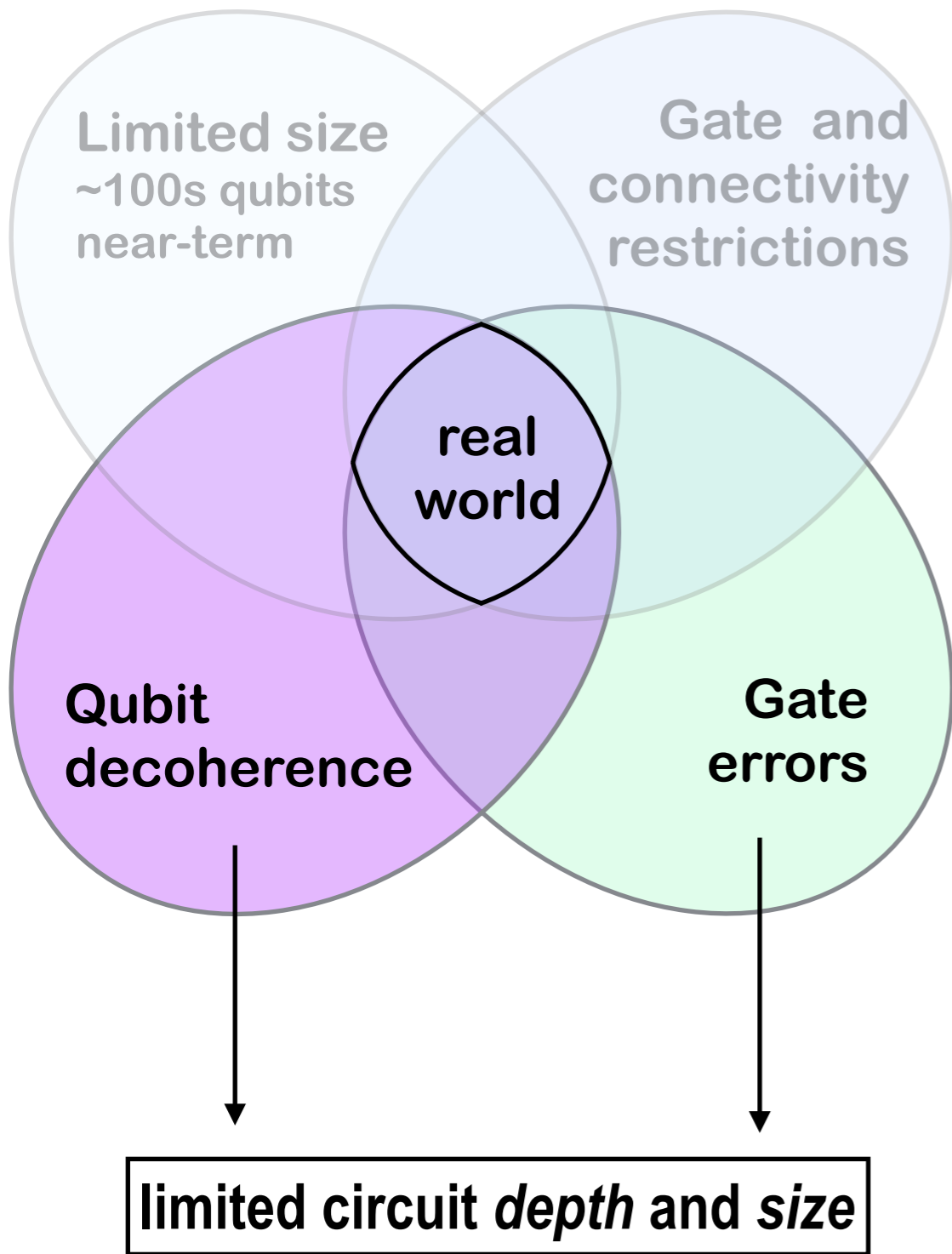
Banana for scale



Also one qubit is *much* more expensive than a banana...

