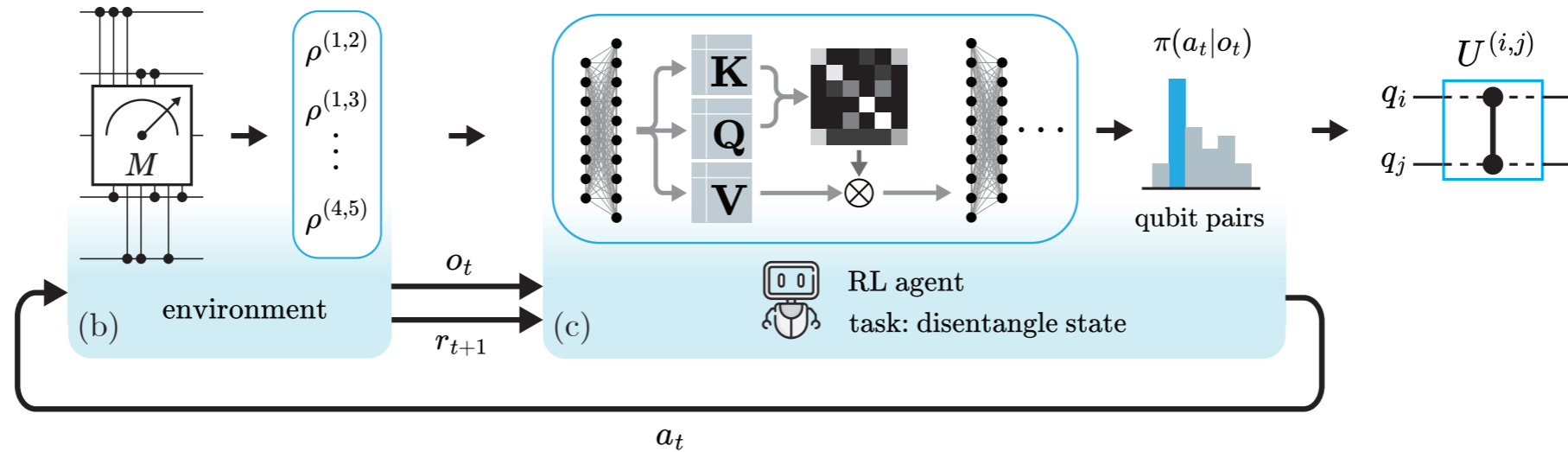
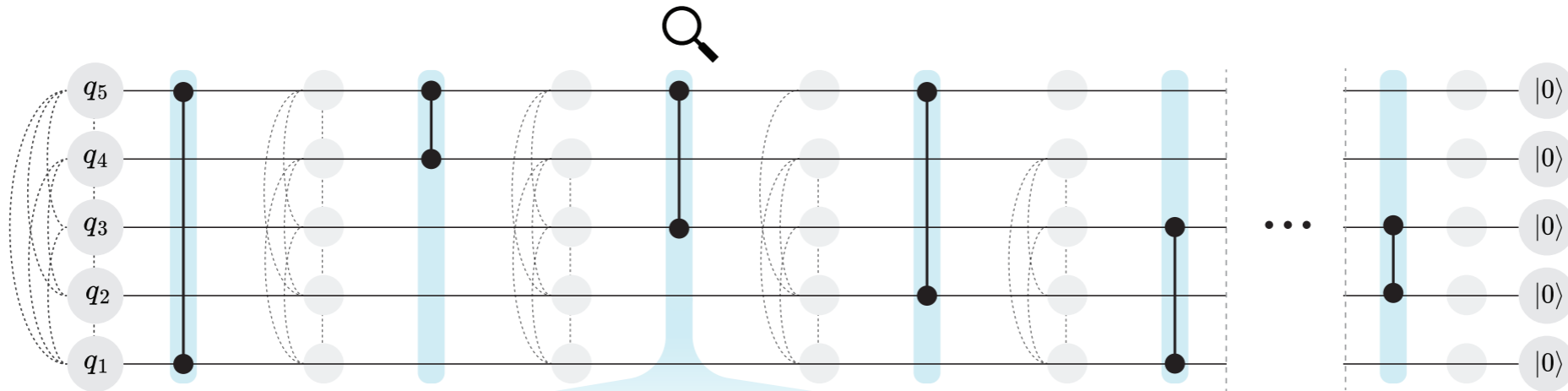




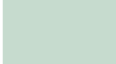

# Deep Reinforcement Learning for Quantum Technology





# RL related talks in ML4QT2

## 2<sup>nd</sup> Workshop "Machine Learning for Quantum Technology" - SCHEDULE

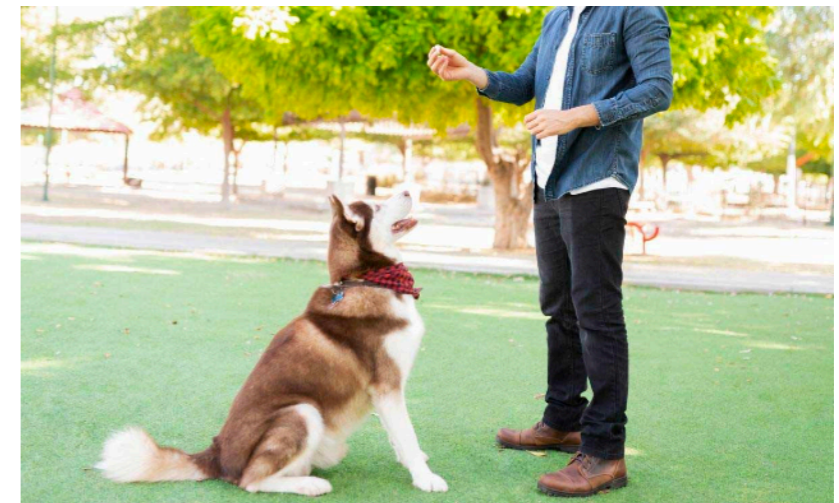
	WEDNESDAY, NOV. 6	THURSDAY, NOV. 7	FRIDAY, NOV. 8
9:00 - 9:30	<b>Registration + opening remarks</b>	<b>Registration</b>	<b>Registration</b>
9:30 - 10:05	Eliška Greplová <i>Autonomous Quantum Control in the age of AI</i>	Monika Aidelsburger <i>Quantum many-body systems under the microscope</i>	Vedran Dunjko <i>Provable exponential quantum advantages in learning from classical data</i>
10:05 - 10:40	Anton Frisk Kockum <i>Quantum state and process tomography with machine learning and gradient descent</i>	Simon Trebst <i>Decoding many-body teleportation</i>	Hans Briegel <i>Towards explainable AI in quantum science</i>
10:40 - 10:55	Martin Gärttner <i>Machine learning assisted quantum simulator readout</i>	Yue Ban <i>Neural-network-assisted parameter estimation for quantum detection</i>	Chenfeng Cao <i>Unveiling quantum phase transitions from traps in variational quantum algorithms</i>
10:55 - 11:20	<b>Coffee break</b>		
11:20 - 11:55	Christopher Eichler <i>Realizing a reinforcement learning agent for real-time quantum feedback</i>	Giuseppe Carleo <i>Neural quantum states for many-body electronic structure and dynamics</i>	Johannes Bausch <i>Machine Learning for Fault-Tolerant Quantum Computation</i>
11:55 - 12:30	Annabelle Bohrdt <i>Trying to solve quantum many-body problems with neural networks</i>	Markus Schmitt <i>(Neural) network representations of many-body wave functions</i>	Evert van Nieuwenburg <i>RL and RL for quantum systems</i>
12:30 - 12:45	Maximilian Prüfer <i>Physics-inspired machine learning models and optimal control for quantum experiments</i>	Dario Poletti <i>Paths towards time evolution with larger neural-network quantum states</i>	Matias Bilkis <i>Automatic re-calibration of quantum devices by RL</i>
12:45 - 13:00	Petr Zapletal <i>Error-tolerant quantum convolutional neural networks for symmetry-protected topological phases</i>	Gorka Muñoz-Gil <i>Representation learning reaches the lab: let machines act!</i>	Clara Wanjura <i>Quantum Equilibrium Propagation for efficient training of quantum systems based on Onsager reciprocity</i>
13:00 - 14:30	<b>Lunch break</b>		
14:30 - 15:05	Marín Bukov <i>Reinforcement learning transmon-qubit entangling gates</i>	Roger Melko <i>Language Models for Quantum Simulation</i>	Jonas Schuff <i>Autonomous tuning of spin qubits</i>
15:05 - 15:40	Volodymyr Sivak <i>Calibration of decoders for quantum error correction using multi-agent reinforcement learning</i>	Markus Heyl <i>Solving 2D quantum matter with neural quantum states</i>	Christof Weitenberg <i>Machine learning and ultracold quantum gases</i>
15:40 - 15:55	Maciej Koch-Janusz <i>Analyzing and constructing efficient data encoding quantum circuits</i>	Cristian Bonato <i>Learning the dynamics of Markovian open quantum systems from experimental data</i>	<b>Closing remarks</b>
15:55 - 16:20	<b>Coffee break</b>		
16:20 - 16:35	Bijita Sarma <i>Fast Hardware-efficient Quantum Gate Design using Optimal Control with Reinforcement Learning Ansatz</i>	Akash Kundu <i>Program synthesis-driven quantum architecture search for optimal quantum circuit design in variational quantum algorithms</i>	<b>LEGEND</b>  Invited talk (30'+5' Q&A)  Contributed talk (12'+3' Q&A)
16:35 - 17:10	Mats Granath <i>Graph neural network based decoders for quantum error correcting codes</i>	Mario Krenn <i>Towards an Artificial Muse for new ideas in Science</i>	
17:10 - 18:00	<b>Poster flash talks (1' each) + poster setup</b>	<b>Poster flash talks (1' each) + poster setup</b>	
From 18:00	<b>Poster session A (including dinner)</b>	<b>Poster session B (including dinner)</b>	



# Outline

## Part 1

- Reinforcement learning (RL) in quantum physics
  - RL as a branch of machine learning
- Applications of RL
  - hallmark applications of RL
  - applications in quantum technologies
- RL framework in a nutshell
  - environment, states, actions, rewards
  - RL algorithms





# Outline

## Part 2

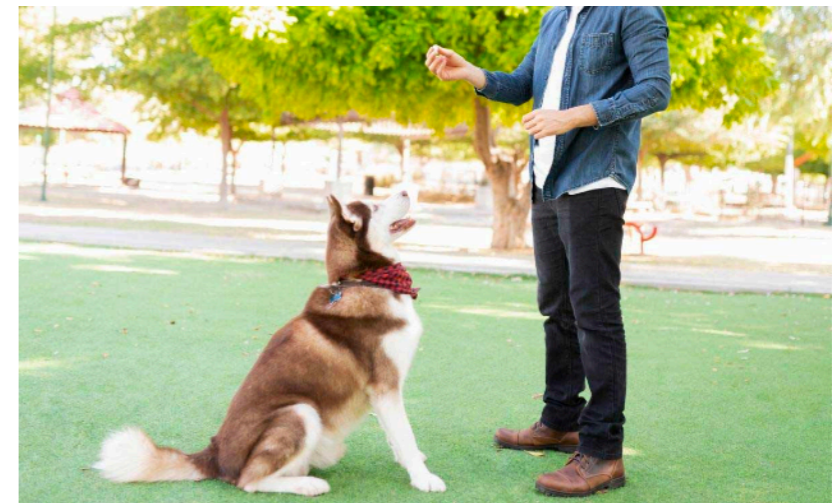
- RL for qubit state preparation
  - effect of noise (measurement shot noise, coherent, incoherent noise)
- experimentally friendly RL framework
  - partially observable environments
  - environment, states, actions, rewards



# Outline

## Part 1

- Reinforcement learning (RL) in quantum physics
  - RL as a branch of machine learning
- Applications of RL
  - hallmark applications of RL
  - applications in quantum technologies
- RL framework in a nutshell
  - environment, states, actions, rewards
  - RL algorithms



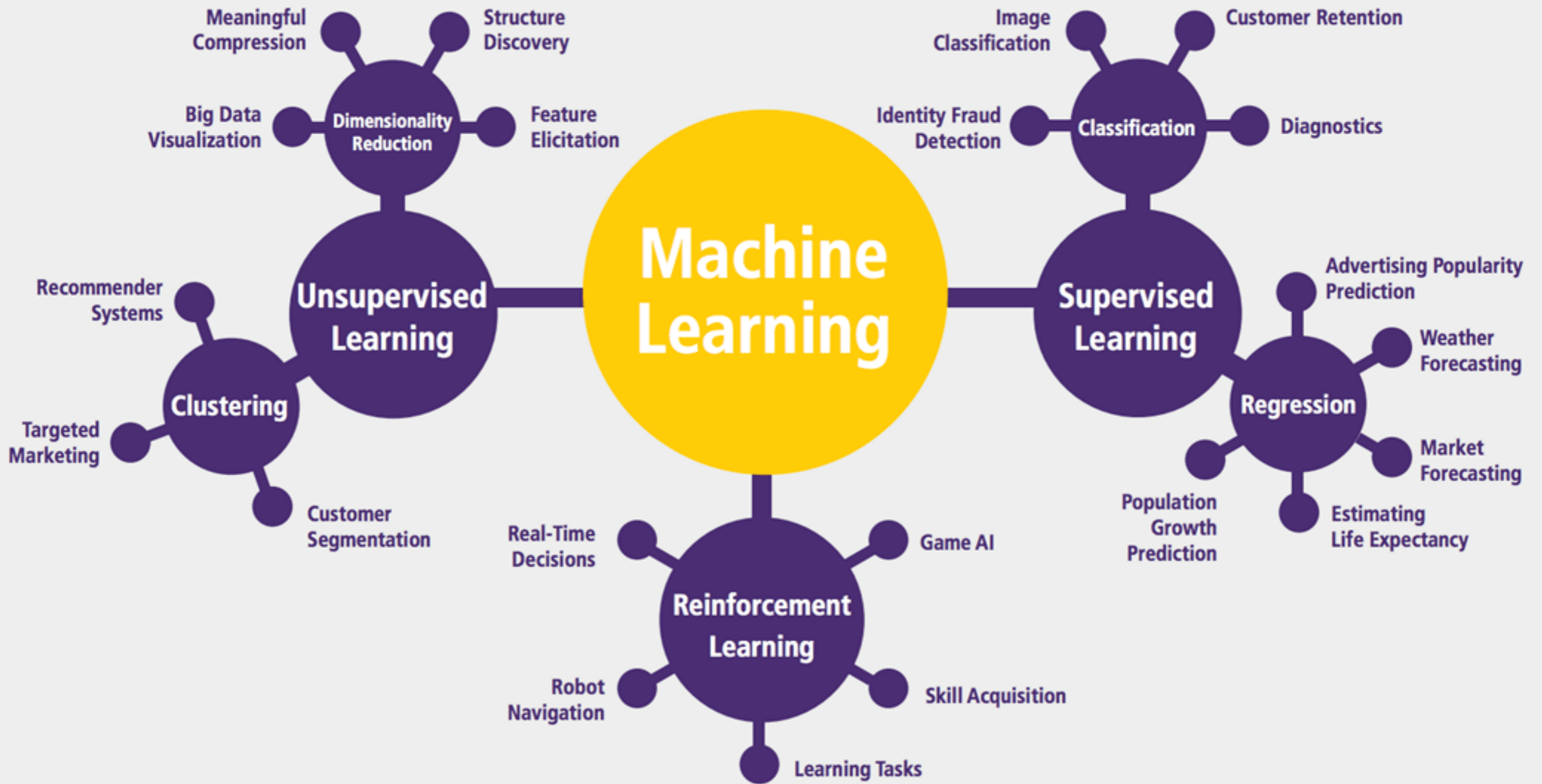


image: Priya Pareek

# Supervised Learning

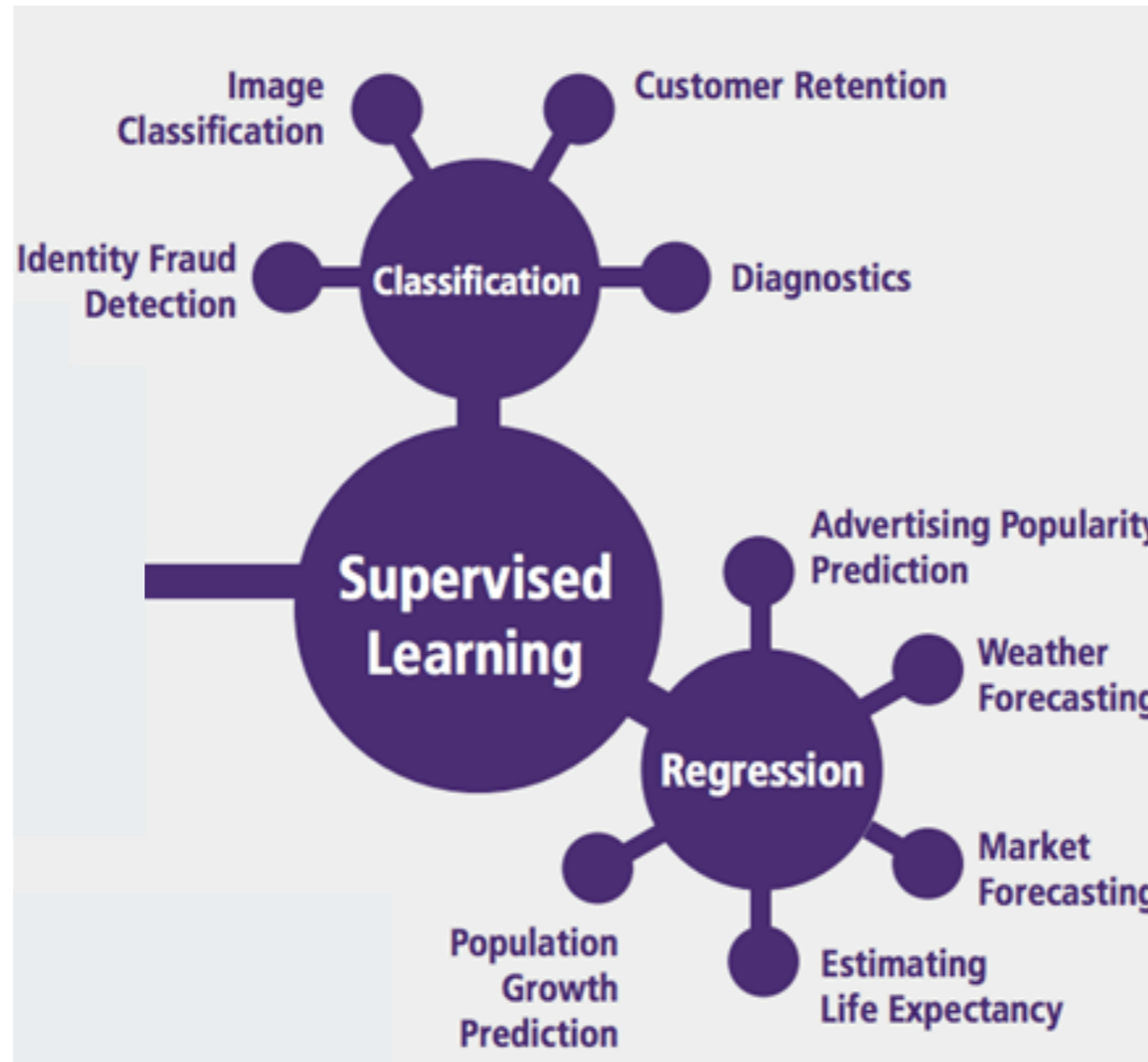
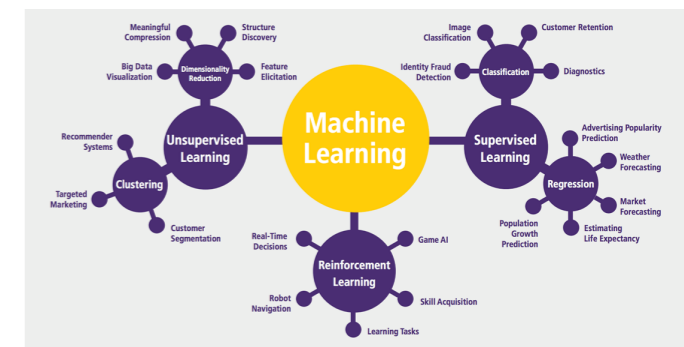
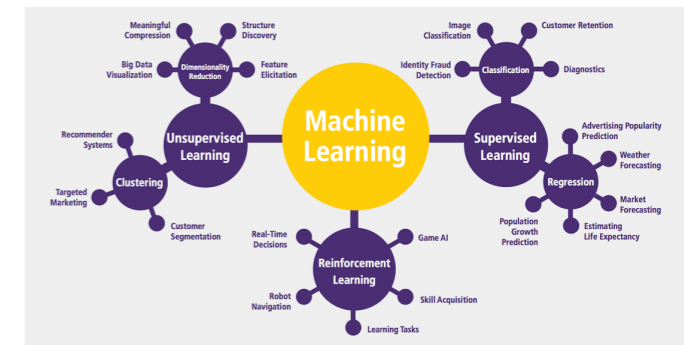
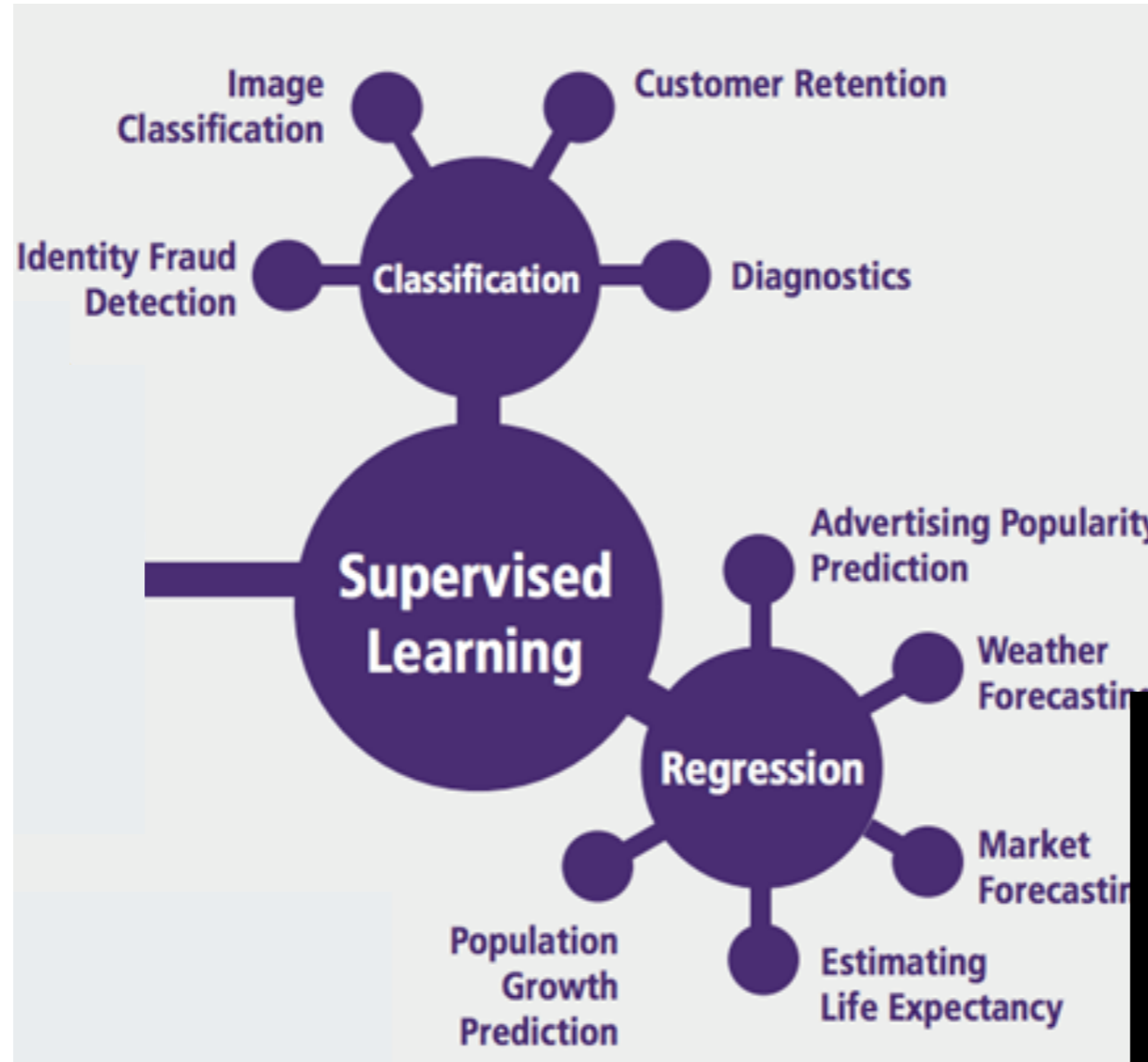


image: Priya Pareek



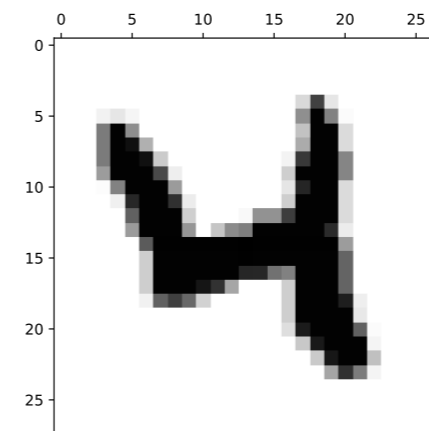
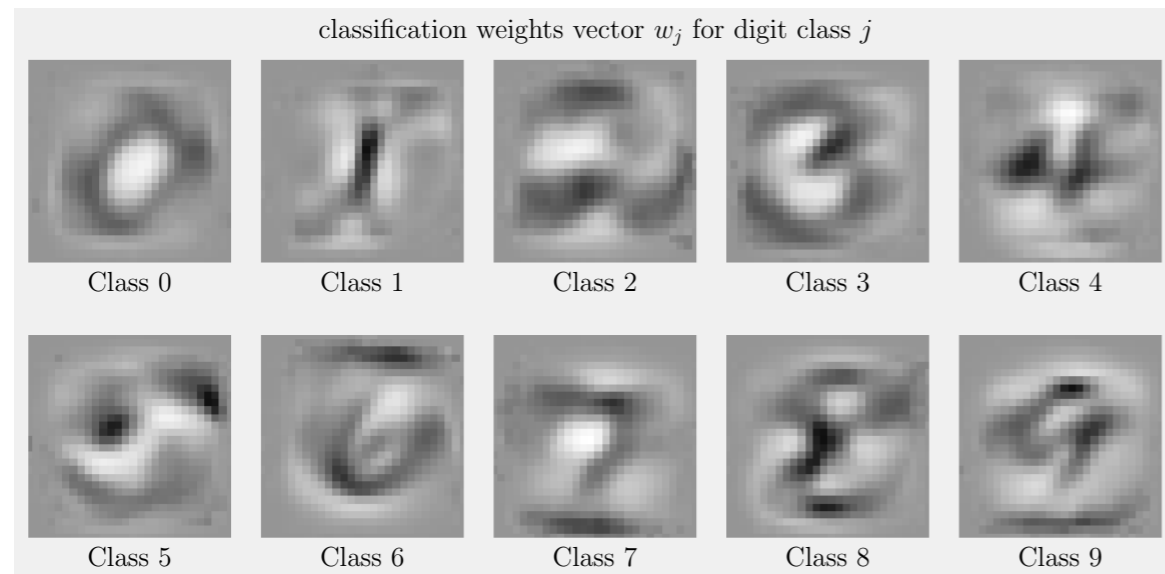
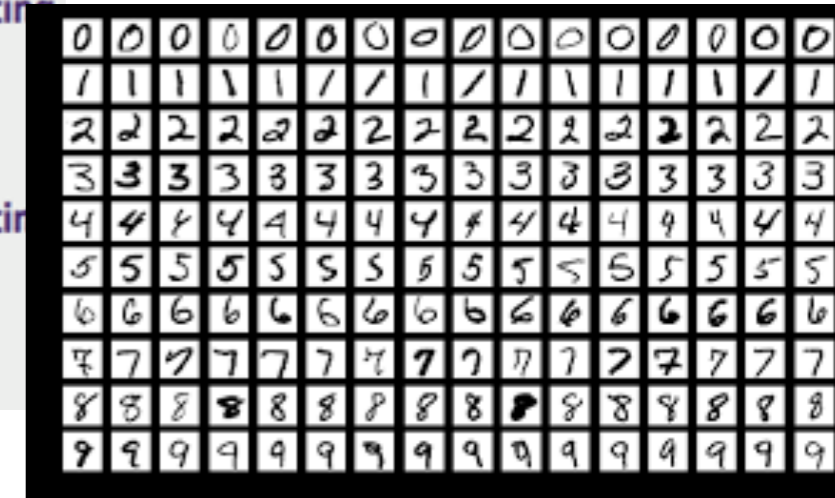
- ▶ learning from examples (labeled data)

# Supervised Learning



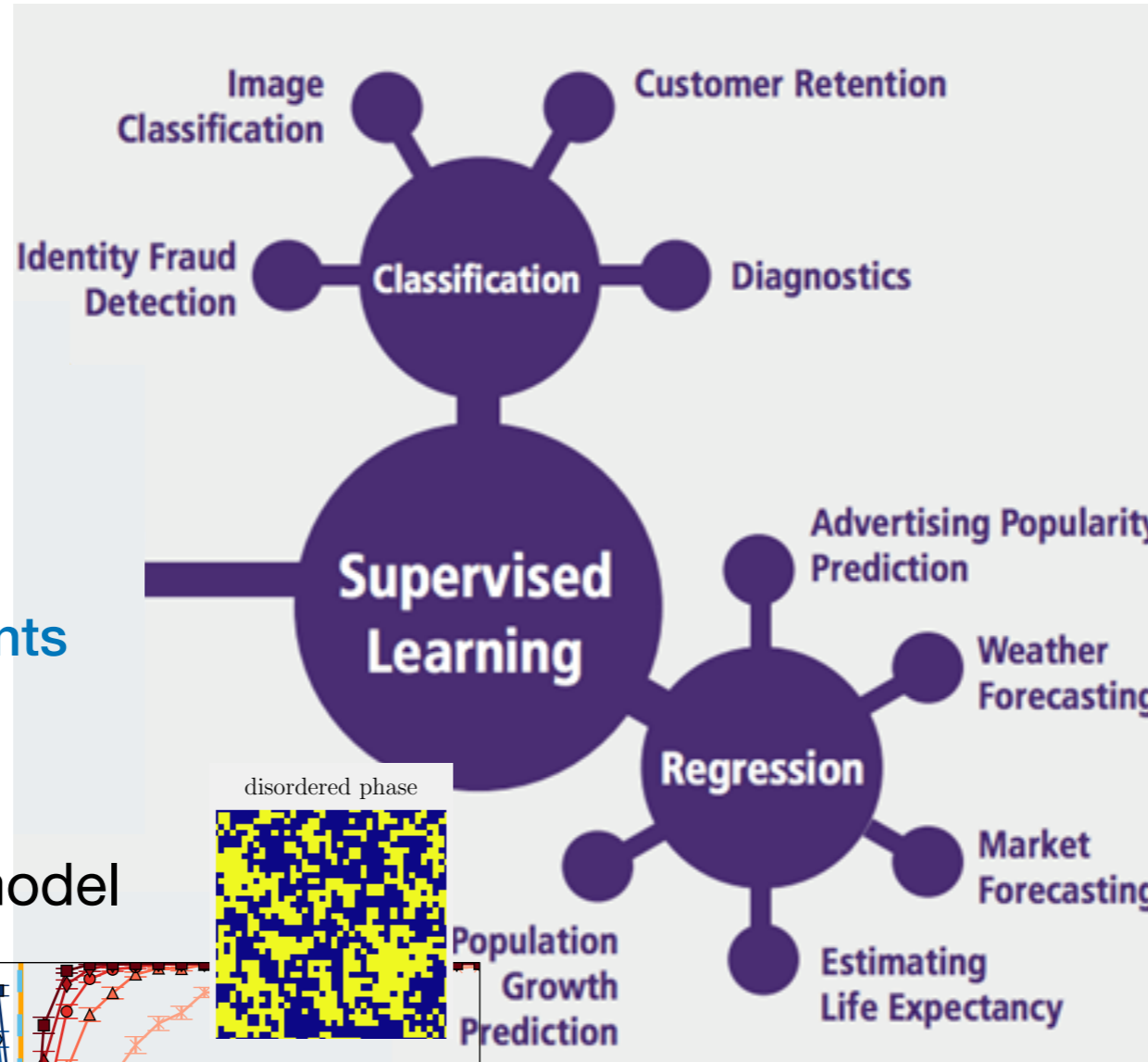
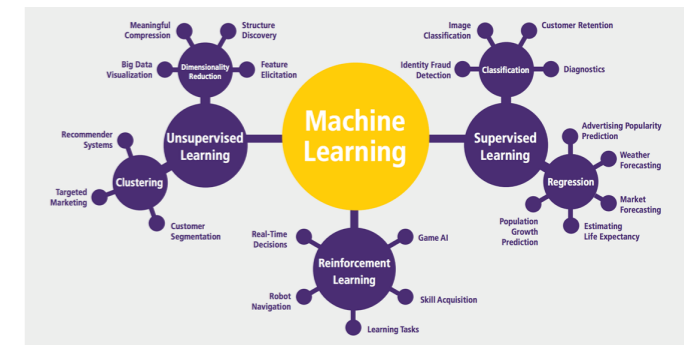
- learning from examples (labeled data)

MNIST





# Supervised Learning



learning from examples (labeled data)

