

The collage consists of five square images arranged in a grid. The top-left image shows an AMD AARTIX processor. The top-right image shows a server rack with multiple drives. The bottom-left image shows a close-up of a microchip with a blue light pattern. The bottom-center image shows a robotic arm interacting with a chip. The bottom-right image shows two Intel processors, one standing upright and one lying flat.

Christof Teuscher teuscher@pdx.edu www.teuscher-lab.com/teaching

Goal: get you started (if you are stuck)

1. Pick an interesting code/workload. E.g., see examples:
 1. QR code recognition
 2. Reinforcement learning (FrozenLake problem)
 3. DeepSpeech (speech-to-text engine, pre-trained)
2. Analyze and benchmark the code to identify possible bottlenecks.
3. Determine total execution time = baseline for HW/SW comparison.

teuscherLab
teuscher-lab.com

 Portland State
UNIVERSITY

Christof Teuscher teuscher@pdx.edu www.teuscher-lab.com/teaching

Scale the problems up to understand better

1. Recognize 1,000, 10,000, 100,000 QR codes
2. Enlarge the FrozenLake grid and place and add a fixed percentage of holes in the ice (e.g., 5%): 5x5, 10x10, 50x50, 100x100, 1,000x1,000 ...

Legends:

- A thick and safe layer of ice which you can walk over
- A hole in the ice. Fall here and you're dead
- The frisbee. Be the hero and get it back
- You

teuscherLab
teuscher-lab.com

Portland State
UNIVERSITY

 Christof Teuscher teuscher@pdx.edu	 www.teuscher-lab.com/teaching
<h1>Total execution time (1)</h1>	
<pre>import time start_time = time.time() # Code to be timed def my_function(): result = 0 for i in range(1000000): result += i return result my_function() end_time = time.time() execution_time = end_time - start_time print("Execution time: {execution_time} seconds")</pre>	<pre>Python import time start_time = time.perf_counter() # Code to be timed def my_function(): result = 0 for i in range(1000000): result += i return result my_function() end_time = time.perf_counter() execution_time = end_time - start_time print("Execution time: {execution_time} seconds")</pre>



Christof Teuscher teuscher@pdx.edu www.teuscher-lab.com/teaching

Total execution time (2)

```
Python

import timeit

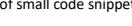
# Code to be timed as a string
code_to_time = """
def my_function():
    result = 0
    for i in range(1000000):
        result += i
    return result
my_function()
"""

execution_time = timeit.timeit(code_to_time, number=1)

print("Execution time: {} seconds".format(execution_time))
```

The `timeit` module is specifically designed for measuring the execution time of small code snippets.

teuscher Lab
teuscher-lab.com



Christof Teuscher teuscher@gdx.edu www.teuscher-lab.com/teaching

Total execution time (3)

```

lbootisec S10_HM_for_AI_and_ML_2025/python/q1 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 5.48999999998884e-07 seconds
Execution time 2: 5.48999999998884e-07 seconds
Execution time 3: 6.22000000004573e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q2 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 2.1457672117164025e-06 seconds
Execution time 2: 2.1457672117164025e-06 seconds
Execution time 3: 2.1457672117164025e-06 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q3 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 4.889999999799995e-07 seconds
Execution time 2: 4.889999999799995e-07 seconds
Execution time 3: 4.889999999799995e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q4 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 2.543999999799986e-06 seconds
Execution time 2: 2.543999999799986e-06 seconds
Execution time 3: 2.543999999799986e-06 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q5 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 9.53671464025e-07 seconds
Execution time 2: 9.53671464025e-07 seconds
Execution time 3: 9.53671464025e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q6 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 5.140000000078095e-07 seconds
Execution time 2: 5.140000000078095e-07 seconds
Execution time 3: 5.140000000078095e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q7 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 2.4000000000000002e-07 seconds
Execution time 2: 2.4000000000000002e-07 seconds
Execution time 3: 2.4000000000000002e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q8 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 5.359999999995429e-07 seconds
Execution time 2: 5.359999999995429e-07 seconds
Execution time 3: 5.359999999995429e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q9 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 1.987340432015e-06 seconds
Execution time 2: 1.987340432015e-06 seconds
Execution time 3: 1.987340432015e-06 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q10 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 5.4099999999904392e-07 seconds
Execution time 2: 5.4099999999904392e-07 seconds
Execution time 3: 5.4099999999904392e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q11 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 4.070000000083972e-07 seconds
Execution time 2: 4.070000000083972e-07 seconds
Execution time 3: 4.070000000083972e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q12 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 5.9509999999916032e-07 seconds
Execution time 2: 5.9509999999916032e-07 seconds
Execution time 3: 5.9509999999916032e-07 seconds
lbootisec S10_HM_for_AI_and_ML_2025/python/q13 cteuschechk python3 qr2_profile.py qr-code-21x21.png
Execution time 1: 1.192492495507812e-06 seconds
Execution time 2: 1.192492495507812e-06 seconds
Execution time 3: 4.520999999977778e-07 seconds