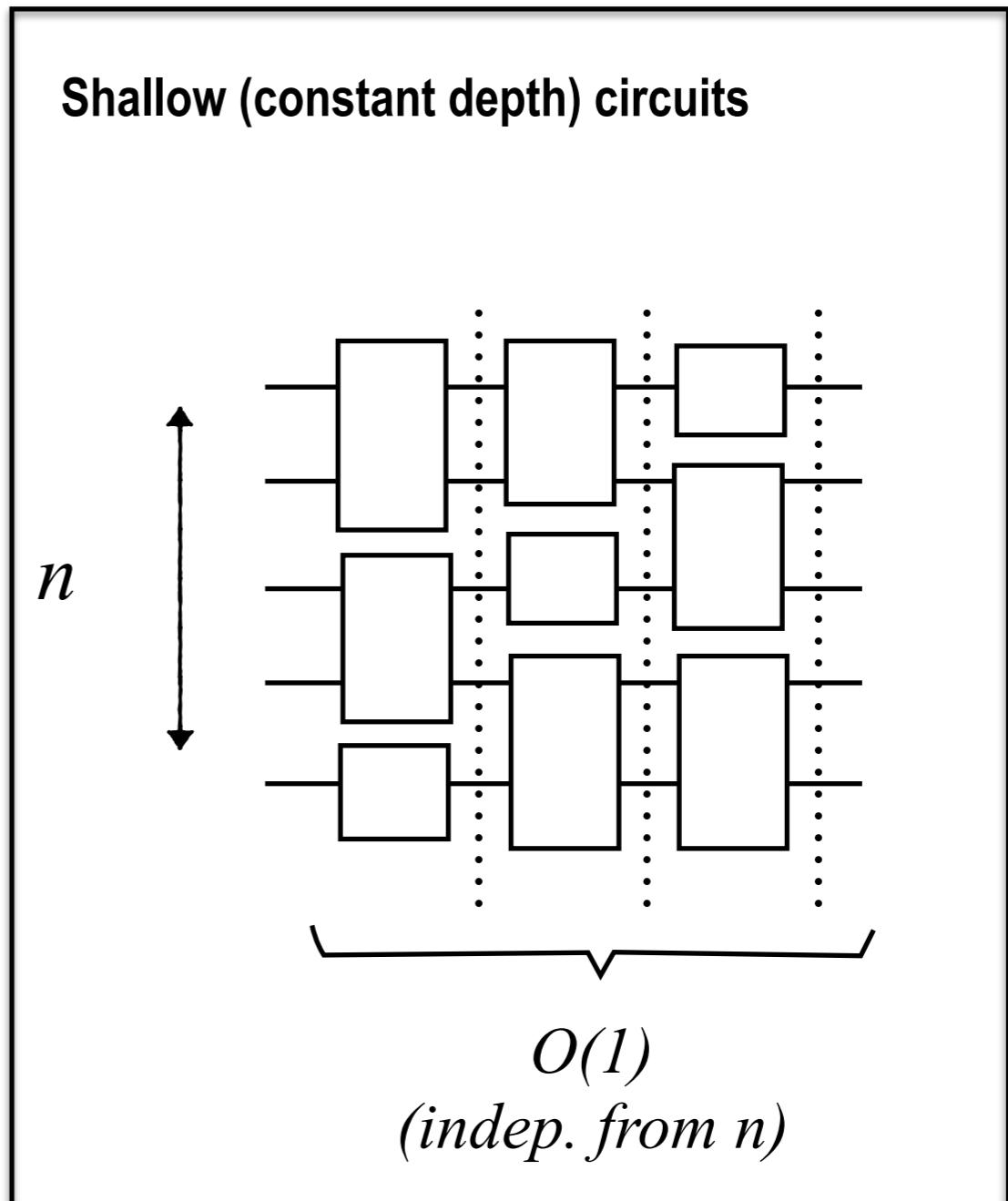


1) Limitations of real-world QCs

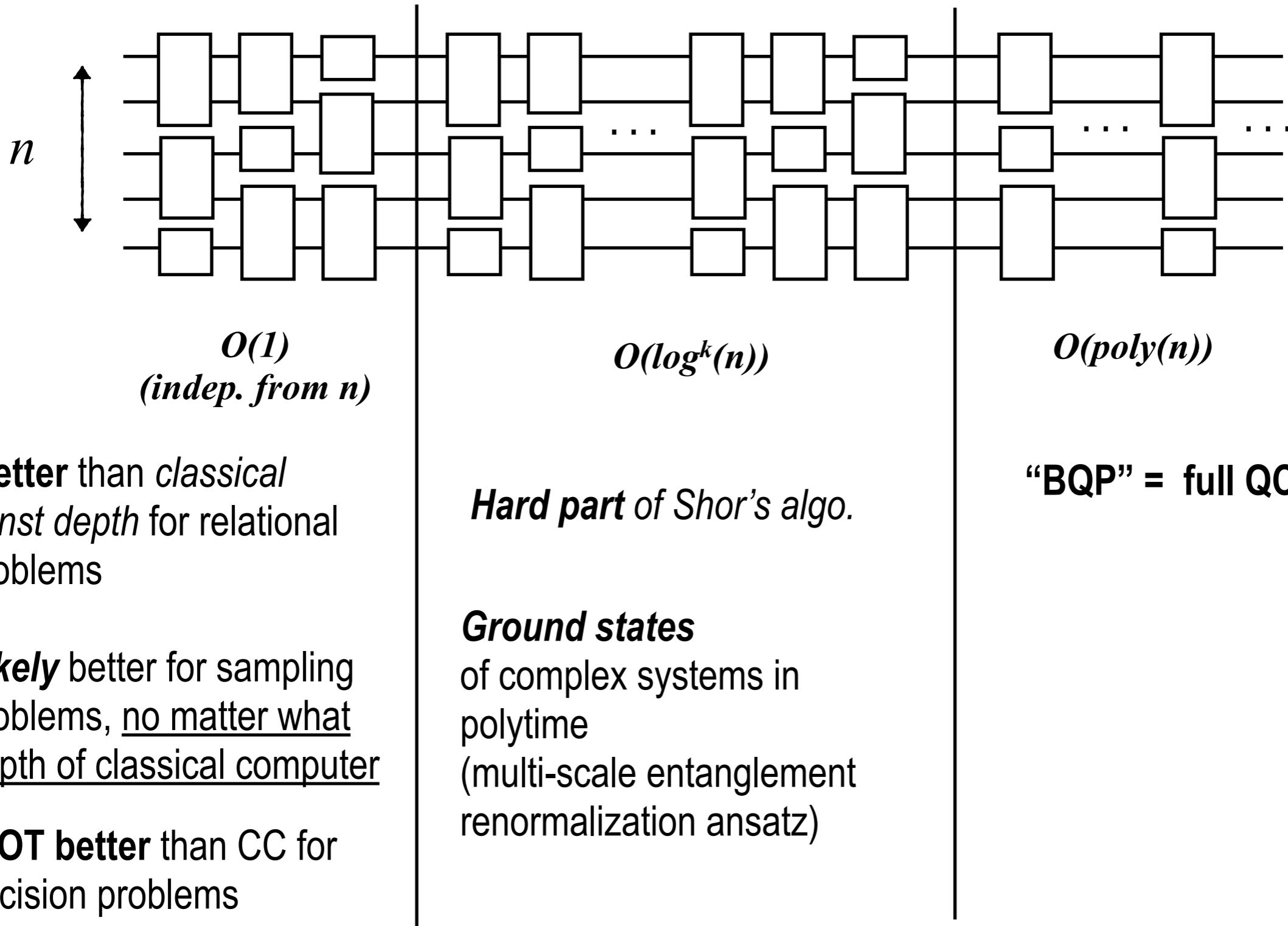


recall: qubits have a life-time (half-life).... up to *ms*
gates can be dozens of ns

Not yet in the same system... nowadays
whatever you can do in 10 - 100 (parallel) gate times



Tangent: Quantum depth complexity



-**better** than *classical*
const depth for relational
problems

-**likely** better for sampling
problems, no matter what
depth of classical computer

-**NOT better** than CC for
decision problems

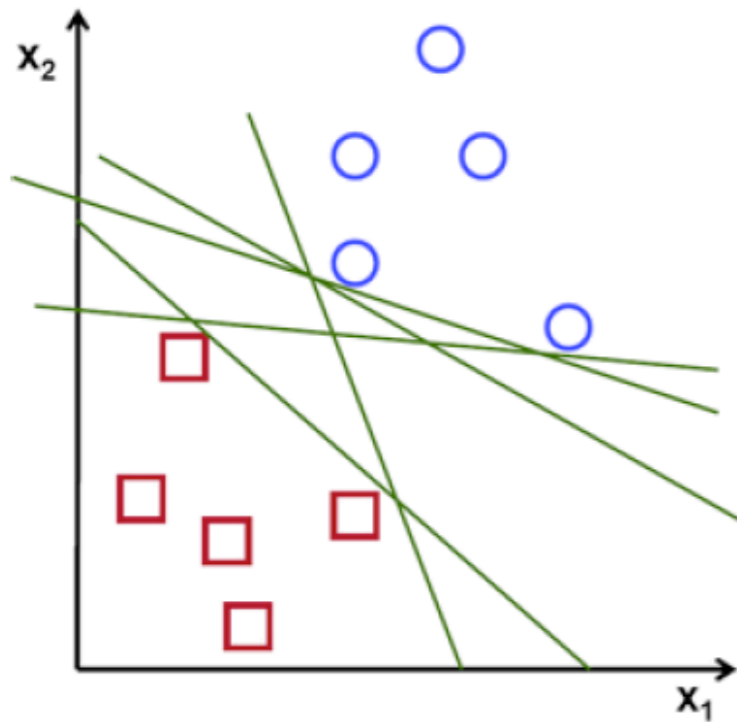
Hard part of Shor's algo.

Ground states
of complex systems in
polytime
(multi-scale entanglement
renormalization ansatz)

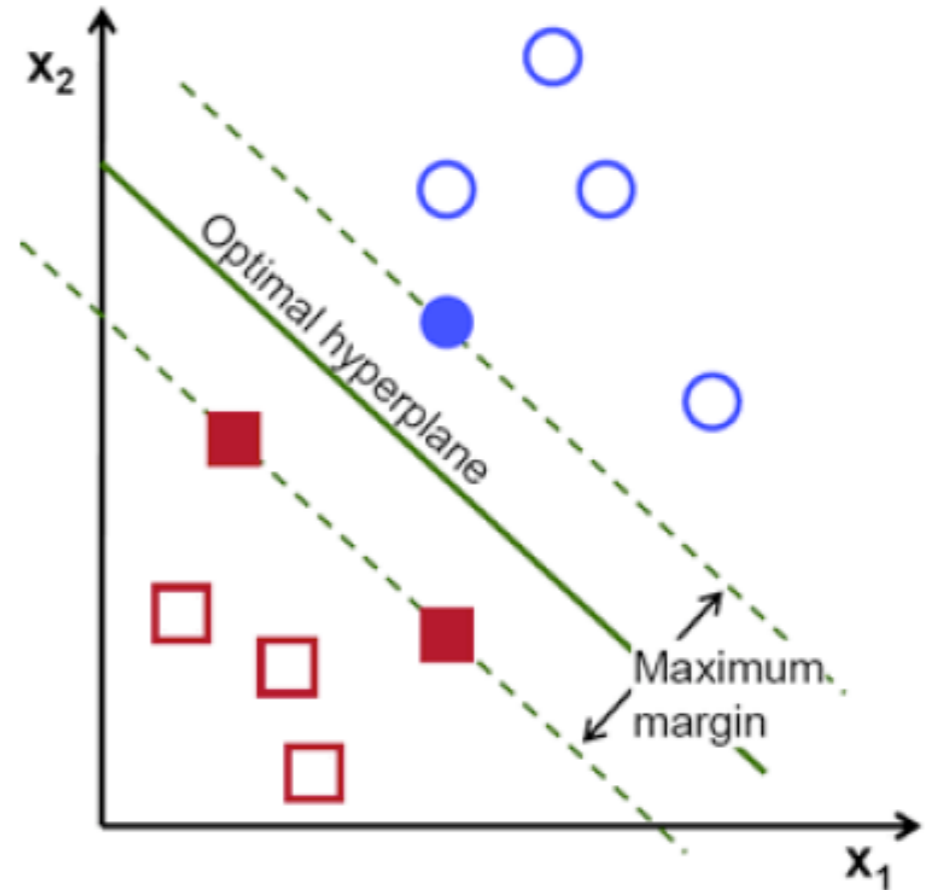
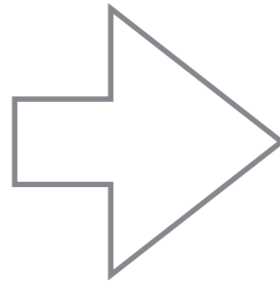
“BQP” = full QC

