

Institutsprojekt – A Distributed Q-Learning Approach for Adaptive Sleep Modes in 5G Networks

Im Rahmen des im 4. Semester des Studiengangs Elektrotechnik, Informationstechnik und Technische Informatik, obligatorischen Institutprojekts hatten wir die Möglichkeit uns genauer mit dem Thema „A Distributed Q-Learning Approach for Adaptive Sleep Modes in 5G Networks“ am Institut für Intelligente Regelungssysteme der RWTH auseinander zu setzen. Das Projekt wurde von uns, einem sechsköpfigen Team bearbeitet und von Oberingenieur Dr. Ing. Michael Reyer und von Universitätsprofessor Dr.-Ing. Christian Ebenbauer betreut. Es basiert auf einem gleichnamigen Paper von Ali El-Amine, Mauricio Iturralde, Hussein Al Haj Hassan and Loutfi Nuaymi. Im Folgenden schildern wir unsere Erwartungen, Erfahrungen, Schwierigkeiten sowie Erkenntnisse, während des Projektes.

Ziele welche wir uns für das Institutsprojekt setzten

Zu Beginn sei zu erwähnen, dass wir uns zusammen mit unserem Betreuer Dr. Michael Reyer darauf verständigt haben, in einer größeren Gruppe (sechs Personen) zu arbeiten. Grund hierfür war es, dass eines der für uns wichtigsten Ziele es war Einblicke in die Arbeit als Ingenieure zu erlangen. Schon während unseres bisherigen Studiums stellte sich heraus, dass man auf Hilfe und Unterstützung von Kommilitonen angewiesen ist. Diese Erkenntnis lässt sich auch auf den Arbeitsalltag von Ingenieuren und Ingenieurinnen übertragen. Auch hier ist es essenziell, das Arbeiten in einem Team zu erlernen, unterschiedliche Erfahrungen zu machen sowie Schwierigkeiten als auch Vorteile festzustellen und zu erlernen. Daher standen wir alle dem Vorschlag, in einer größeren Gruppe zu arbeiten, positiv gegenüber.

Des Weiteren hielten wir es schon zu Beginn des Projekts für wichtig, nicht lediglich einen Programmcode für ein gegebenes Problem zu schreiben, sondern setzten auch einen Schwerpunkt darauf, das gegebene Problem wirklich zu verstehen, um uns mit der Problemstellung intensiv beschäftigen zu können. Wir planten daher schon anfangs Zeit dafür ein das gegebene Paper zu durchdringen, Sekundärliteratur zu recherchieren als auch ausführlich das Paper im kollektiv zu besprechen und unterschiedliche Lösungsansätze zu diskutieren.

Ferner war es für uns von größter Wichtigkeit, nicht nur die Problemstellung ausführlich zu verstehen, sondern ebenfalls uns intensiv mit dem maschinellen Lernen, speziell dem Q-Learning zu beschäftigen.

Zudem war natürlich die Implementation einer Simulationsumgebung in MATLAB ausschlaggebend, um unsere Überlegungen und Thesen zu testen, zu überprüfen und zu dokumentieren. Hier sahen wir es ebenfalls als Ziel an, auch abseits vom Thema unsere Fähigkeiten in MATLAB zu verbessern, da wir es als wichtigen Kenntnis für das spätere Arbeitsleben als Ingenieur ansehen.

Abschließend war es uns auch wichtig, im Rahmen des Projekts viel Spaß miteinander zu haben und neue Freundschaften zu knüpfen, was sich im Nachhinein auch bestätigte.

Problemstellung

Jede Basisstation hat einen gewissen Energieverbrauch, welcher reduziert werden soll, möglichst ohne Performanz zu verlieren. Wir machen uns zunutze, dass es auch Zeitpunkte gibt, an denen eine Basisstation keinen Consumer bedienen muss. Während diesen Perioden ist es somit nicht nötig dauerhaft im Ruhemodus zu bleiben, da dies Energie verbraucht. Wir führen daher sogenannte Schlafmodi ein, in welche jede Basisstation wechseln kann, um den Energieverbrauch runterzuschrauben. Unsere Simulation modelliert drei Schlafmodi, wobei die Energieersparnis vom ersten zum dritten Schlafmodus zunimmt. Auf der anderen Seite benötigen die Basisstationen nun eine gewisse Zeit, um aus diesem Energiesparmodus wieder in den aktiven Zustand zu wechseln und Consumer zu bedienen. Dies führt zu einer Zeitverzögerung der Bedienung (Delay).