```

teuscherLab teuscher-lab.com

Portland State
UNIVERSITY

Christof Teuscher teuscher@pdx.edu

Total execution time (4)

```
# Execution time
start_time = time.time()
start_time = time.perf_counter()
start_time = time.default_timer()

repeats = 10
if __name__ == "__main__":
    import sys
    for i in range(repeats):
        if len(sys.argv) == 2:
            print("Usage: python qr_recognizer.py <image_path>")
            sys.exit(1)
        image_path = sys.argv[1]
        recognizer = QRCodeRecognizer()
        try:
            result = recognizer.recognize(image_path)
            print(result['code_recognition_result'])
            print(result)
            except Exception as e:
                print(f"Error: {e}")
        end_time = time.time()
        end_time = time.perf_counter()
        execution_time = end_time - start_time
        execution_time = time.default_timer() - start_time
        print(f"Execution time 1: {execution_time} seconds")
        print(f"Execution time 2: {execution_time} seconds")
        print(f"Execution time 3: {execution_time} seconds")
```

teuscher•Lab teuscher-lab.com

www.teuscher-lab.com/teaching

Total execution time (4)

Execution time 1: 1.8990371227264404 seconds
Execution time 2: 1.899044882999997 seconds
Execution time 3: 1.8990454230000002 seconds

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

Code analysis (1)

Static Code Analysis Tools

1. Using PyCallGraph

PyCallGraph is a tool that creates call graph visualizations for Python applications:

```
python
# Install with: pip install pycallgraph graphviz
from pycallgraph import PyCallGraph
from pycallgraph.output import GraphvizOutput

def analyze_qr_code():
    with PyCallGraph(output=GraphvizOutput()):
        recognizer = QRCodeRecognizer()
        result = recognizer.recognize("path_to_qr_image.png")
        print(result)

analyze_qr_code()
```

2. Using pyan for Static Analysis

Pyan analyzes Python code and generates approximate call graphs:

```
python
# Install with: pip install pyan
# Usage from command line
$ pyan qr_recognizer.py --dot > qr_graph.dot
$ dot -Tpng qr_graph.dot -o qr_graph.png
```

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

Code analysis (2)

Runtime Analysis Tools

3. Using cProfile with Visualization

```
python
import cProfile
import pstats
from pstats import SortKey

# Profile the execution
cProfile.run('recognizer.recognize("path_to_qr_image.png")')

# Analyze the results
p = pstats.Stats('qr_stats')
p.strip_dirs().sort_stats(SortKey.CUMULATIVE).print_stats(10)

# For visualization, you can use gprof2dot
# pip install gprof2dot
# gprof2dot -f pstats qr_stats | dot -Tpng -o qr_profile.png
```

4. Using SnakeViz for Interactive Visualization

```
python
# Install with: pip install snakeviz
# Run profiling
import cProfile
cProfile.run('recognizer.recognize("path_to_qr_image.png")', 'qr_stats'

# Then from command line:
# snakeviz qr_stats
```

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

Code analysis (3)