1.2) Support vector machines

$$D = \{(x_i, y_i)\}_i \quad x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$$



*separating hyperplanes
(linear classifier, not SVM)*



SVM: max-margin hyperplanes

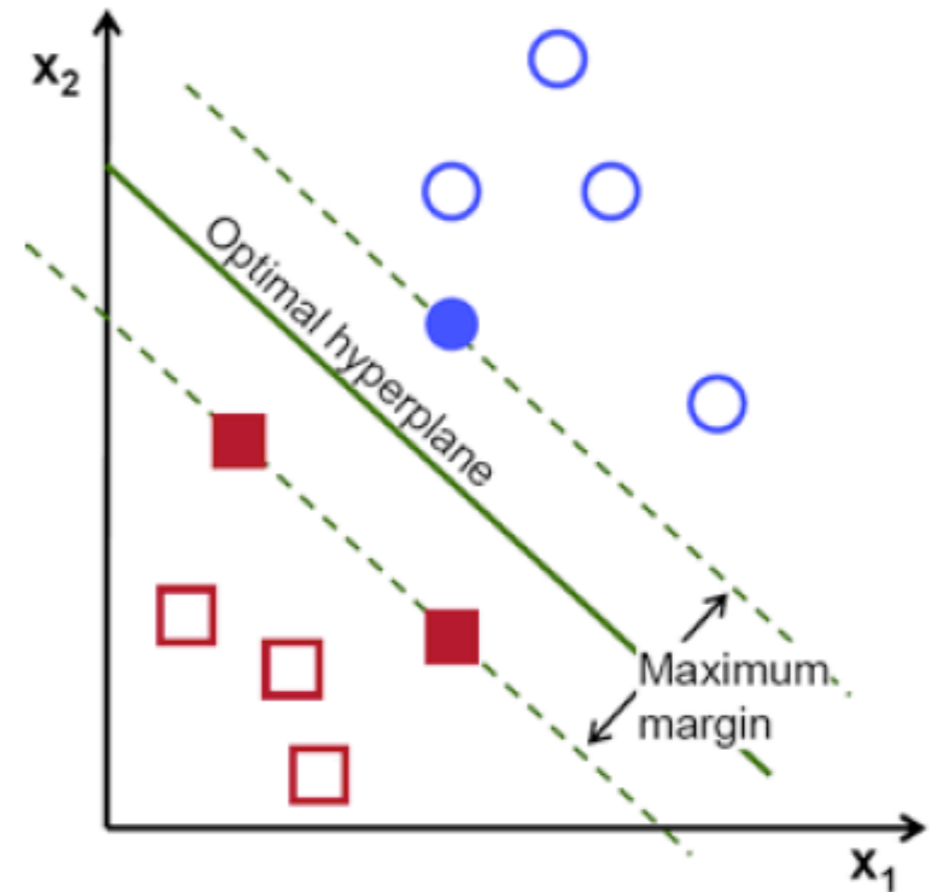
$$D = \{(x_i, y_i)\}_i \quad x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$$

Quadratic problem:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

such that $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, N.$

$$\arg \max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$



SVM: max-margin hyperplanes

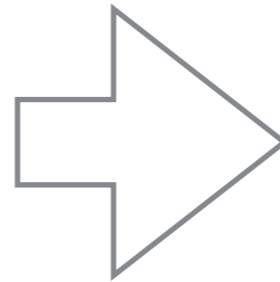
Note: defined on the basis of
“*support vectors*”

Primal problem:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

such that $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, N.$

$$\arg \max_{\mathbf{w}, b} \min_{i \in \{1, \dots, N\}} \frac{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}{\|\mathbf{w}\|}$$



Dual problem:

$$\arg \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i)^\top \mathbf{x}_j,$$

such that $\alpha_i \geq 0, \quad \text{for } i = 0, \dots, N,$

$$\text{and } \sum_{i=1}^N \alpha_i y_i = 0.$$

$$\mathbf{w} = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i.$$

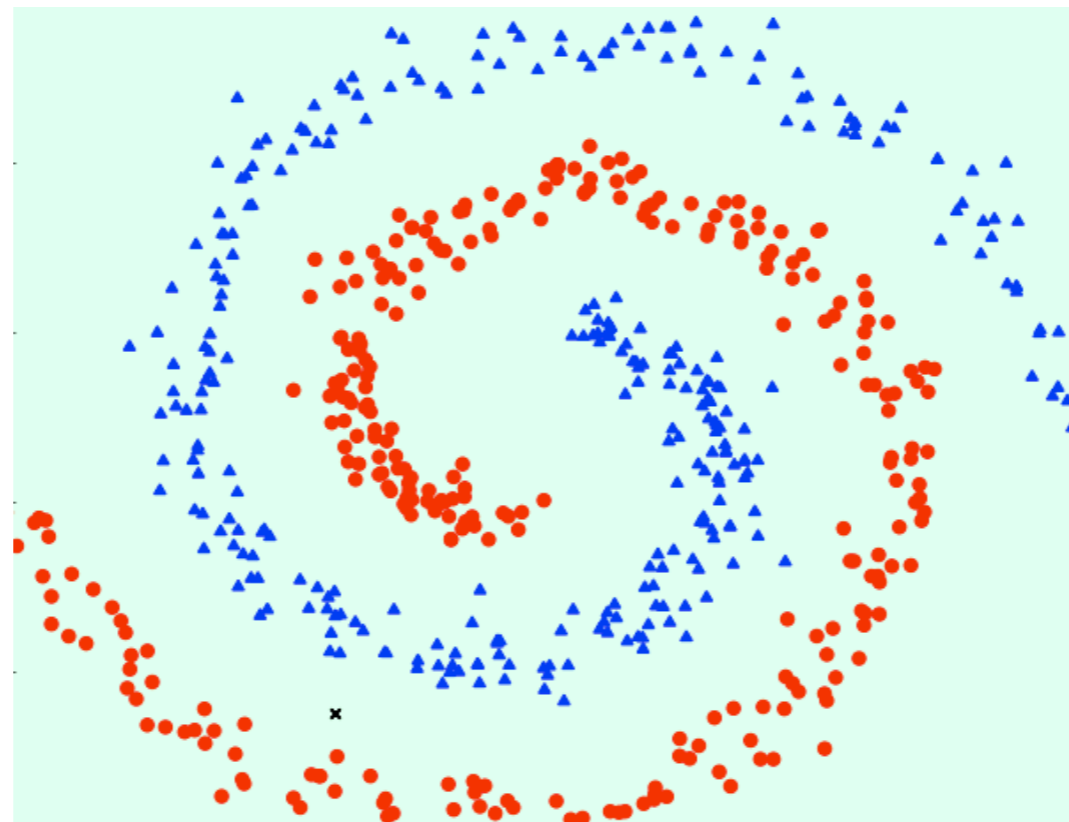
Why bother with dual problem? Representation in *terms of datapoints*