- ▶ MNIST
- ▶ recognizing hand-written digits

- ▶ classify phases of matter
- ▶ determine critical points from data

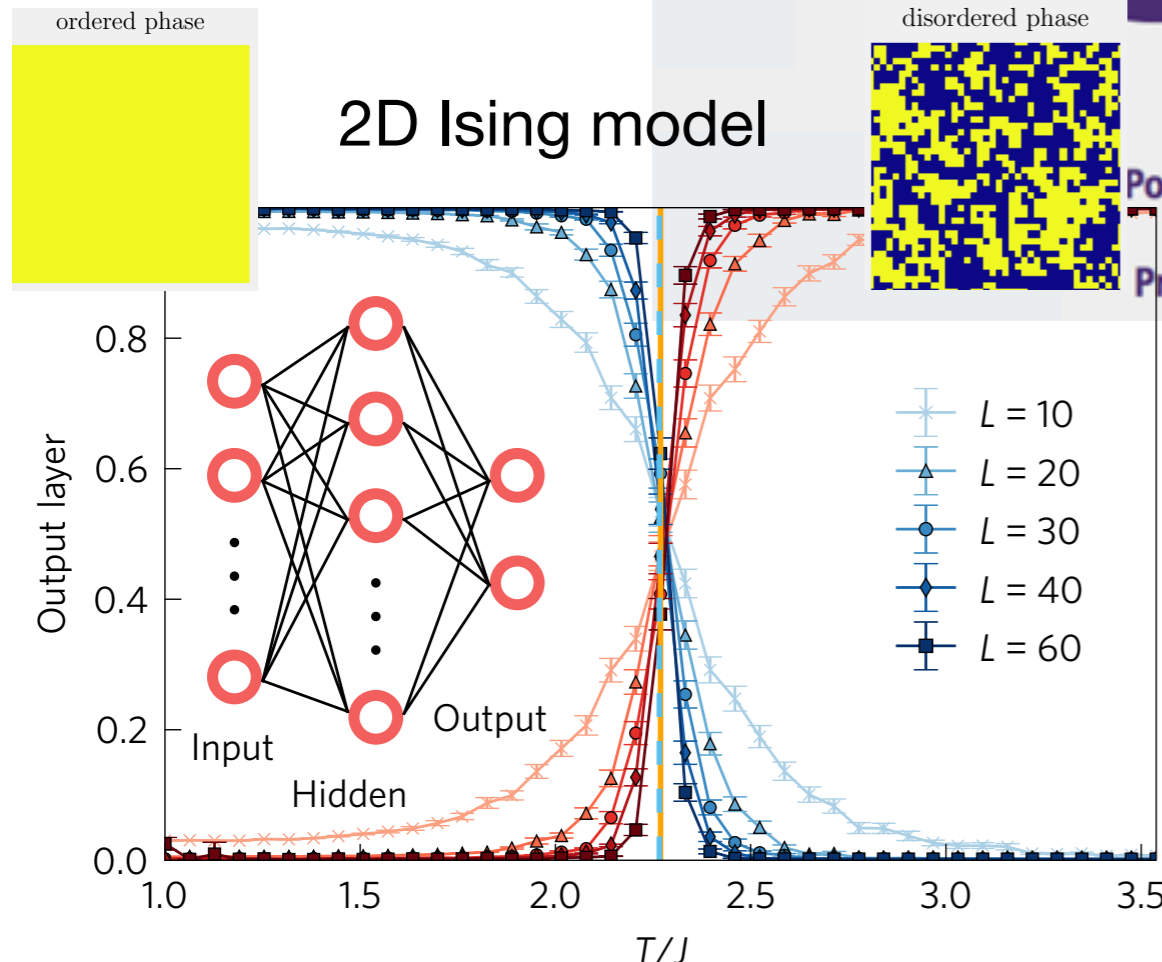
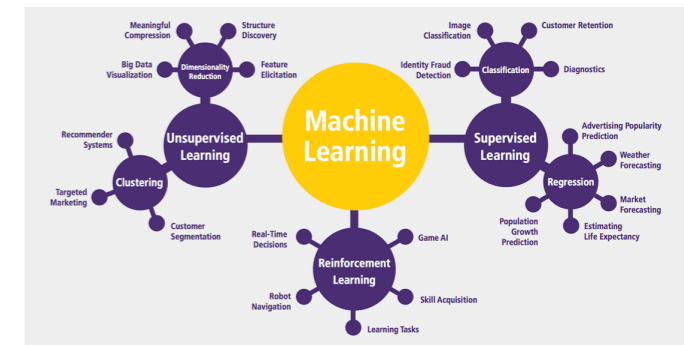


image: Priya Pareek

# Unsupervised Learning



▶ learning the distribution that generated data

- ▶ compose music
- ▶ draw paintings
- ▶ write text (Chat GPT)

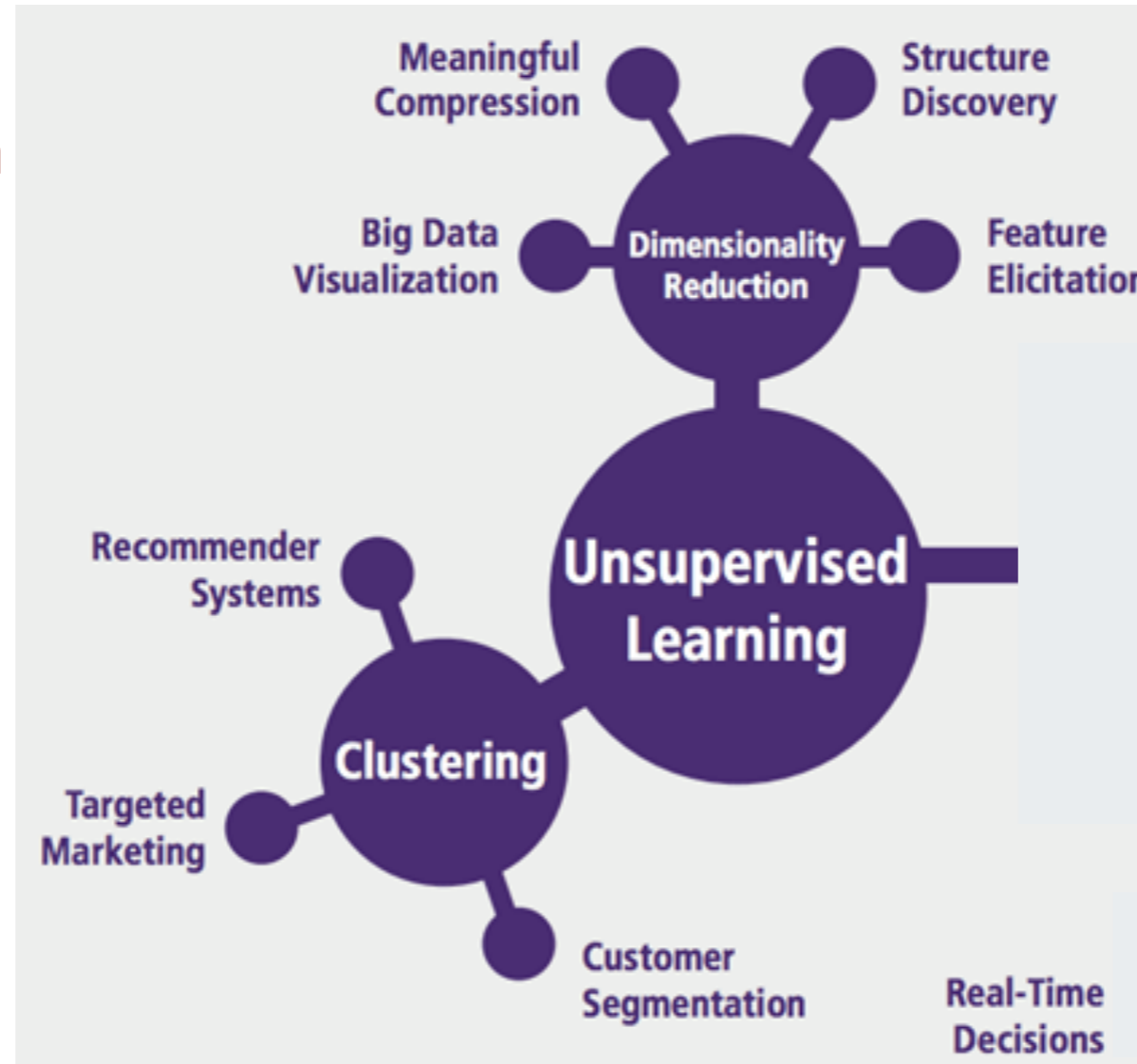


image: Priya Pareek

# Unsupervised Learning

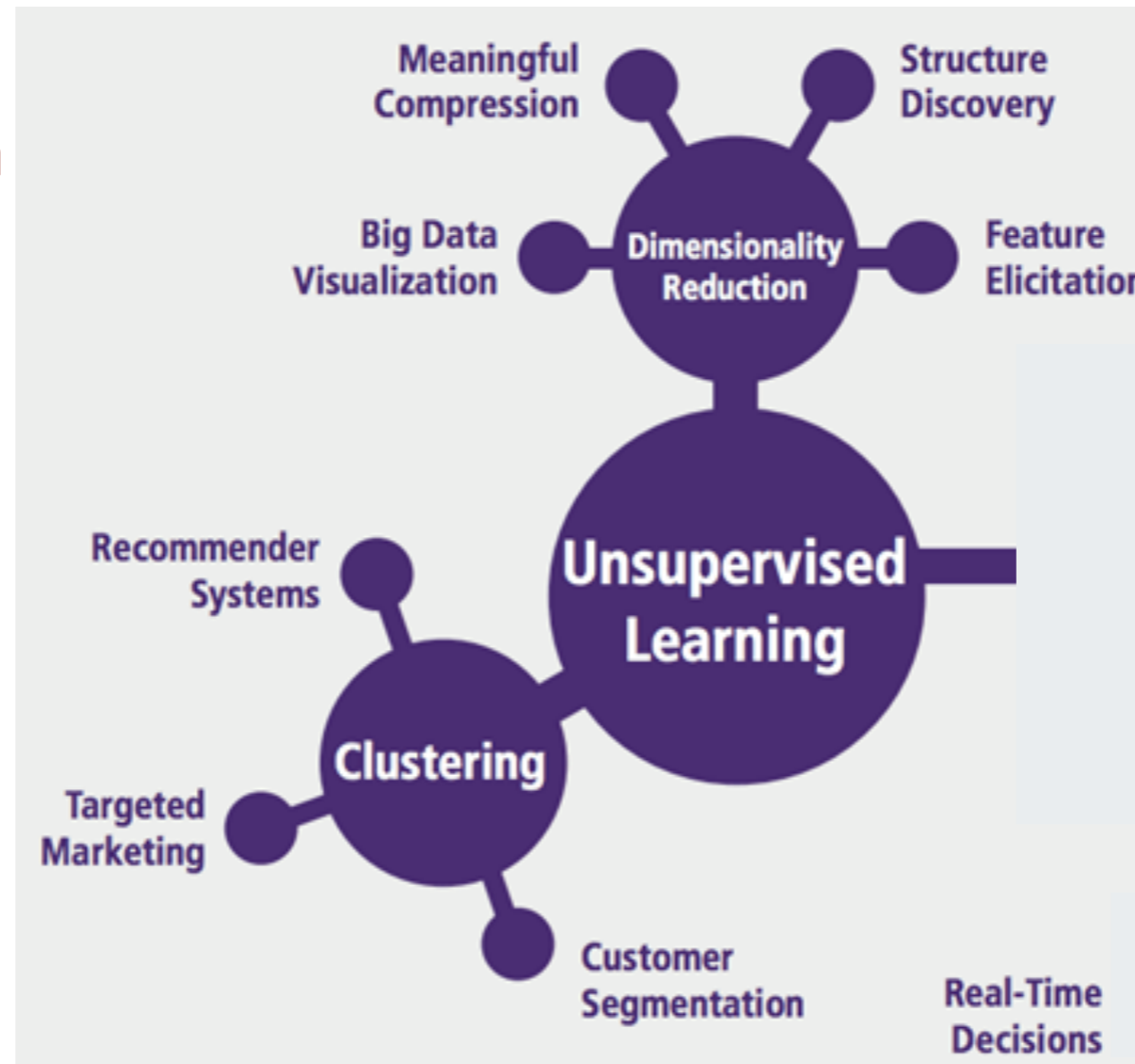
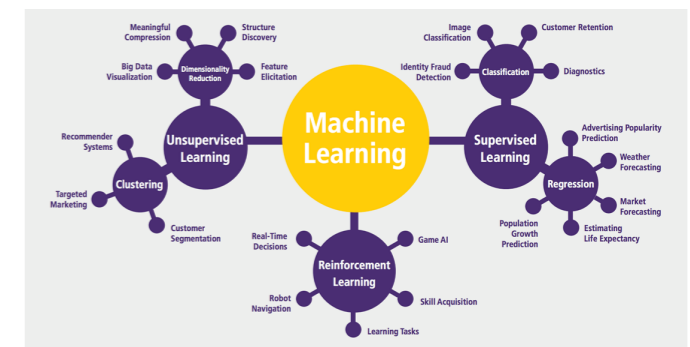
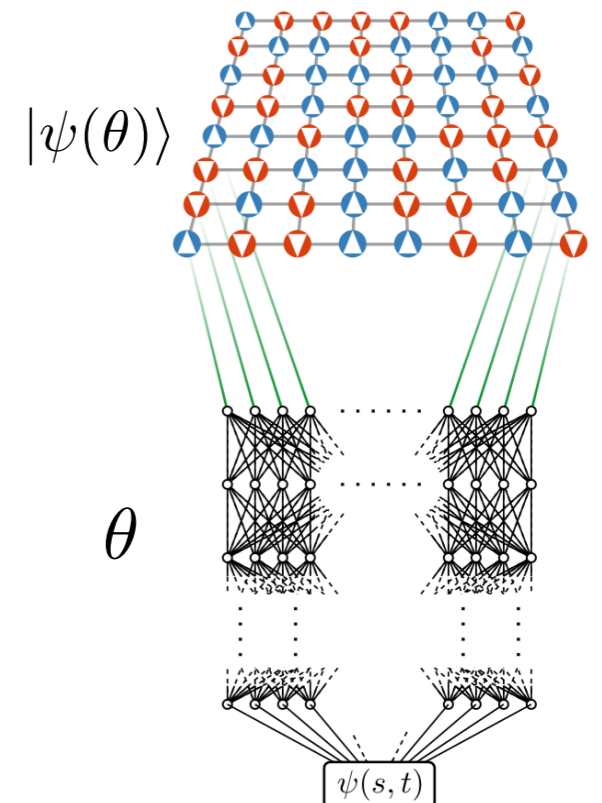


image: Priya Pareek

## Neural Quantum States

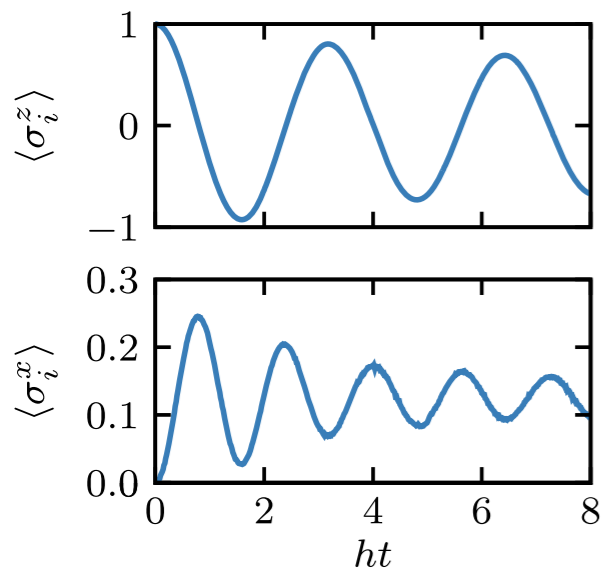
$$E(\theta) = \langle \psi(\theta) | H | \psi(\theta) \rangle$$

$$\partial_{\theta} E(\theta) \stackrel{!}{=} 0$$

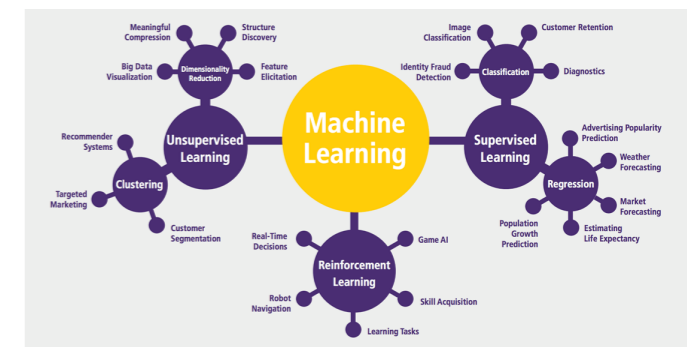


learning the distribution that generated data

- compose music
- draw paintings
- write text (Chat GPT)



# Reinforcement Learning



- learning from experience

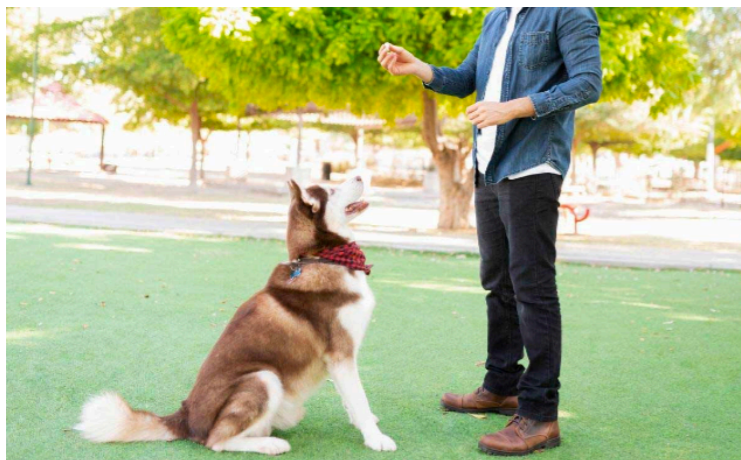


image: Canine Journal

RL entails interactive dynamics

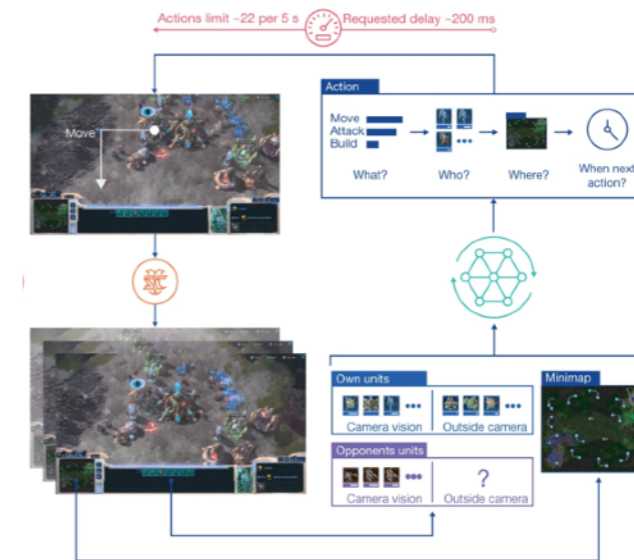
# What is reinforcement learning used for?

## Mastering the game of Go with deep neural networks and tree search



Silver, et. al, Nature 529 484–489 (2016)

## Mastering video games (StarCraft II)



Vinyals, et. al, Nature 350 (2019)

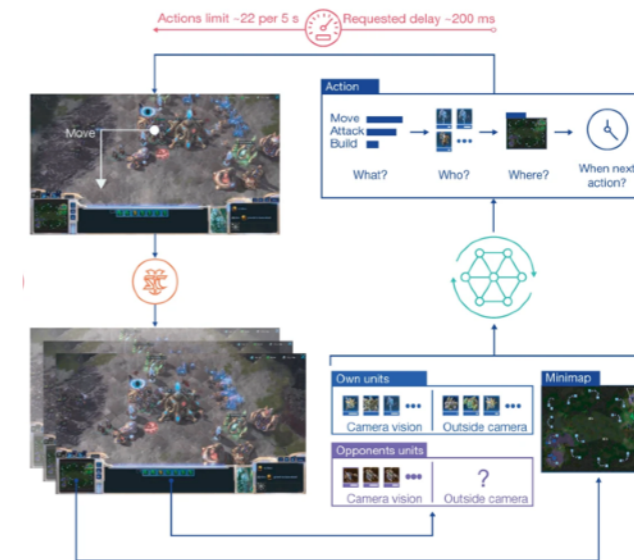
# What is reinforcement learning used for?

## Mastering the game of Go with deep neural networks and tree search



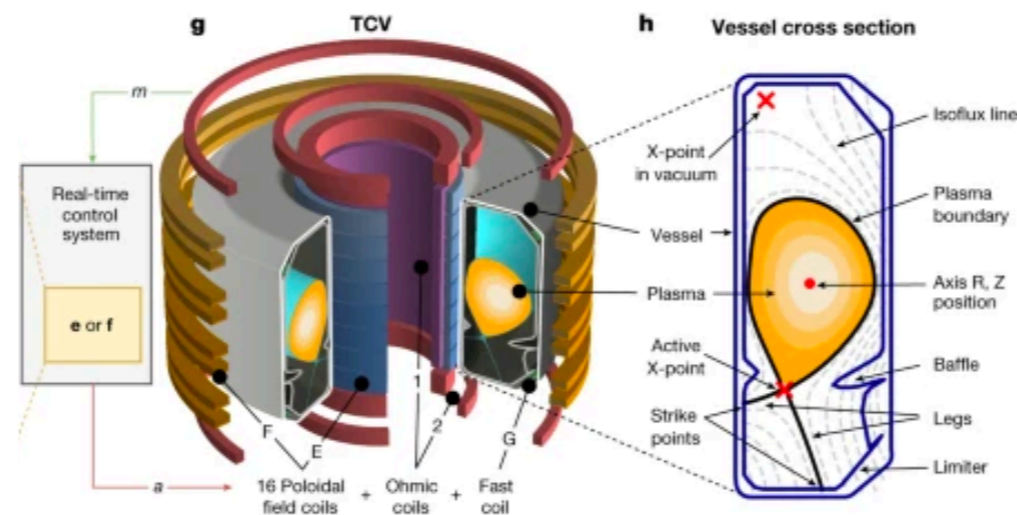
Silver, et. al, Nature 529 484–489 (2016)

## Mastering video games (StarCraft II)



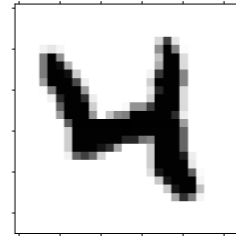
Vinyals, et. al, Nature 350 (2019)

## Magnetic control of tokamak plasmas thru deep RL

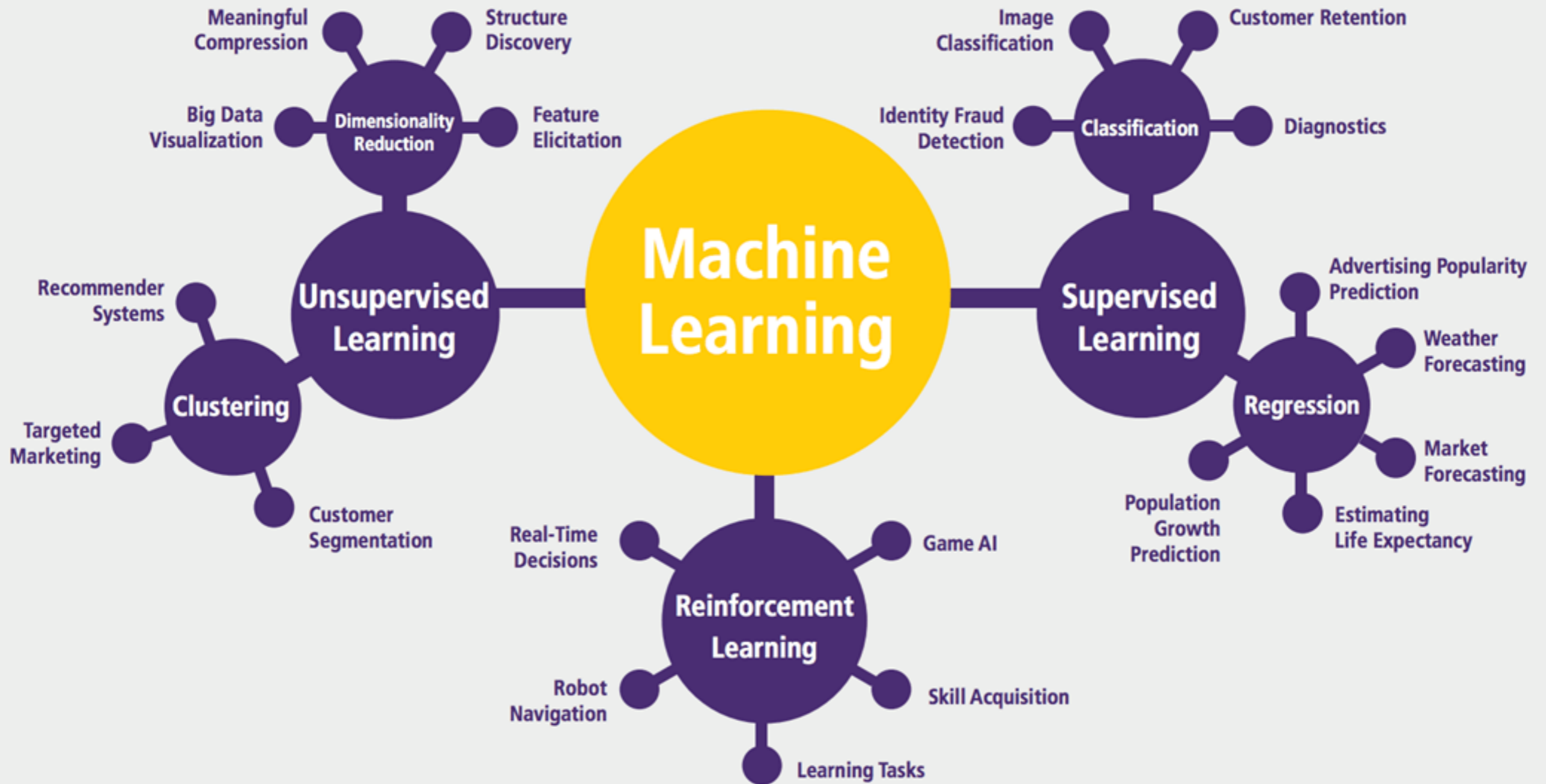


Degrave, et. al, Nature 602 414–419 (2022)

▶ learning distribution that generates data



▶ learning from examples (labeled data)



▶ learning from experience

image: Priya Pareek

▶ learning distribution that generates data



chemistry, 2024

▶ learning from examples (labeled data)

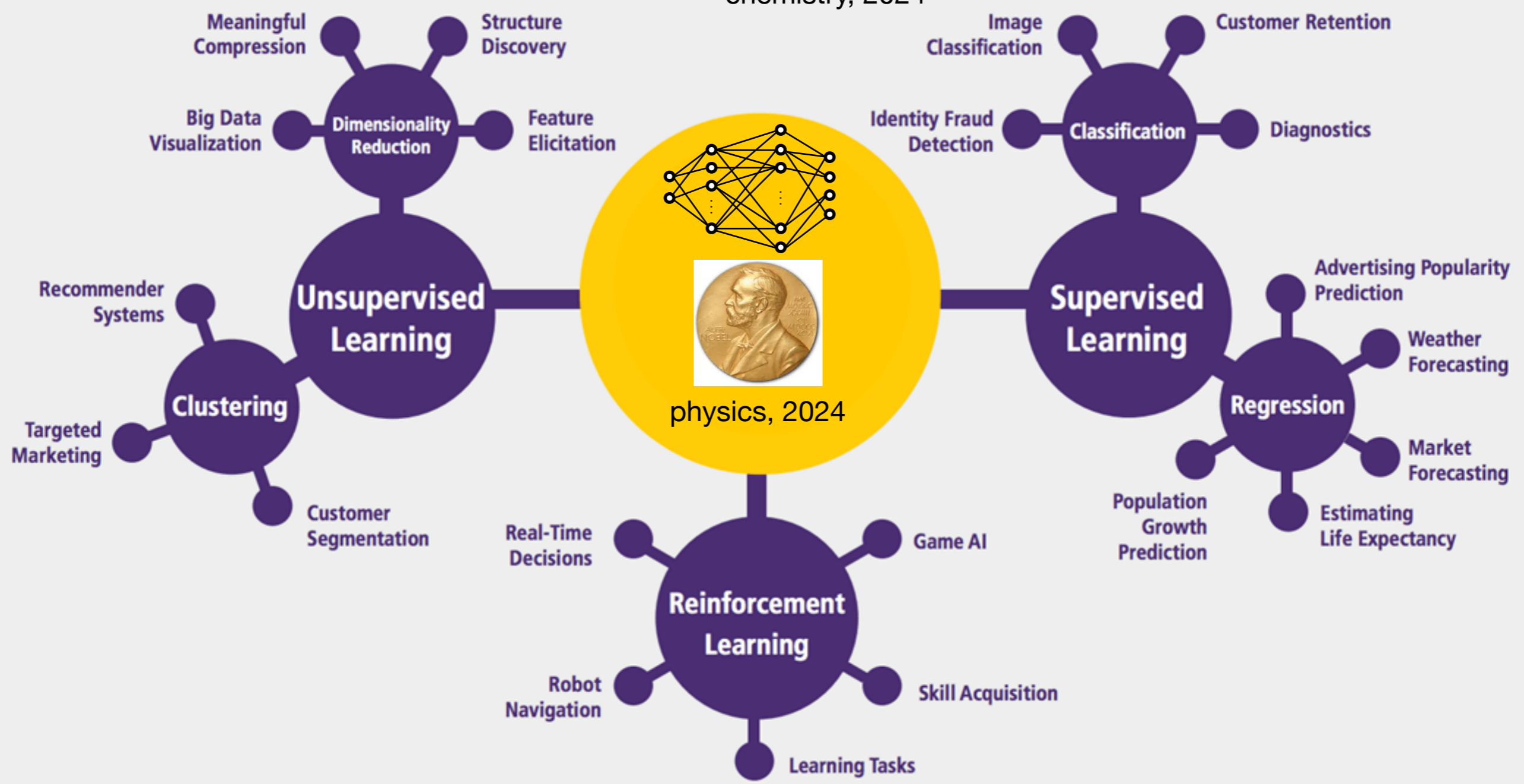


image: Priya Pareek

▶ learning from experience





# Outline

## Part 1

### ✓ Reinforcement learning (RL) in quantum physics

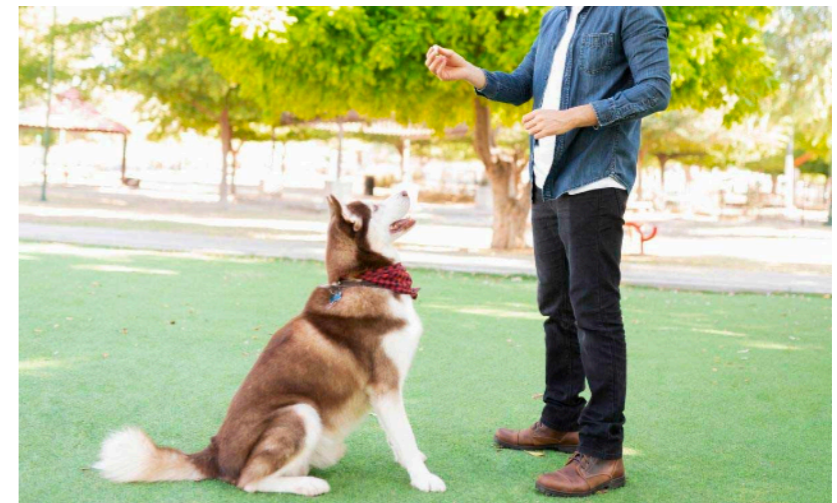
- RL as a branch of machine learning

### • Applications of RL

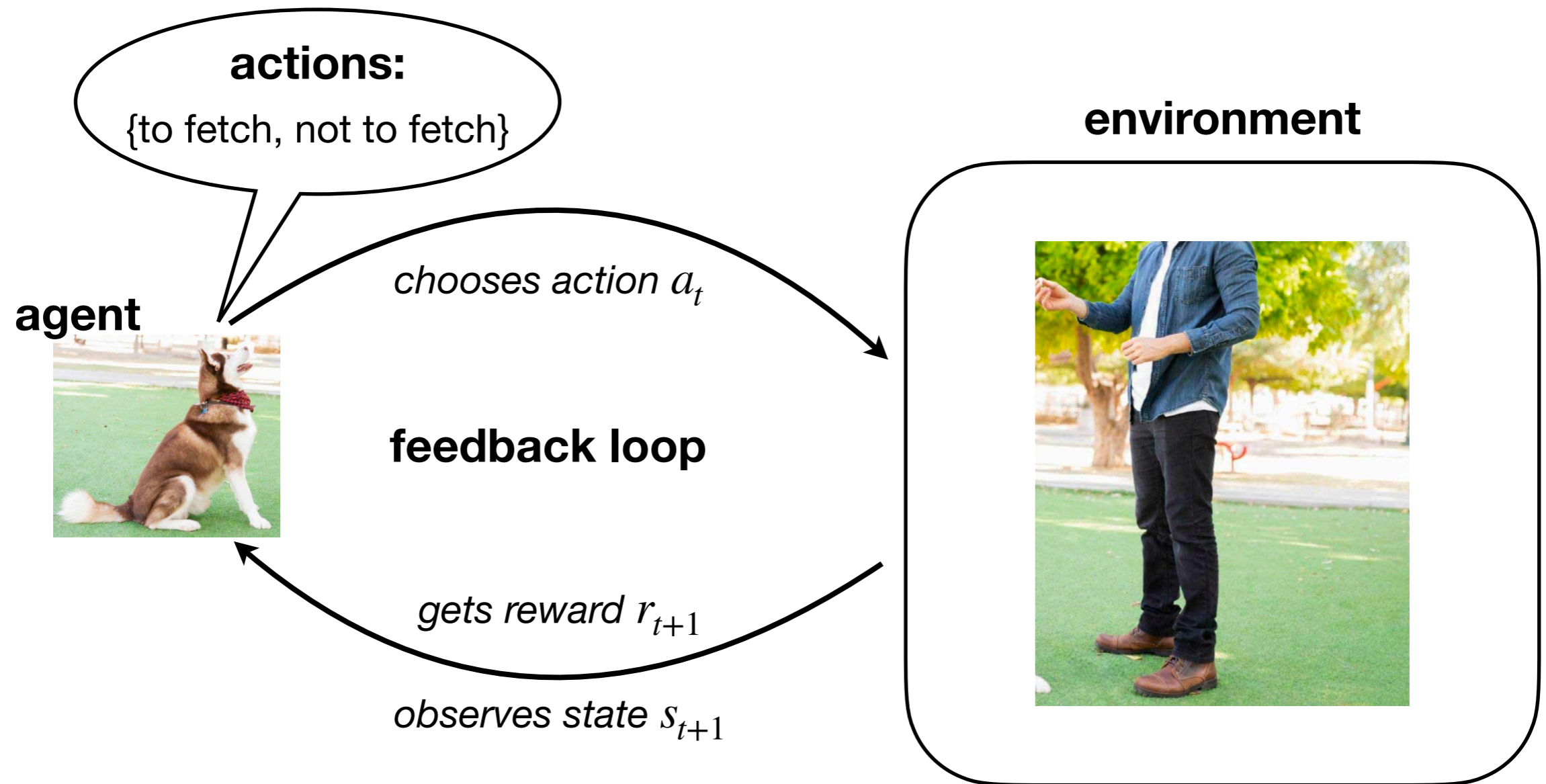
- hallmark applications of RL
- applications in quantum technologies