Es stellt sich somit die Frage welchen Modus jede Basisstation wählen soll, um bei möglichst guter Performanz (geringer Delay) möglichst wenig Energie zu verbrauchen.

Maschinelles Lernen (Q-Learning)

Einen Verbesserungsansatz, welchen wir im Rahmen des Projekts untersucht haben, um das gegebene Problem zu optimieren, ist die Optimierung mithilfe von maschinellem Lernen, genauer Q-Learning.

Maschinelles Lernen bzw. Reinforcement Learning ist ein Konzept zur Verbesserung von Algorithmen. Hierbei wird ein Agent ohne Kenntnisse seiner Umgebung und den Einfluss bzw. die Folgen seiner Aktionen mit dem Problem konfrontiert. Dieser Agent soll im Folgenden sein Handeln darauf trainieren, eine möglichst große Belohnung zu erhalten.

Zu erwähnen sei, dass es mehrere Arten gibt, das Reinforcement Learning System zu trainieren, wir uns im Rahmen des Projekts schwerpunktmäßig jedoch auf einen fundamentalen, relativ simplen Algorithmus konzentriert haben, dem Q-Learning.

Ziel des Q-Learning Ansatzes ist es eine optimale Policy zu erhalten. Hierunter versteht man das gelernte Verhalten des Agenten, welches im sagt, welche Aktion er in einem bestimmten Zustand in einer bestimmten Umgebung ausführen soll. Diese Policy wird mittels sogenannter Q-Values (Formel (2)) in einer Matrix gespeichert, welche die zu erwartenden Belohnungen darstellen. Der Index m steht hierbei für die den jeweiligen Agenten (hier: Basisstation) und der Index t stellt den Zeitschritt dar.

Diese Matrix ist so aufgebaut, dass jede Zeile für eine bestimmte Beobachtung steht und jede Spalte für eine auszuführende Aktion. Diese Matrix bzw. ihre Q-Values werden während der Durchführung der Simulation immerzu aktualisiert, um, wie bereits oben erwähnt, eine optimale Lösung eines gegebenen Problems zu erlangen. Die Aktualisierung geschieht wie in Formel (2) zu sehen, anhand des vorherigen Q-Values und den neu gemachten Erfahrungen bei der Erkundung der Umgebung des Agenten. Hierbei regelt der Parameter Alpha ($\alpha \in (0,1)$; Learning Rate), inwieweit bzw. wie schnell der Algorithmus aus der aktuellen Aktion Schlüsse zieht und der Diskontinuitätsfaktor Gamma ($\gamma \in [0,1)$) regelt, wie stark die Zukunft in unsere Bewertung einfließen soll. Je höher der Wert des Parameters desto größer ist der Einfluss der zukünftigen Werte.

Zu Beginn wird unsere Q-Matrix als Null-Matrix initialisiert und der Algorithmus springt in eine Schleife, welche bei dem Zeitparameter t gleich Null startet. Es wird gemäß Formel (1) eine aktuelle Aktion ausgewählt. Dies geschieht zunächst anhand Formel (1).

$$a_m^t = \begin{cases} \arg \max_a Q(s_m^t, a), & \text{wenn } y > \epsilon \\ \text{rand}(\mathcal{A}), & \text{sonst} \end{cases}, m = 1, \dots, M. \quad (1)$$

Nachdem nun die aktuelle Aktion a_m^t ermittelt wurde, wird diese ausgeführt. Es wird somit anhand Formel (2) einem Tupel s_m^t und a_m^t zum Zeitpunkt t ein zukünftiger Zustand s_m^{t+1} ermittelt. Das kleine m steht hier für die Basisstation. Auf die Bedeutung von r_m^t wird weiter unten genauer eingegangen.

Schließlich werden die Q-Values anhand Formel (2) aktualisiert.

$$Q(s_m^t, a_m^t) = Q(s_m^t, a_m^t) + \alpha \left[r_m^t + \gamma \max_a Q(s_m^{t+1}, a) - Q(s_m^t, a_m^t) \right] \quad (2)$$

Des Weiteren hat der Agent beim Q-Learning stets zwei Möglichkeiten des Lernens. Entweder folgt unser Agent den bereits erkundeten und für besten identifizierten Pfad weiter oder er beschreitet einen neuen Weg, um entweder noch nicht optimalen oder sich ändernden Bedingungen folgen zu können. Dieser Ansatz nennt sich Exploration.