- sparser evaluation
- only inner products matter
- was handy for *quantum tricks*

$$(\mathbf{w}^*)^\top \mathbf{x} + b^* = \left(\sum_{i=1}^N \alpha_i y_i (\mathbf{x}_i)^\top \mathbf{x} \right) + b^*.$$

$$\alpha_i \alpha_j y_i y_j (\mathbf{x}_i)^\top \mathbf{x}_j,$$

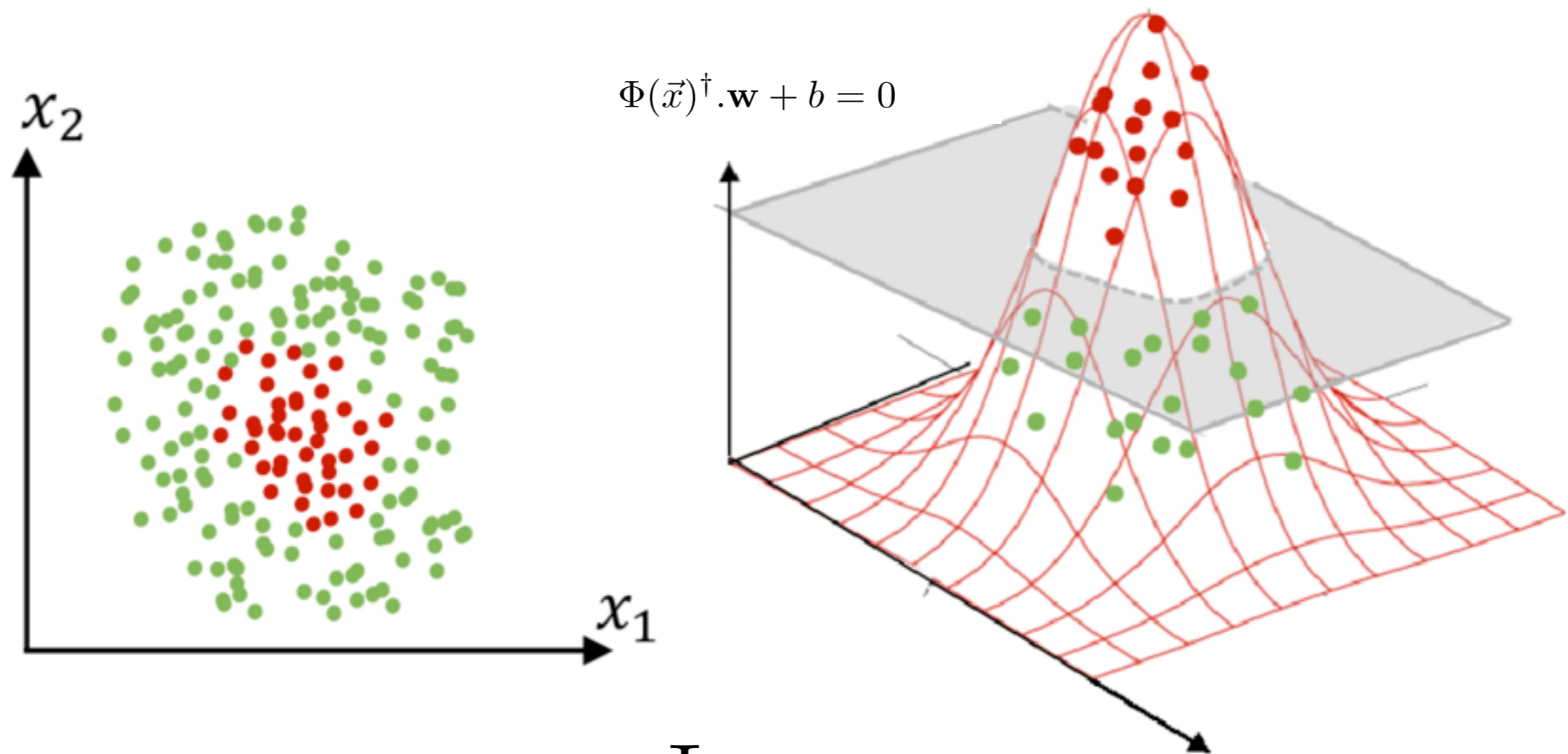
*Why one should actually bother with SVMs:
when data is NOT linearly separable*



Non-separable datasets?

-slack variables (this lead to QSVM - type 1)

-feature mapping and the kernel trick



$$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D \quad \vec{x} \xrightarrow{\Phi} \Phi(\vec{x})$$

c.f.: Cover's theorem...

The kernel trick:

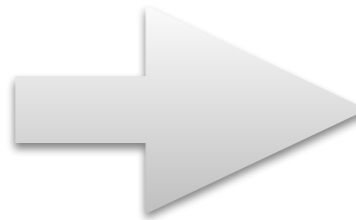
one can “train” and evaluate SVM classifiers in rich feature spaces without ever mapping data-points into said spaces. They can even be infinite dimensional

The kernel trick

Recall... only inner products matter:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (\phi = \Phi \dots)$$

$$\arg \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i)^{\top} \mathbf{x}_j,$$



$$\arg \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

kernels can sometimes be evaluated (much) more efficiently directly:

E.g. (stupidly)

$$(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \mapsto \phi(\mathbf{x}) = (x_1 x_1 \quad x_1 x_2 \quad x_1 x_3 \quad x_2 x_1 \quad x_2 x_2 \quad x_2 x_3 \quad x_3 x_1 \quad x_3 x_2 \quad x_3 x_3)^{\top}$$

$$\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \sum_{i=1}^d \sum_{j=1}^d x_i z_i x_j z_j \quad \text{Runtime for } \phi(\mathbf{x}): \mathcal{O}(d^2)$$

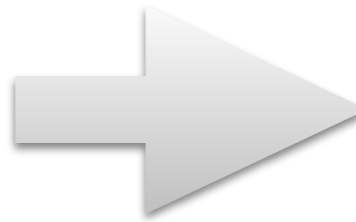
c.f. Mercer's theorem

The kernel trick

Recall... only inner products matter:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (\phi = \Phi \dots)$$

$$\arg \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i)^{\top} \mathbf{x}_j,$$



$$\arg \max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

$$\phi(\mathbf{x}) = (x_1 x_1 \quad x_1 x_2 \quad x_1 x_3 \quad x_2 x_1 \quad x_2 x_2 \quad x_2 x_3 \quad x_3 x_1 \quad x_3 x_2 \quad x_3 x_3)^{\top}$$

reverse-engineered:
$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^{\top} \mathbf{z})^2 = \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{i=1}^d x_i z_i \right) = \sum_{i=1}^d \sum_{j=1}^d x_i z_i x_j z_j = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle.$$

Directly:

Let $\mathbf{x} = (x_1, \dots, x_d)^{\top}$, $\mathbf{z} = (z_1, \dots, z_d)^{\top}$ and

$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^{\top} \mathbf{z})^2.$$

Runtime: $\mathcal{O}(d)$.