### • RL framework in a nutshell

- environment, states, actions, rewards
- RL algorithms



# Reinforcement learning (RL) in a nutshell



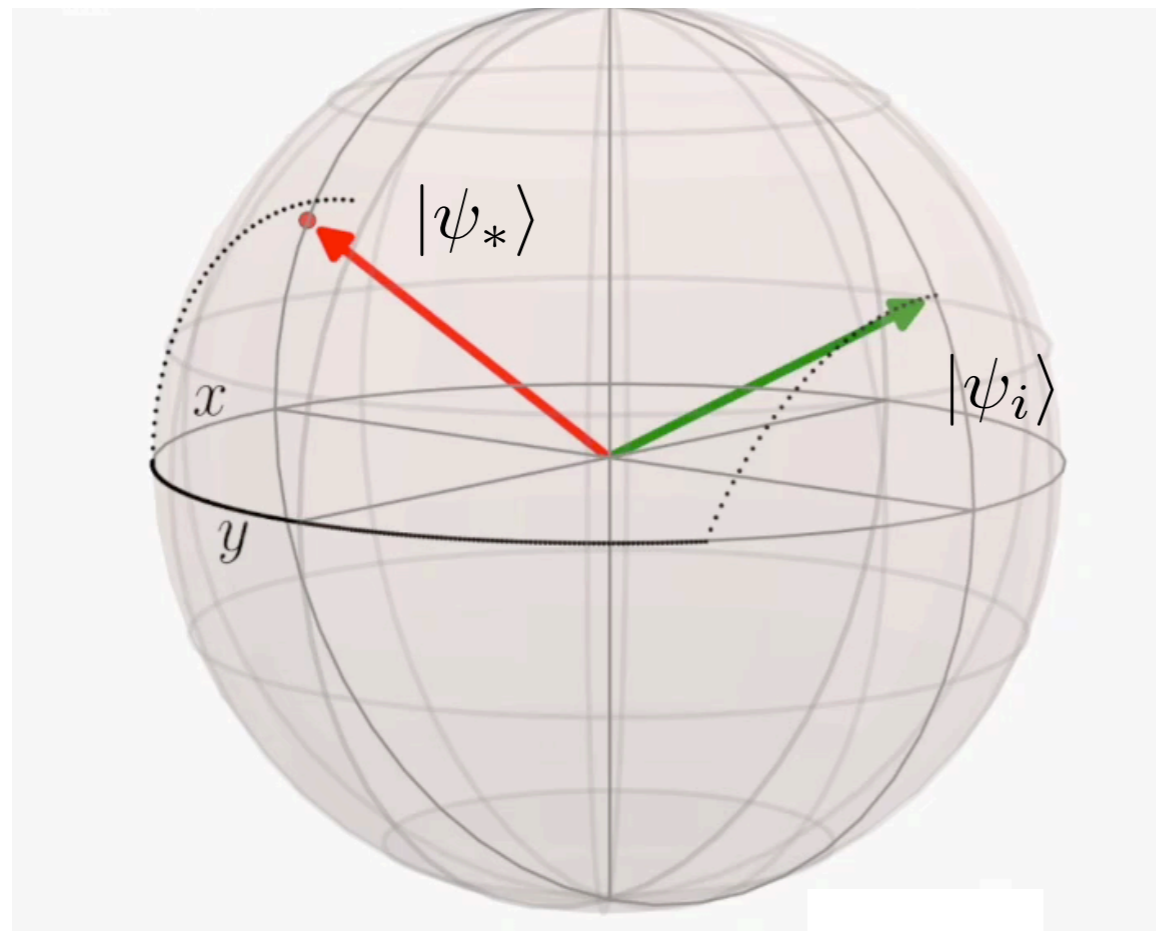
# Applications of RL in Quantum Physics

- quantum control

$$H(t) = -\frac{1}{2} (Z + h_x(t)X)$$

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

**Bloch sphere**



**state preparation = fidelity optimization**

$$F_h(t) = |\langle \psi_* | \psi(t) \rangle|^2$$

$$|\psi(t)\rangle = \mathcal{T} e^{-i \int_0^t dt' H(t')} |\psi_i\rangle$$

$$h_x(t) = ?$$

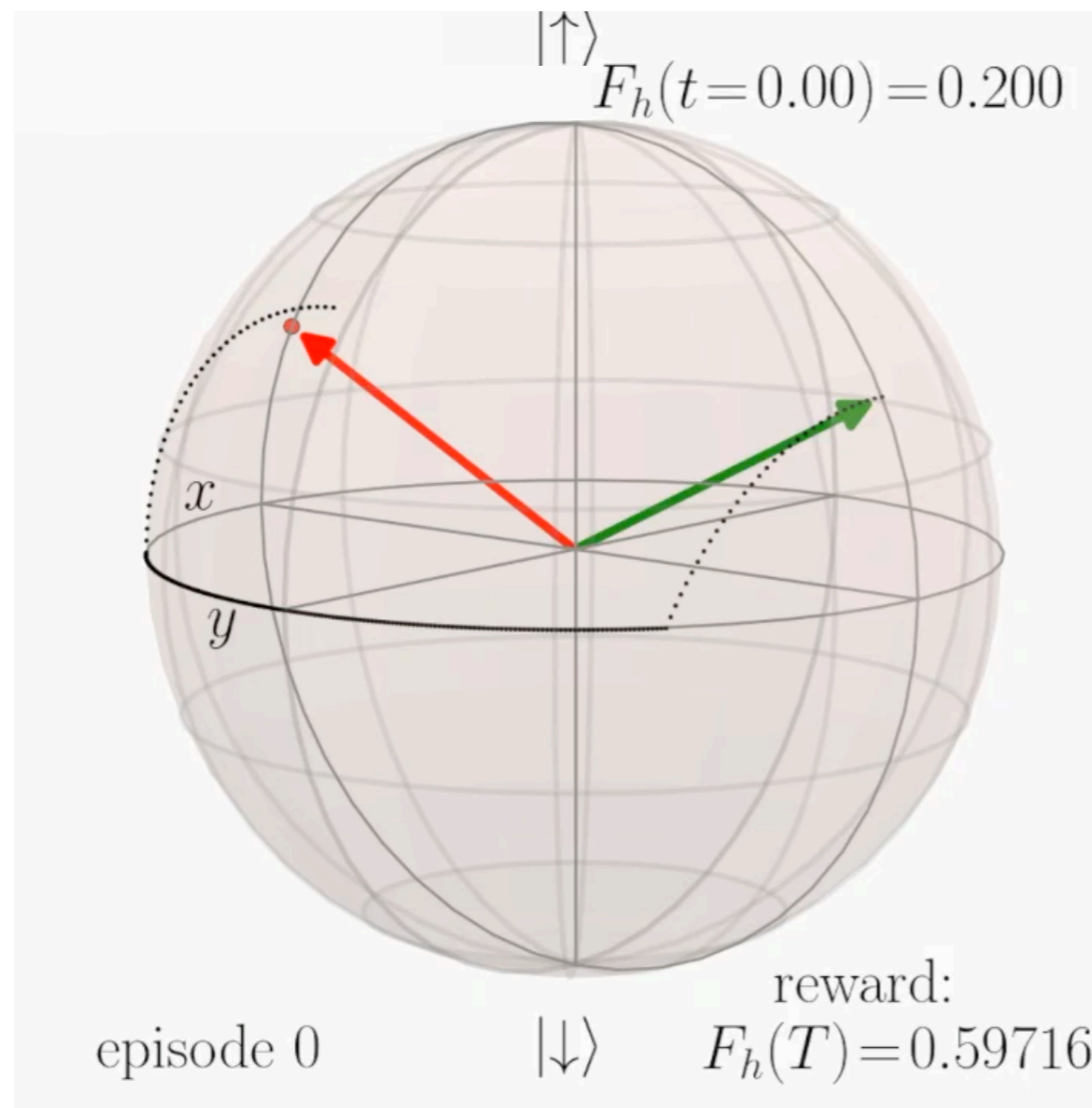
# Applications of RL in Quantum Physics

- quantum control

$$H(t) = -\frac{1}{2} (Z + h_x(t)X)$$

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

## Bloch sphere



state preparation = fidelity optimization

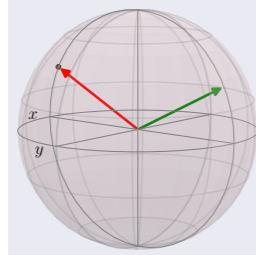
$$F_h(t) = |\langle \psi_* | \psi(t) \rangle|^2$$

$$|\psi(t)\rangle = \mathcal{T} e^{-i \int_0^t dt' H(t')} |\psi_i\rangle$$

$$h_x(t) = ?$$

# Applications of RL in Quantum Physics

- quantum control



MB et al, PRX 8 031086 (2018)

Niu et al, npj 5 33 (2019)

Sivak et al, PRX 12, 011059 (2022)

Gispén et al, MSML (2021)

Reuer, Nat Comm 14 7138 (2023)

Yao et al, PRX 11 (3), 031070 (2021)

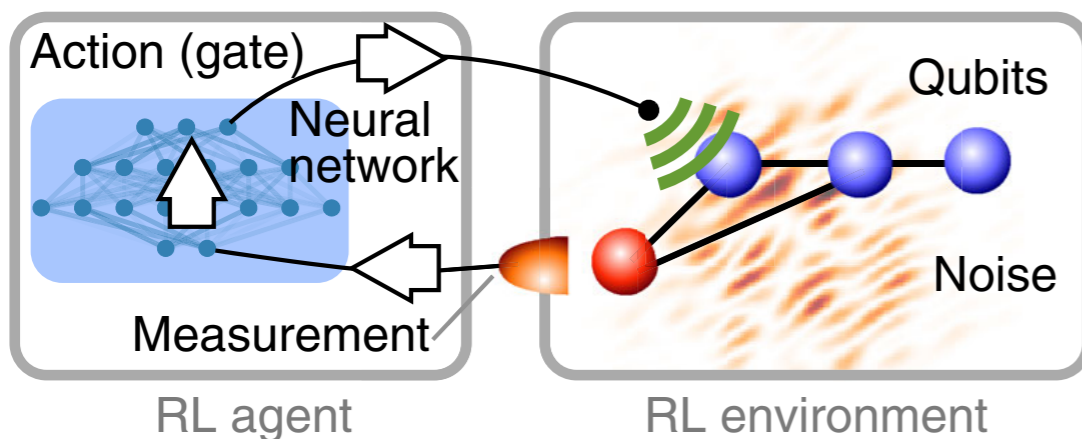
Porotti, Comm Phys 2 (2019)

Dalgaard et al, npj 6 6 (2020)

+ many more

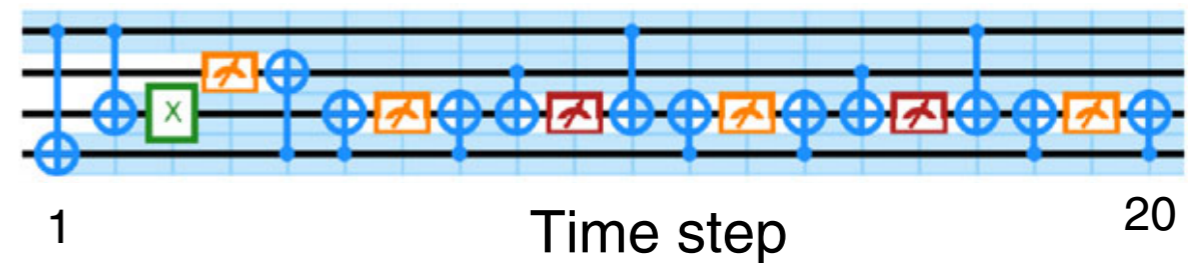
- quantum error correction

- ▶ **task:** find error correcting code that protects qubits from decoherence



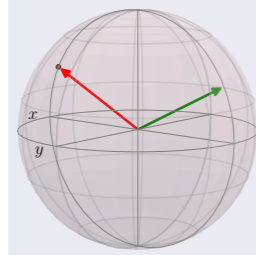
Fössel et al, PRX 8 031086 (2018)

Olle et al, arXiv:2311.04750



# Applications of RL in Quantum Physics

- quantum control



MB et al, PRX 8 031086 (2018)

Niu et al, npj 5 33 (2019)

Sivak et al, PRX 12, 011059 (2022)

Gispen et al, MSML (2021)

Reuer, Nat Comm 14 7138 (2023)

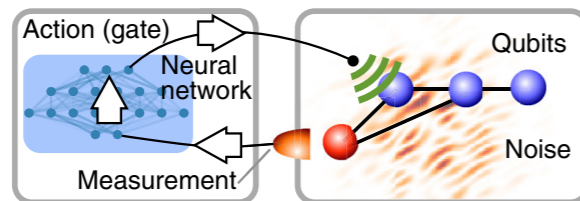
Yao et al, PRX 11 (3), 031070 (2021)

Porotti, Comm Phys 2 (2019)

Dalgaard et al, npj 6 6 (2020)

+ many more

- quantum error correction



Fössel et al, PRX 8 031086 (2018)

Andreasson et al, Quantum 3 183 (2019)

Sweke et al, ML Sci Tech 2 025005 (2020)

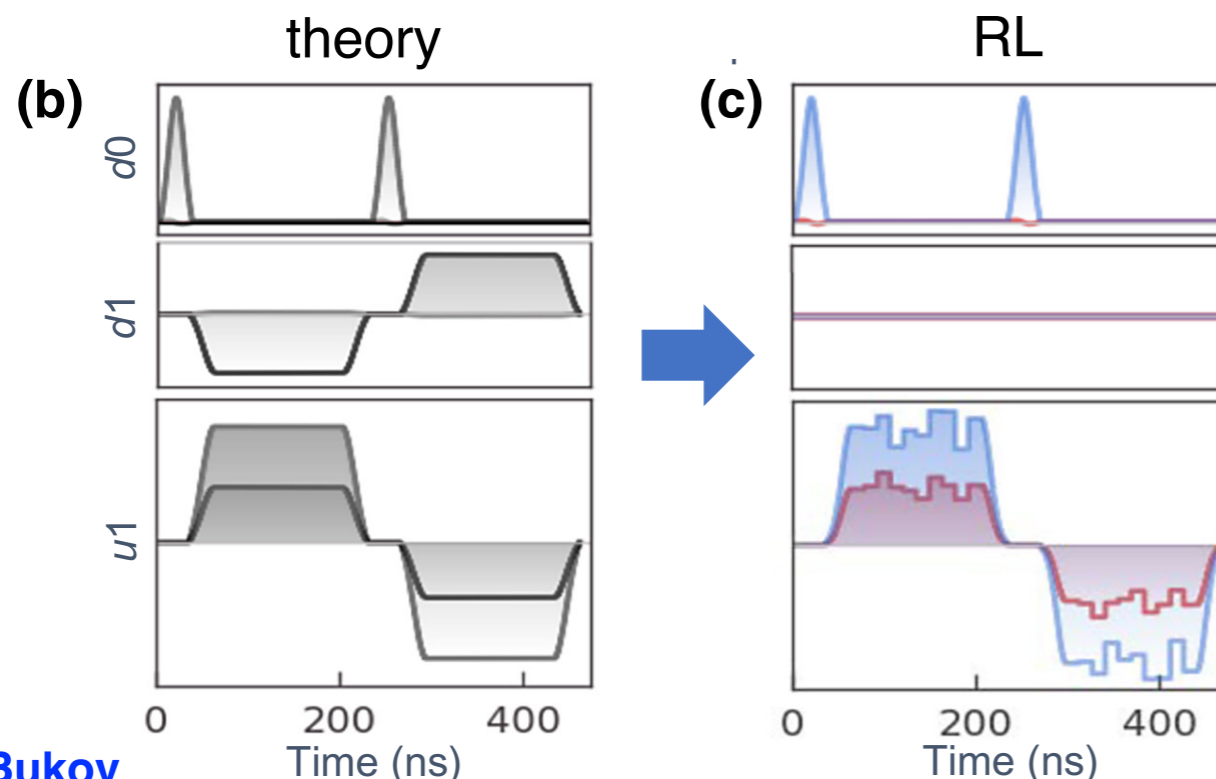
Sivak et al, Nature 616 50-55 (2023)

Olle et al, arXiv:2311.04750

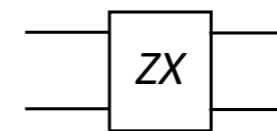
+ many more

- quantum gate design

- task:** find high-fidelity pulses that emulate gates on a quantum computer



Baum et al, PRX Quantum 2, 040324 (2021)



cross resonance gate

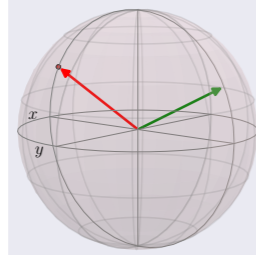
— IBM  $(1.41 \pm 0.06) \times 10^{-2}$

— DRL  $(6.50 \pm 0.61) \times 10^{-3}$

infidelity / error rate

# Applications of RL in Quantum Physics

## ● quantum control



MB et al, PRX 8 031086 (2018)

Niu et al, npj 5 33 (2019)

Sivak et al, PRX 12, 011059 (2022)

Gispen et al, MSML (2021)

Reuer, Nat Comm 14 7138 (2023)

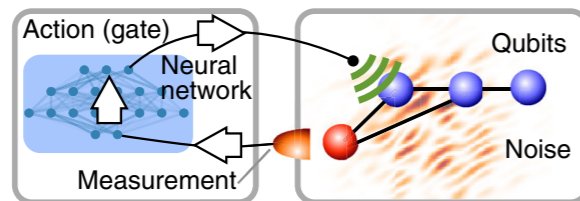
Yao et al, PRX 11 (3), 031070 (2021)

Porotti, Comm Phys 2 (2019)

Dalgaard et al, npj 6 6 (2020)

+ many more

## ● quantum error correction



Fössel et al, PRX 8 031086 (2018)

Andreasson et al, Quantum 3 183 (2019)

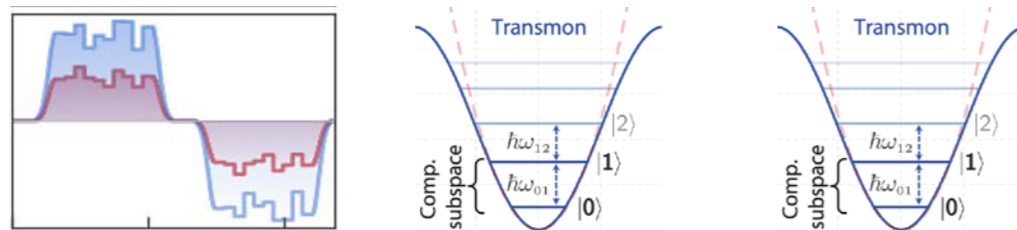
Sweke et al, ML Sci Tech 2 025005 (2020)

Sivak et al, Nature 616 50-55 (2023)

Olle et al, arXiv:2311.04750

+ many more

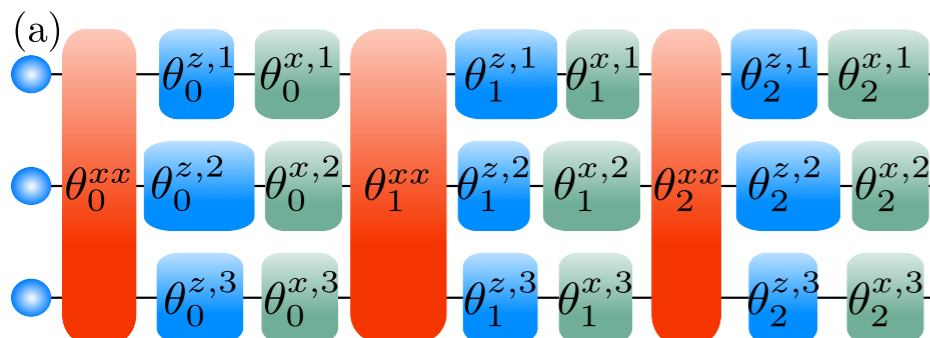
## ● quantum gate design



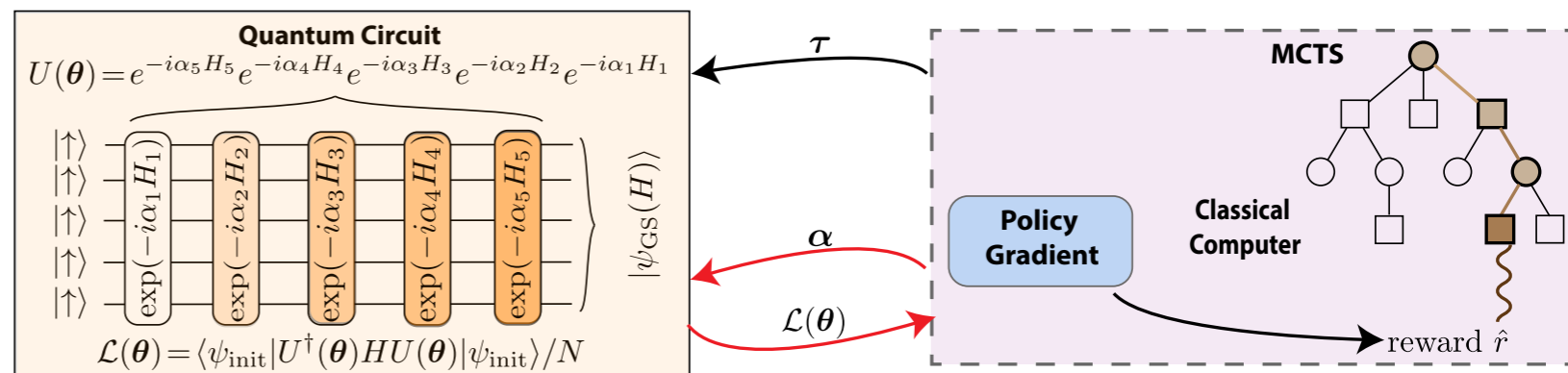
Baum et al, PRX Quantum 2, 040324 (2021)

Nguyen et al, ML Sci & Tech, 5, 025066 (2024)

## ● quantum circuit design and synthesis



Bolens et al, PRL 2021



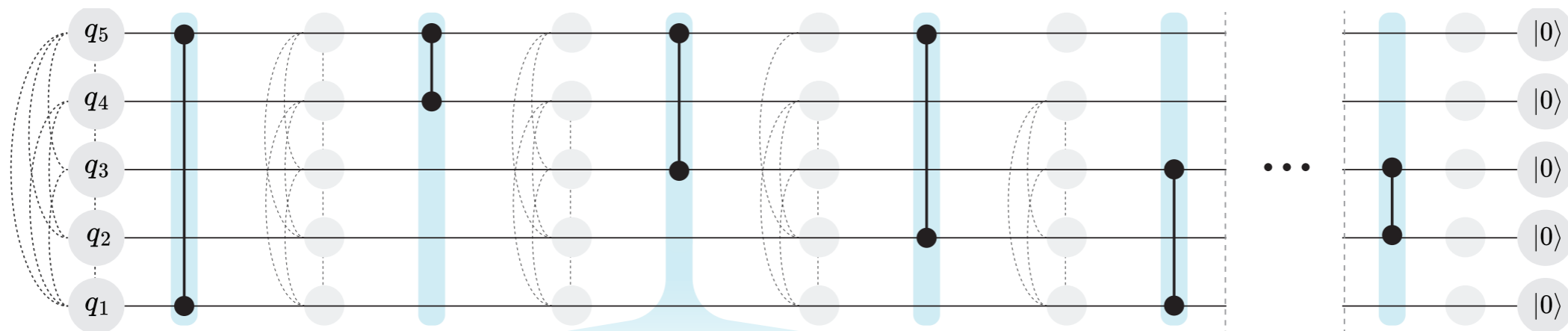
Yao et al, MSML 2019, 2021, 2022

+ many more

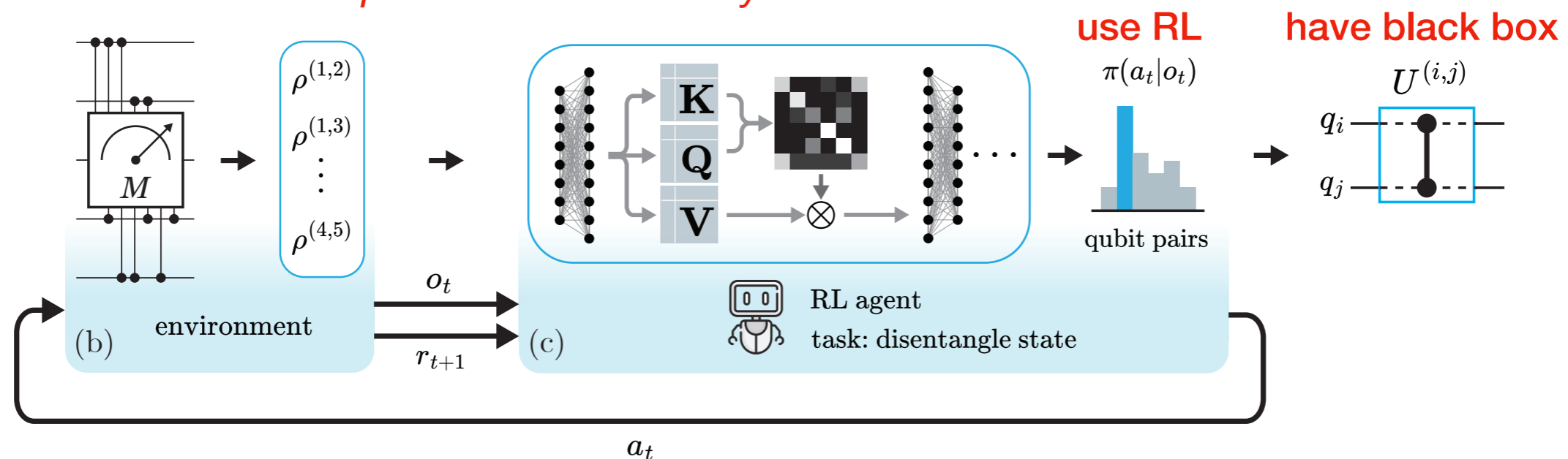
# Applications of RL in Quantum Physics

- quantum control
- quantum gate design
- quantum error correction
- quantum circuit design and synthesis

► construct disentangling circuits



→ access to all 2-qubit reduced density matrices

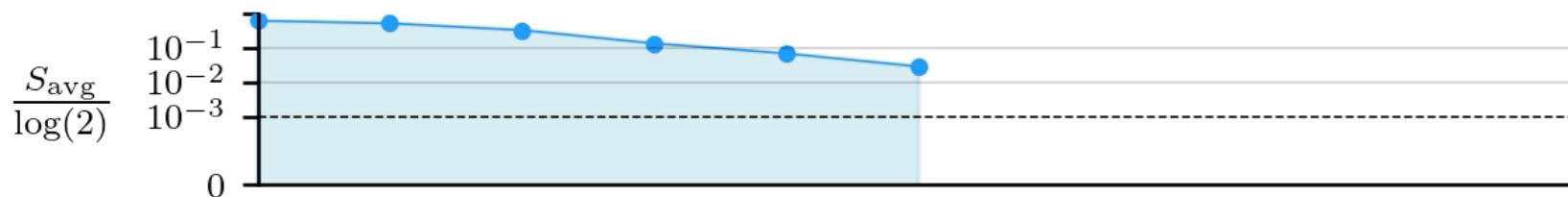
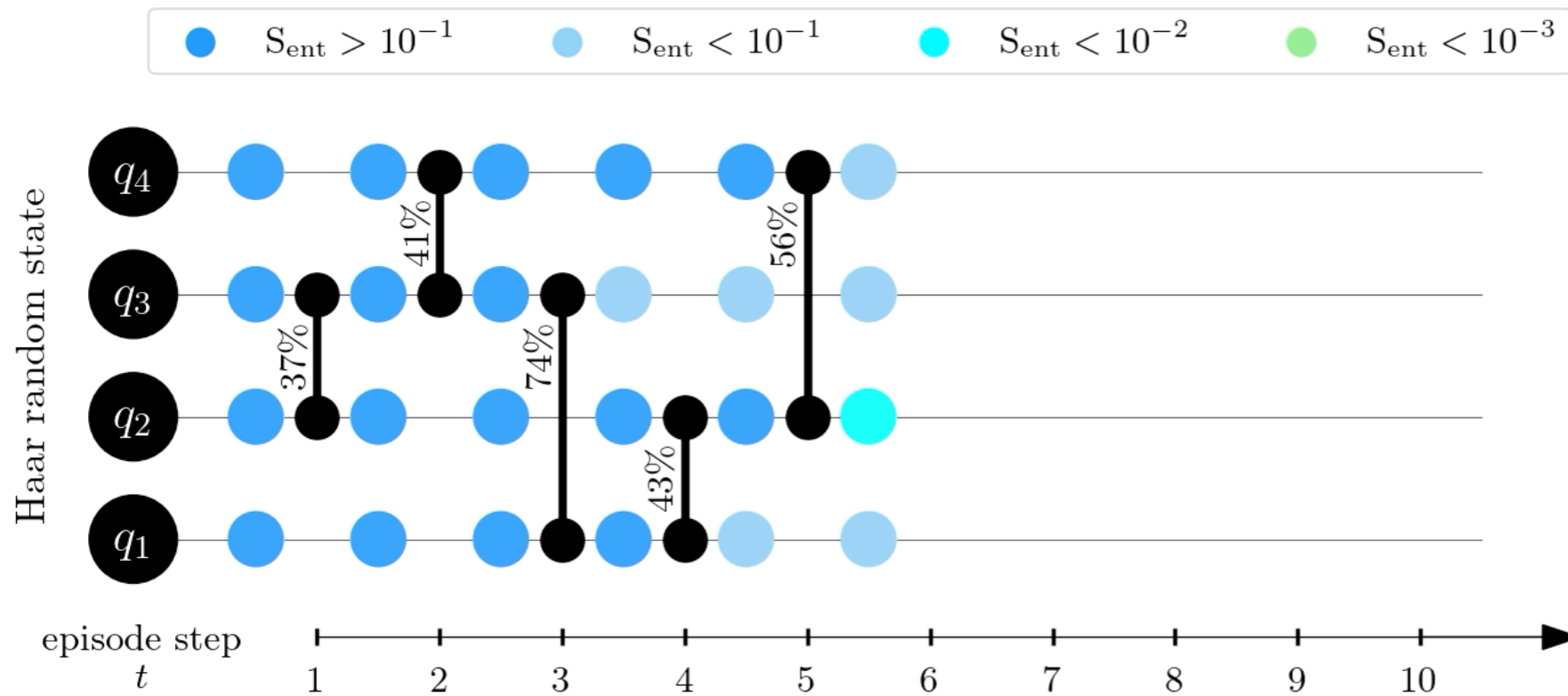




# Applications of RL in Quantum Physics

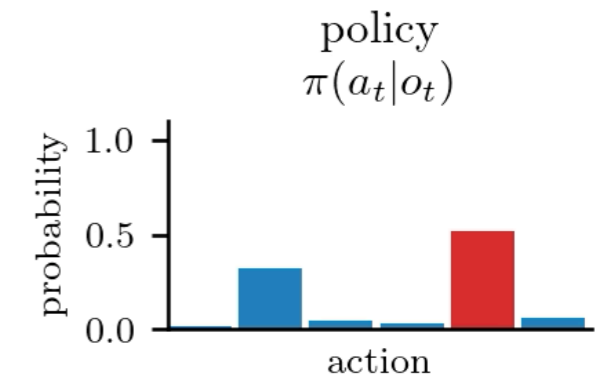
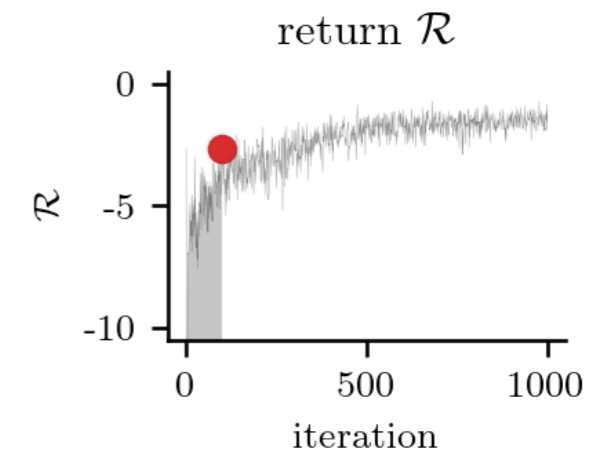
- quantum control
- quantum gate design
- quantum error correction
- quantum circuit design and synthesis
  - construct disentangling circuits

then it starts minimizing the overall number of gates applied..



$$S_{\text{avg}} = \frac{1}{L} \sum_{j=1}^L S_{\text{ent}}[\rho^{(j)}]$$

training iteration = 100

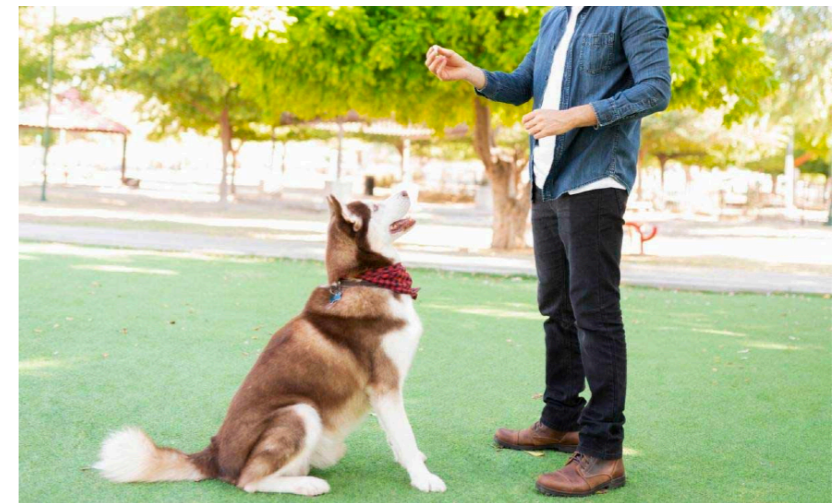




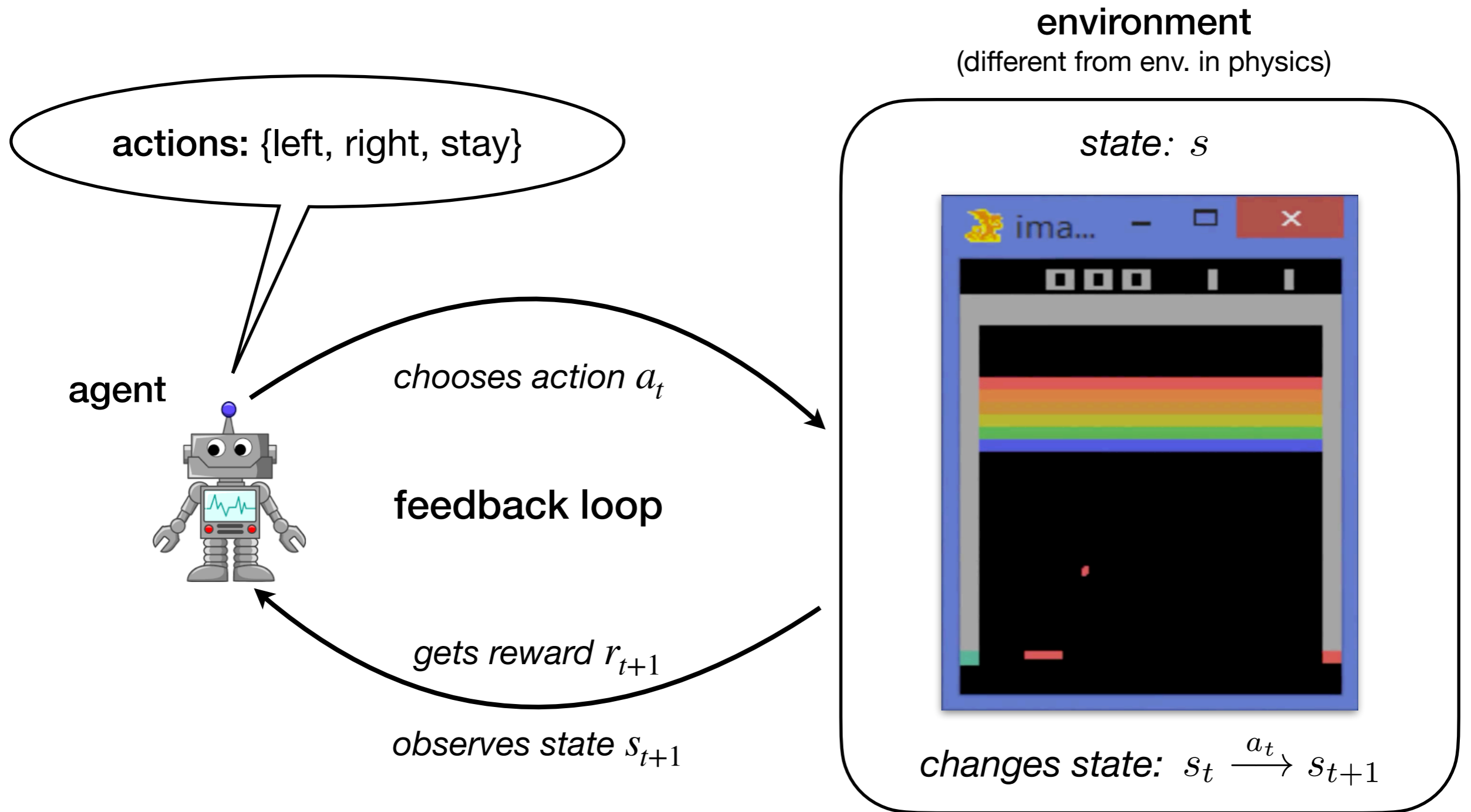
# Outline

## Part 1

- ✓ Reinforcement learning (RL) in quantum physics
  - RL as a branch of machine learning
  