Der Trade-Off zwischen diesen beiden Möglichkeiten ist in unserem Fall mittels der „epsilon-greedy“-Methode implementiert. Hierbei justiert ein Parameter Epsilon ($\epsilon \in [0,1]$) den Trade-Off zwischen dem weiteren Handeln nach der verfolgten Strategie und einer zufälligen Aktion.

Unser Epsilon Parameter gibt somit an, mit welcher Wahrscheinlichkeit der Algorithmus eine zufällige Aktion wählt oder mit welcher Wahrscheinlichkeit er einen bereits beschrittenen Pfad weiter folgt (1). Dabei wird $y \in [0,1]$ jedes Mal, gleichverteilt zufällig, neu bestimmt. Somit sollte ein Anteil Epsilon der Aktionen zufällig, also explorativ, sein. Die rand()-Funktion haben wir ebenfalls als gleichverteilte Zufallsfunktion implementiert.

Zustandsraum:

Der Zustandsraum entspricht den verschiedenen Zuständen der Basisstation während ihrer Inaktivitätsphase. Er kann sich in einem der verschiedenen Schlafmodi SM_1, SM_2 oder SM_3 oder in keinem davon befinden (Ruhezustand).

$$\mathcal{S} = \{\text{idle}, SM_1, SM_2, SM_3\}$$

Aktionsraum:

Der Aktionsraum enthält die möglichen Entscheidungen, die an jedem Entscheidungspunkt zu treffen sind. Die Aktion besteht darin, n_i zu entscheiden, wie viel Zyklen die Basisstation in der folgenden Schlafmodusstufe bzw. Ruhezustand verbleiben kann. Wir bezeichnen mit N_i die Menge der möglichen Werte, die n_i annehmen kann, für $i \in \{0, 1, 2, 3\}$.

$$\mathcal{A} = N_3 \times N_2 \times N_1 \times N_0$$

Zur Diskretisierung setzen wir die Menge N_i wie folgendes:

$$N_i = \{n_i = \lfloor x_i \cdot n \rfloor\}$$

mit $x_i = \frac{\tau}{\theta_i}$, τ ist die durchschnittliche Wartezeit zwischen zwei Nutzern, θ_i ist die Dauer einer Schlafzyklus von SM_i , und n nehmen wir aus den folgenden Werten:

$$n \in \{0.01, 0.02, 0.03, \dots, 0.09, 0.1, 0.2, 0.3, \dots, 0.9, 1, 2, 3, \dots, 9\}$$

Gestaltung der Simulation

Orientiert am Paper und in Absprache mit unserem Betreuer Herrn Reyer haben wir uns auf folgende Struktur und Funktionsweise unserer Simulation geeinigt.

Wir verwenden eine zweidimensionale Karte, auf welcher Consumer und Basisstationen platziert werden. Diese werden jeweils durch ein der Klasse spezifisches Icon angezeigt und geben Auskunft über deren Status.

Die Basisstationen bekommen im Vorhinein eine Position zugewiesen und sind in jeder unserer Simulationen gleich. Sie haben eine feste Bandbreite, welche gleich auf die Consumer aufgeteilt wird, dabei wird die Beanspruchung in Prozent gerechnet.

Die Schlafmodi werden wie im Paper angegeben übernommen. Somit haben wir drei verschiedene Schlafmodi, welche die auf der folgenden Abbildung beschriebenen Eigenschaften besitzen. Der Vierte bleibt ungenutzt, da er, wie im Paper auch festgestellt wurde, zu lang ist, als dass die Basisstation noch als aktiv von den Consumern angesehen werden kann. Außerdem hat dieser nur wenig zusätzlichen Energiegewinn.

Sleep level	Deactivation duration	Minimum sleep duration	Activation duration
idle	0 s	0.1 s	0 s
SM₁	35.5 μ s	71 μ s	35.5 μ s
SM₂	0.5 ms	1 ms	0.5 ms
SM₃	5 ms	10 ms	5 ms

Tabelle 1

Hierbei entscheidet dann der Agent, falls die Basisstation im Ruhezustand ist, wie sie fortfahren sollte anhand des Q-Learning Algorithmus.

Die Anzahl der Consumer wird vorher festgelegt und diese anschließend über Ort und Zeit gleichverteilt. Beim Erscheinen werden sie dann automatisch der Basisstation mit kleinstem Abstand zugeordnet.