Yay, quadratic speedup

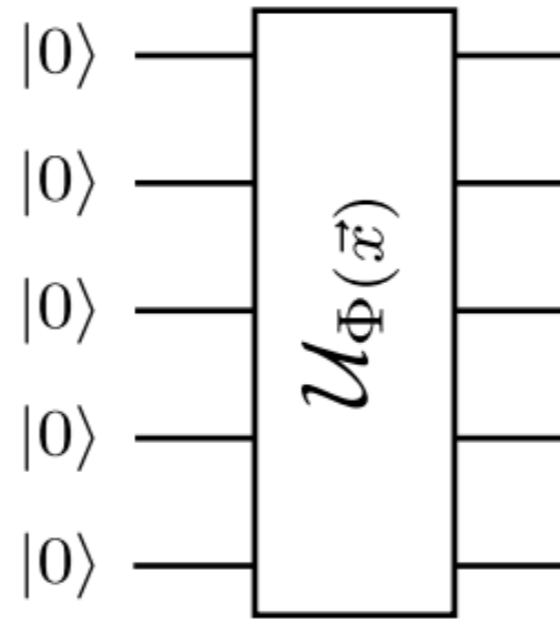
c.f. Mercer's theorem

The kernel trick:

one can “train” and evaluate SVM classifiers in rich feature spaces without ever mapping data-points into said spaces. They can even be infinite dimensional

Feature maps matter,
and sometimes kernels are *not* efficiently computable...

Feature maps matter,
and sometimes kernels are *not* efficiently computable...



$$\vec{x} \mapsto \mathcal{U}_{\Phi}(\vec{x})|0\rangle = |\Phi(\vec{x})\rangle$$

$$U_{\Phi(\vec{x})} = \exp \left(i \sum_{S \subseteq [n]} \phi_S(\vec{x}) \prod_{i \in S} Z_i \right)$$

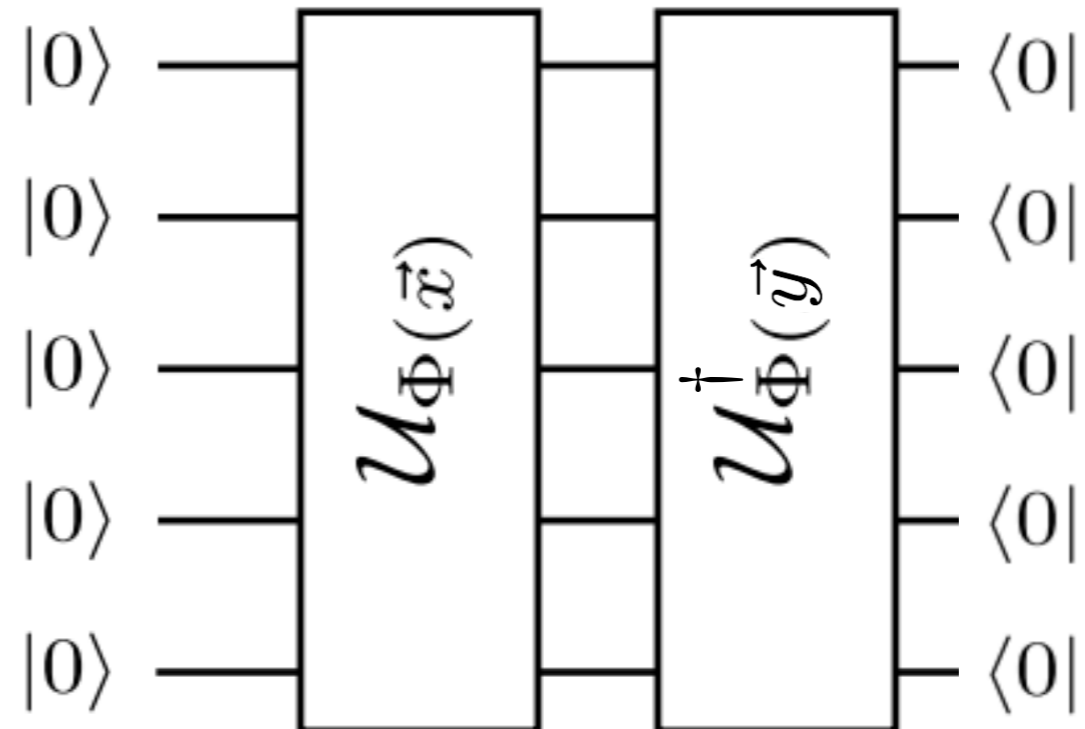
$$\phi_{\{i\}}(\vec{x}) = x_i \text{ and } \phi_{\{1,2\}}(\vec{x}) = (\pi - x_1)(\pi - x_2)$$

Feature maps matter,
and sometimes kernels are *not* efficiently computable...

Kernel!

$$|\langle \Phi(\vec{y}) | \Phi(\vec{x}) \rangle|^2$$

Can be hard to compute.



Do this quantumly
(recall QC is good for inner products)

$$U_{\Phi(\vec{x})} = \exp \left(i \sum_{S \subseteq [n]} \phi_S(\vec{x}) \prod_{i \in S} Z_i \right)$$

$$\phi_{\{i\}}(\vec{x}) = x_i \text{ and } \phi_{\{1,2\}}(\vec{x}) = (\pi - x_1)(\pi - x_2)$$

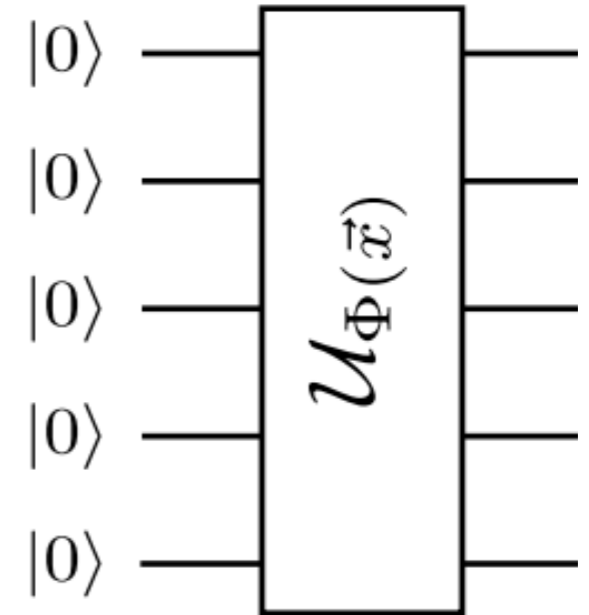
Feature maps matter,
and sometimes kernels are *not* efficiently computable...

$$U_{\Phi(\vec{x})} = \exp \left(i \sum_{S \subseteq [n]} \phi_S(\vec{x}) \prod_{i \in S} Z_i \right)$$

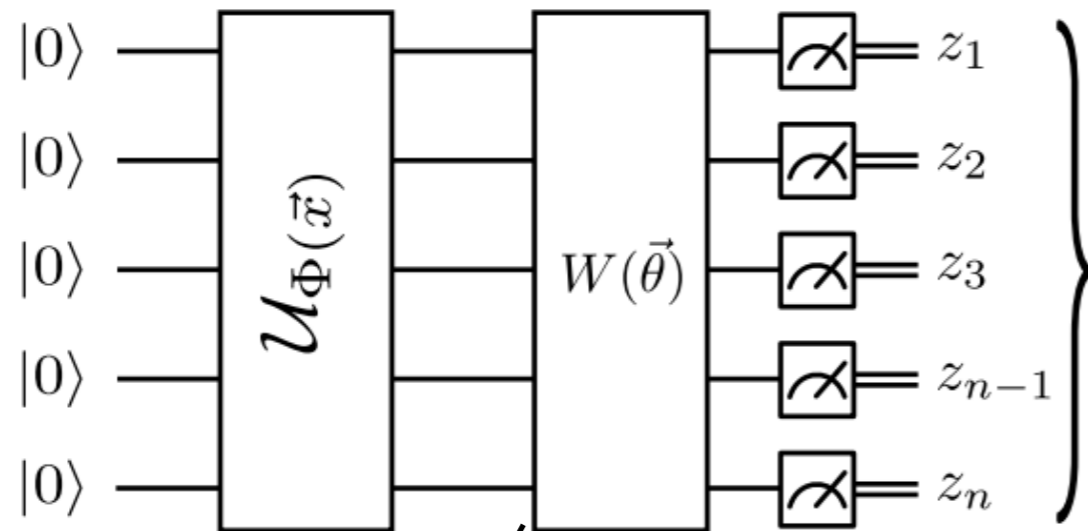
$$\phi_{\{i\}}(\vec{x}) = x_i \text{ and } \phi_{\{1,2\}}(\vec{x}) = (\pi - x_1)(\pi - x_2)$$