- ✓ Applications of RL
  - hallmark applications of RL
  - applications in quantum technologies
  
- **RL framework in a nutshell**
  - environment, states, actions, rewards
  - RL algorithms



# Reinforcement Learning (RL) formalism



Mnih et al., Nature 518 (2015) [Google DeepMind]

# RL in a Nutshell

- RL formalism

- ▶ action space:  $\mathcal{A} = \{\text{left, stay, right}\}$
- ▶ state space:  $\mathcal{S}$  pixelized image of the screen
- ▶ reward function:  $r = \text{score}$



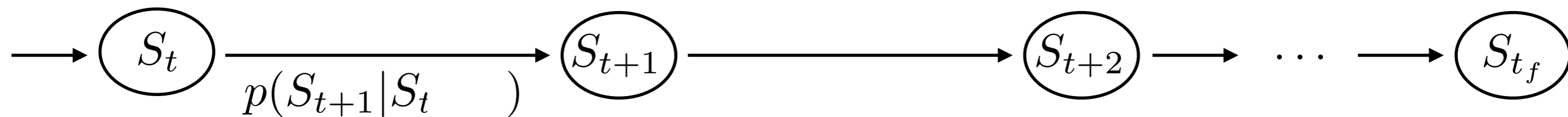
# RL in a Nutshell

- RL formalism

- ▶ action space:  $\mathcal{A} = \{\text{left, stay, right}\}$
- ▶ state space:  $\mathcal{S}$  pixelized image of the screen
- ▶ reward function:  $r = \text{score}$



- RL episode is a Markov decision process



- ▶ transition probability:  $p(S_{t+1} | S_t, \quad)$

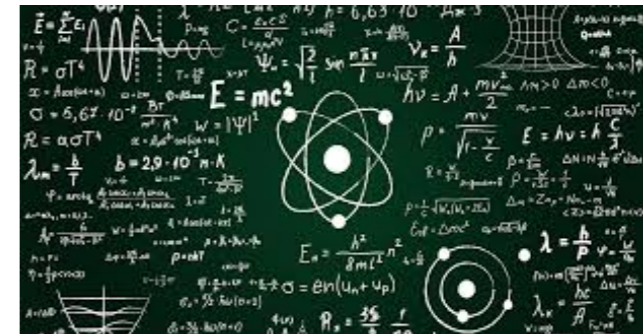


image: Kardashev Scale Wiki

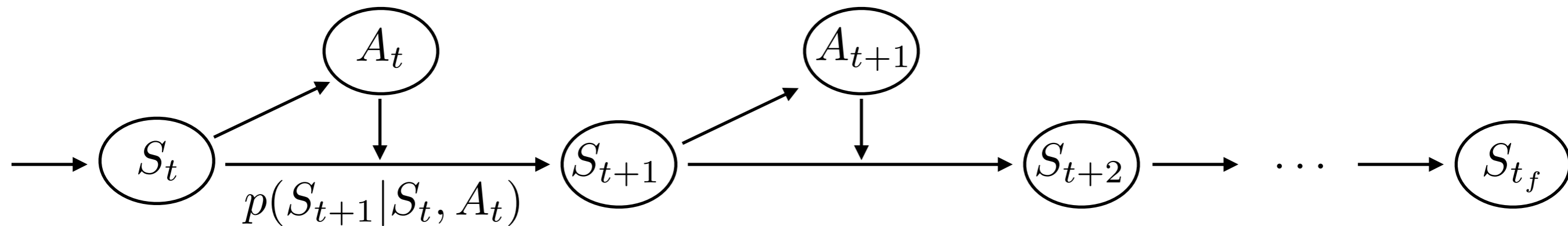
# RL in a Nutshell

- RL formalism

- ▶ action space:  $\mathcal{A} = \{\text{left, stay, right}\}$
- ▶ state space:  $\mathcal{S}$  pixelized image of the screen
- ▶ reward function:  $r = \text{score}$



- RL episode is a Markov **decision** process



- ▶ transition probability:  $p(S_{t+1} | S_t, A_t)$



image: Kardashev Scale Wiki

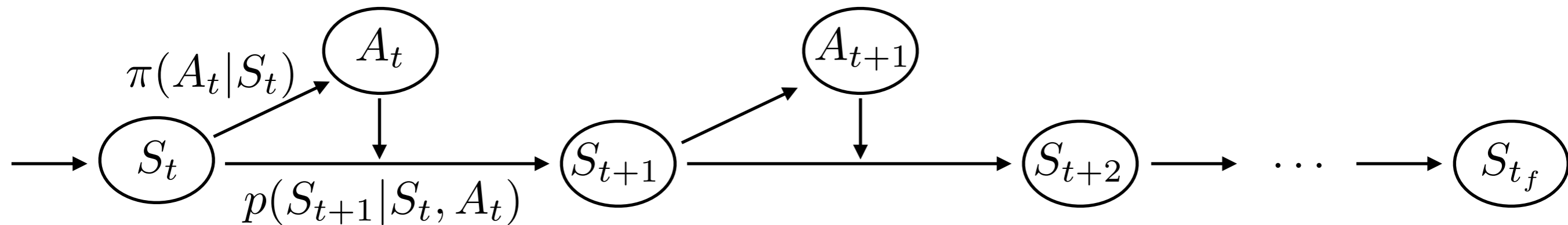
# RL in a Nutshell

- RL formalism

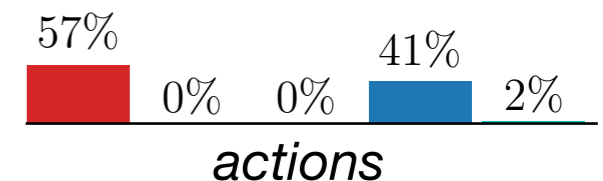
- ▶ action space:  $\mathcal{A} = \{\text{left, stay, right}\}$
- ▶ state space:  $\mathcal{S}$  pixelized image of the screen
- ▶ reward function:  $r = \text{score}$



- RL episode is a Markov **decision** process



- ▶ transition probability:  $p(S_{t+1} | S_t, A_t)$
- ▶ policy:  $\pi(A_t, S_t)$  — probability to take action  $A_t$  in the state  $S_t$



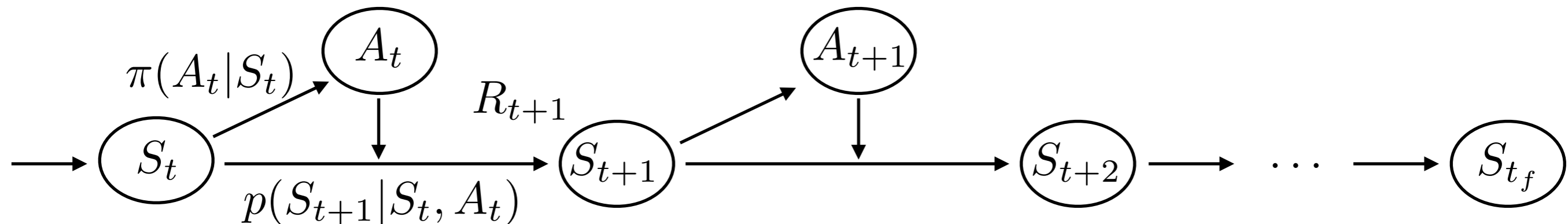
# RL in a Nutshell

- RL formalism

- ▶ action space:  $\mathcal{A} = \{\text{left, stay, right}\}$
- ▶ state space:  $\mathcal{S}$  pixelized image of the screen
- ▶ reward function:  $r = \text{score}$



- RL episode is a Markov **decision** process



- RL **objective**: *find policy* which maximizes the total *expected return*

$$J = \mathbb{E}_{a \sim \pi(a|s)} [R_{t+1} + \dots + R_{t_f} | S_0 = s]$$





# Outline

## Part 1

### ✓ Reinforcement learning (RL) in quantum physics

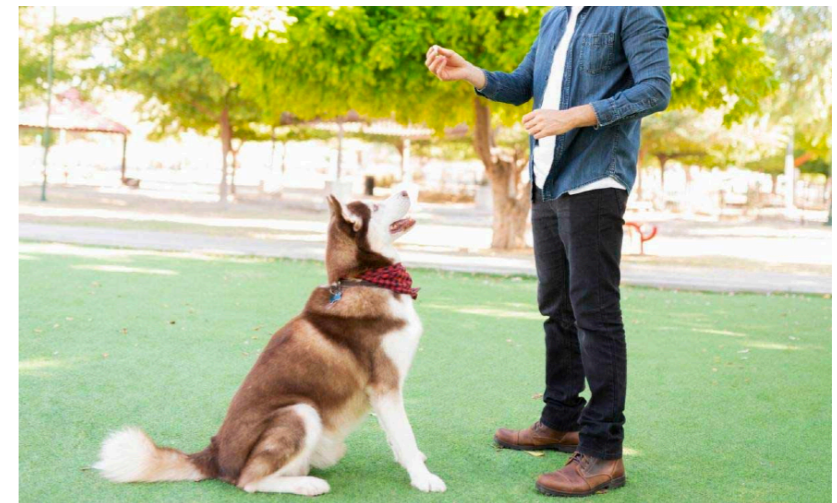
- RL as a branch of machine learning

### ✓ Applications of RL

- hallmark applications of RL
- applications in quantum technologies

### • RL framework in a nutshell

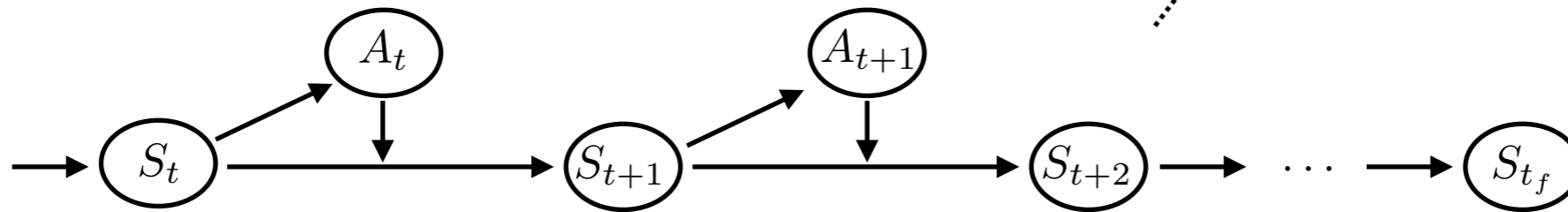
- environment, states, actions, rewards
- RL algorithms



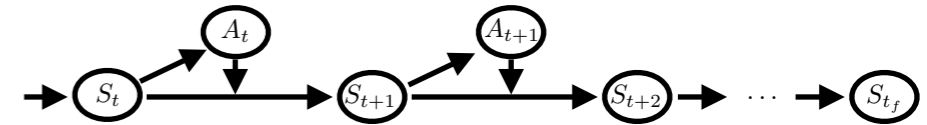
# Policy Gradient

- *in practice: evaluate* RL objective by sampling trajectories

$$J = \mathbb{E}_{a \sim \pi(a|s)} \left[ \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \sum_t R_t(\tau_j)$$



# Policy Gradient



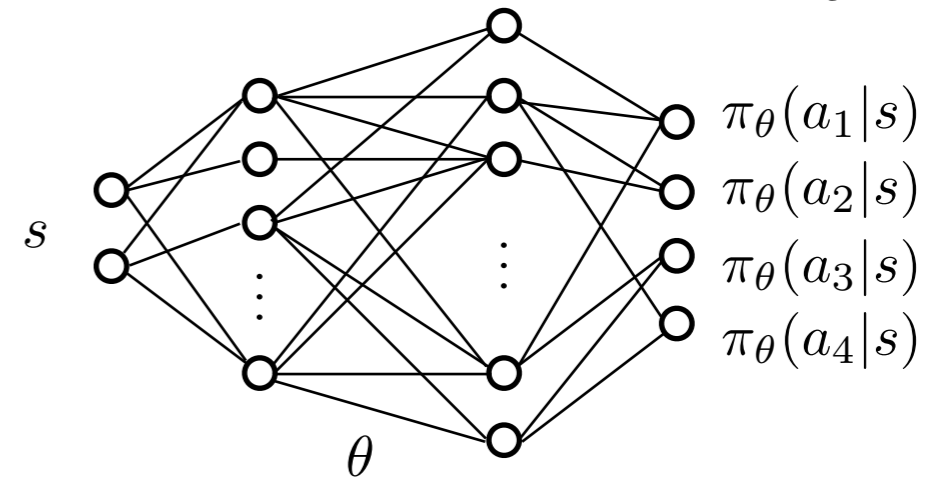
- *in practice: evaluate* RL objective by sampling trajectories

$$J = \mathbb{E}_{a \sim \pi(a|s)} \left[ \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \sum_t R_t(\tau_j)$$

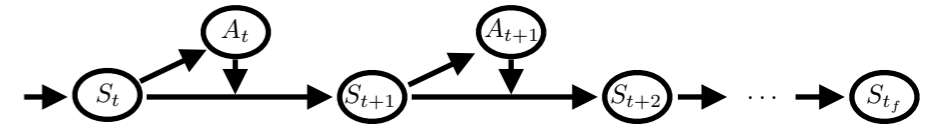
57%	0%	0%	41%	2%
$a_1$	$a_2$	$a_3$	$a_4$	...

- **improve policy**

→ parametrize policy  $\pi(a|s) \approx \pi_\theta(a|s)$



# Policy Gradient



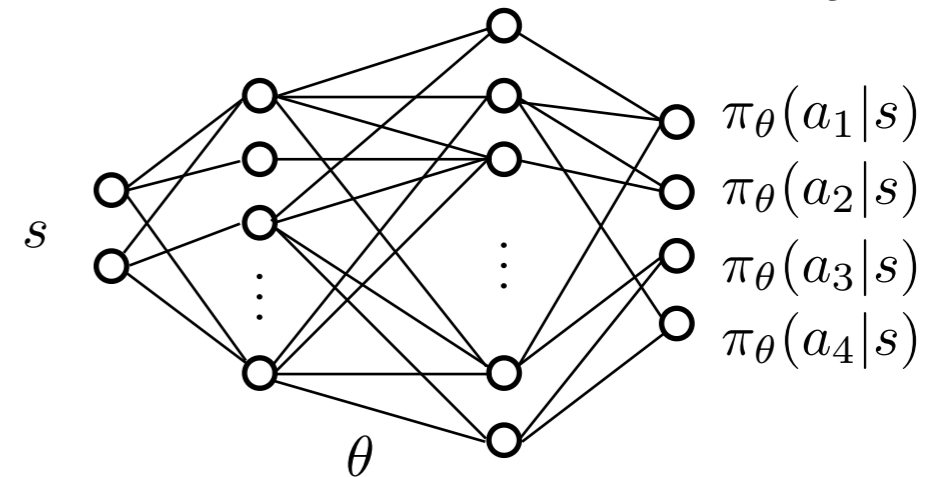
- *in practice: evaluate* RL objective by sampling trajectories

$$J = \mathbb{E}_{a \sim \pi(a|s)} \left[ \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \sum_t R_t(\tau_j)$$

A bar chart showing the distribution of actions  $a_1, a_2, a_3, a_4, \dots$  sampled from the policy. The probabilities are:  $a_1$  (57%),  $a_2$  (0%),  $a_3$  (0%),  $a_4$  (41%), and others (2%).

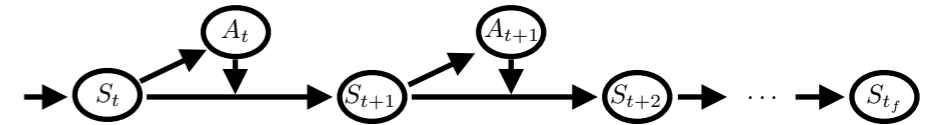
- **improve policy**

- parametrize policy  $\pi(a|s) \approx \pi_\theta(a|s)$
- compute gradients



$$\nabla_\theta J(\theta) = \mathbb{E}_{a \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$

# Policy Gradient



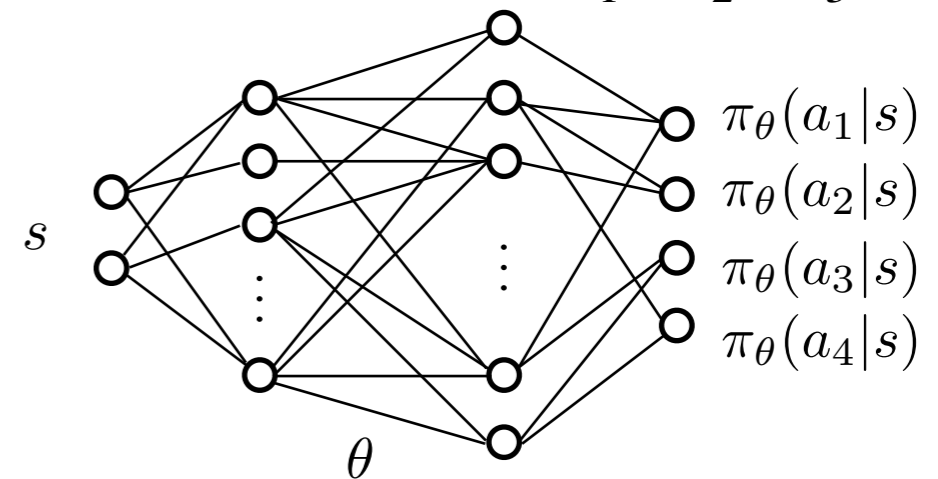
- *in practice: evaluate* RL objective by sampling trajectories

$$J = \mathbb{E}_{a \sim \pi(a|s)} \left[ \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \sum_t R_t(\tau_j)$$

A bar chart showing the distribution of actions  $a_1, a_2, a_3, a_4, \dots$  sampled from the policy. The probabilities are:  $a_1$  (57%),  $a_2$  (0%),  $a_3$  (0%),  $a_4$  (41%), and others (2%).

- **improve policy**

- parametrize policy  $\pi(a|s) \approx \pi_\theta(a|s)$
- compute gradients



$$\nabla_\theta J(\theta) = \mathbb{E}_{a \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$

- update parameters  $\theta$  to maximize return

$$\theta_{\text{new}} = \theta_{\text{old}} + \alpha \nabla_\theta J(\theta)$$

*gradient ascent*

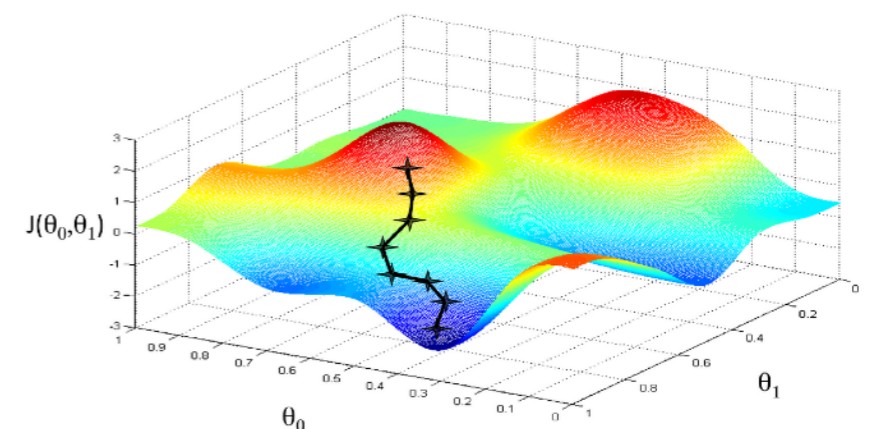
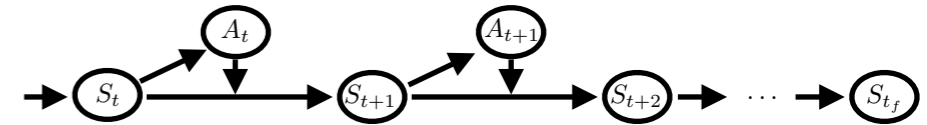


image: medium.com

# Policy Gradient



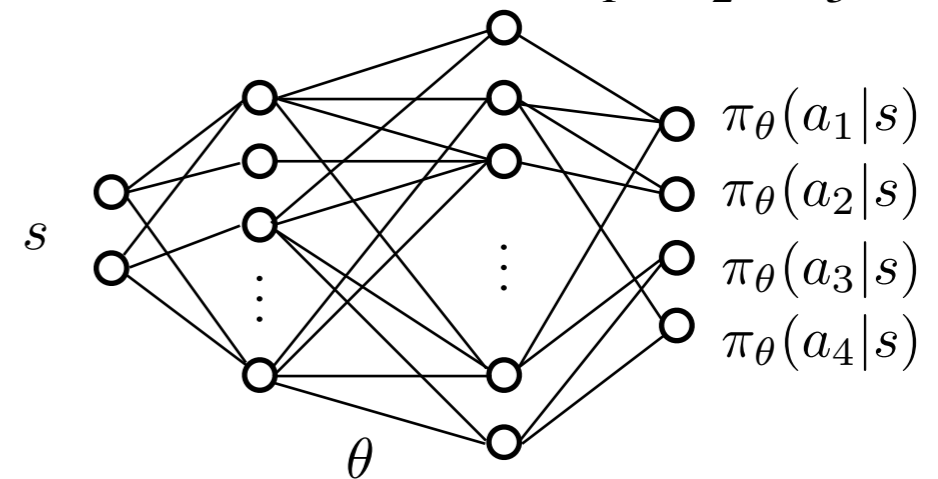
- *in practice: evaluate* RL objective by sampling trajectories

$$J = \mathbb{E}_{a \sim \pi(a|s)} \left[ \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \sum_t R_t(\tau_j)$$

A bar chart showing the distribution of actions  $a_1, a_2, a_3, a_4, \dots$  sampled from the policy. The probabilities are:  $a_1$  (57%),  $a_2$  (0%),  $a_3$  (0%),  $a_4$  (41%), and others (2%).

- **improve policy**

- parametrize policy  $\pi(a|s) \approx \pi_\theta(a|s)$
- compute gradients



$$\nabla_\theta J(\theta) = \mathbb{E}_{a \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta \sum_t R_t \mid S_0 = s \right] \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$

- update parameters  $\theta$  to maximize return

$$\theta_{\text{new}} = \theta_{\text{old}} + \alpha \nabla_\theta J(\theta)$$

- pseudo-loss function  $\tilde{J}(\theta)$   $\nabla_\theta \tilde{J}(\theta) = \nabla_\theta J_{\text{MC}}(\theta)$

$$\tilde{J}(\theta) = \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$

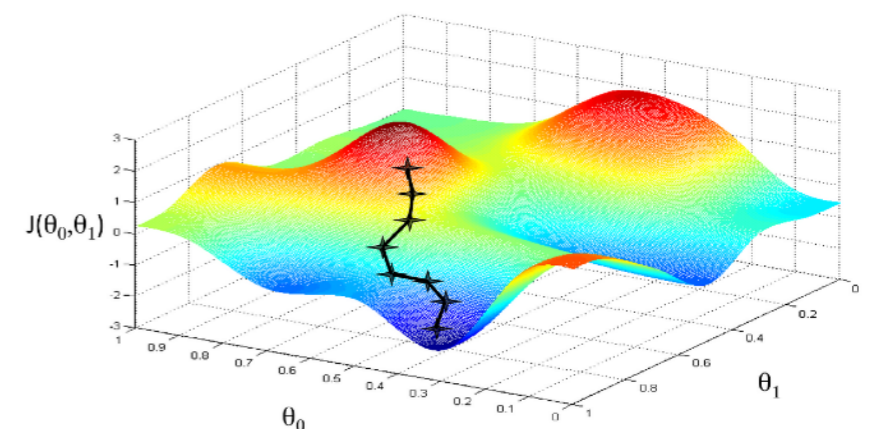


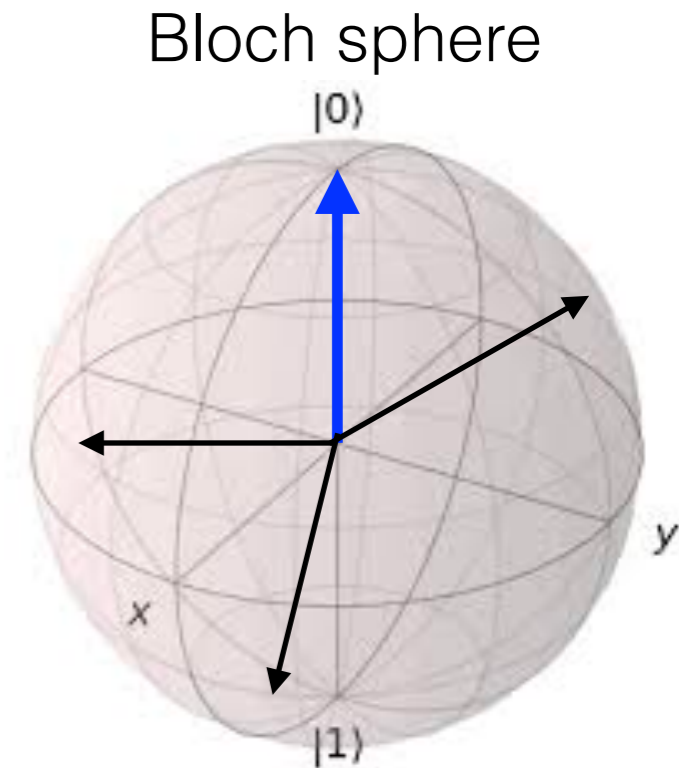
image: medium.com

# Example: RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$



# RL for Quantum State Initialization

→ **quantum control:** two-level system

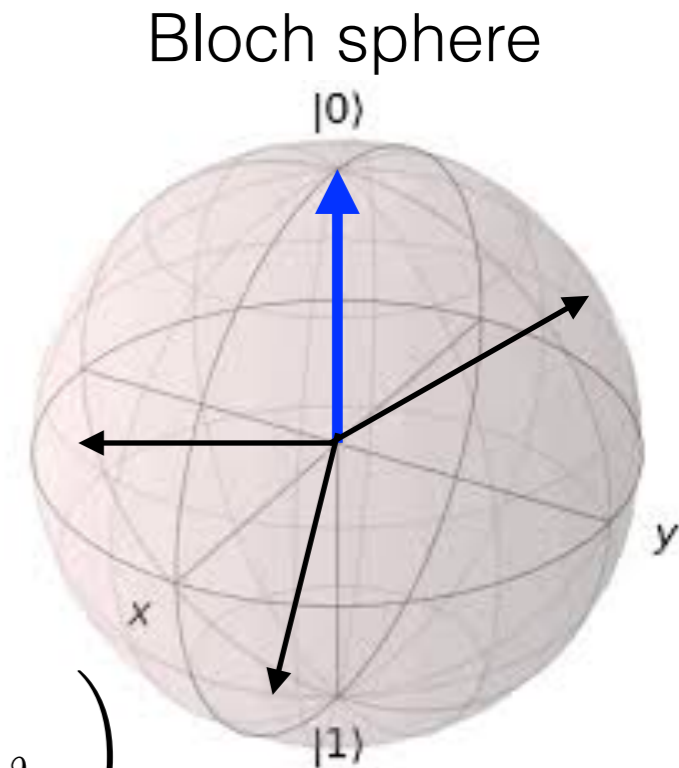
- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$





# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

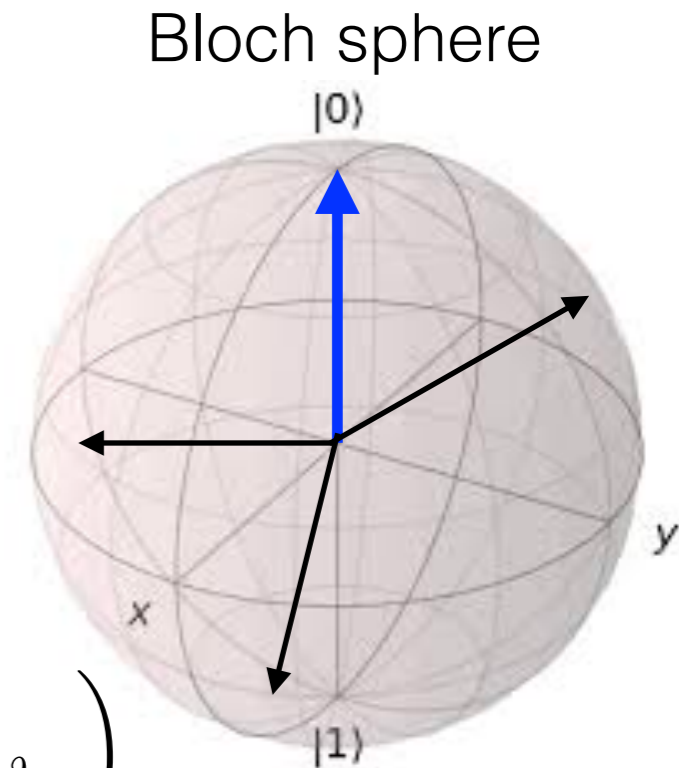
$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$   $|\psi'\rangle = U_\alpha|\psi\rangle$   $s \mapsto s'$



# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

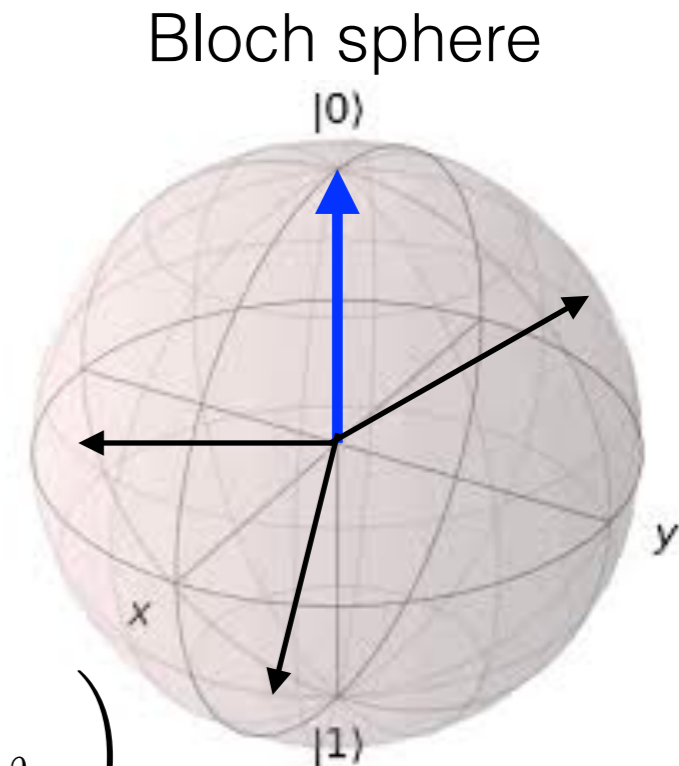
$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$        $|\psi'\rangle = U_\alpha|\psi\rangle$        $s \mapsto s'$
- define reward:  $r_t = |\langle 0|\psi_t\rangle|^2$



# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

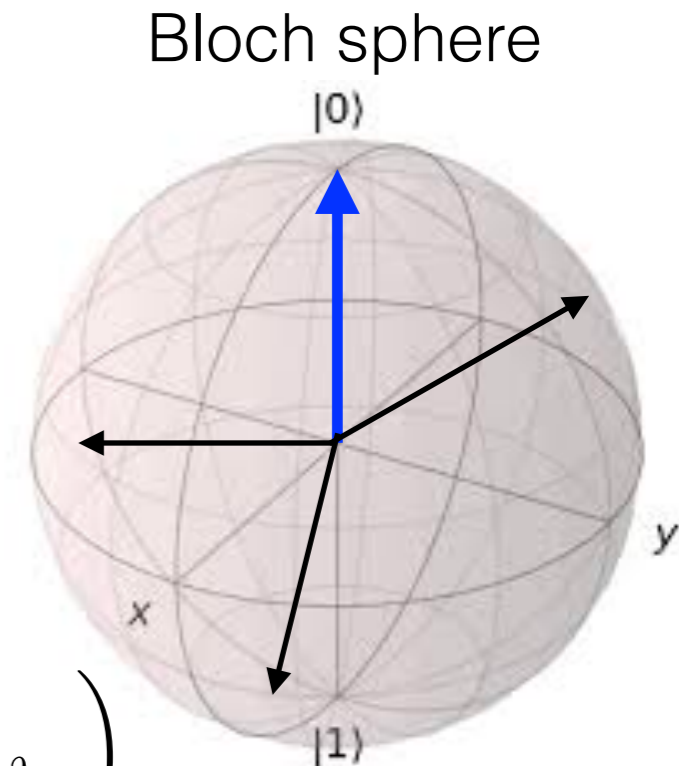
- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$   $|\psi'\rangle = U_\alpha|\psi\rangle$   $s \mapsto s'$

- define reward:  $r_t = |\langle 0|\psi_t\rangle|^2$

- choose **RL algorithm:** policy gradient

simulate trial trajectories  $\tau_j$

$$\nabla_\theta J(\theta) \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$



# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

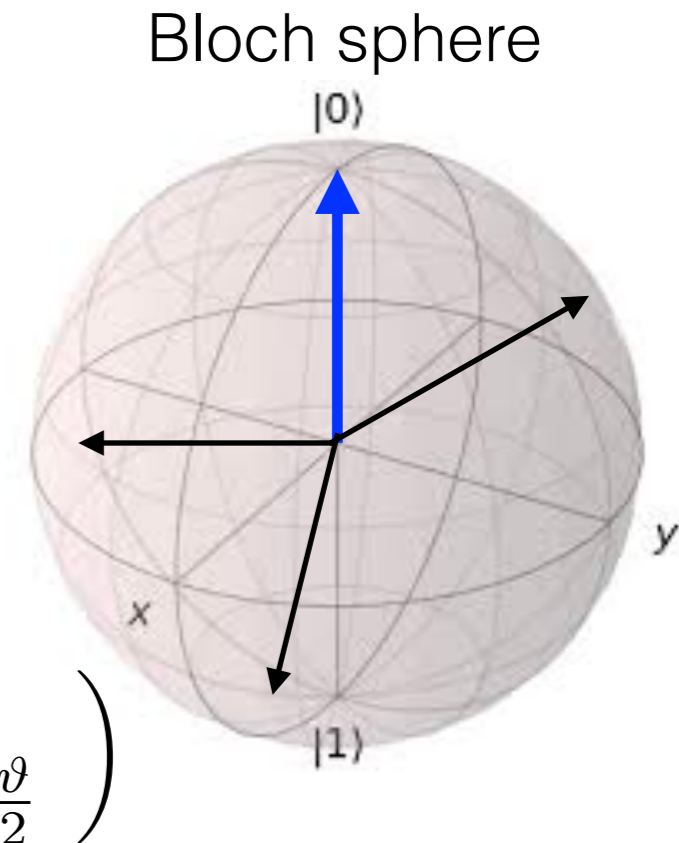
$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$        $|\psi'\rangle = U_\alpha|\psi\rangle$        $s \mapsto s'$
- define reward:  $r_t = |\langle 0|\psi_t\rangle|^2$



**problem:** *continuous* state space has infinitely many configurations

# RL with Function Approximation

→ **problem:** state space has exponentially/continuously many configurations  $|\mathcal{A}|^{N_T}$

- can we estimate values of not yet encountered states?

# RL with Function Approximation

→ **problem:** state space has exponentially/continuously many configurations  $|\mathcal{A}|^{N_T}$

- can we estimate values of not yet encountered states?