Jeder Benutzer fragt eine gewisse Menge Daten an, welche dann von der Basisstation mit der lieferbaren Datenrate übertragen werden. Die Übertragungsrate wird aber nicht nur durch die Menge der Consumer pro Basisstation skaliert, sondern auch durch die Übertragungsqualität, welche mit Hilfe der Signal-to-Interference-plus-Noise Ratio (SINR) berechnet wird.

$$\text{SINR}_m(k) = \frac{P_m^{\text{Tx}} h_m(k)}{\sigma^2 + \sum_{m' \in \mathcal{S}, m' \neq m} P_{m'} h_{m'}(k)} \quad (3)$$

Hierbei ist P die Übertragungsleistung der jeweiligen Basisstation, bei uns sind diese alle gleich, $h_m(k)$ ist die Kanalverstärkung, welche den Pfadverlust und „shadowing-effect“ repräsentiert. Zuletzt ist σ^2 die additive Rauschleistung.

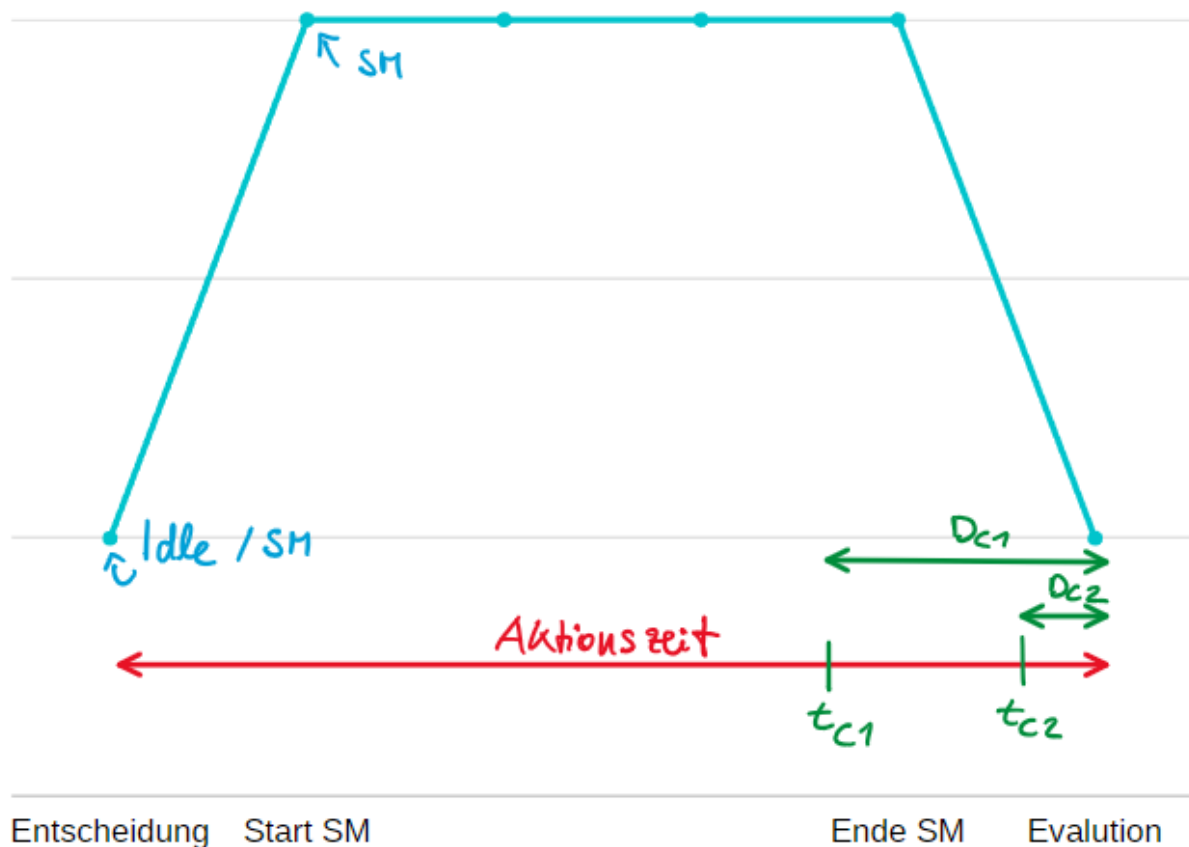
Somit kommen wir zur gelieferten Übertragungsrate durch das Shannon-Hartley Theorem, welches wie folgt lautet:

$$R_m(k) = \alpha \times W \times \log_2(1 + \text{SINR}_m(k)) \quad (4)$$

dabei repräsentiert W die Bandbreite und α skaliert diese auf den Anteil, der genutzt wird, um die Daten zu übertragen.