$$e^{i\phi_{\{l,m\}}(\vec{x})} Z_l Z_m = \begin{array}{c} \bullet \text{---} \text{---} \bullet \\ | \quad \quad | \\ \oplus \text{---} \boxed{Z_\phi} \text{---} \oplus \end{array}$$

$$\mathcal{U}_\Phi = H^{\otimes n} U_\Phi H^{\otimes n} U_\Phi \cdots H^{\otimes n} U_\Phi$$

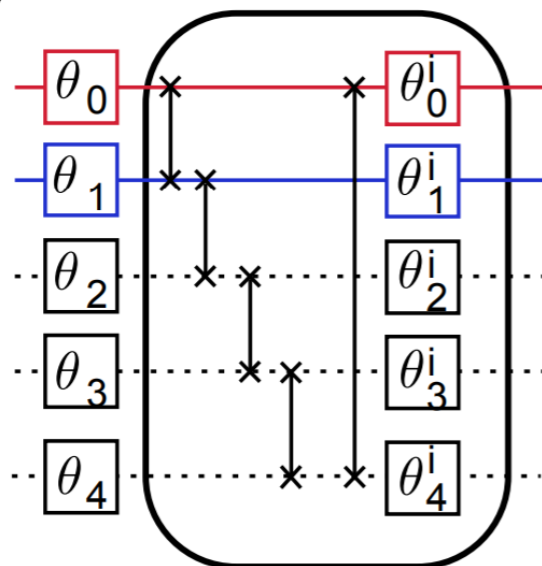


But there is also the fully quantum version:



$f(z)$

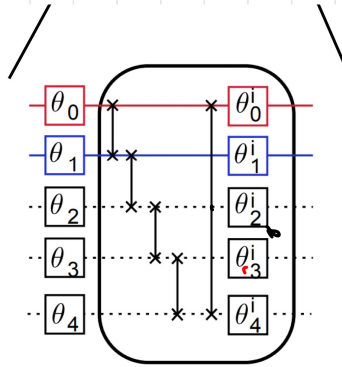
$$f(z) : \{0, 1\}^n \rightarrow \{-1, 1\}$$



repeat l - times

$$U(\theta_{m,t}) = e^{i\frac{1}{2}\theta_{m,t}^z Z_m} e^{i\frac{1}{2}\theta_{m,t}^y Y_m}$$

By the way... CIRCUITS OF THIS TYPE



WHICH DEVIATE FROM THE DISCRETE GATESET $\{H, \pi/8, CNOT\}$

BUT UTILIZE (A NUMBER OF) CONTINUOUS PARAMETER ELEMENTS ARE CALLED

PARAMETRIZED OR VARIATIONAL CIRCUITS

THEY ARE EXPERIMENTALLY WELL-MOTIVATED

$$\exp(i H \Delta t) \rightarrow U(\theta) = \exp(i H \theta) \dots$$

How does it output a label?

$$\text{label}(\vec{y}) = \tilde{m}(\vec{x}) = \text{sign}(\underbrace{\langle \Phi(\vec{x}) | W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) | \Phi(\vec{x}) \rangle}_{\text{involves running circuit many times}} + b)$$

involves running circuit many times

How does it output a label?

$$\text{label}(\vec{y}) = \tilde{m}(\vec{x}) = \text{sign}(\langle \Phi(\vec{x}) | W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) | \Phi(\vec{x}) \rangle + b)$$

How does it learn?

Optimize θ to minimize some loss/error/empirical risk on dataset

Involves evaluation of label function many times...

How does it output a label?

$$\text{label}(\vec{y}) = \tilde{m}(\vec{x}) = \text{sign}(\langle \Phi(\vec{x}) | W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) | \Phi(\vec{x}) \rangle + b)$$

How does it learn?

Optimize θ to minimize some loss/error/empirical risk on dataset

What does it *do*?

$$w_\alpha(\vec{\theta}) = \text{tr} \left[W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) P_\alpha \right]$$

$$\Phi_\alpha(\vec{x}) = \langle \Phi(\vec{x}) | P_\alpha | \Phi(\vec{x}) \rangle$$

$$\tilde{m}(x) = \text{sign} \left(2^{-n} \sum_\alpha w_\alpha(\vec{\theta}) \Phi_\alpha(\vec{x}) + b \right)$$

How does it output a label?

$$\text{label}(\vec{y}) = \tilde{m}(\vec{x}) = \text{sign}(\langle \Phi(\vec{x}) | W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) | \Phi(\vec{x}) \rangle + b)$$

How does it learn?

Optimize θ to minimize some loss/error/empirical risk on dataset

What does it *do*?

$$w_\alpha(\vec{\theta}) = \text{tr} [W^\dagger(\vec{\theta}) \mathbf{f} W(\vec{\theta}) P_\alpha]$$

$$\Phi_\alpha(\vec{x}) = \langle \Phi(\vec{x}) | P_\alpha | \Phi(\vec{x}) \rangle$$

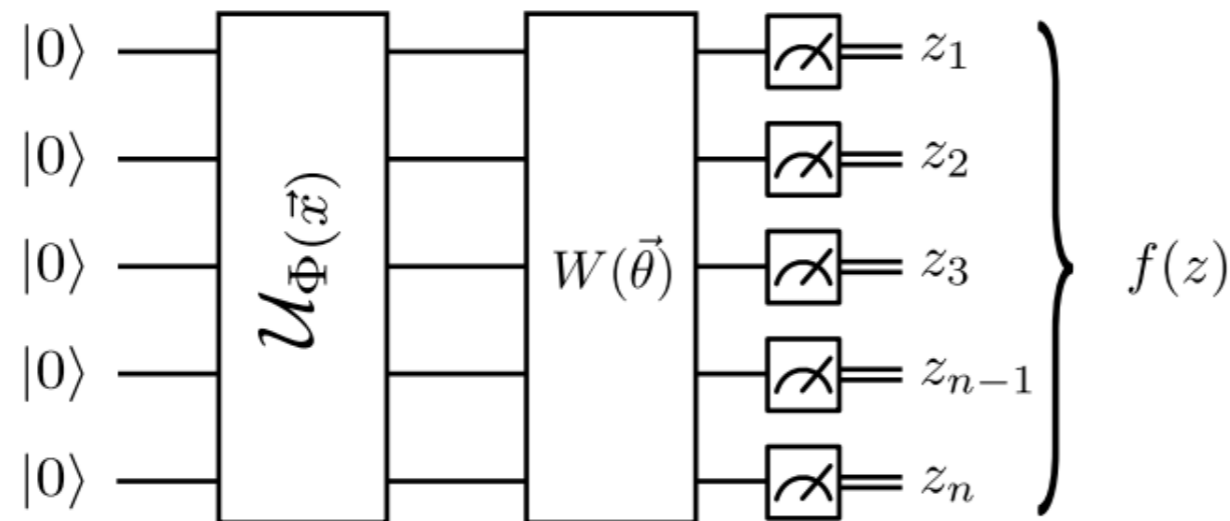
$$\tilde{m}(x) = \text{sign} \left(2^{-n} \sum_\alpha w_\alpha(\vec{\theta}) \Phi_\alpha(\vec{x}) + b \right)$$

-limitations on the model
come into play here...
-not *all hyperplanes*
reachable...