→ YES, via inter- & extrapolation: parametrize the Q-function/policy

$$\pi(a|s) \rightarrow \pi_{\theta}(a|s)$$

variational parameters  $\theta$



# RL with Function Approximation

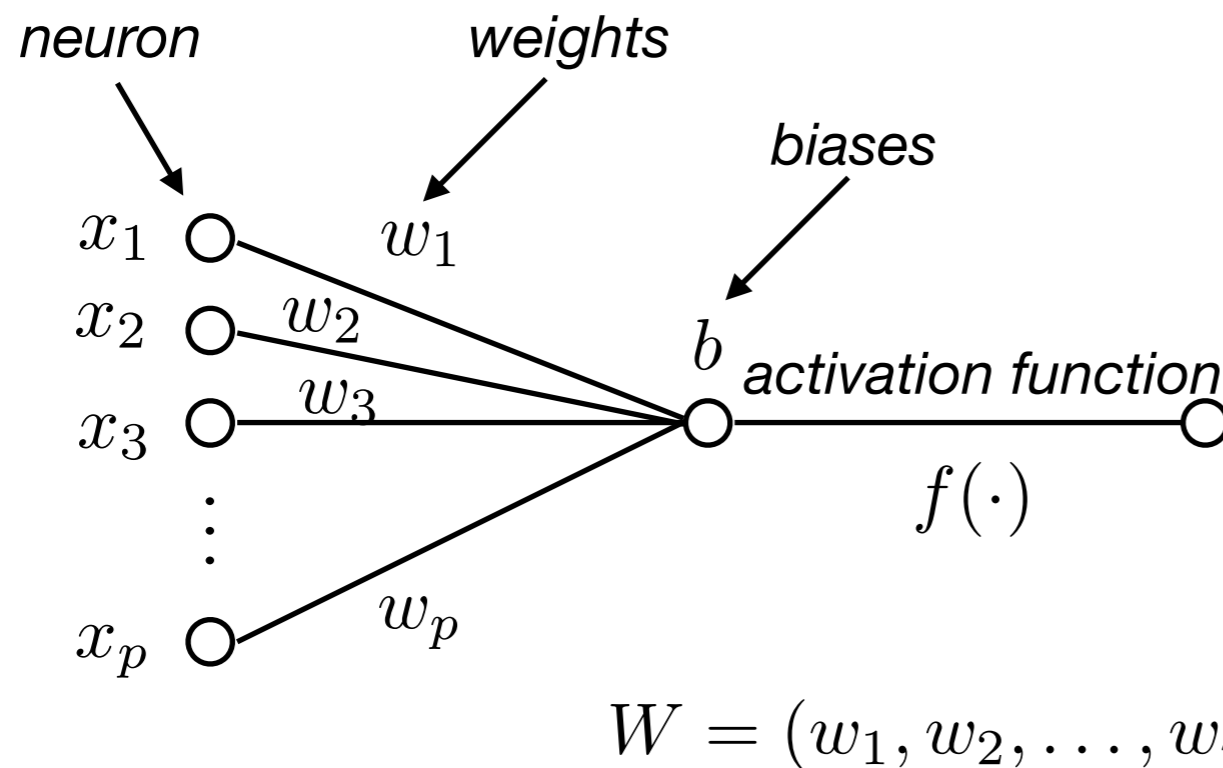
→ **problem:** state space has exponentially/continuously many configurations  $|\mathcal{A}|^{N_T}$

- can we estimate values of not yet encountered states?

→ YES, via inter- & extrapolation: parametrize the Q-function/policy

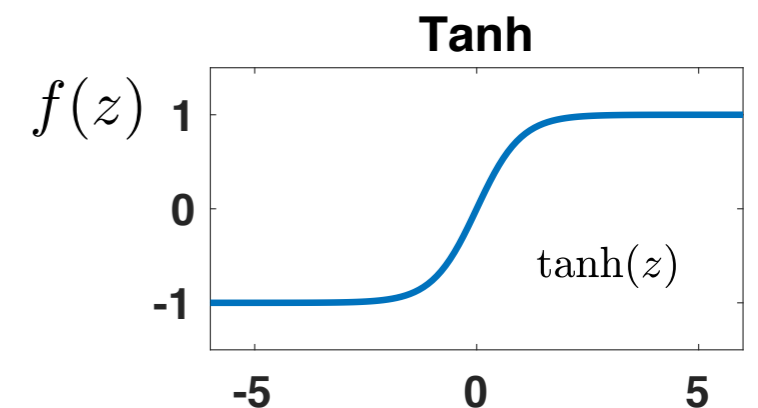
$$\pi(a|s) \rightarrow \pi_{\theta}(a|s)$$

- typical approach: use deep neural network (**Deep RL**)



$$y = f(W \cdot x + b)$$

$$\theta = \{W, b\}$$



# RL with Function Approximation

→ **problem:** state space has exponentially/continuously many configurations  $|\mathcal{A}|^{N_T}$

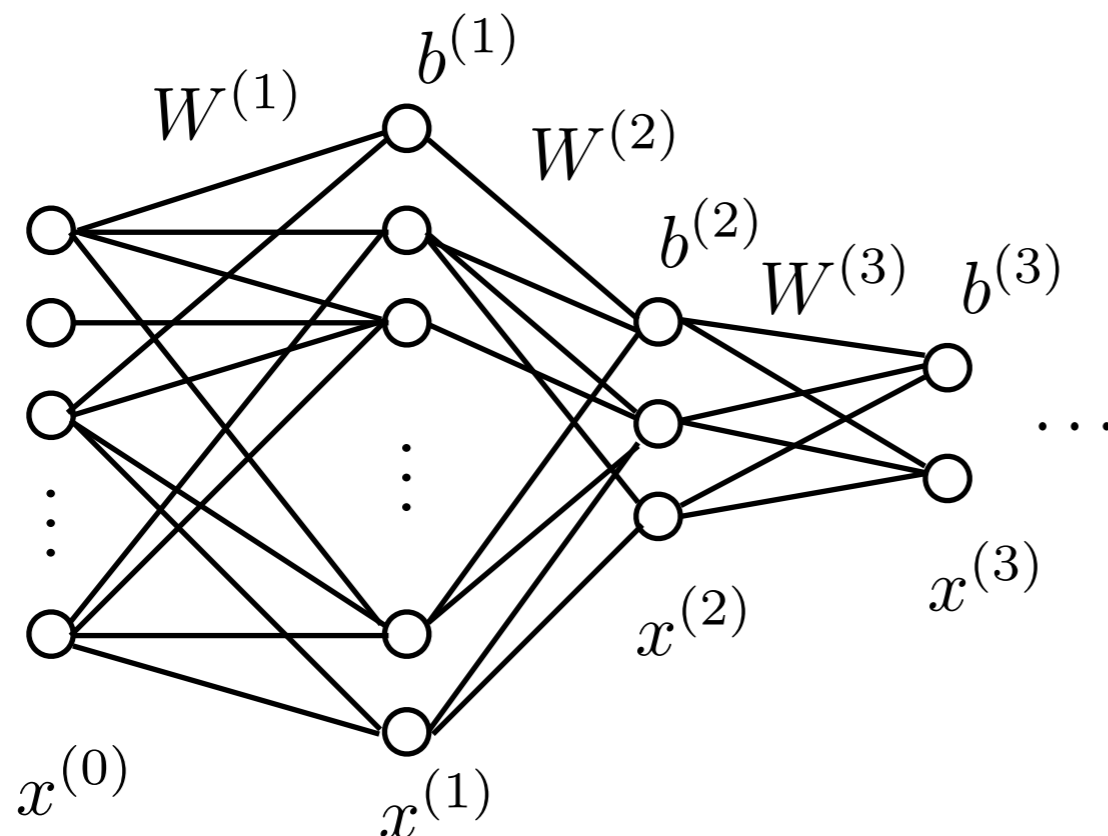
- can we estimate values of not yet encountered states?

→ YES, via inter- & extrapolation: parametrize the Q-function/policy

$$\pi(a|s) \rightarrow \pi_{\theta}(a|s)$$

$$\theta = \{W, b\}$$

- typical approach: use deep neural network (**Deep RL**)



$x^{(0)}$  input layer

$$x_i^{(1)} = f^{(1)} \left( W_{ij}^{(1)} x_j^{(0)} + b_i^{(1)} \right)$$

$$x_i^{(2)} = f^{(2)} \left( W_{ij}^{(2)} x_j^{(1)} + b_i^{(2)} \right)$$

$b_i^{(l)}$  : bias vector of layer  $l$

$W_{ij}^{(l)}$  : weight matrix of layer  $l$

$f^{(l)}$  : activation function of layer  $l$



# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

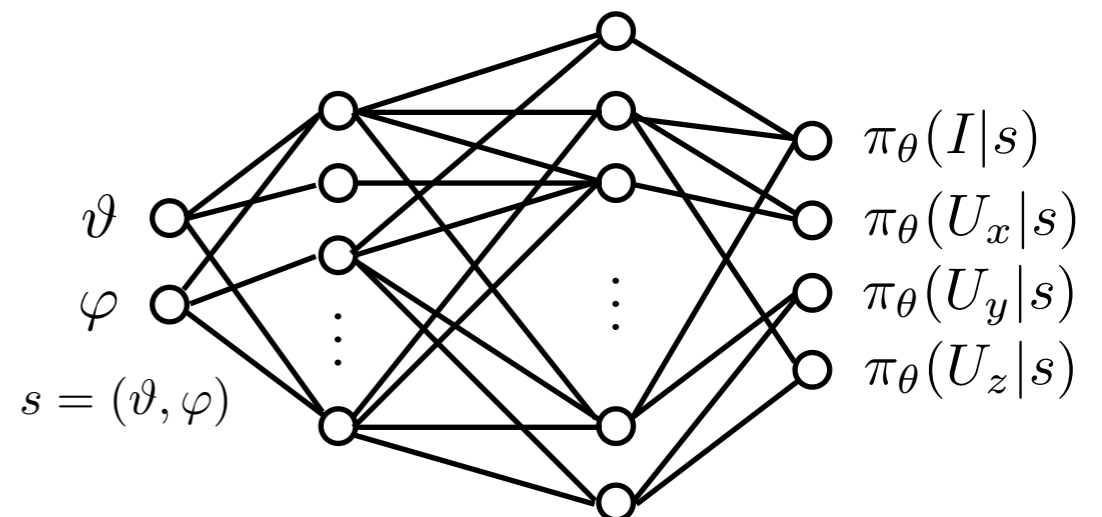
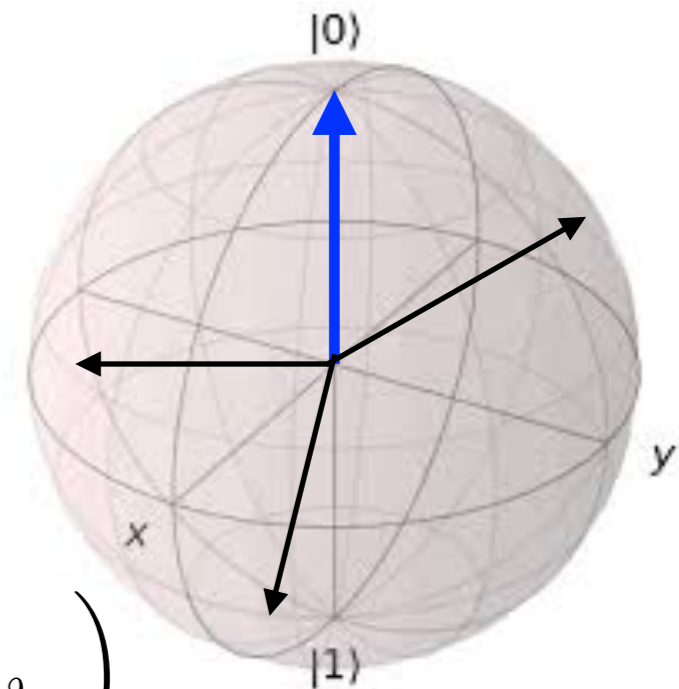
- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$

- define reward:  $r_t = |\langle 0 | \psi_t \rangle|^2$

- choose **RL algorithm:** policy gradient

$$\pi_\theta(a|s)$$

Bloch sphere



# RL for Quantum State Initialization

→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos \frac{\vartheta}{2} \\ e^{i\varphi} \sin \frac{\vartheta}{2} \end{pmatrix}$

RL state space:  $\mathcal{S} = \{(\vartheta, \varphi) | \vartheta \in [0, \pi], \varphi \in [0, 2\pi)\}$

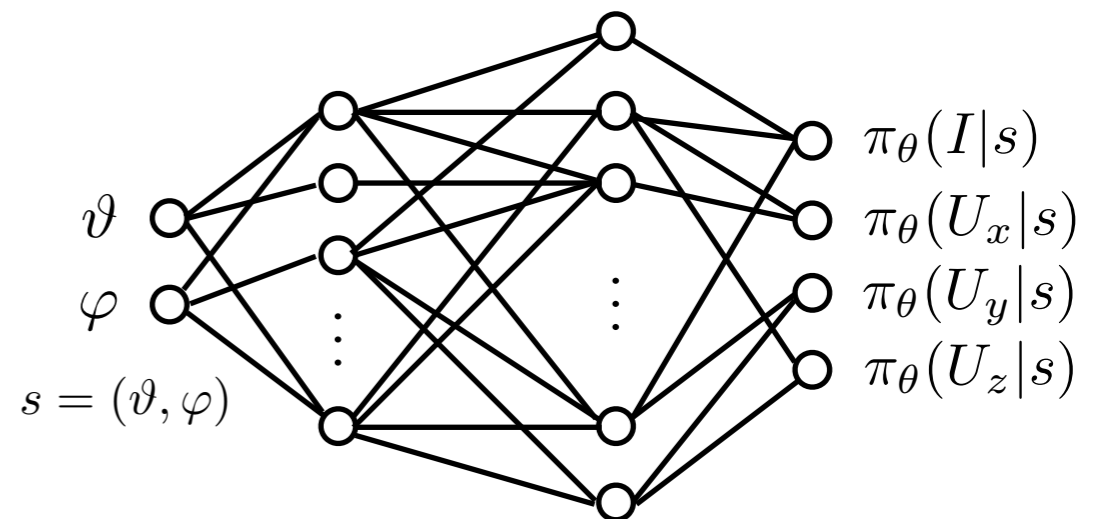
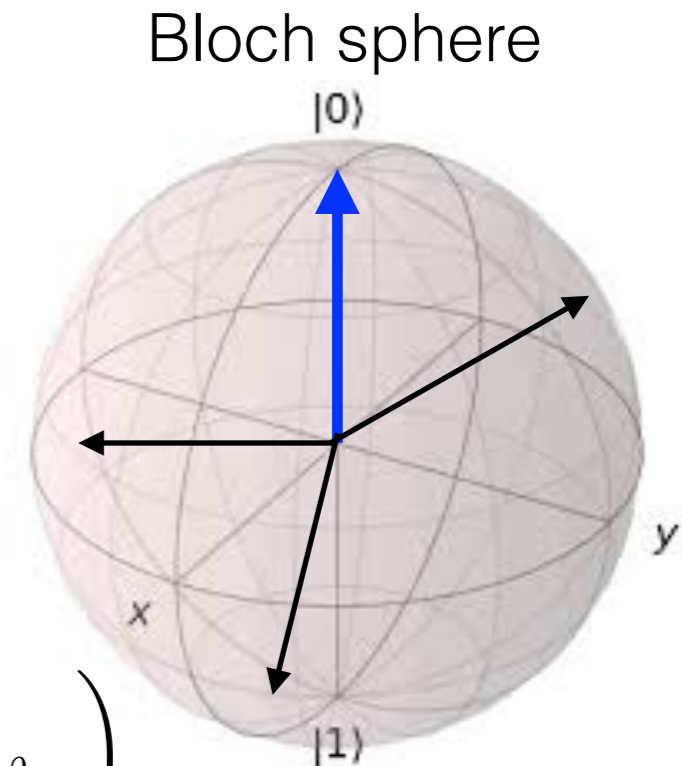
- define **RL actions:**  $\mathcal{A} = \{I, U_\alpha | \alpha = x, y, z\}$

- define reward:  $r_t = |\langle 0 | \psi_t \rangle|^2$

- choose **RL algorithm:** policy gradient

simulate trial trajectories  $\tau_j$

$$\nabla_\theta J(\theta) \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$



# RL for Quantum State Initialization

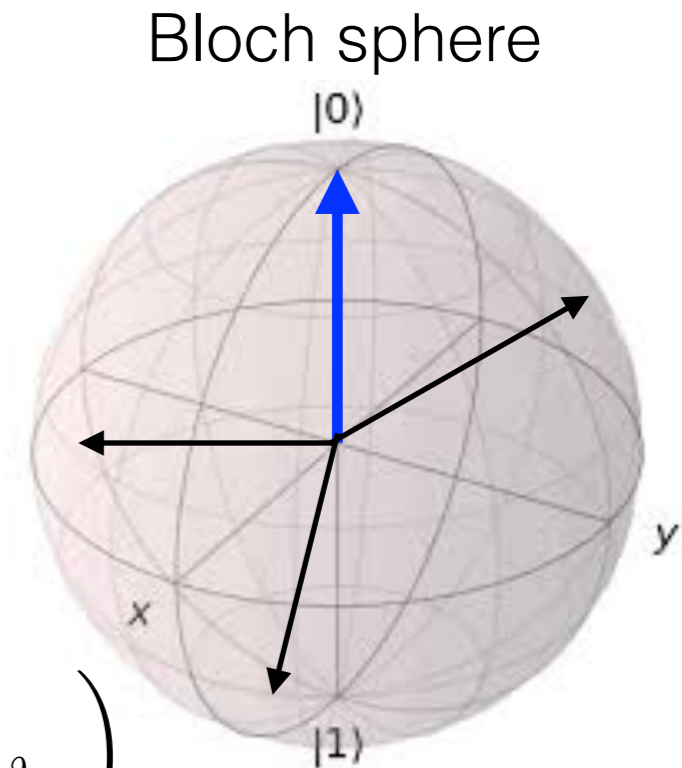
→ **quantum control:** two-level system

- **task:** prepare  $|0\rangle$  using infinitesimal rotations

$$U_\alpha = \exp\left(-i\frac{\delta t}{2}\sigma^\alpha\right) \quad \alpha = x, y, z$$

- define **RL states:**

parametrization of state of system:  $|\psi\rangle = \begin{pmatrix} \cos\frac{\vartheta}{2} \\ e^{i\varphi}\sin\frac{\vartheta}{2} \end{pmatrix}$



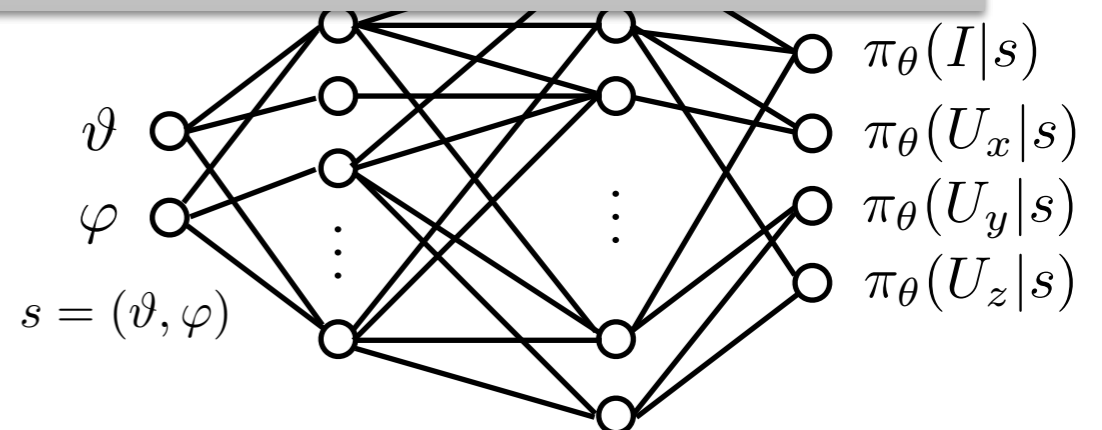
Check out Jupyter notebook for how this works in practice!

[https://github.com/mgbukov/RL\\_quantum](https://github.com/mgbukov/RL_quantum)



- choose **RL algorithm:** policy gradient  
simulate trial trajectories  $\tau_j$

$$\nabla_\theta J(\theta) \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \nabla_\theta \log \pi_\theta(\tau_j) \sum_t R_t(\tau_j)$$



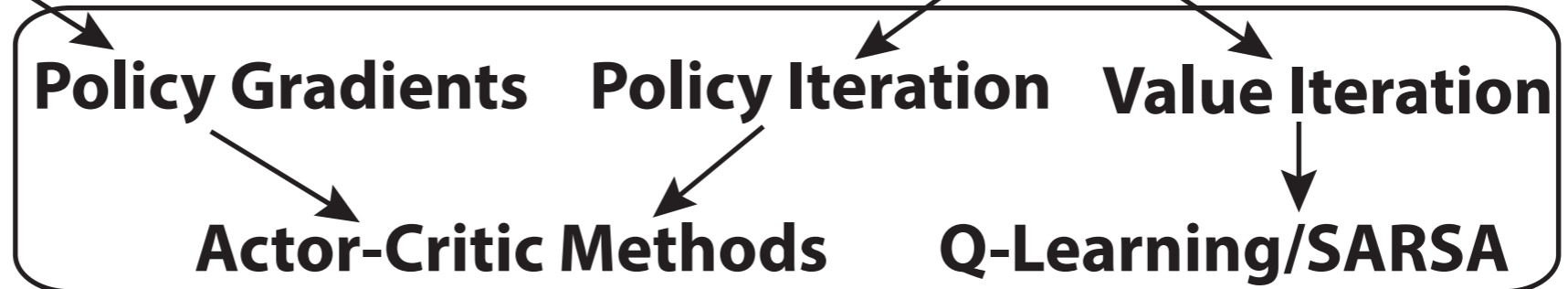
# What other RL Algorithms are there?

→ overview of RL algorithms

*Policy Optimization*

↓  
Evolutionary  
Methods

*Dynamic Programming*



**Policy Gradients**

**Policy Iteration**

**Value Iteration**

**Actor-Critic Methods**

**Q-Learning/SARSA**

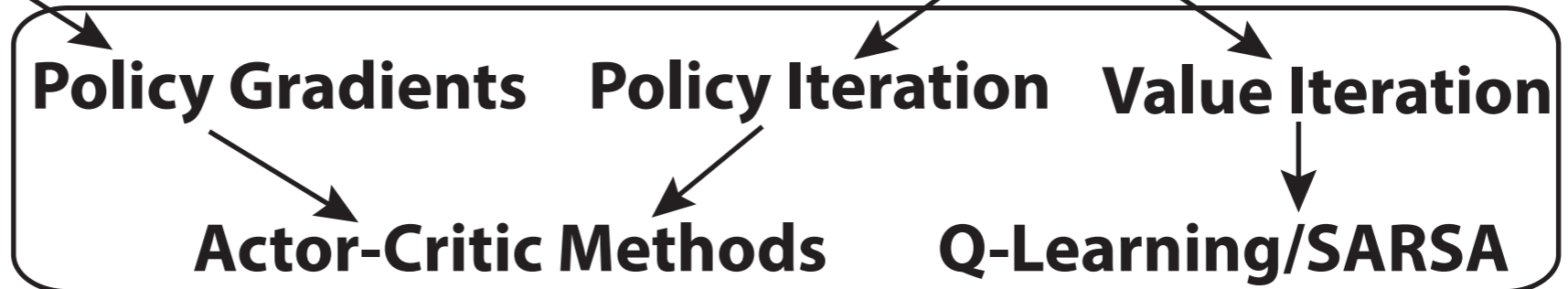
# What other RL Algorithms are there?

→ overview of RL algorithms

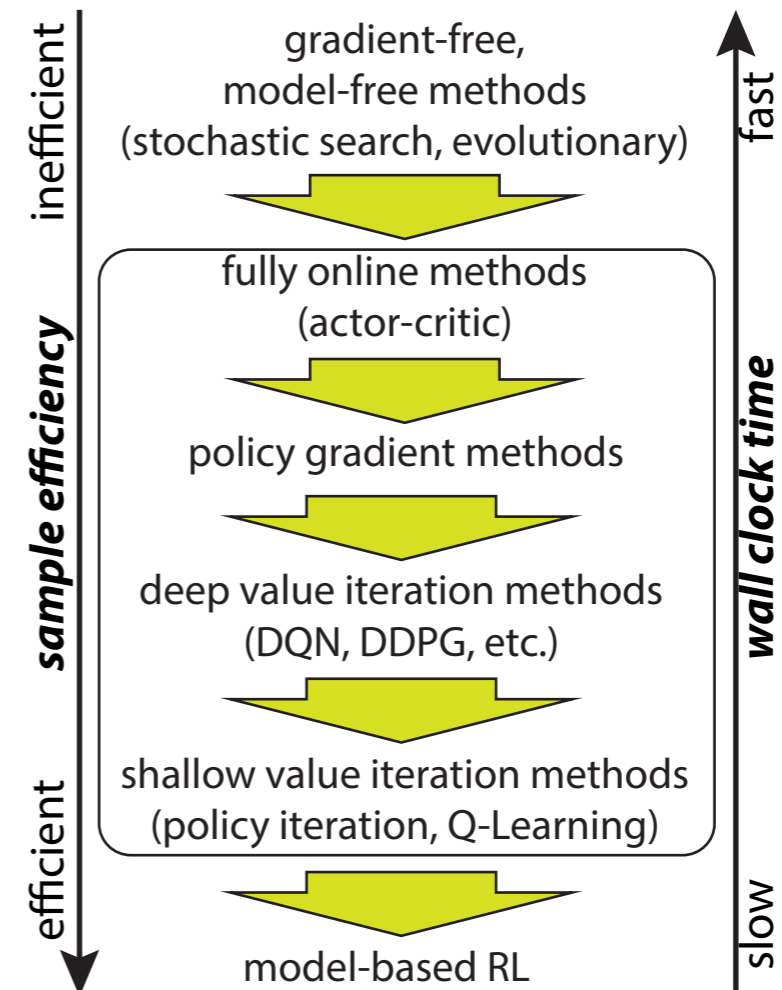
*Policy Optimization*

*Dynamic Programming*

Evolutionary  
Methods



→ which algorithm to use?



# Value function methods



## → Value Iteration methods

- value function: **expected** total return under the policy  $\pi(a|s)$  from state  $s$

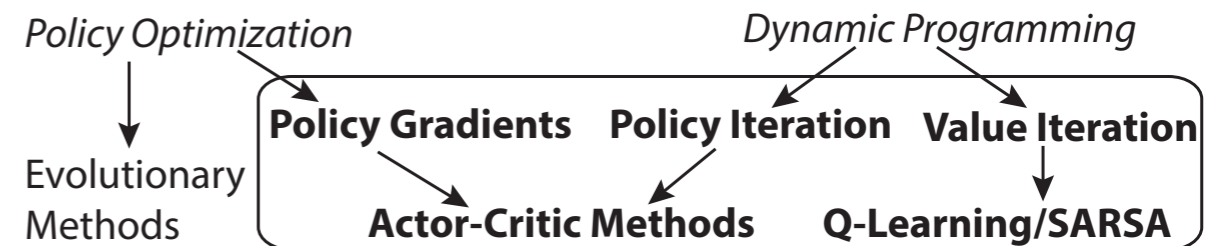
$$v_{\pi}(s) = \mathbb{E}_{a \sim \pi(a|s)} [G_t | S_t = s]$$

$$G_t = R_{t+1} + G_{t+1}$$

**problem:** cannot reconstruct the policy



# Value function methods



## → Value Iteration methods

- value function: **expected** total return under the policy  $\pi(a|s)$  from state  $s$

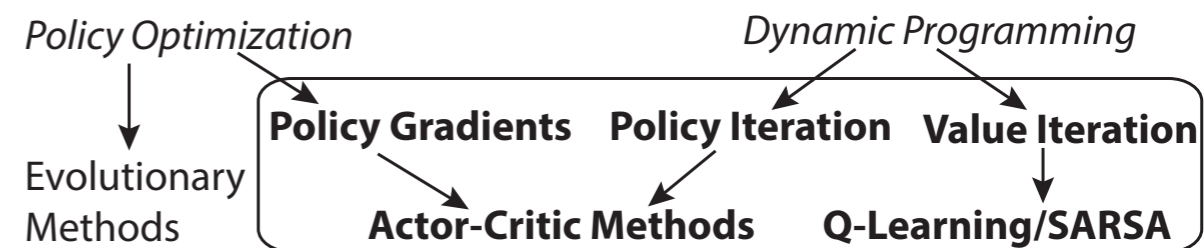
$$v_{\pi}(s) = \mathbb{E}_{a \sim \pi(a|s)} [G_t | S_t = s] \qquad G_t = R_{t+1} + \gamma v_{\pi}(S_{t+1})$$

- action-value (or Q-) function: **expected** total return under the policy  $\pi(a|s)$  starting from state  $s$  and taking action  $a$ :

$$Q_{\pi}(s, a) = \mathbb{E}_{a \sim \pi(a|s)} [G_t | S_t = s, A_t = a]$$



# Value function methods



## → Value Iteration methods

- value function: **expected** total return under the policy  $\pi(a|s)$  from state  $s$

$$v_{\pi}(s) = \mathbb{E}_{a \sim \pi(a|s)} [G_t | S_t = s] \qquad G_t = R_{t+1} + \gamma v_{\pi}(S_{t+1})$$

- action-value (or Q-) function: **expected** total return under the policy  $\pi(a|s)$  starting from state  $s$  and taking action  $a$ :

$$Q_{\pi}(s, a) = \mathbb{E}_{a \sim \pi(a|s)} [G_t | S_t = s, A_t = a]$$

- ## → optimal action-value function:
- $$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$
- $$\pi_*(a|s) = \operatorname{argmax}_a Q_*(s, a)$$

Bellman's equation: 
$$Q_*(s, a) = \sum_{s'} p(s'|s, a) \left[ r(s, s', a) + \gamma \max_{a'} Q_*(s', a') \right]$$



# Meaning of Q-function

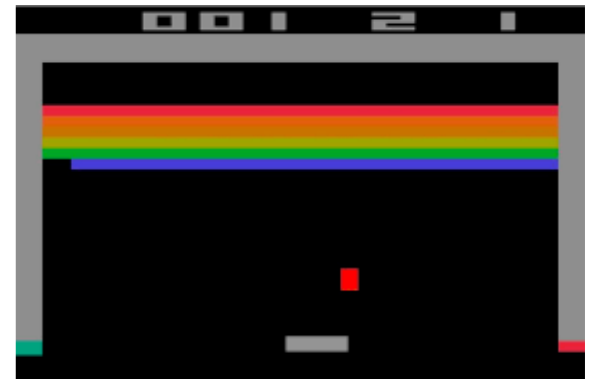
→ assign a value to each state

- first step deterministic, then follow policy

$$Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$$



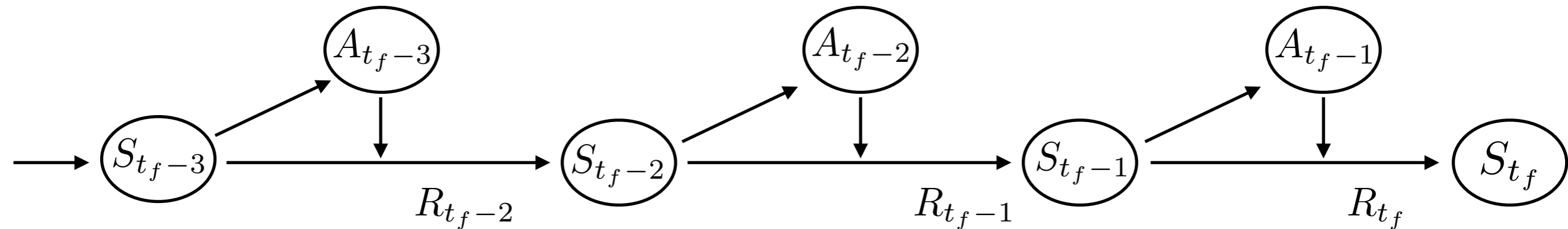
# Meaning of Q-function



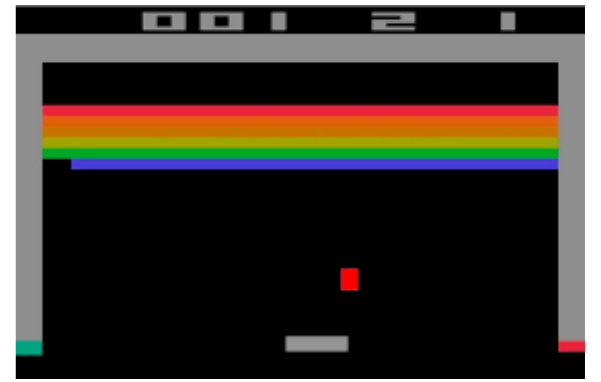
→ assign a value to each state

- first step deterministic, then follow policy

$$Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$$



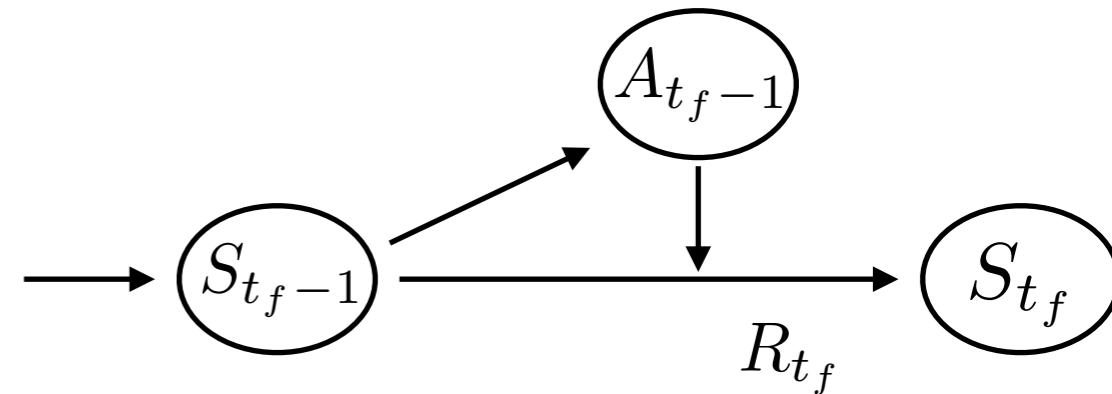
# Meaning of Q-function



→ assign a value to each state

- first step deterministic, then follow policy

$$Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$$



$$Q(S_{t_f-1}, A_{t_f-1}) = R_{t_f}$$

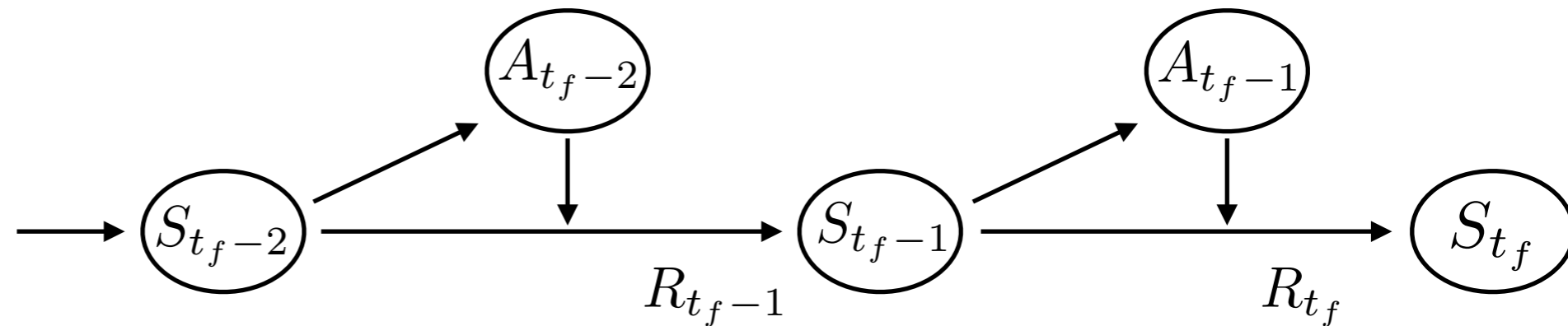
# Meaning of Q-function



→ assign a value to each state

- first step deterministic, then follow policy

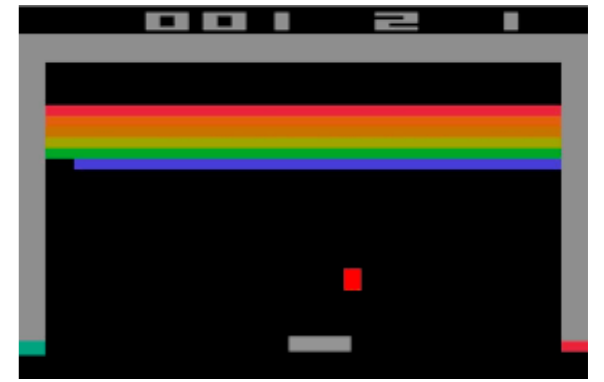
$$Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$$



$$Q(S_{t_f-1}, A_{t_f-1}) = R_{t_f}$$

$$Q(S_{t_f-2}, A_{t_f-2}) = R_{t_f-1} + R_{t_f}$$

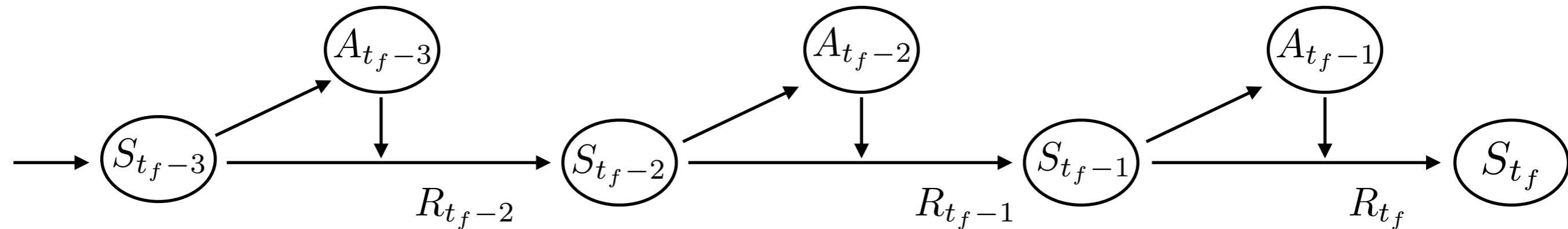
# Meaning of Q-function



→ assign a value to each state

- first step deterministic, then follow policy

$$Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$$



$$Q(S_{t_f-1}, A_{t_f-1}) = R_{t_f}$$

$$Q(S_{t_f-2}, A_{t_f-2}) = R_{t_f-1} + R_{t_f}$$

$$Q(S_{t_f-3}, A_{t_f-3}) = R_{t_f-2} + R_{t_f-1} + R_{t_f}$$

# Q-Learning

→ assign a value to each state  $Q(s, a) = \mathbb{E}_{a \sim \pi}[R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$