Performance Visualization

6. Visualize Performance Bottlenecks

To visualize execution time of different components:

```
python
import time
import numpy as np
import matplotlib.pyplot as plt

def benchmark_nr_components(image_path, iterations=10):
    recognizer = QRCodeRecognizer()
    components = {}

    # Benchmark image loading
    start = time.time()
    for _ in range(iterations):
        img = recognizer.load_image(image_path)
    components["load_image"] = (time.time() - start) / iterations

    # Benchmark thresholding
    start = time.time()
```

Using Specialized Tools

For more advanced analysis, consider these specialized tools:

- Doxxygen** - Can generate call graphs if properly configured
- CodeViz** - Works with C/C++ extensions but can be used on some Python projects
- Memray** - For memory profiling with visualization capabilities
- Py-Spy** - For sampling profiler that can generate flame graphs

Each of these methods offers different insights into your code's structure and performance. Combining a few of them will give you the most comprehensive understanding of your QR code recognizer's data and control flow.

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

Christof Teuscher teuscher@pdx.edu

cProfile and SnakeViz

1. pip install snakeviz
2. python3 -m cProfile -o output.prof qr_recognizer.py qr.png
3. snakeviz output.prof

teuscher•Lab teuscher-lab.com

www.teuscher-lab.com/teaching

cProfile and SnakeViz

teuscher•Lab teuscher-lab.com

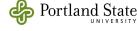
Christof Teuscher teuscher@pdx.edu

www.teuscher-lab.com/teaching

teuscher•Lab teuscher-lab.com

Portland State UNIVERSITY

teuscher Lab
teuscher-lab.com



```
1. kernprof -l q_learning_profiling.py
2. python3 -m line_profiler -rmt "q_learning_profiling.py.lprof"
```

```
#!/usr/bin/python3
# @profile
def q_learning(episodes):
    n = 8
    # loop through best path for each episode
    while(episodes < episodes):
        episode = end
        if episode == 1:
            # get reward and add to array for plot
            reward = self.State.getReward()
            self.plot(reward)
            self.plot_reward.append(reward)

            # assign reward to each Q_value in state
            for i in range(n):
                for a in self.actions:
                    self.Q[i][a] = round(reward,3)

            #reset state
            self.State.state = self.State.state_idn

            #set rewards to zero and iterate to next episode
            self.rewards = 0
        else:
            to_arbitrary_low_value = compare net state actions
            next_state = next_state, next_action and current reward
            next_state, next_action = self.getAction()
            self.State.state = next_state
            reward = self.State.getReward()
            self.plot(reward)
            self.rewards += reward
```

```

Christof Teuscher · teuscher@pdx.edu · www.teuscher-lab.com/teaching

Total Lines: 1,9371
File: a_learning_profiling.py
Function: Q_Learning at line 141

Line #   Hits    Time  Per Hit % Time Line Contents
=====  ======  ======  ======  ======  ======
141
142
143      1      1.0      1.0     0.0
144
145 94289 24463.0     0.3     1.3
146
147 94288 22656.0     0.2     1.2
148
149 18000 18360.0     1.4     0.7
150 18000 3836.0      0.4     0.2
151 18000 8266.0      0.6     0.3
152
153
154 18000 3839.0     0.3     0.2
155 58000 12719.0     0.3     0.7
156 48000 26152.0     0.7     1.4
157
158
159 18000 9741.0     1.0     0.5
160 18000 2886.0      0.4     0.2
161
162
163 18000 2983.0     0.3     0.2
164 18000 3693.0      0.6     0.2

                                         #profile
def Q_Learning(self,episodes):
    for e in range(episodes):
        #iterate through best path for each episode
        while(x < episodes):
            #check if state is end
            if self.isEnd:
                #get current reward and add to array for plot
                self.rewards += self.currentReward
                self.rewards -= self.lastReward
                self.plot_reward(self.rewards)

                #get state, assign reward to each Q_value in state
                i,j = self.State.state
                for a in self.actions:
                    self.new_Q[(i,j,a)] = round(reward,3)

            #reset state
            self.State = State()
            self.isEnd = self.State.isEnd

            #set rewards to zero and iterate to next episode
            self.rewards = 0
            x+=1