-not maximal margin
attained!

**The group with this project
will clarify this in report.**

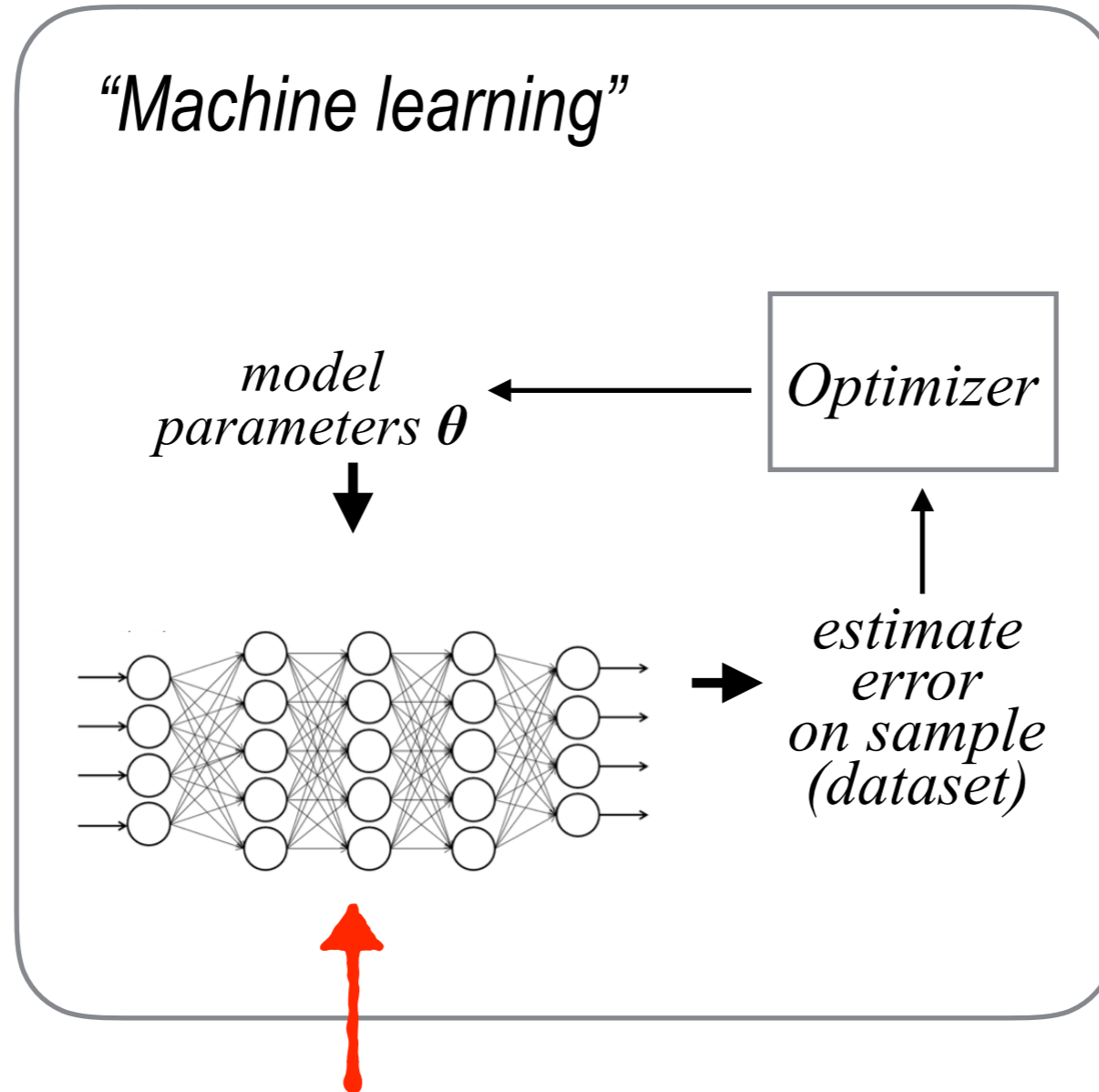
Summary:



- estimate probability of -1/1
- take estimate of expected value
- use this to label
- loop optimizing θ on dataset

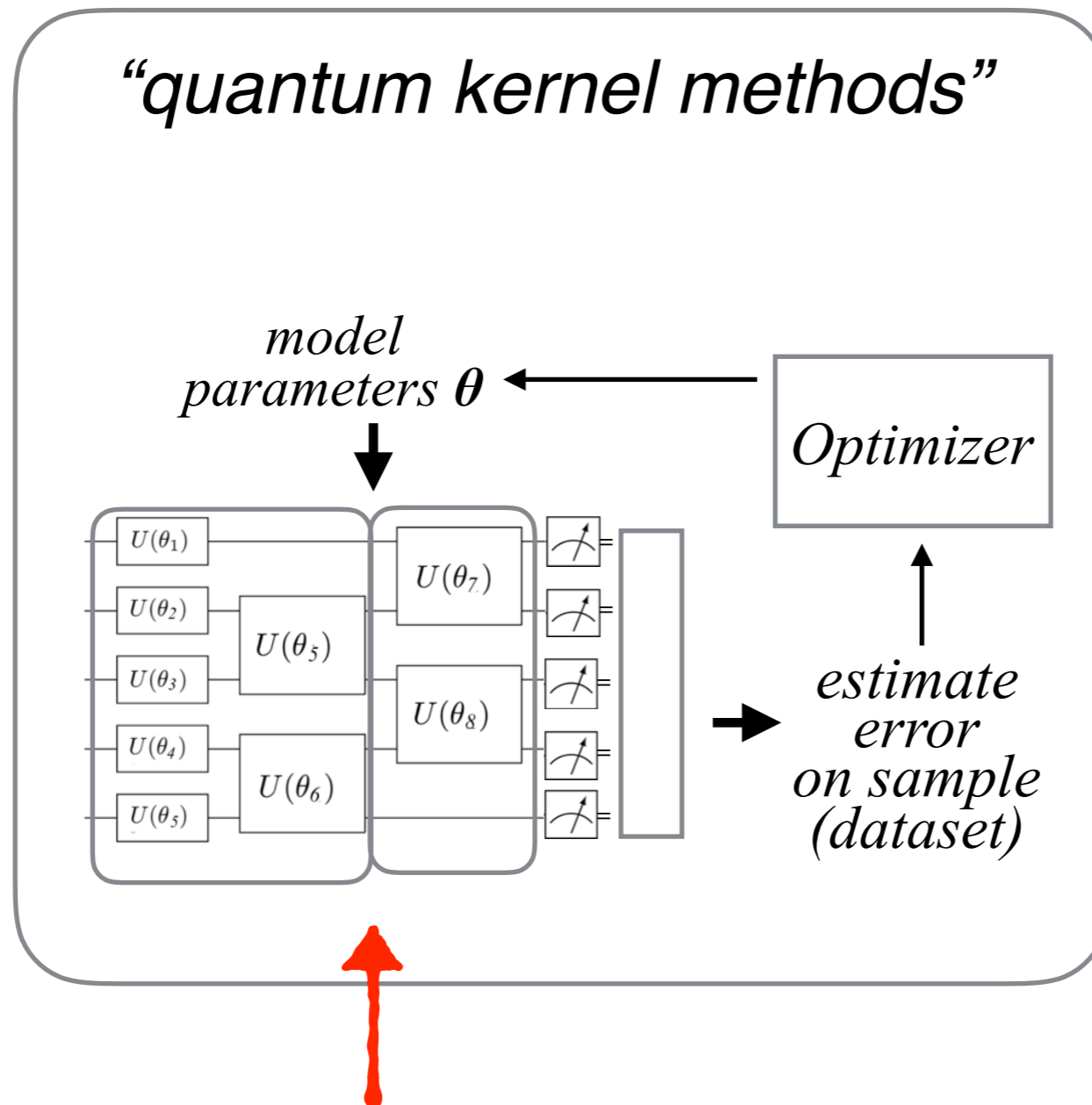
These are basic ideas, with (some) steps omitted.
The group with this project will report this precisely.