- first step deterministic, then follow policy

$Q(s, a)$	$A_1$	$A_2$	$A_3$
$S_1$	1.5	0.5	2.1
$S_2$	1.3	4.4	0.2
$S_3$	0.9	0.9	3.3

# Q-Learning

→ assign a value to each state  $Q(s, a) = \mathbb{E}_{a \sim \pi}[R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$

- first step deterministic, then follow policy

$Q(s, a)$	$A_1$	$A_2$	$A_3$
$S_1$	1.5	0.5	2.1
$S_2$	1.3	4.4	0.2
$S_3$	0.9	0.9	3.3

$$A = \operatorname{argmax}_a Q(S, a)$$

$$\pi(S_1) = A_3$$

$$\pi(S_2) = A_2$$

$$\pi(S_3) = A_3$$

- take action which maximizes the Q-value at each step

# Q-Learning

→ assign a value to each state  $Q(s, a) = \mathbb{E}_{a \sim \pi} [R_{t+1} + \dots + R_{t_f} | S_0 = s, A_0 = a]$

- first step deterministic, then follow policy

$Q(s, a)$	$A_1$	$A_2$	$A_3$
$S_1$	1.5	0.5	2.1
$S_2$	1.3	4.4	0.2
$S_3$	0.9	0.9	3.3

$$A = \operatorname{argmax}_a Q(S, a)$$

$$\pi(S_1) = A_3$$

$$\pi(S_2) = A_2$$

$$\pi(S_3) = A_3$$

- take action which maximizes the Q-value at each step

→ iterate following two steps until convergence

- error according to definition  $\delta_t = Q(S_t, A_t) + R_{t+1} - \max_a Q(S_{t+1}, a)$

- update current function  $Q_{\text{new}} = Q_{\text{old}} + \alpha \delta_t$   $\alpha \in [0, 1)$



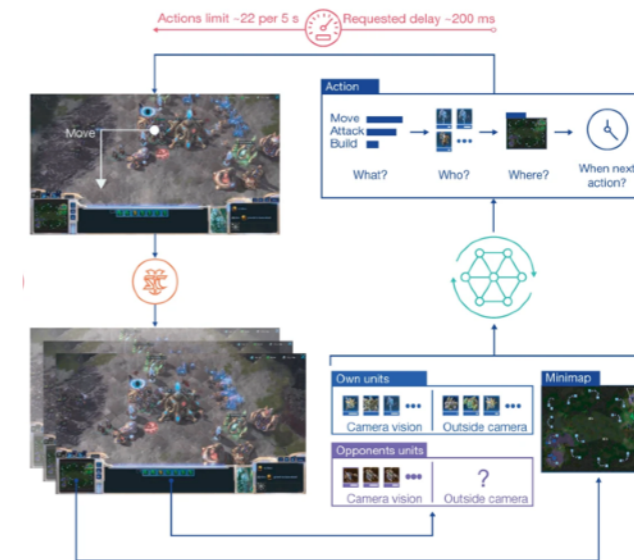
# What is reinforcement learning used for?

## Mastering the game of Go with deep neural networks and tree search



Silver, et. al, Nature 529 484–489 (2016)

## Mastering video games (StarCraft II)



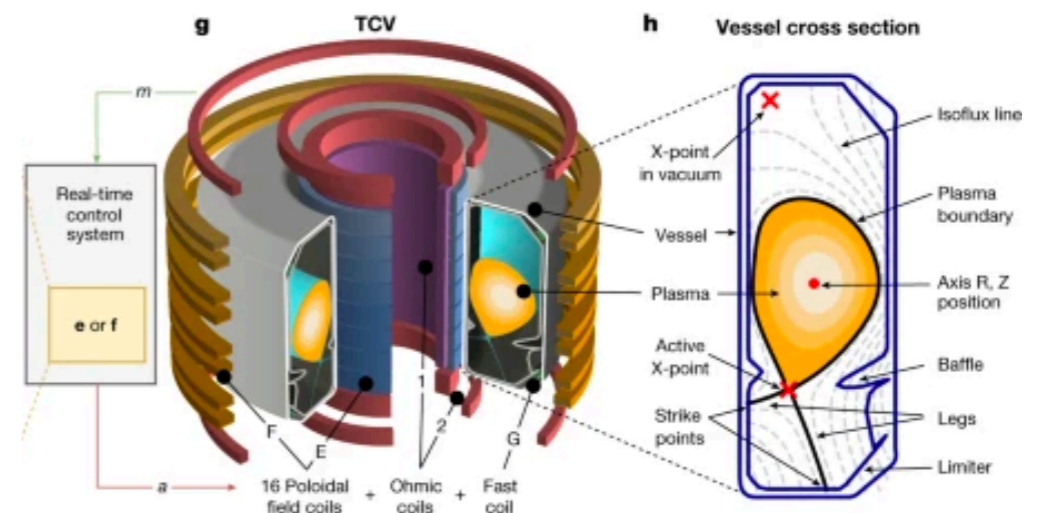
Vinyals, et. al, Nature 350 (2019)

## Atari games



Mnih et al., Nature 518 (2015) [Google DeepMind]

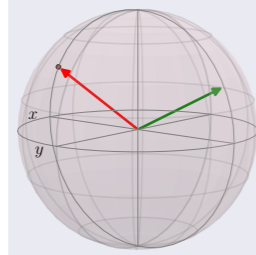
## Magnetic control of tokamak plasmas thru deep RL



Degrave, et. al, Nature 602 414–419 (2022)

# Applications of RL in Quantum Physics

## ● quantum control



MB et al, PRX 8 031086 (2018)

Niu et al, npj 5 33 (2019)

Sivak et al, PRX 12, 011059 (2022)

Gispen et al, MSML (2021)

Reuer, Nat Comm 14 7138 (2023)

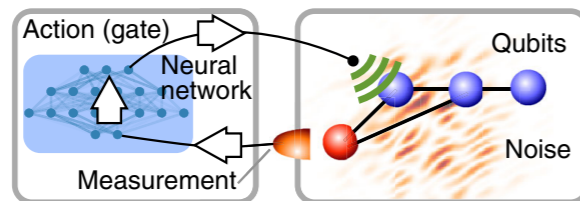
Yao et al, PRX 11 (3), 031070 (2021)

Porotti, Comm Phys 2 (2019)

Dalgaard et al, npj 6 6 (2020)

+ many more

## ● quantum error correction



Fössel et al, PRX 8 031086 (2018)

Andreasson et al, Quantum 3 183 (2019)

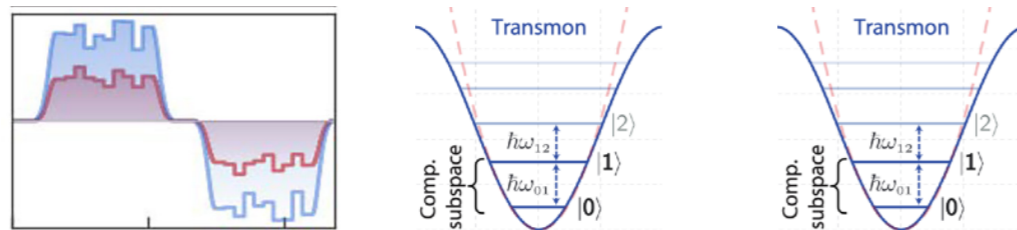
Sweke et al, ML Sci Tech 2 025005 (2020)

Sivak et al, Nature 616 50-55 (2023)

Olle et al, arXiv:2311.04750

+ many more

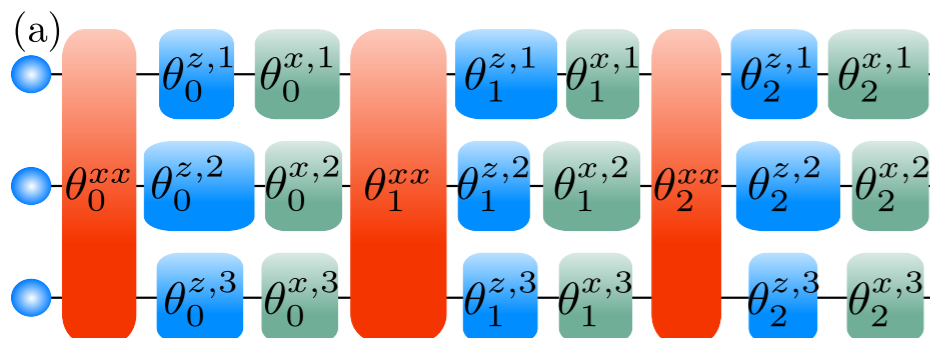
## ● quantum gate design



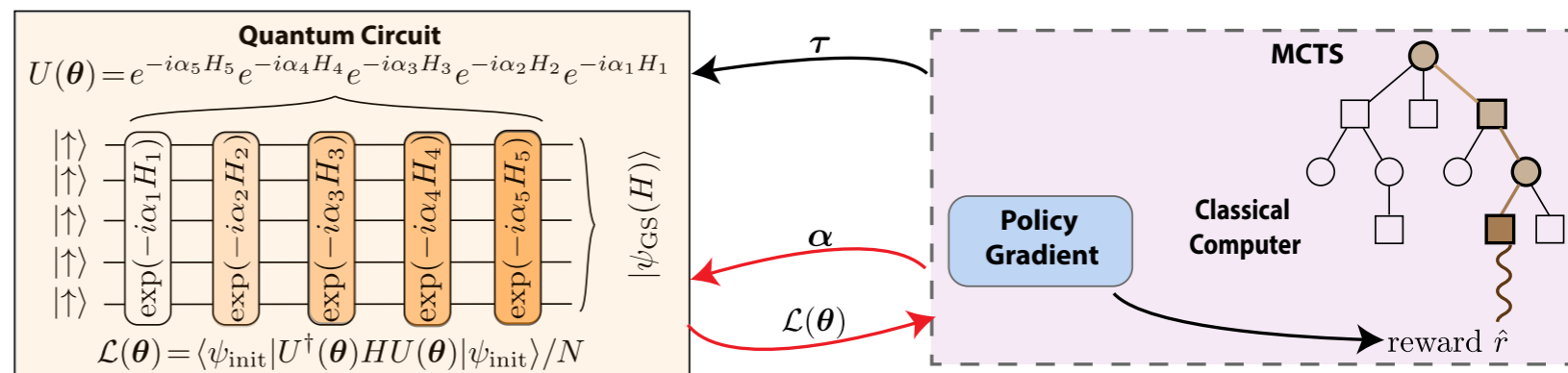
Baum et al, PRX Quantum 2, 040324 (2021)

Nguyen et al, ML Sci & Tech, 5, 025066 (2024)

## ● quantum circuit design and synthesis



Bolens et al, PRL 2021



Yao et al, MSML 2019, 2021, 2022

+ many more

# What advantages does RL offer?

- model-free: requires no pre-knowledge of the controlled physical system
  - unknown sources of noise (e.g., quantum computing)
  - not fully known Hamiltonian (solid state materials, superconducting qubits, etc.)
- adaptive: transfers acquired knowledge that might allow us to identify connections between unrelated phenomena
- interactive: designed for feedback control
  - quantum feedback control
- autonomous: provides novel insights into automating complex manipulation protocols
  - experiments



# Outline

## Part 2

- RL for qubit state preparation
  - effect of noise (measurement shot noise, coherent, incoherent noise)
- experimentally friendly RL framework
  - partially observable environments
  - environment, states, actions, rewards



# Outline

## Part 2

- RL for qubit state preparation
  - effect of noise (measurement shot noise, coherent, incoherent noise)
- experimentally friendly RL framework
  - partially observable environments
  - environment, states, actions, rewards

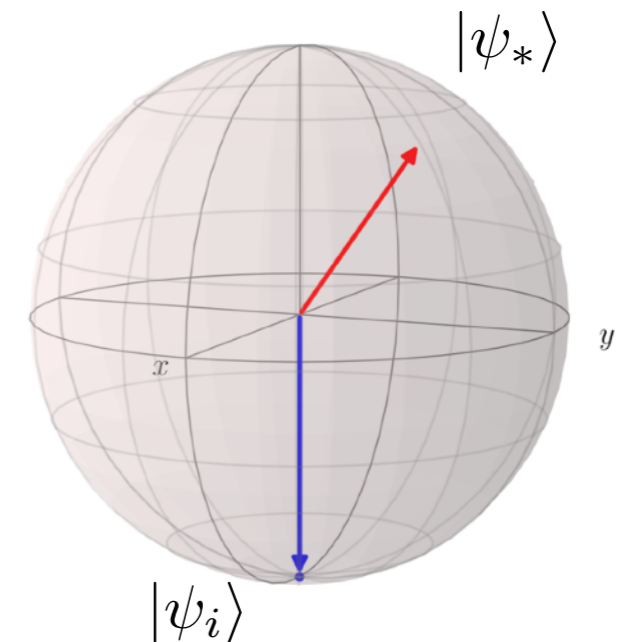
# RL frameworks for quantum systems

- **rewards obtained from (projective) measurements**
  - can be given only once, or else need to restart episode
- **quantum states cannot be observed**
  - extracting data from quantum system changes its state
- **quantum data is binary (qubits)**
  - minimal amount of information, difficult to learn
- **quantum data is probabilistic**
  - Heisenberg uncertainty

# RL for qubit control

→ goal: prepare state of a single qubit

- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$



# RL for qubit control

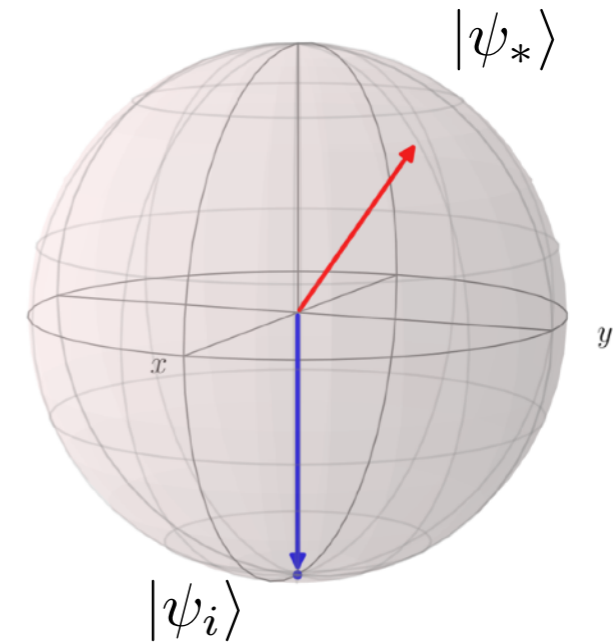
→ goal: prepare state of a single qubit

- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$

→ use quantum gates:

$$U_{\text{ctrl}}(\alpha, \beta, \gamma) = e^{-i\gamma\sigma^z/2}e^{-i\beta\sigma^y/2}e^{-i\alpha\sigma^x/2}$$

- need to determine optimal angles  $\alpha, \beta, \gamma$
- implements arbitrary rotations





# RL for qubit control

→ goal: prepare state of a single qubit

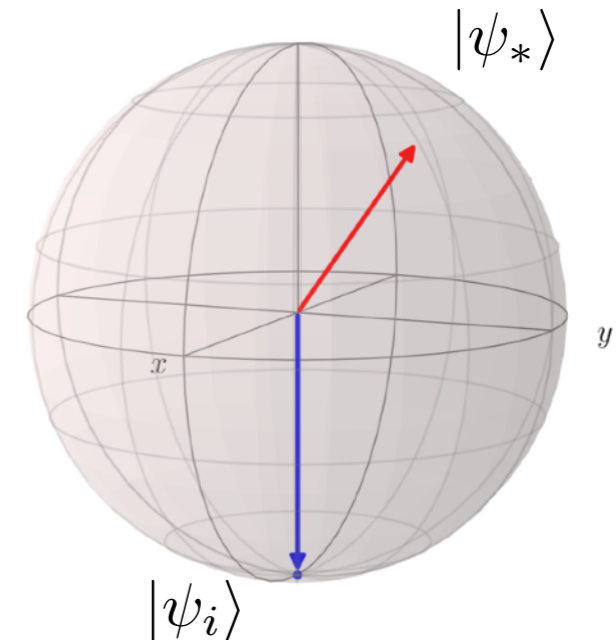
- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$

→ use quantum gates:

$$U_{\text{ctrl}}(\alpha, \beta, \gamma) = e^{-i\gamma\sigma^z/2}e^{-i\beta\sigma^y/2}e^{-i\alpha\sigma^x/2}$$

- need to determine optimal angles  $\alpha, \beta, \gamma$
- implements arbitrary rotations

→ solution: maximize the fidelity  $F(\alpha, \beta, \gamma) = |\langle\psi_*|U_{\text{ctrl}}(\alpha, \beta, \gamma)|\psi_i\rangle|^2$



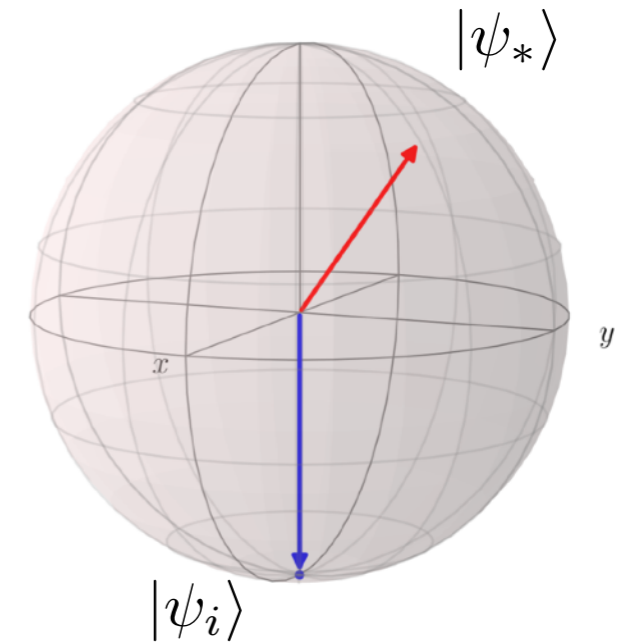
# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$

issues:

→ how do we compute  $F = |\langle\psi_*|U_{\text{ctrl}}|\psi_i\rangle|^2$  ?



# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$

issues:

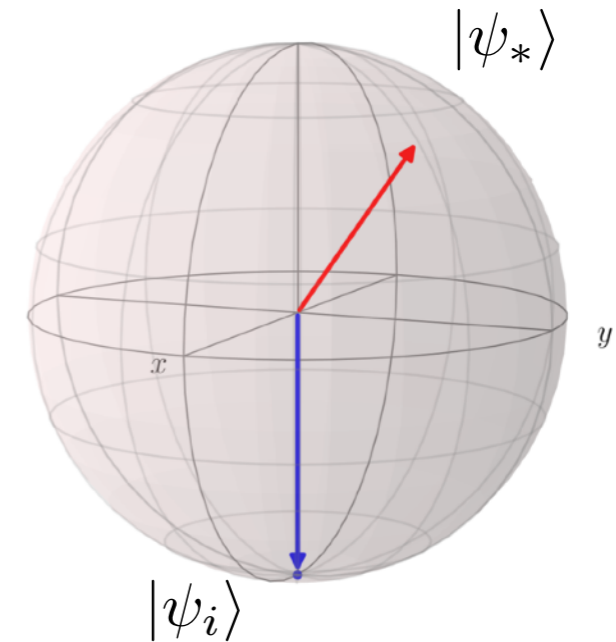
→ how do we compute  $F = |\langle\psi_*|U_{\text{ctrl}}|\psi_i\rangle|^2$  ?

- need to measure state in target direction  $|\psi_*\rangle$

find operator whose eigenstate is  $|\psi_*\rangle$ :  $\sigma_* = \hat{n}_* \cdot \vec{\sigma}$

$$\hat{n}_* = \hat{n}_*(\theta_*, \varphi_*)$$

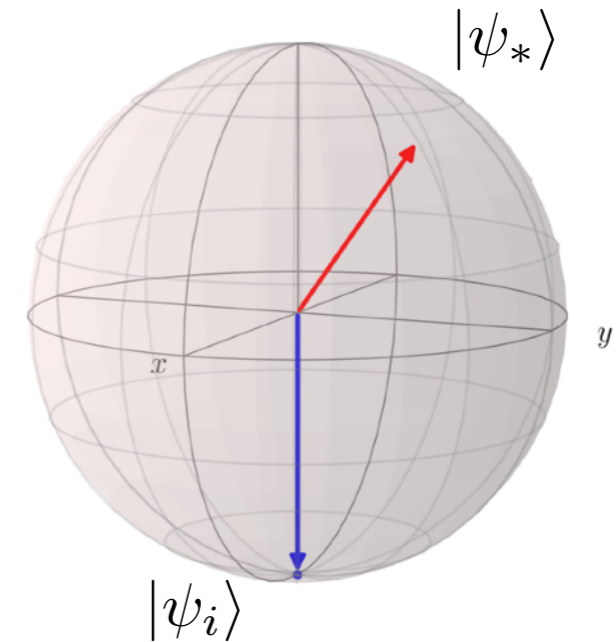
eigenvectors:  $\sigma_*|\psi_*\rangle = +|\psi_*\rangle$ ,  $\sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$



# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

- initial state:  $|\psi_i\rangle = |1\rangle$
- target state:  $|\psi_*\rangle = \cos\frac{\pi}{8}|0\rangle + \sin\frac{\pi}{8}e^{i\frac{\pi}{3}}|1\rangle$



## issues:

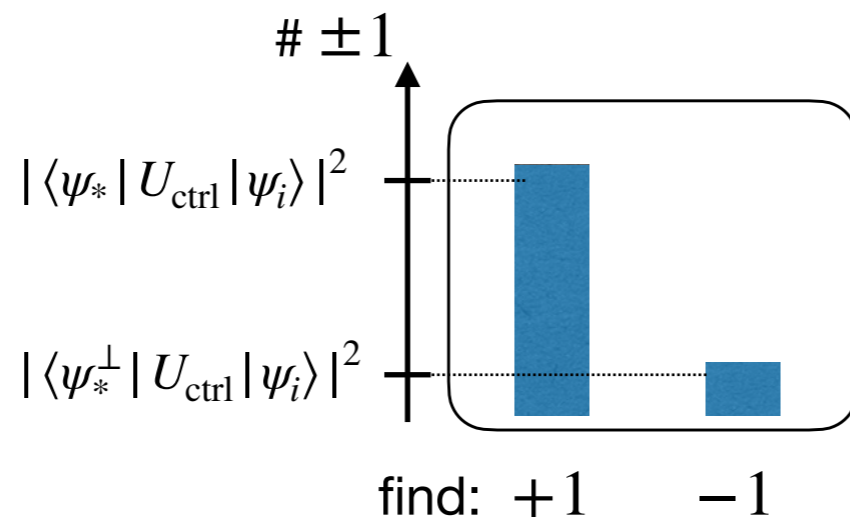
→ how do we compute  $F = |\langle\psi_*|U_{\text{ctrl}}|\psi_i\rangle|^2$  ?

- need to measure state in target direction  $|\psi_*\rangle$

find operator whose eigenstate is  $|\psi_*\rangle$ :  $\sigma_* = \hat{n}_* \cdot \vec{\sigma}$

$$\hat{n}_* = \hat{n}_*(\theta_*, \varphi_*)$$

eigenvectors:  $\sigma_*|\psi_*\rangle = +|\psi_*\rangle$ ,  $\sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$



measurement shot noise!

discrete measurement outcome

# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

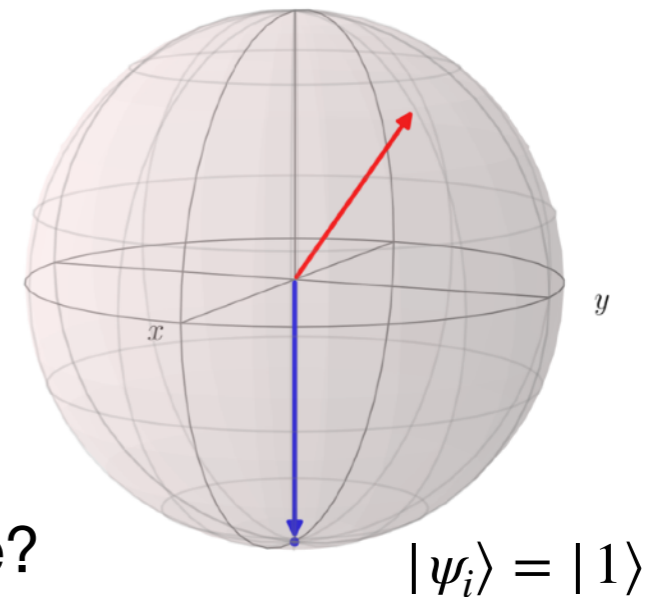
## issues:

→ how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

▸ measurement shot noise

→ how do we measure  $\sigma_* = \hat{n}_* \cdot \vec{\sigma}$  on a quantum device?

$$|\psi_*\rangle = \cos \frac{\pi}{8} |0\rangle + \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |1\rangle$$



# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

## issues:

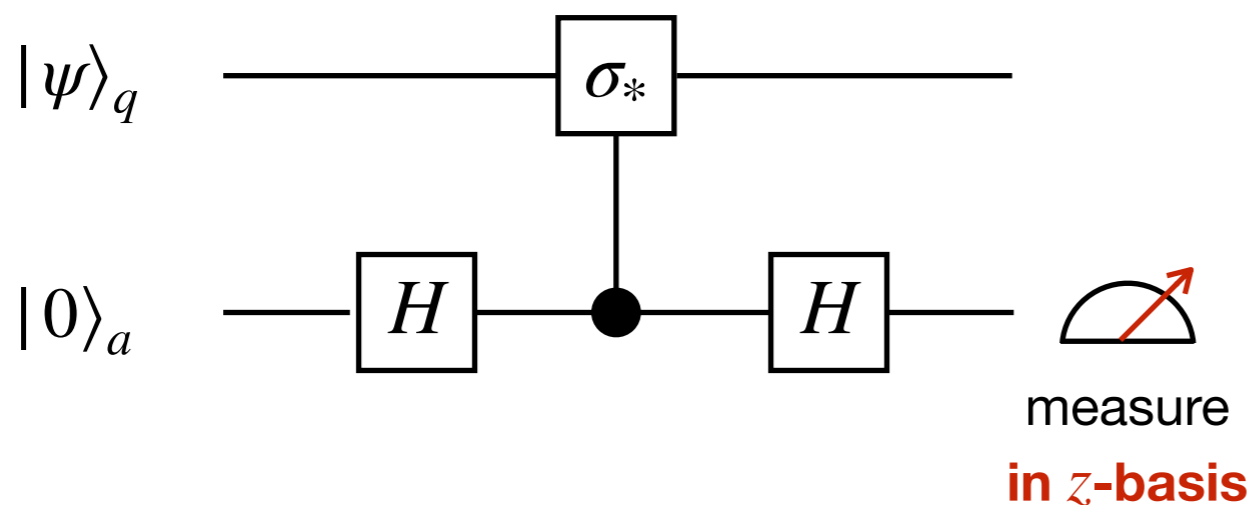
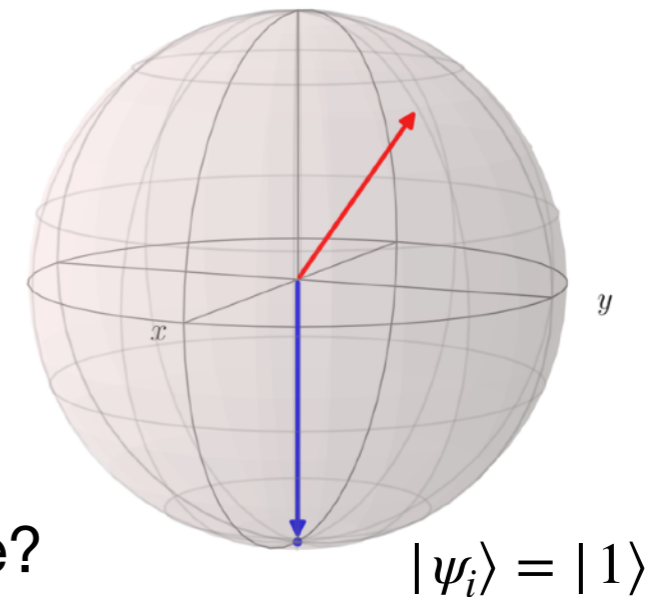
→ how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

▸ measurement shot noise

→ how do we measure  $\sigma_* = \hat{n}_* \cdot \vec{\sigma}$  on a quantum device?

• need ancilla  $|\psi\rangle = |\psi\rangle_q |0\rangle_a$

$$|\psi_*\rangle = \cos \frac{\pi}{8} |0\rangle + \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |1\rangle$$



$$\text{Hadamard gate: } H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$\text{controlled-}\sigma_* \text{ gate: } C\sigma_* = \begin{cases} \sigma_{*q} & \text{if ancilla is in } |1\rangle_a \\ 1 & \text{if ancilla is in } |0\rangle_a \end{cases}$$

# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

issues:

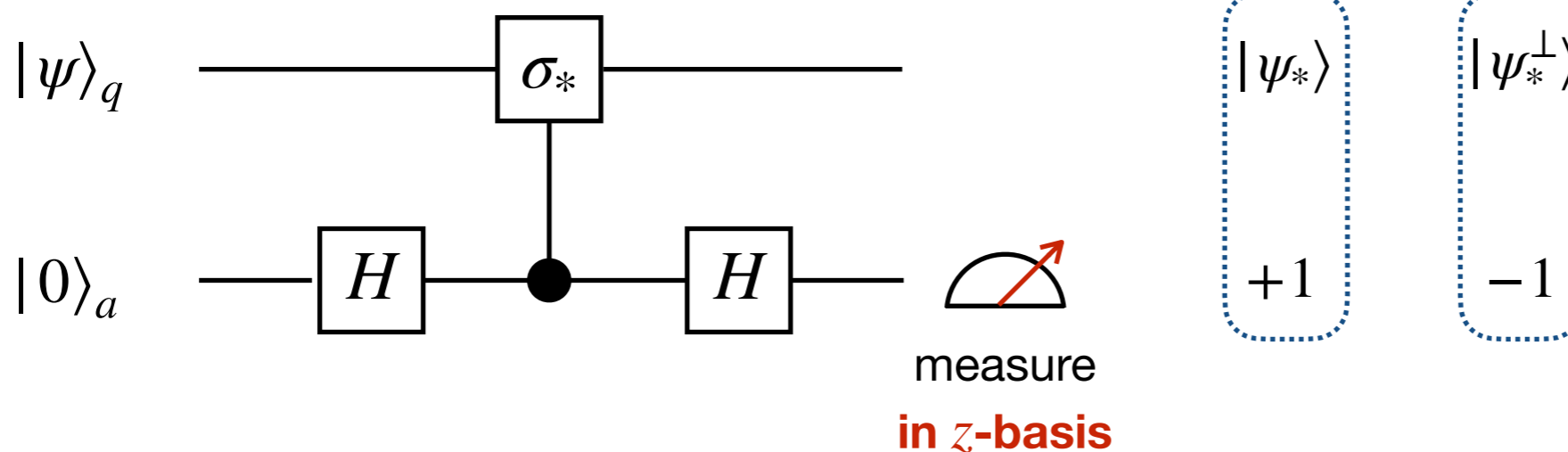
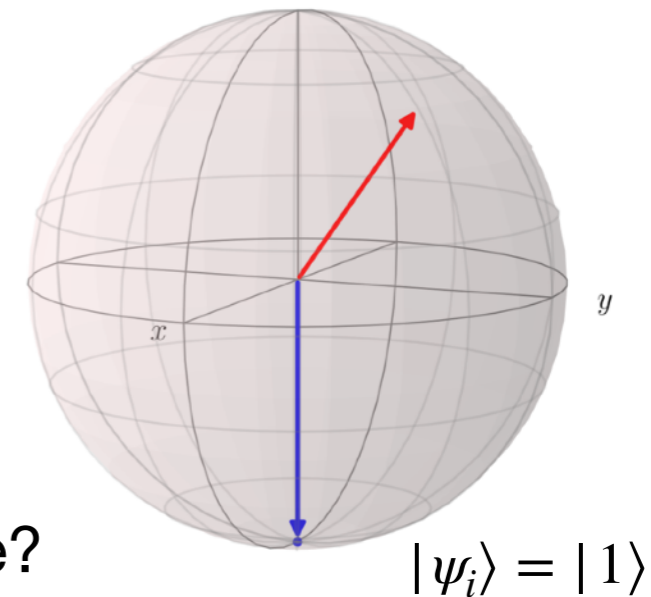
→ how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

▸ measurement shot noise

→ how do we measure  $\sigma_* = \hat{n}_* \cdot \vec{\sigma}$  on a quantum device?