```

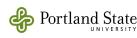
teuscher Lab
teuscher-lab.com



 teuscher  Lab
teuscher-lab.com

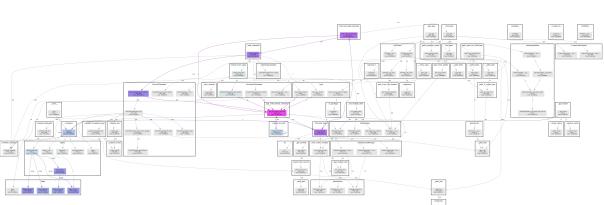
 Portland State
UNIVERSITY

teuscher Lab
teuscher-lab.com

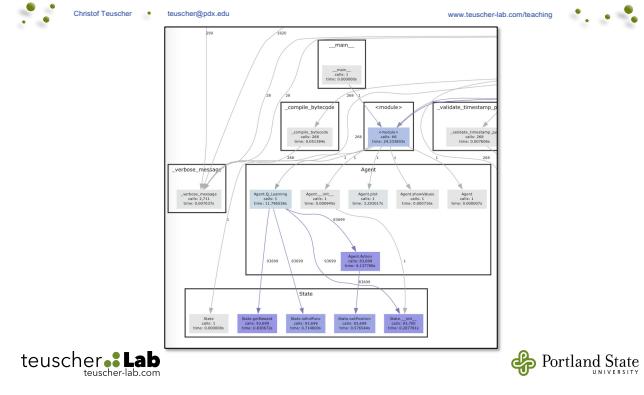


Pycallgraph2: Dynamic Call Graph Generation

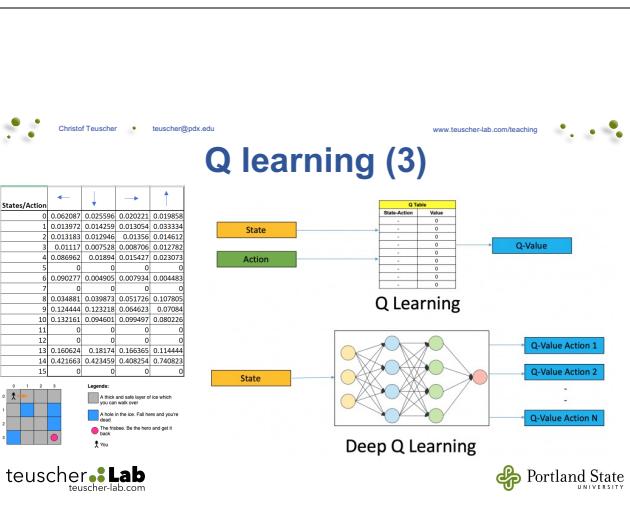
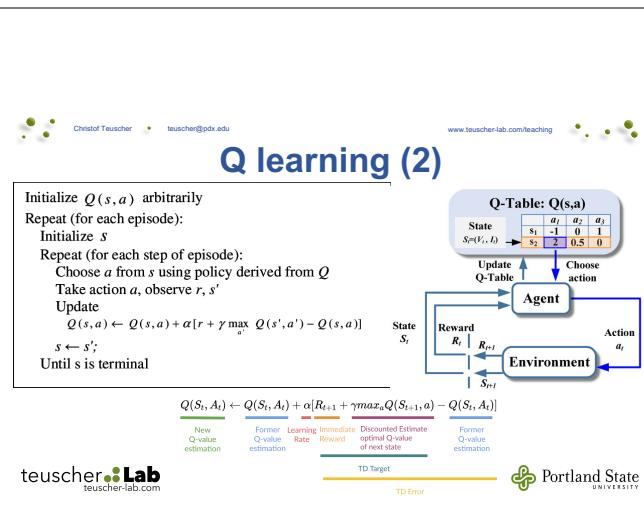
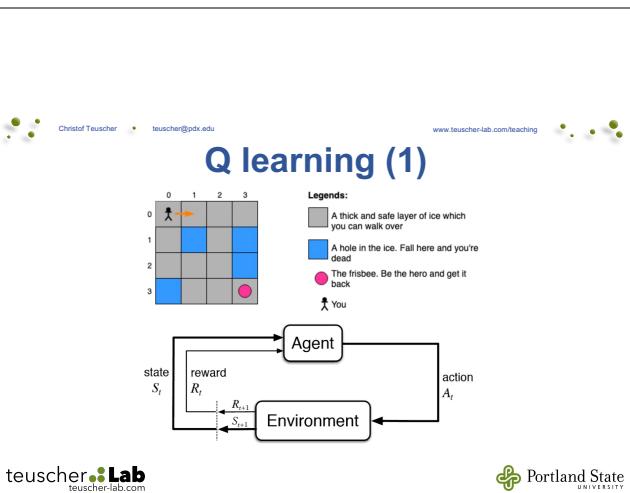
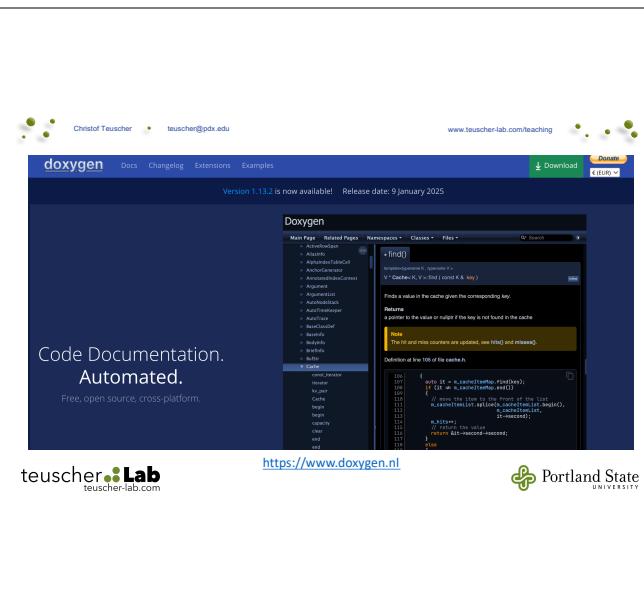
1. Requires `pycallgraph2` and `graphviz`
 2. `pycallgraph graphviz --output-file output.png a_learning_orig.py`