Note this is much like training NNs or other general models



*family of functions.
if it's "good", we can generalize well*

But you train a “quantum” network, without backprop, ofc.

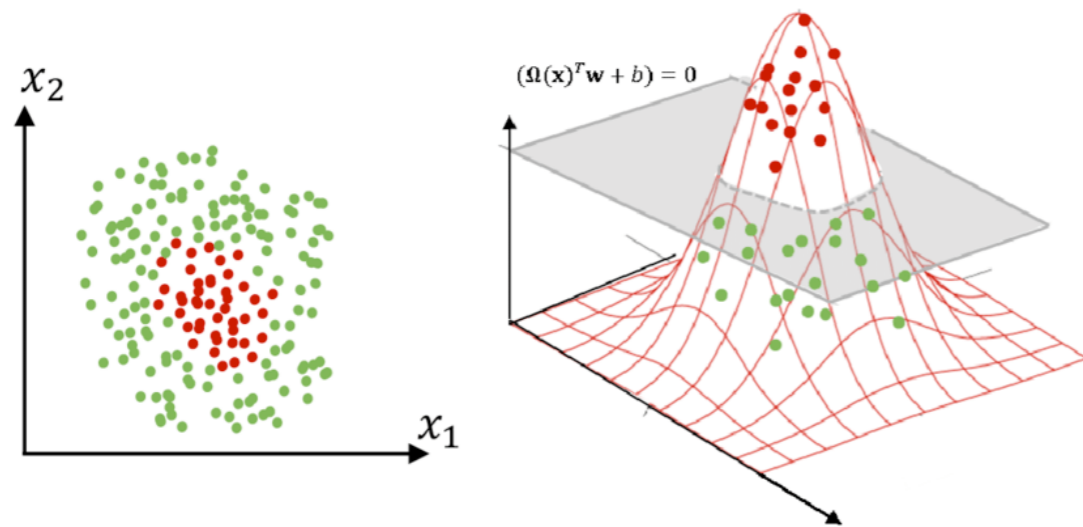


How about “shallow quantum circuits”?

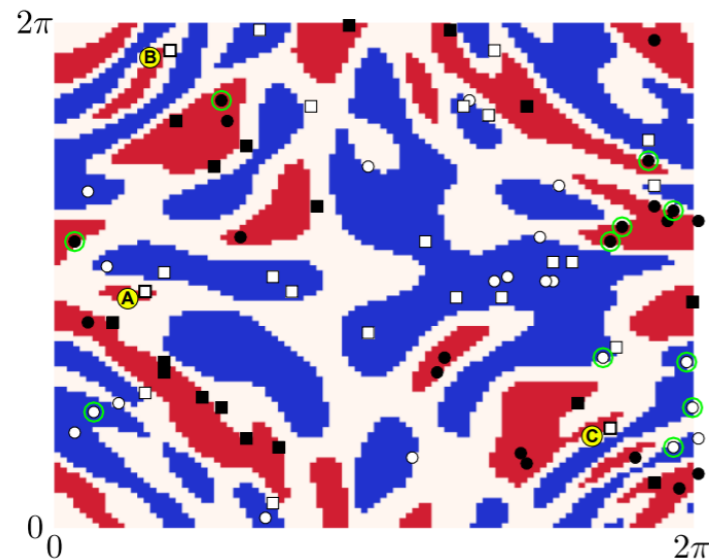
-instead neural network, train a QC!

-related to ideas from q. condensed-matter physics (VQE)

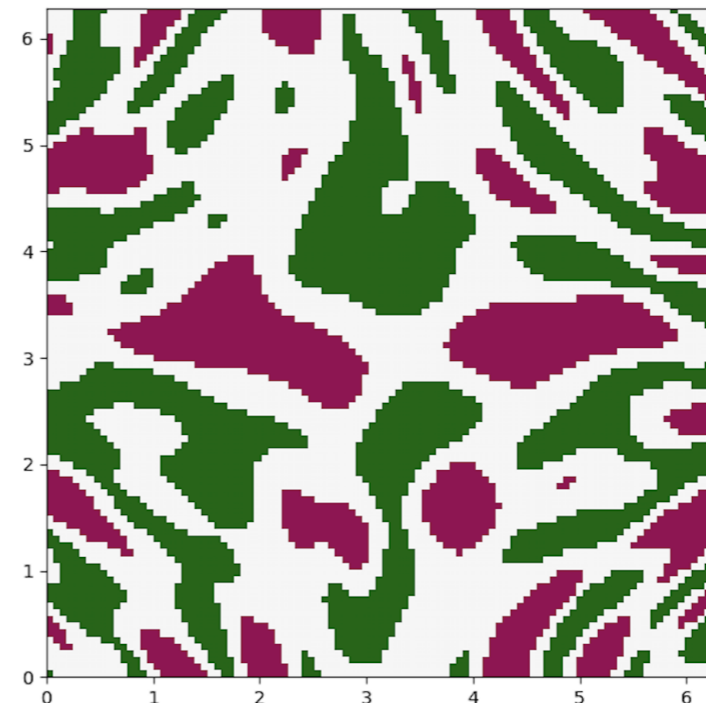
BUT it can be interpreted as SVM
So what does it do?



Two slices of quantum kernels:



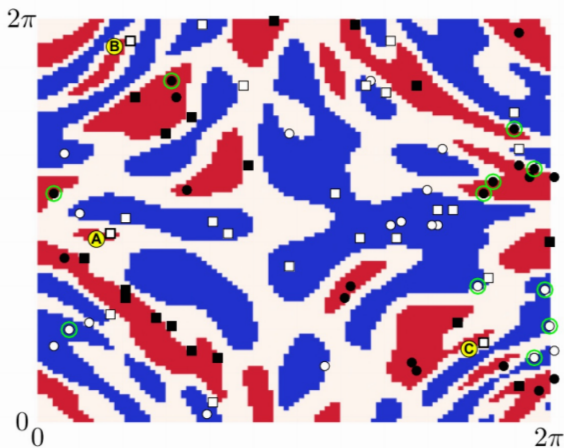
ORIGINAL PAPER



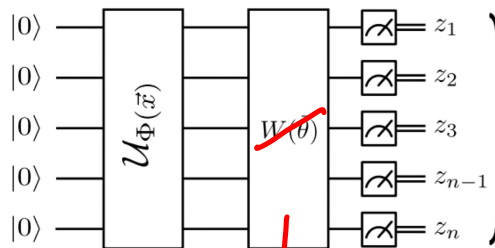
MSc THESIS MARDIROSIAN (LIACS)

- PERFORMANCE OF SUCH Q-KERNELS STUDIED IN A NUMBER OF WORKS

- IN ORIGINAL PAPER: GENERALIZATION PERFORMANCE ON ARTIFICIAL DATASETS



100% correct classification...



$f(z)$ →

Point labels
generated in
the same way
as classification

RANDOMLY CHOSEN

UNITARY (BUT IMPLEMENTABLE!)

PROBABLY = CHOOSE θ RANDOMLY ..

Cost functions and optimization?

Noise tolerance!

Advantages?

Two models, back-to-back?

That's the question...

Supervised learning with quantum enhanced feature spaces

Havlicek, Córcoles, Temme, Harrow, Kandala, Chow, Gambetta

Nature. vol. 567, pp. 209-212 (2019)