• need ancilla  $|\psi\rangle = |\psi\rangle_q |0\rangle_a$

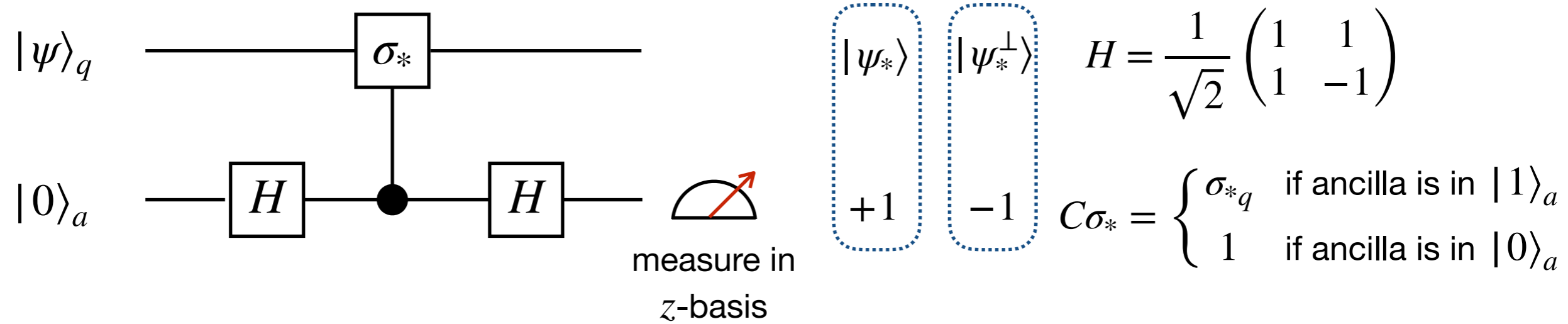
$$|\psi_*\rangle = \cos \frac{\pi}{8} |0\rangle + \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |1\rangle$$



$$\text{Hadamard gate: } H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$\text{controlled-}\sigma_* \text{ gate: } C\sigma_* = \begin{cases} \sigma_{*q} & \text{if ancilla is in } |1\rangle_a \\ 1 & \text{if ancilla is in } |0\rangle_a \end{cases}$$

# Measuring arbitrary operator





# Measuring arbitrary operator

$$|\psi\rangle_q \text{ --- } \boxed{\sigma_*}$$



$$|\psi\rangle_q = a |\psi_*\rangle_q + b |\psi_*^\perp\rangle_q$$

$$\sigma_* |\psi_*\rangle = + |\psi_*\rangle, \quad \sigma_* |\psi_*^\perp\rangle = - |\psi_*^\perp\rangle$$

$$\sigma^z |0\rangle = + |0\rangle, \quad \sigma^z |1\rangle = - |1\rangle$$

# Measuring arbitrary operator

$|\psi\rangle_q$  —

$|0\rangle_a$  —



$$|\psi\rangle_q = a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q$$

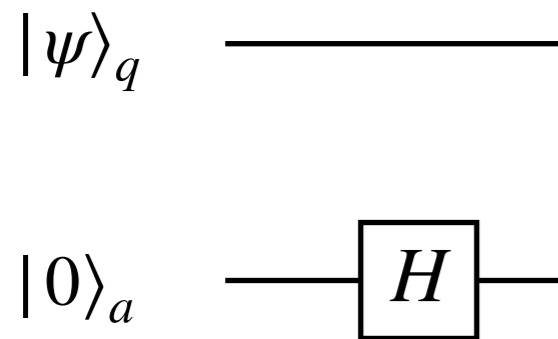
$$\sigma_*|\psi_*\rangle = +|\psi_*\rangle, \quad \sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$$



$$|\psi\rangle_q|0\rangle_a = a|\psi_*\rangle_q|0\rangle_a + b|\psi_*^\perp\rangle_q|0\rangle_a$$

$$\sigma^z|0\rangle = +|0\rangle, \quad \sigma^z|1\rangle = -|1\rangle$$

# Measuring arbitrary operator



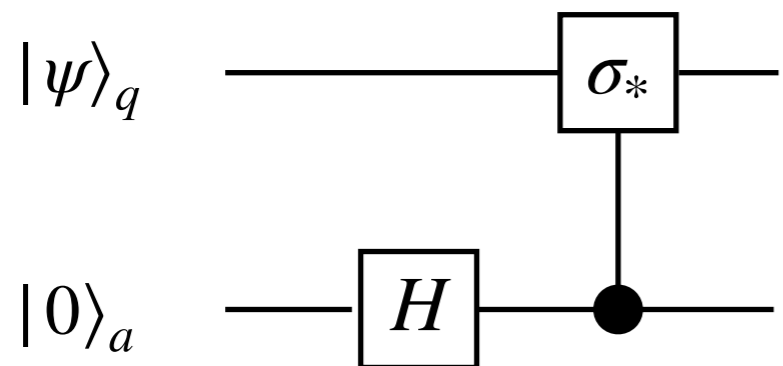
$$H = \frac{1}{\sqrt{2}} \begin{matrix} & |0\rangle & |1\rangle \\ \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} & |0\rangle \\ & |1\rangle \end{matrix}$$

$$\rightarrow |\psi\rangle_q = a |\psi_*\rangle_q + b |\psi_*^\perp\rangle_q \quad \sigma_* |\psi_*\rangle = + |\psi_*\rangle, \quad \sigma_* |\psi_*^\perp\rangle = - |\psi_*^\perp\rangle$$

$$\rightarrow |\psi\rangle_q |0\rangle_a = a |\psi_*\rangle_q |0\rangle_a + b |\psi_*^\perp\rangle_q |0\rangle_a \quad \sigma^z |0\rangle = + |0\rangle, \quad \sigma^z |1\rangle = - |1\rangle$$

$$\begin{aligned} \rightarrow H_a |\psi\rangle_q |0\rangle_a &= \frac{a}{\sqrt{2}} |\psi_*\rangle_q (|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}} |\psi_*^\perp\rangle_q (|0\rangle_a + |1\rangle_a) \\ &= \frac{1}{\sqrt{2}} (a |\psi_*\rangle_q + b |\psi_*^\perp\rangle_q) |0\rangle_a + \frac{1}{\sqrt{2}} (a |\psi_*\rangle_q + b |\psi_*^\perp\rangle_q) |1\rangle_a \end{aligned}$$

# Measuring arbitrary operator



$$H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{matrix} |0\rangle \\ |1\rangle \end{matrix}$$

$$C\sigma_* = \begin{cases} \sigma_{*q} & \text{if ancilla is in } |1\rangle_a \\ 1 & \text{if ancilla is in } |0\rangle_a \end{cases}$$

$$\rightarrow |\psi\rangle_q = a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q \quad \sigma_*|\psi_*\rangle = +|\psi_*\rangle, \quad \sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$$

$$\sigma^z|0\rangle = +|0\rangle, \quad \sigma^z|1\rangle = -|1\rangle$$

$$\rightarrow |\psi\rangle_q |0\rangle_a = a|\psi_*\rangle_q |0\rangle_a + b|\psi_*^\perp\rangle_q |0\rangle_a$$

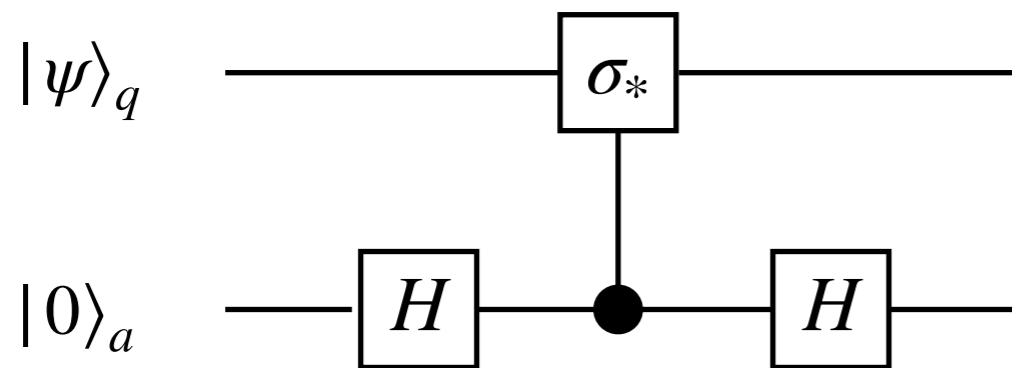
$$\rightarrow H_a |\psi\rangle_q |0\rangle_a = \frac{a}{\sqrt{2}} |\psi_*\rangle_q (|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}} |\psi_*^\perp\rangle_q (|0\rangle_a + |1\rangle_a)$$

$$= \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |0\rangle_a + \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |1\rangle_a$$

$$\rightarrow C\sigma_* H_a |\psi\rangle_q |0\rangle_a = \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |0\rangle_a + \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q - b|\psi_*^\perp\rangle_q) |1\rangle_a$$

$$= \frac{a}{\sqrt{2}} |\psi_*\rangle_q (|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}} |\psi_*^\perp\rangle_q (|0\rangle_a - |1\rangle_a)$$

# Measuring arbitrary operator



$$H = \frac{1}{\sqrt{2}} \begin{matrix} & |0\rangle & |1\rangle \\ \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} & |0\rangle \\ & |1\rangle \end{matrix}$$

$$C\sigma_* = \begin{cases} \sigma_{*q} & \text{if ancilla is in } |1\rangle_a \\ 1 & \text{if ancilla is in } |0\rangle_a \end{cases}$$

$$\rightarrow |\psi\rangle_q = a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q \quad \sigma_*|\psi_*\rangle = +|\psi_*\rangle, \quad \sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$$

$$\sigma^z|0\rangle = +|0\rangle, \quad \sigma^z|1\rangle = -|1\rangle$$

$$\rightarrow |\psi\rangle_q |0\rangle_a = a|\psi_*\rangle_q |0\rangle_a + b|\psi_*^\perp\rangle_q |0\rangle_a$$

$$\rightarrow H_a |\psi\rangle_q |0\rangle_a = \frac{a}{\sqrt{2}} |\psi_*\rangle_q (|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}} |\psi_*^\perp\rangle_q (|0\rangle_a + |1\rangle_a)$$

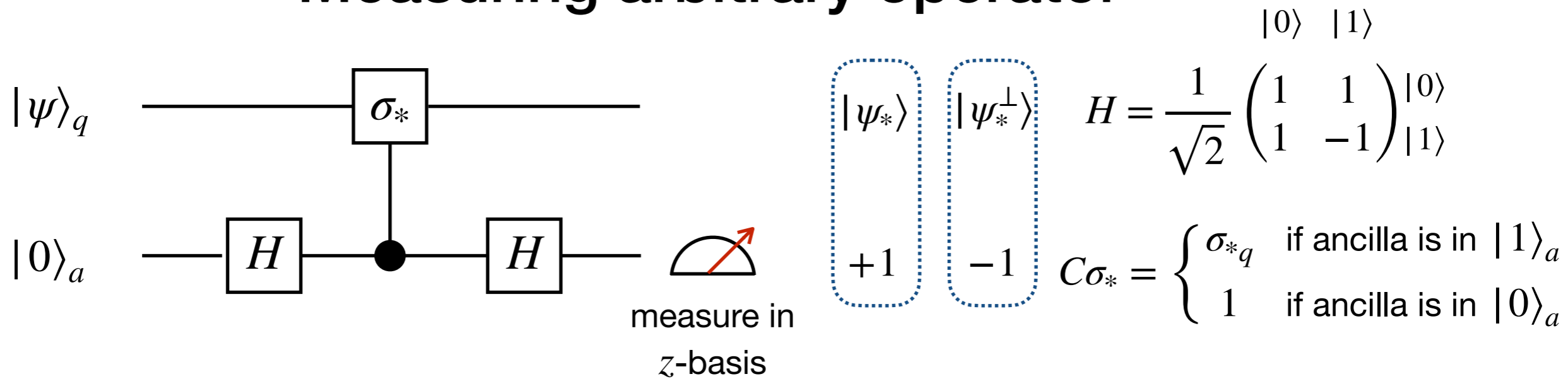
$$= \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |0\rangle_a + \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |1\rangle_a$$

$$\rightarrow C\sigma_* H_a |\psi\rangle_q |0\rangle_a = \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q) |0\rangle_a + \frac{1}{\sqrt{2}} (a|\psi_*\rangle_q - b|\psi_*^\perp\rangle_q) |1\rangle_a$$

$$= \frac{a}{\sqrt{2}} |\psi_*\rangle_q (|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}} |\psi_*^\perp\rangle_q (|0\rangle_a - |1\rangle_a)$$

$$\rightarrow H_a C\sigma_* H_a |\psi\rangle_q |0\rangle_a = a|\psi_*\rangle_q |0\rangle_a + b|\psi_*^\perp\rangle_q |1\rangle_a$$

# Measuring arbitrary operator



→  $|\psi\rangle_q = a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q$        $\sigma_*|\psi_*\rangle = +|\psi_*\rangle, \quad \sigma_*|\psi_*^\perp\rangle = -|\psi_*^\perp\rangle$

→  $|\psi\rangle_q|0\rangle_a = a|\psi_*\rangle_q|0\rangle_a + b|\psi_*^\perp\rangle_q|0\rangle_a$        $\sigma^z|0\rangle = +|0\rangle, \quad \sigma^z|1\rangle = -|1\rangle$

→  $H_a|\psi\rangle_q|0\rangle_a = \frac{a}{\sqrt{2}}|\psi_*\rangle_q(|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}}|\psi_*^\perp\rangle_q(|0\rangle_a + |1\rangle_a)$

$= \frac{1}{\sqrt{2}}(a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q)|0\rangle_a + \frac{1}{\sqrt{2}}(a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q)|1\rangle_a$

→  $C\sigma_*H_a|\psi\rangle_q|0\rangle_a = \frac{1}{\sqrt{2}}(a|\psi_*\rangle_q + b|\psi_*^\perp\rangle_q)|0\rangle_a + \frac{1}{\sqrt{2}}(a|\psi_*\rangle_q - b|\psi_*^\perp\rangle_q)|1\rangle_a$

$= \frac{a}{\sqrt{2}}|\psi_*\rangle_q(|0\rangle_a + |1\rangle_a) + \frac{b}{\sqrt{2}}|\psi_*^\perp\rangle_q(|0\rangle_a - |1\rangle_a)$

→  $H_aC\sigma_*H_a|\psi\rangle_q|0\rangle_a = a|\psi_*\rangle_q|0\rangle_a + b|\psi_*^\perp\rangle_q|1\rangle_a$

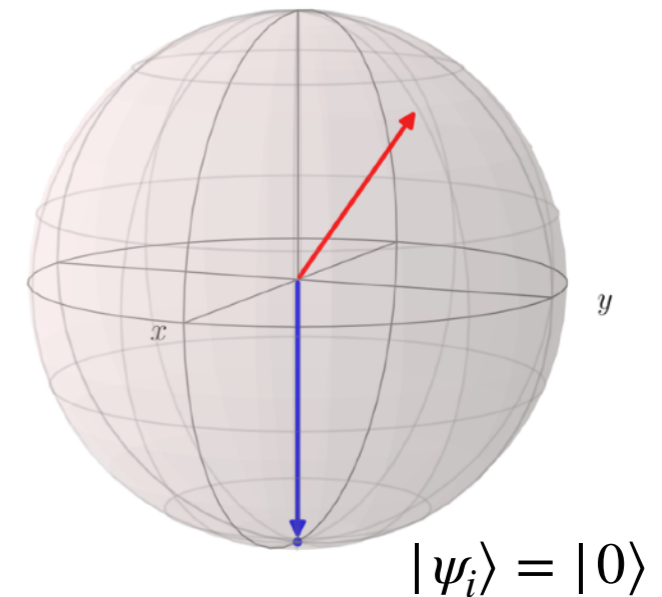
# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?
2. qubit + ancilla apparatus can be noisy
  - coherent noise

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



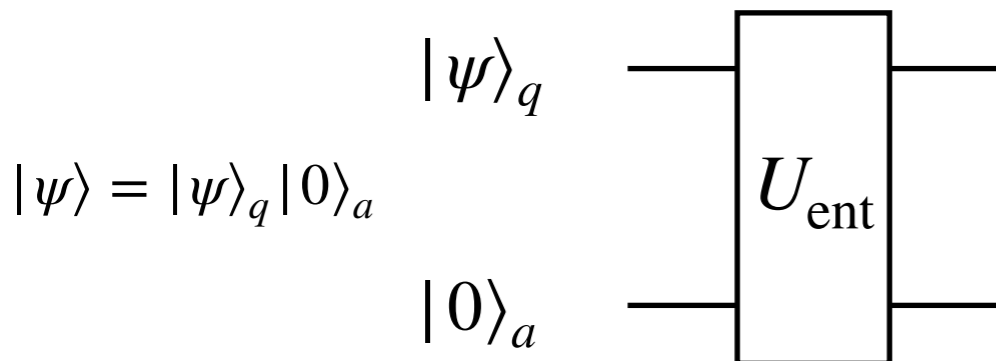
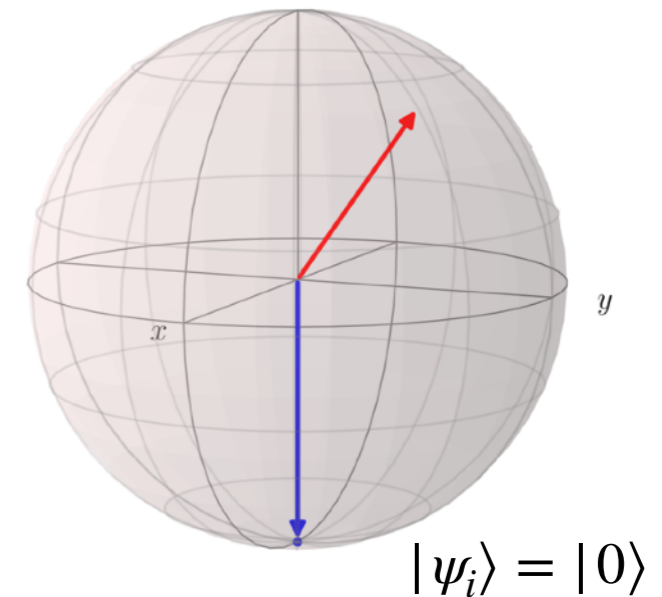
# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?
2. qubit + ancilla apparatus can be noisy
  - coherent noise

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



$$U_{\text{ent}} = e^{-ia\sigma_a^x \sigma_q^x} e^{-ib\sigma_a^y \sigma_q^y} e^{-ic\sigma_a^z \sigma_q^z}$$

noise strength set by  $a, b, c$

qubit and ancilla can get entangled



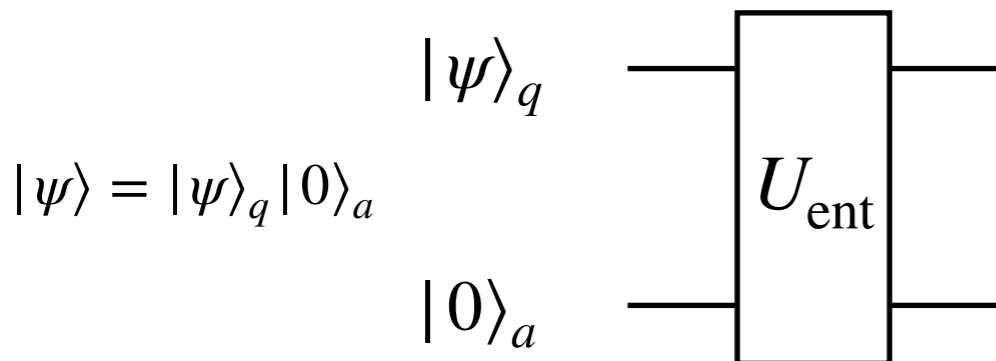
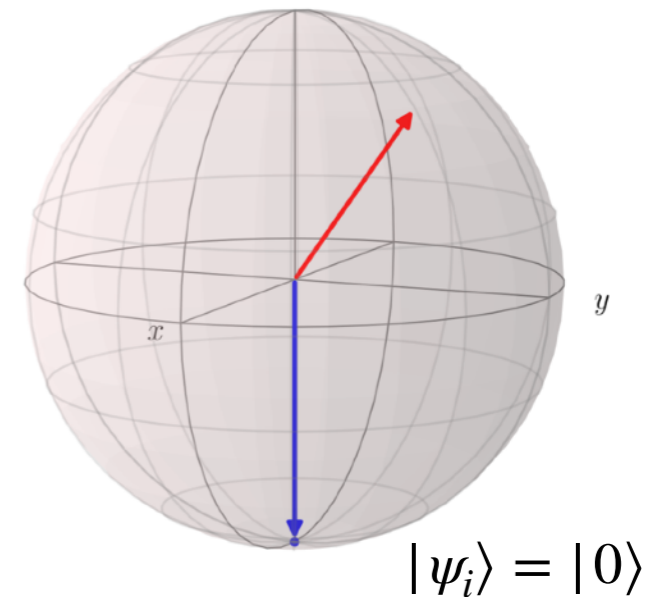
# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?
2. qubit + ancilla apparatus can be noisy
  - coherent noise

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



$$U_{\text{ent}} = e^{-ia\sigma_a^x \sigma_q^x} e^{-ib\sigma_a^y \sigma_q^y} e^{-ic\sigma_a^z \sigma_q^z}$$

noise strength set by  $a, b, c$

qubit and ancilla can get entangled

how do we get rid of this entanglement?

measure ancilla in  $z$ -basis to reset it!

# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

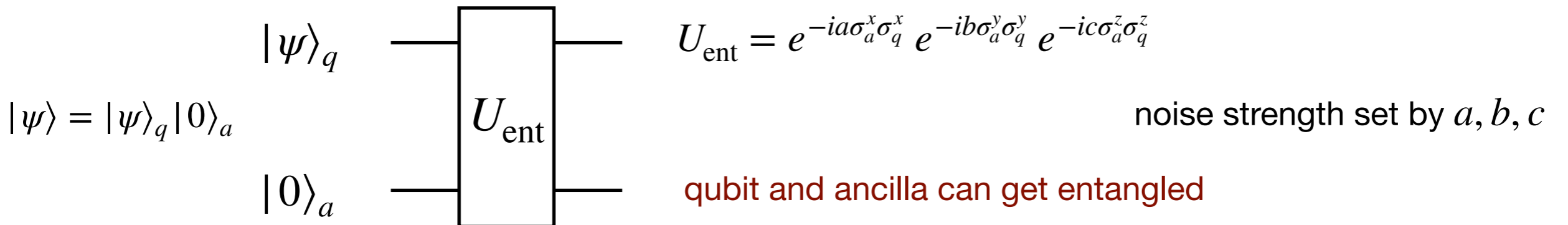
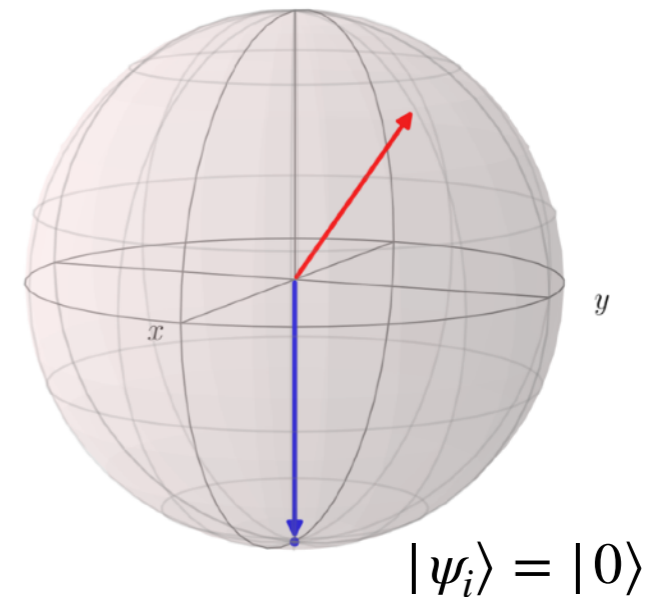
## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

2. qubit + ancilla apparatus can be noisy

- coherent noise

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



how do we get rid of this entanglement?

measure ancilla in  $z$ -basis to reset it!

$$|\psi\rangle \longrightarrow \begin{cases} \frac{P_a^z |\psi\rangle}{\sqrt{p}}, & \text{with prob. } p = |\langle \psi | P_a^z | \psi \rangle|^2 \quad \& \text{ measurement outcome } -1 \\ \frac{(1 - P_a^z) |\psi\rangle}{\sqrt{1-p}}, & \text{with prob. } 1 - p \quad \& \text{ measurement outcome } +1 \end{cases}$$

$$P_a^z = 1_q \otimes \frac{1}{2}(1 - \sigma^z)_a$$

# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

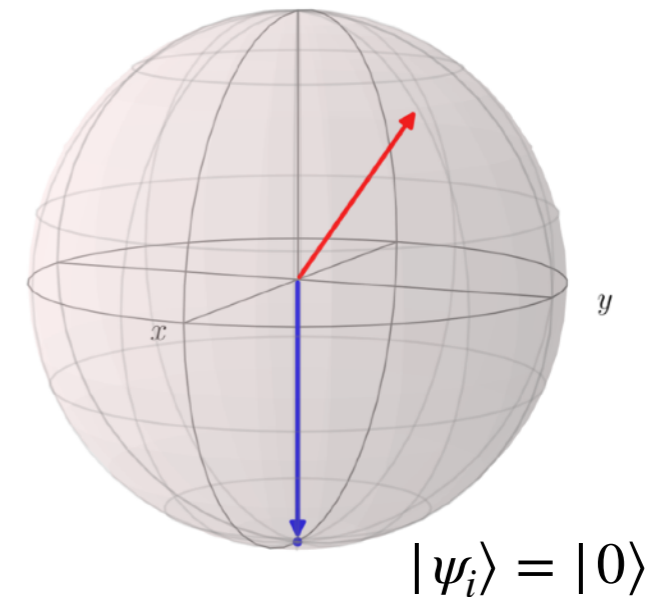
## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

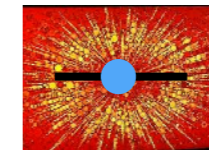
2. qubit + ancilla apparatus can be noisy

- coherent noise
- incoherent noise: spontaneous decay of qubit state

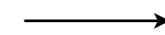
$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



photon  $\gamma$



qubit



qubit

# RL for *noisy* qubit control

→ goal: prepare state of a single qubit

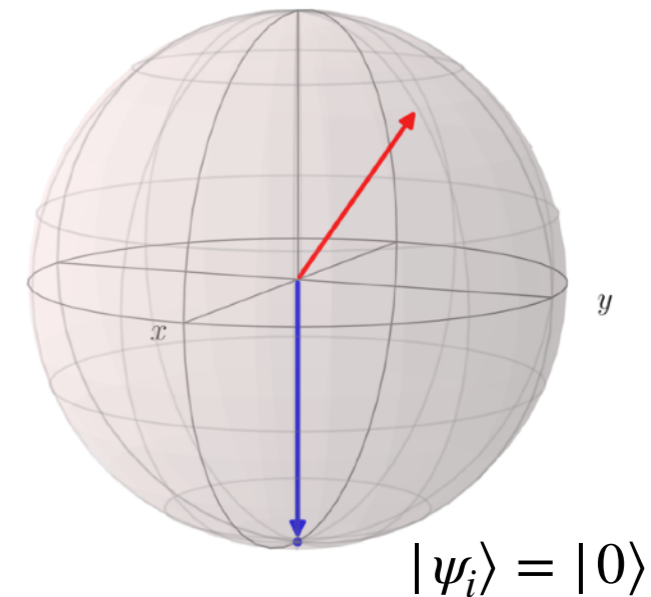
## issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?

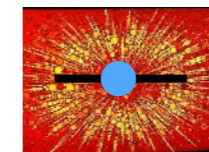
2. qubit + ancilla apparatus can be noisy

- coherent noise
- incoherent noise: spontaneous decay of qubit state

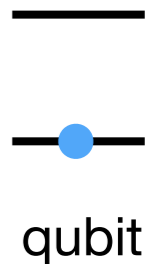
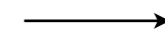
$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$



photon  $\gamma$



qubit



qubit

$$|\psi\rangle_q \longrightarrow \begin{cases} \frac{P_q^z |\psi\rangle}{\sqrt{|\langle \psi | P_z | \psi \rangle_q|}}, & \text{with prob. } p_{\text{emit}} \\ |\psi\rangle_q, & \text{with prob. } 1 - p_{\text{emit}} \end{cases}$$

$$P_q^z = \frac{1}{2}(1 - \sigma^z)_q$$

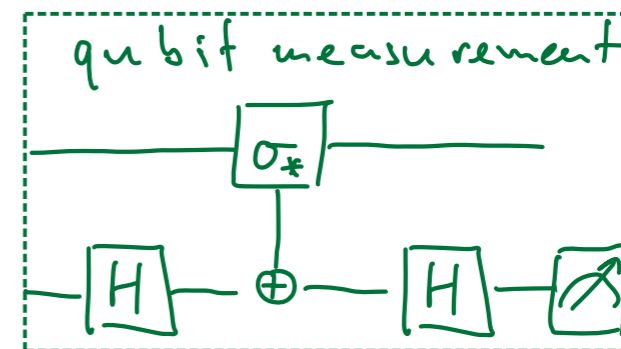
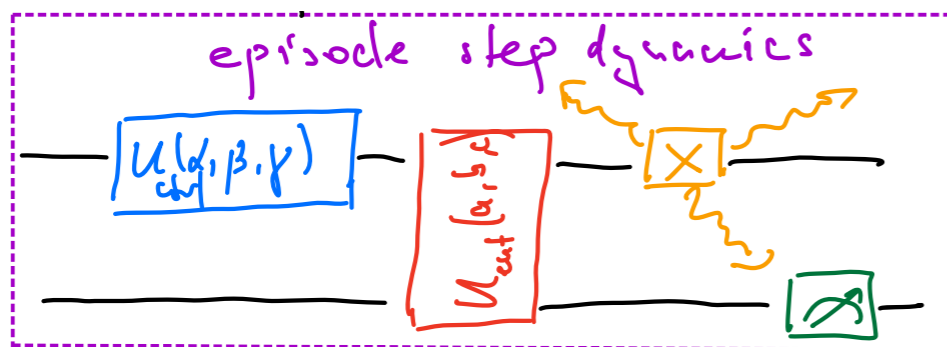
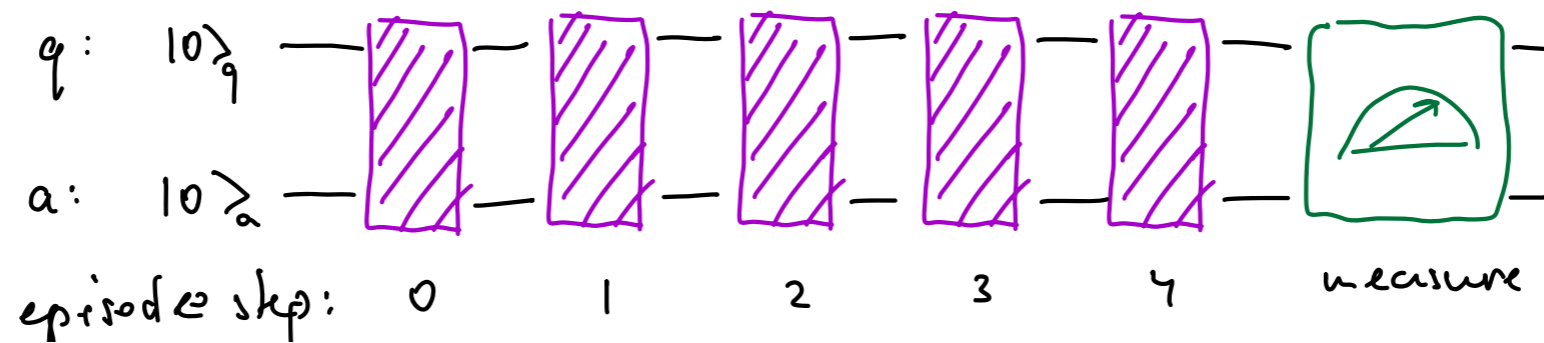
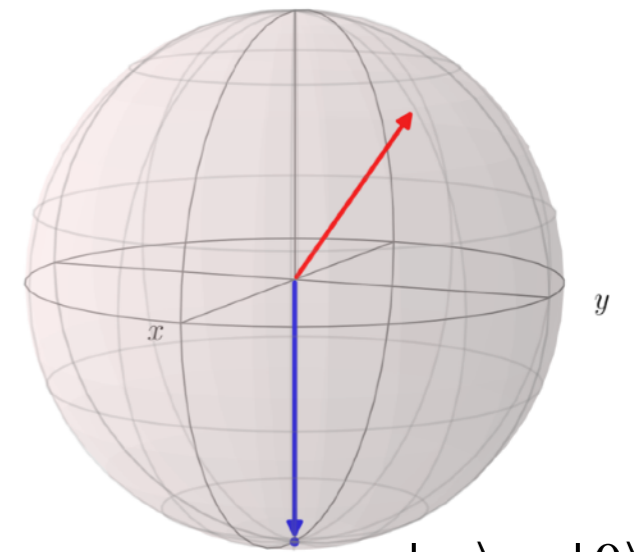
# RL for noisy qubit control

→ goal: prepare state of a single qubit

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$

issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?
2. qubit + ancilla apparatus can be noisy



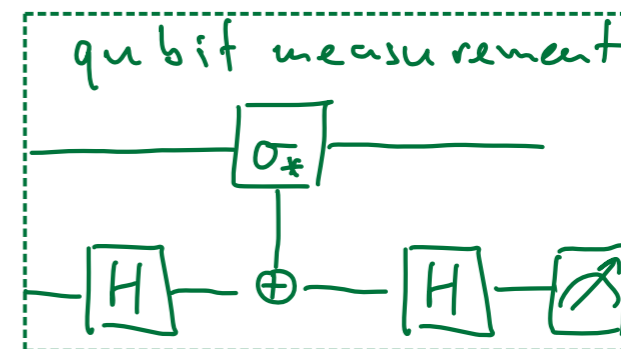
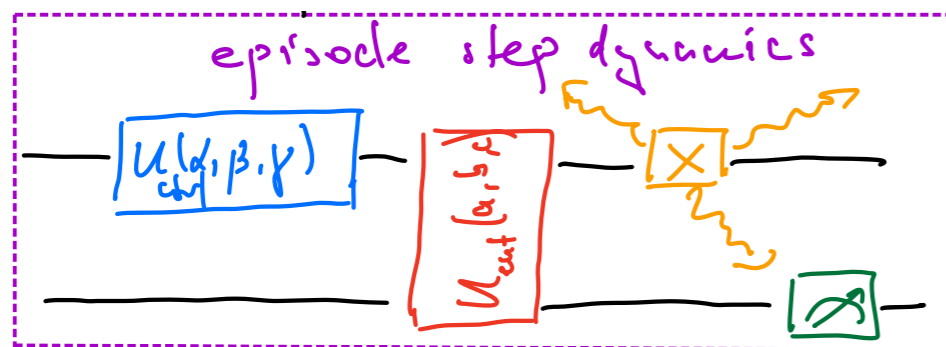
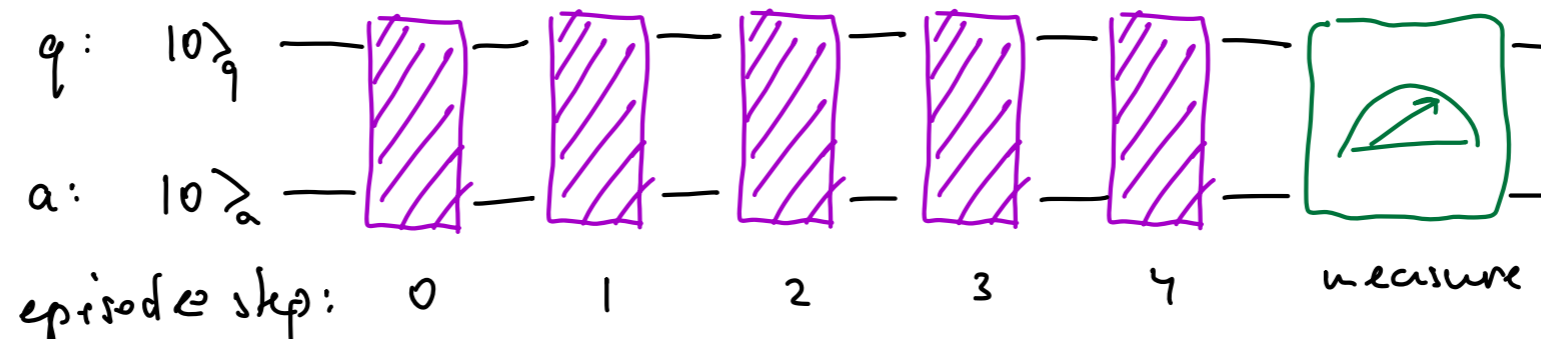
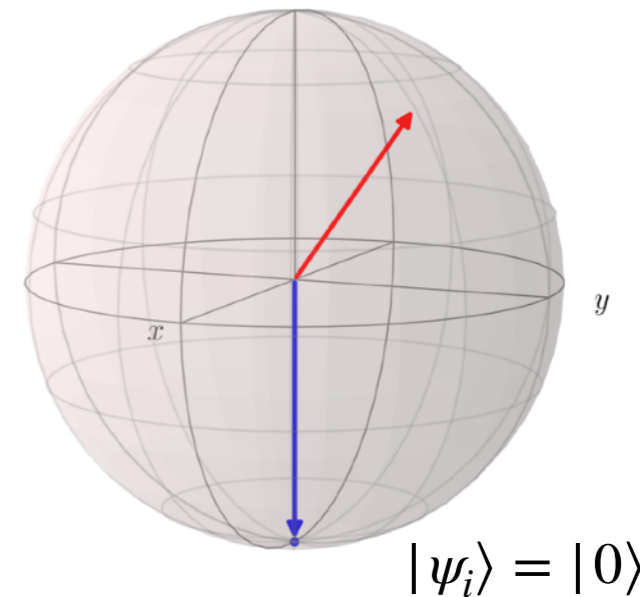
# RL for noisy qubit control

→ goal: prepare state of a single qubit

$$|\psi_*\rangle = \sin \frac{\pi}{8} e^{i\frac{\pi}{3}} |0\rangle + \cos \frac{\pi}{8} |1\rangle$$

issues:

1. how do we compute  $F = |\langle \psi_* | U_{\text{ctrl}} | \psi_i \rangle|^2$  ?
2. qubit + ancilla apparatus can be noisy



3. accessible data: is all binary!!