teuscherLab
teuscher-lab.com



ncalls	file	fn	ttime	percall	cumtime	percall	filename:lineno(function)
1			3.438	3.438	3.743	3.743	<@<built-in method matplotlib.backends._macosx.show>
46/44			0.7926	0.01801	0.7994	0.01817	<@<built-in method _mp_create_dynamic>
1			0.4153	0.4153	0.9049	0.9049	q_learning.py 141(Q_Learning)
83522			0.1256	1.954e-06	0.2928	3.506e-06	q_learning.py 111(Action)
292			0.09859	0.0003376	0.09859	0.0003376	<@<method 'read' of '_io.BufferedReader' objects>
1			0.07742	0.07742	0.1346	0.1346	backend_bases.py 2673(create_with_canvas)
7178			0.06007	8.359e-06	0.107	1.491e-05	inspect.py 626(cleanups)
8403			0.05586	6.647e-06	0.1172	1.394e-05	<@<method 'choice' of 'numpy.random.mtrand.RandomState' objects>
123364			0.055	4.451e-07	0.055	4.451e-07	<@<built-in method builtins.round>
292			0.04718	0.0001616	0.04718	0.0001616	<@<built-in method marshal.loads>
115			0.04713	0.0004098	0.04713	0.0004098	<@<method 'set_message' of 'matplotlib.backends._macosx.NavigationTools'
93522			0.04602	4.912e-07	0.04602	4.912e-07	q_learning.py 35(getReward)
7			0.04631	0.006673	0.04541	0.006488	<@<method 'draw_path' of 'matplotlib.backends._backend_agg.RendererAgg'
83522			0.04079	4.884e-07	0.04079	4.884e-07	q_learning.py 53(nextPosition)



Efficient HW Q learning implementation

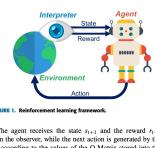


FIGURE 1. Reinforcement learning framework.