- ancilla measurement
- single photon detection
- qubit measurement



# Outline

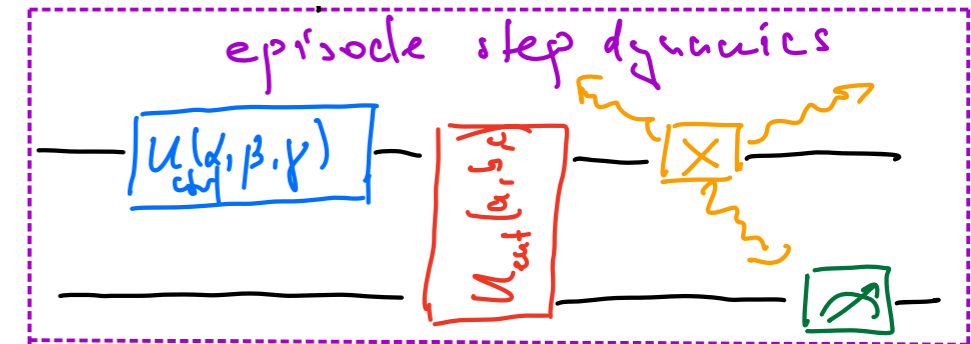
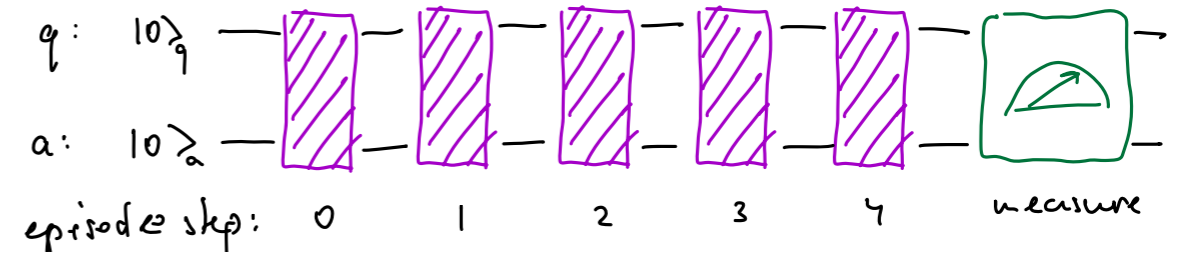
## Part 2

- RL for qubit state preparation
  - effect of noise (measurement shot noise, coherent, incoherent noise)
- experimentally friendly RL framework
  - partially observable environments
  - environment, states, actions, rewards

# RL framework

→ rewards

- qubit measurement output:  $\pm 1$  (binary)





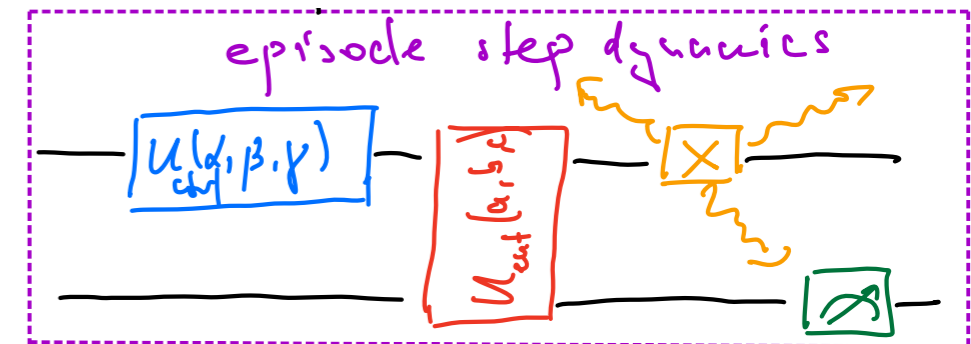
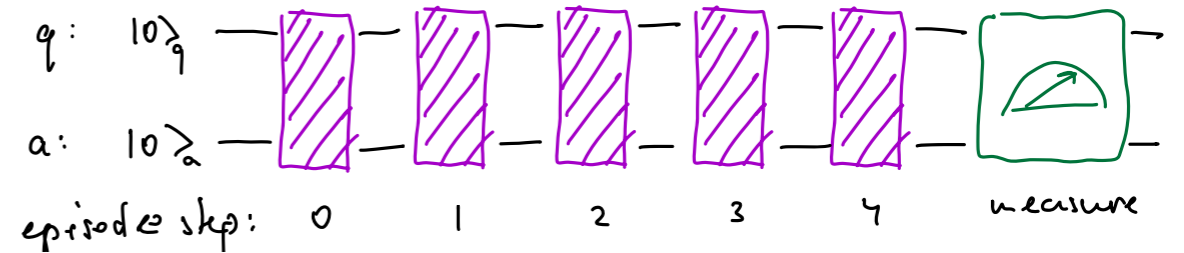
# RL framework

## → rewards

- qubit measurement output:  $\pm 1$  (binary)

## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$  : continuous (!)



# RL framework

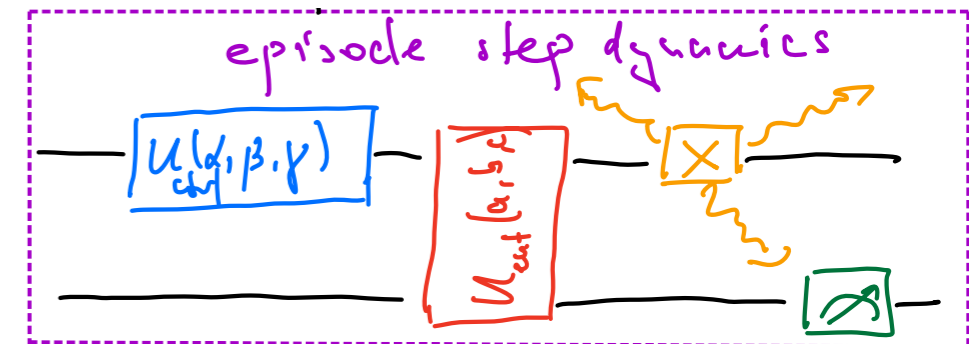
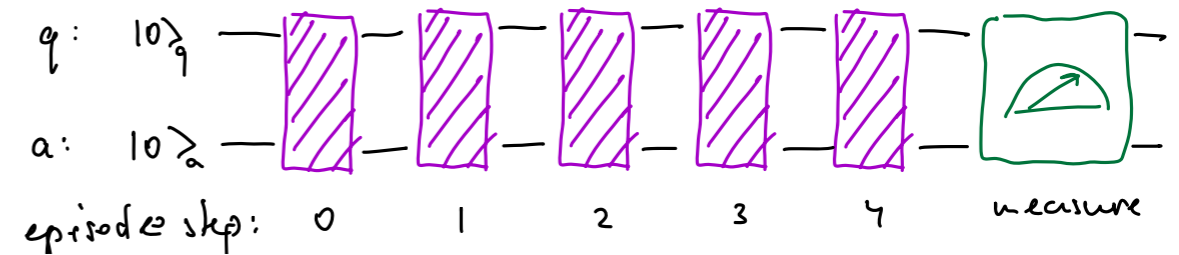
## → rewards

- qubit measurement output:  $\pm 1$  (binary)

## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$  : **continuous (!)**

## → states

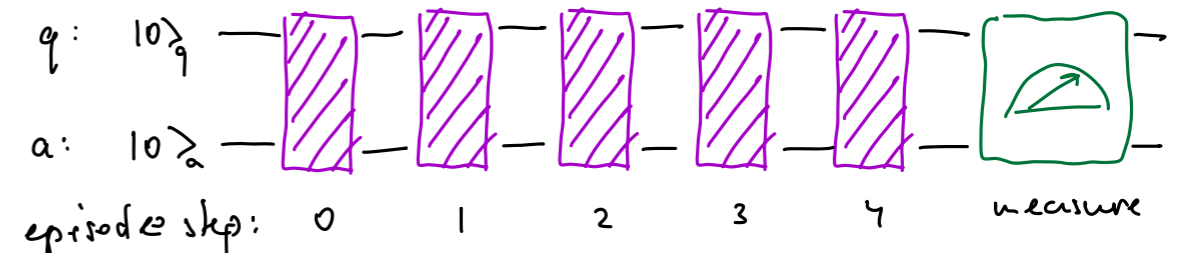


**quantum states cannot be measured/observed!!!**

# RL framework

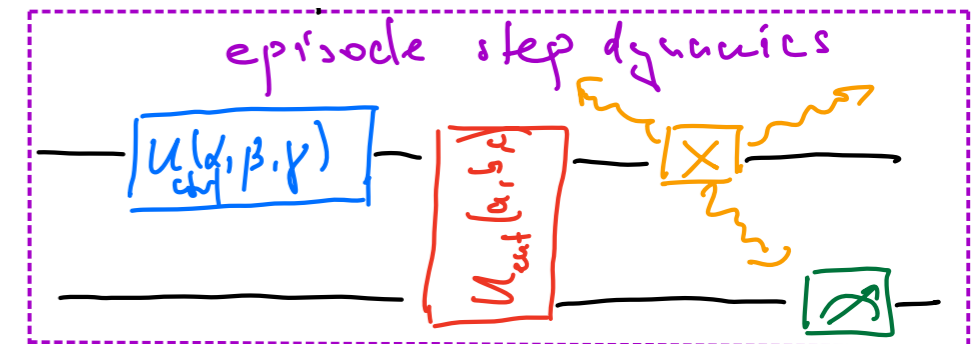
→ rewards

- qubit measurement output:  $\pm 1$  (binary)



→ actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$  : **continuous (!)**



→ states → observations

- wall clock time: one-hot representation of the step number

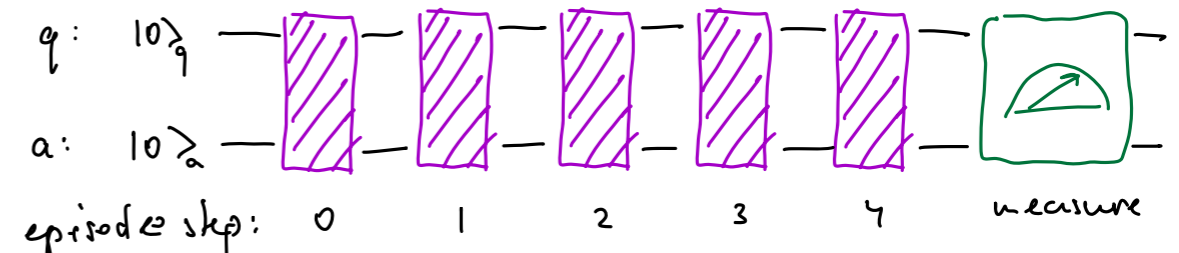
$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

quantum states cannot be measured/observed!!!

# RL framework

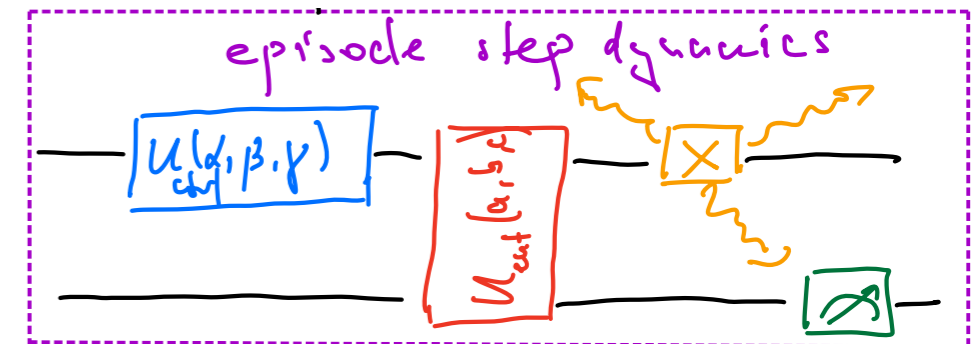
## → rewards

- qubit measurement output:  $\pm 1$  (binary)



## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$ : continuous (!)



## → states → observations

- wall clock time: one-hot representation of the step number
- ancilla measurement output:  $\pm 1$  (binary)

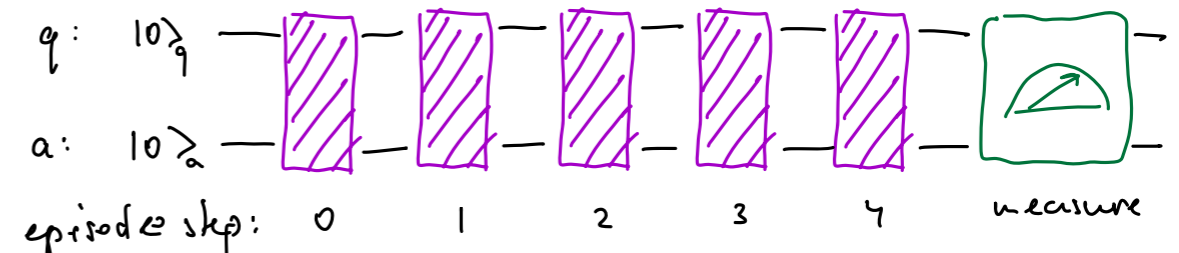
$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \end{pmatrix}$$

quantum states cannot be measured/observed!!!

# RL framework

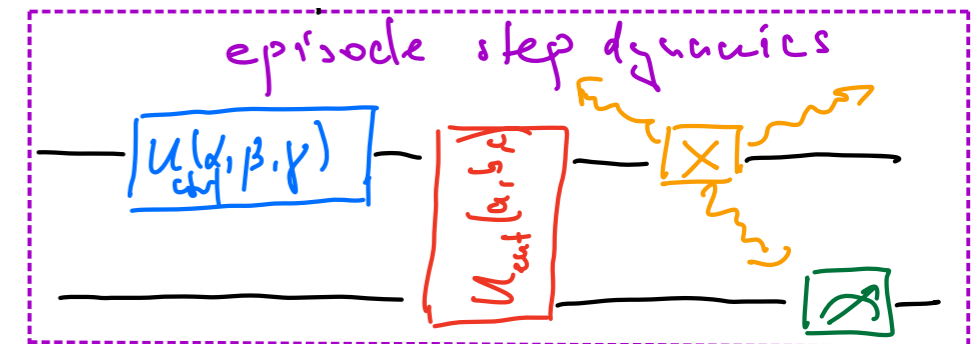
## → rewards

- qubit measurement output:  $\pm 1$  (binary)



## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$  : continuous (!)



## → states → observations

- wall clock time: one-hot representation of the step number
- ancilla measurement output:  $\pm 1$  (binary)
- spontaneously emitted photon detection:  $\pm 1$  (binary)

-1 = 'qubit found in GS'  
(photon detected)

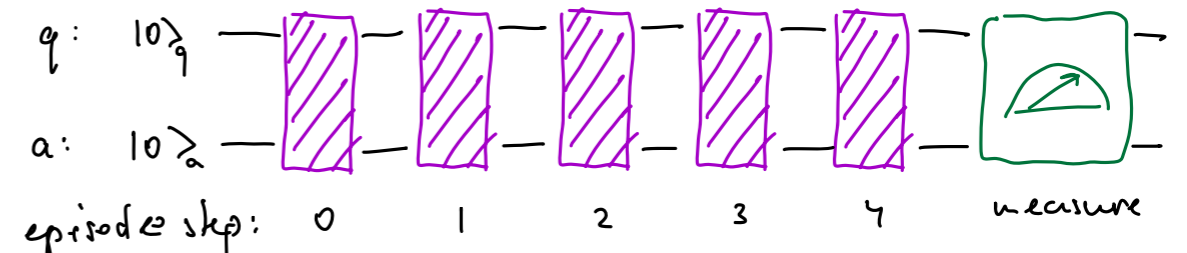
$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$$

**quantum states cannot be measured/observed!!!**

# RL framework

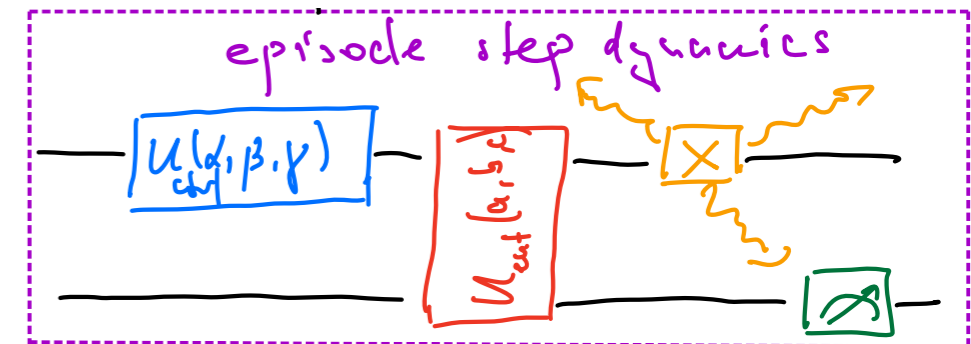
## → rewards

- qubit measurement output:  $\pm 1$  (binary)



## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$  : continuous (!)



## → states → observations

- wall clock time: one-hot representation of the step number
- ancilla measurement output:  $\pm 1$  (binary)
- spontaneously emitted photon detection:  $\pm 1$  (binary)

-1 = 'qubit found in GS'  
(photon detected)

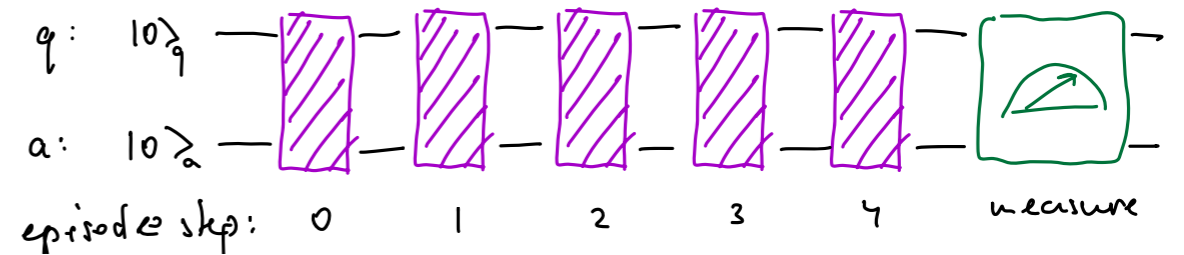
$$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ -1 \\ 1 \\ 1 \end{pmatrix} \quad \longrightarrow \quad \begin{pmatrix} 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \\ 0 & 1 & 1 \end{pmatrix} = o_t \quad \text{RL observation}$$

**quantum states cannot be measured/observed!!!**

# RL framework

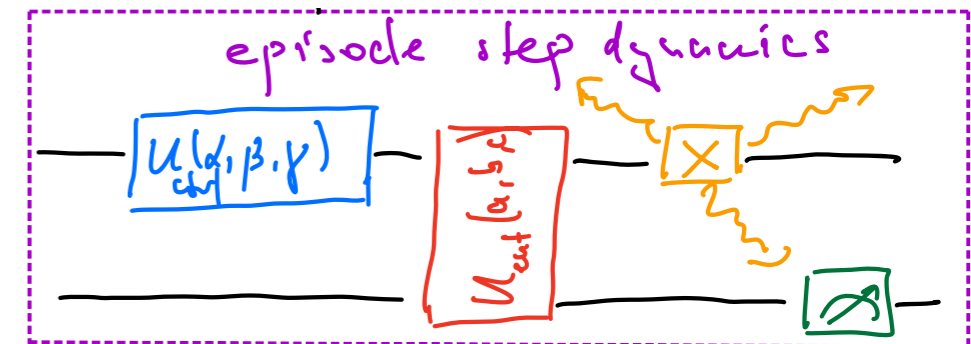
## → rewards

- qubit measurement output:  $\pm 1$  (binary)



## → actions

- angles of control unitary  $U(\alpha, \beta, \gamma)$ : continuous (!)



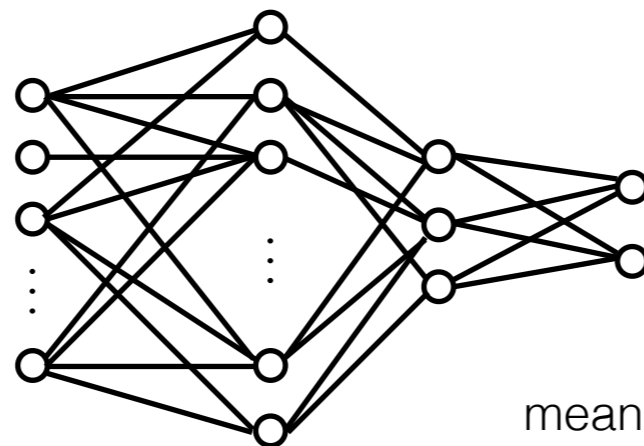
## → states → observations

- wall clock time: one-hot representation of the step number
- ancilla measurement output:  $\pm 1$  (binary)
- spontaneously emitted photon detection:  $\pm 1$  (binary)

## → agent

$$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

observation



neural net

means & variances

$$\begin{aligned} &\mathcal{N}(\mu_\alpha; \sigma_\alpha) \\ &\mathcal{N}(\mu_\beta; \sigma_\beta) \\ &\mathcal{N}(\mu_\gamma; \sigma_\gamma) \end{aligned}$$

sample angles

$$U_{\text{ctrl}}(\alpha, \beta, \gamma)$$



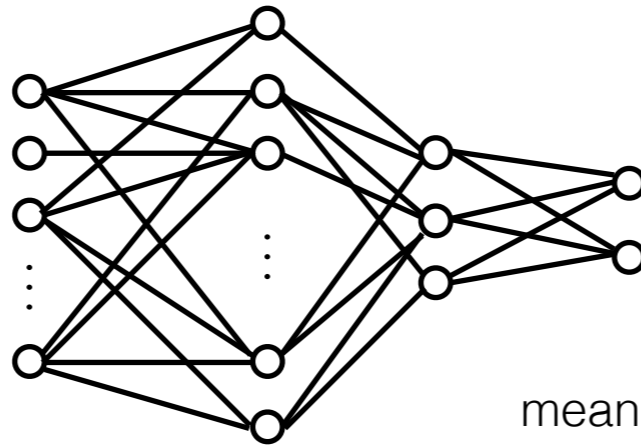
$$(\alpha, \beta, \gamma)$$

# Training curves

→ agent

$$\begin{pmatrix} 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

observation



neural net

$$\begin{matrix} \vec{\mu} \\ \vec{\sigma} \end{matrix}$$

means & variances

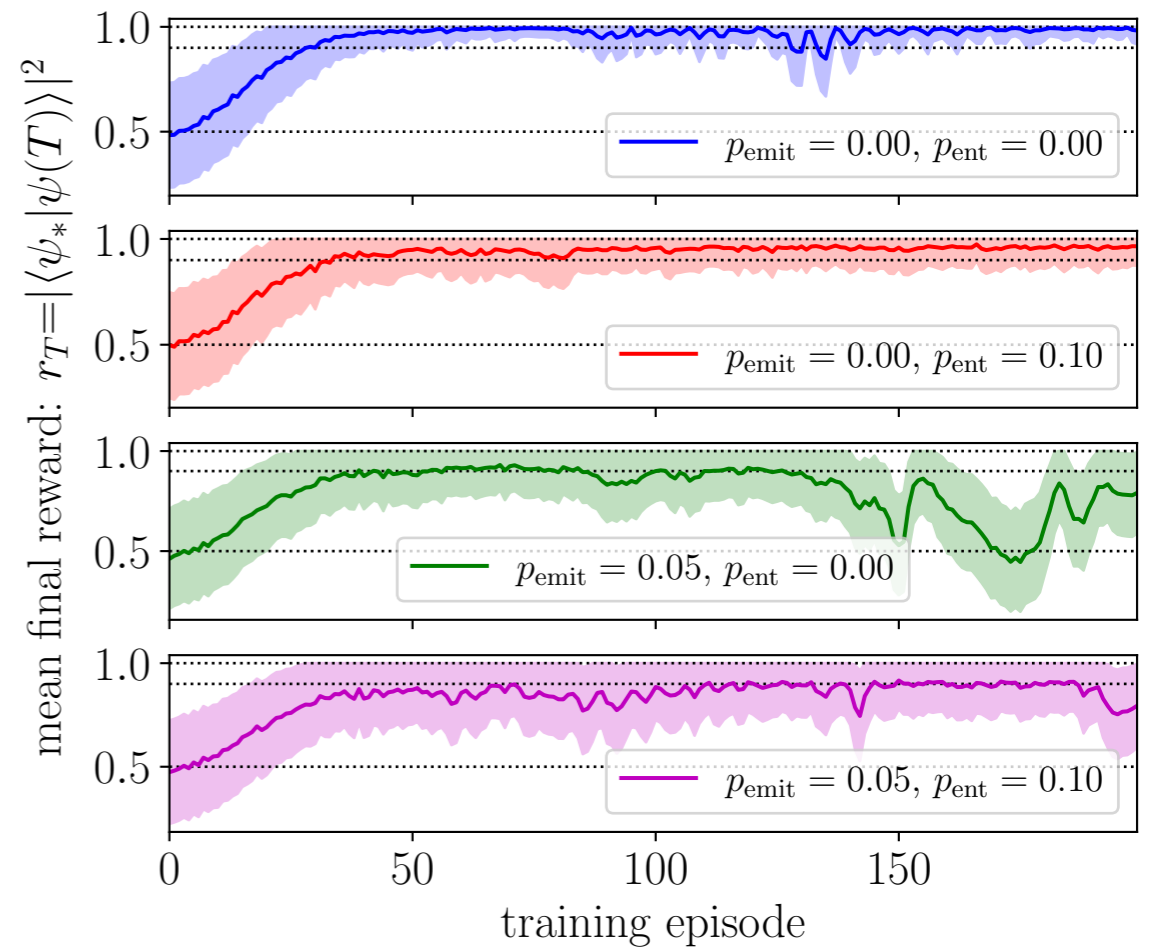
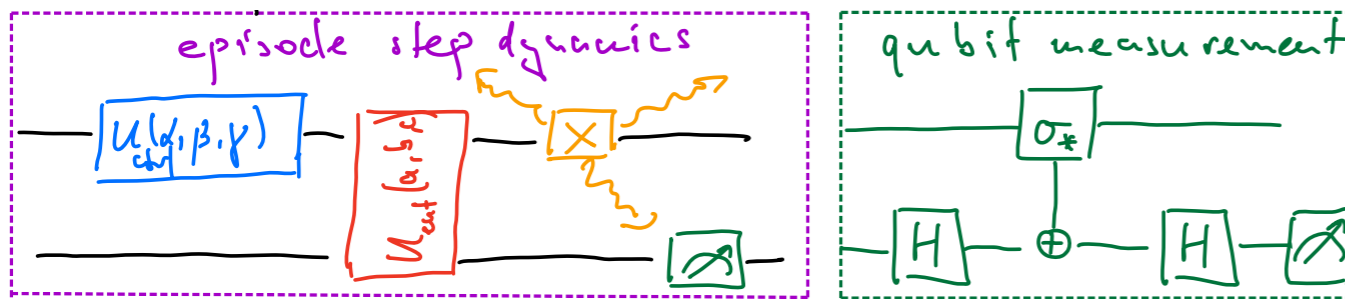
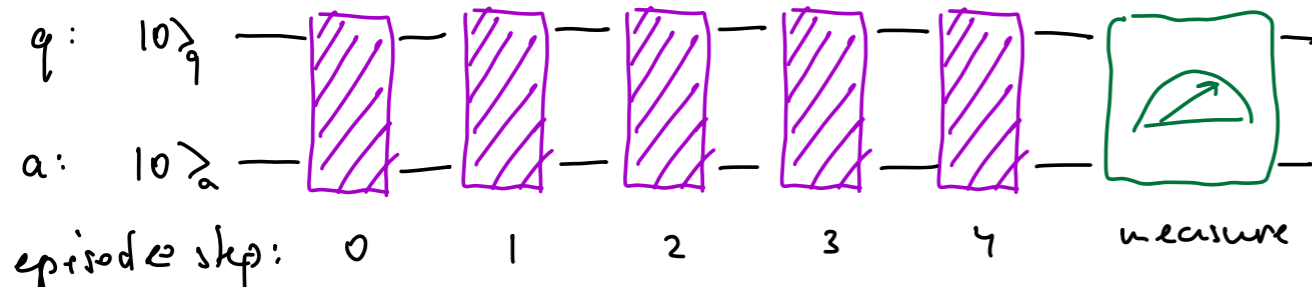


$$\begin{matrix} \mathcal{N}(\mu_\alpha; \sigma_\alpha) \\ \mathcal{N}(\mu_\beta; \sigma_\beta) \\ \mathcal{N}(\mu_\gamma; \sigma_\gamma) \end{matrix}$$

$$\longrightarrow (\alpha, \beta, \gamma)$$

sample angles

$$U_{\text{ctrl}}(\alpha, \beta, \gamma)$$





# Hands-on policy gradient

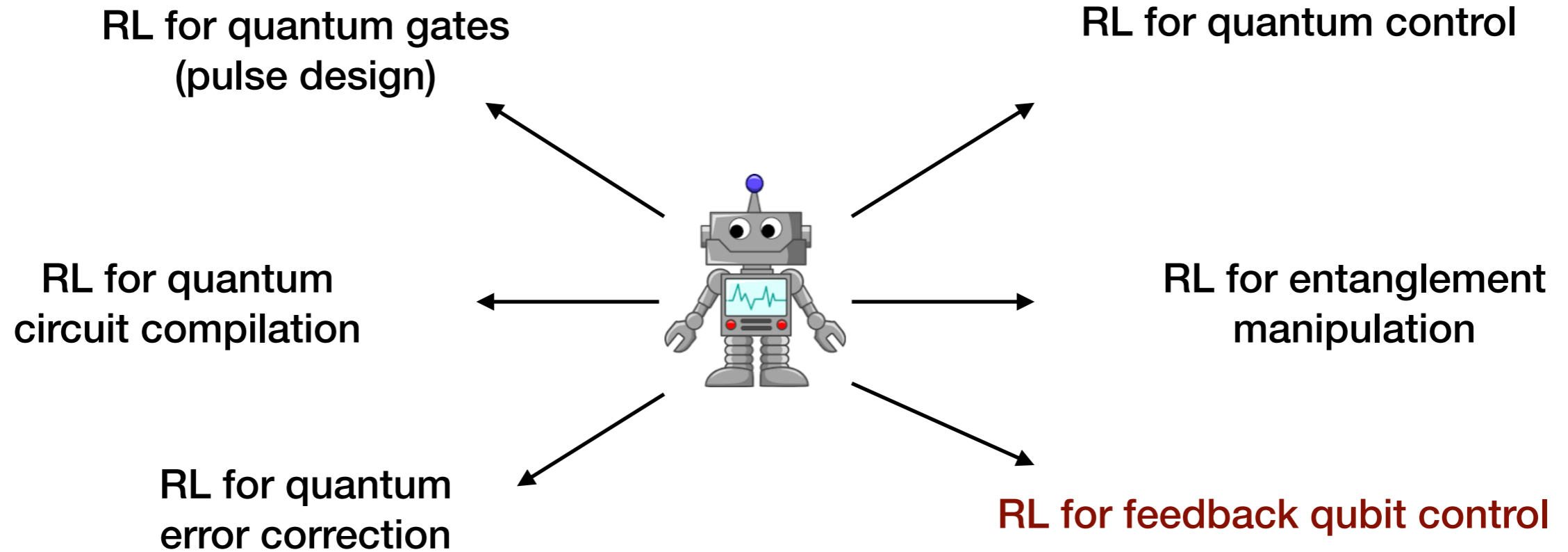
Check out Jupyter notebook for how this works in practice!

[https://github.com/mgbukov/RL\\_quantum](https://github.com/mgbukov/RL_quantum)





# RL in Quantum Physics: an Overview



Check out Jupyter notebook for how this works in practice!

[https://github.com/mgbukov/RL\\_quantum](https://github.com/mgbukov/RL_quantum)



*Thanks for your attention!*

# Useful Literature

M. Nielsen, *Neural Networks and Deep Learning* (online book)

Sutton and Barto, *Reinforcement Learning: an Introduction*, MIT press

S. Levine, *You Tube*, UC Berkeley (videos of lecture course)

M. Bukov, *lecture course*, Sofia University

[http://quantum-dynamics.phys.uni-sofia.bg/teaching/WiSe\\_2020\\_RL\\_class/](http://quantum-dynamics.phys.uni-sofia.bg/teaching/WiSe_2020_RL_class/)

Check out Jupyter notebook for how this works in practice!

[https://github.com/mgbukov/RL\\_quantum](https://github.com/mgbukov/RL_quantum)