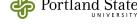
 $z^{-1} = \text{delay registers}$ <https://doi.org/10.1109/ACCESS.2019.2961174>

FIGURE 2. High level architecture of the Q-Learning agent.

 $z^{-1} = \text{delay registers}$

Efficient HW Q learning implementation

Figure 3 shows the Q-Learning accelerator. The Q-Matrix is stored into Z Dual-Port RAMs, named Action RAMs. Consequently, we have one memory block per action. The RAMs are addressed by the current state s_{t+1} and the number of memory locations corresponds to the number of states N . The read address is the next state s_{t+1} , while the write address is the current state s_t . The enable signals for the Action RAMs, generated by a decoder driven by the current action a_t , select the value $Q(s_t, a_t)$ to be updated. The Action RAMs outputs correspond to a row of the Q-Matrix $Q(s_{t+1}, A)$.

The signal $Q(s_t, a_t)$ is obtained by delaying the output of the memory blocks and then selecting the Action RAM through a multiplexer driven by a_t . A MAX block fed by the output of the Action RAMs generates $\max(Q(s_{t+1}, A))$.

The Q-Update block implements the Q-Matrix update equation (1) generating $Q_{\text{new}}(s_t, a_t)$ to be stored into the corresponding Action RAM.

The accelerator can be also used for Deep Q-Learning [23] applications if the Action RAMs are replaced with Neural Network-based approximators.

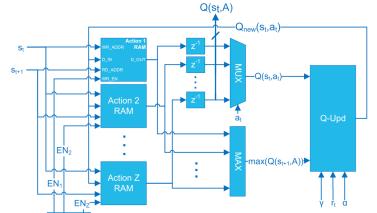
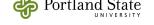
<https://doi.org/10.1109/ACCESS.2019.2961174>

FIGURE 3. Q-Learning accelerator architecture.

Efficient HW Q learning implementation

$$Q_{\text{new}}(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_a Q(s_{t+1}, a) \right) \quad (1)$$

The variables in (1) refer to:

- s_t and s_{t+1} : current and next state of the environment.
- a_t and a_{t+1} : current and next action chosen by the agent (according to its policy).
- γ : discount factor, $\gamma \in [0, 1]$. It defines how much the agent has to take into account long-run rewards instead of immediate ones.
- α : learning rate, $\alpha \in [0, 1]$. It determines how much the newest piece of knowledge has to replace the older one.
- r_t : current reward value.

How many multipliers?

<https://doi.org/10.1109/ACCESS.2019.2961174><https://doi.org/10.1109/ACCESS.2019.2961174>

Efficient HW Q learning implementation

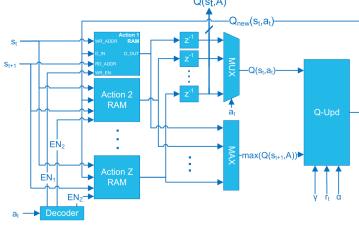
FIGURE 4. MAX block tree architecture for $Z = 6$.

FIGURE 5. Q-matrix updater block architecture.

<https://doi.org/10.1109/ACCESS.2019.2961174>

Efficient HW Q learning implementation

1) APPROXIMATED MULTIPLIERS

The main speed limitation in the update block is the propagation delay of the multipliers. Using a similar approach to [32], it is possible to replace the full multipliers shown in Fig. 5 with approximated multipliers based on barrel shifters [33]. In this way, we are approximating α and γ with a number equal to their nearest power of two (single shifter), or to the nearest sum of powers of two (two or more shifters). Due to the fact that $\alpha, \gamma \in [0, 1]$, only right shifts have been used.

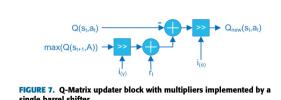


FIGURE 7. Q-Matrix updater block with multipliers implemented by a single barrel shifter.

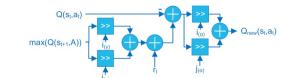


FIGURE 8. Q-Matrix updater block with multipliers implemented by two barrel shifters.

<https://doi.org/10.1109/ACCESS.2019.2961174>