In unserer Simulation haben wir schließlich den Q-Learning-Algorithmus wie oben beschrieben implementiert.

Jede Basisstation wird von einem individuellen, von den anderen Basisstationen unabhängigen Agenten verwaltet. Somit spiegelt sich die Anzahl der Basisstationen im Parameter m wider. Die Aktionen, welche der Agent durchführt, sind in unserer Simulation durch die Wechsel der Schlafmodi repräsentiert, während die Zustände durch den aktuellen Modus dargestellt werden. Je größer der Q-Wert einer Aktion in dem aktuellen Zustand, desto eher wird die Basisstation in diesen Schlafmodus tendieren zu wechseln. Ziel jeder Basisstation ist es den Reward, welcher nach jeder abgeschlossenen Aktion, wie in der folgenden Abbildung beschrieben, berechnet wird, zu maximieren.



Wir definieren die Belohnung, analog wie im Paper, als die gewichtete Summe des Energiegewinns G und der zusätzlichen Verzögerung D , die beide aus dem während einer Episode gewählten Schlafmodus resultieren. (Formel (5))

$$r = (1 - \eta)G - \eta D \quad (5)$$

Anhand des Parameters η wird der Trade-Off zwischen Energieverbrauch und Delay justiert. Der Delay ist hier die Zeit, die ein Consumer warten muss, bis seine Basisstation aus dem Schlafmodus erwacht ist und ihn wieder bedienen kann.

Abschließend sei zu erwähnen, dass wir uns dazu entschieden haben die Parameter, bis auf das Shadowing, für die Simulation analog zu denen aus dem gestellten Paper zu verwenden, siehe Tabelle 2. Ferner setzten wir die Users arrival anders als in der Tabelle angegeben fest.

Parameter	Wert
Basisstation Tx Leistung	45 dBm
Bandbreite	20 MHz
Thermisches Rauschen	-174 dBm/Hz
Pfadverlust	$128.1 + 37.6 \log_{10}(d) \text{ dB}$
Ankunft des Consumers	Gleichverteilt über Ort und Zeit
Datenanfrage	Weibull Verteilung mit Parametern λ und k
Scale Parameter	$\lambda = 441.305$
Shape Parameter	$k = 0.8$

Tabelle 2

Um den zeitlichen Ablauf zu simulieren, jedoch nicht diese Zeit warten zu müssen, entschieden wir uns, diese durch Events zu modellieren. Dies wird durch das triggern von Events realisiert. Ein Event ist dabei das Eintreffen eines Consumers, der Beginn und das Ablaufen seiner Übertragungszeit, aber auch das Aufwachen und Einschlafen einer Basisstation. Diese sorgen dann für eine Neuevaluierung der Situation und passt die Daten, zum Beispiel wann ein Consumer mit seiner Datenübertragung fertig wird, an, falls nötig. Basisstationen überprüfen dabei immer beim Ende einer Übertragung, ob es die Möglichkeit zu schlafen gibt und wenn ja, wie lange.

Der Weg zur Implementierung in MATLAB

Hier wurde uns die Wahl gelassen, welche Umgebung wir genau benutzen wollen, um unsere Simulation zu implementieren. Nachdem wir uns ausgetauscht hatten, wer welche Vorkenntnisse in der Programmierung besitzt und welche Sprache welche Vorteile bietet, wurde uns schnell klar, dass es auf die MATLAB, hinauslaufen wird. Aufgrund der Nutzung von MATLAB in einem Lego-Mindstorm-Projekts im 1. Semester, als auch in anderen Fächern waren wir mit er MATLAB Umgebung schon vertraut.

Ein zusätzlicher Vorteil, welcher die Benutzung von MATLAB unserer Meinung nach hat, ist, dass es sich um einen weit verbreiteten Industriestandard handelt, wodurch sich die Vertiefung/Verbesserung unserer MATLAB-Fähigkeiten doppelt lohnt. Ferner ermöglicht es MATLAB, die Implementierung, als auch durch die Nutzung des App-Designers die Simulation und die GUI zu erstellen und zu verknüpfen.

Simulationsevents

Die Events sind wie folgt klassifiziert. Die Struktur dafür enthält 4 Schlüsselwörter: „name“ für den Namen des Events, „type“ für den Typ, „ind“ für den Index des durchgeführten Objekts (Basisstation und Consumer) sowie „time“ für den Zeitpunkt, für den dieses Event getriggert wird. Alle Events werden anhand ihrer Zeitpunkte in einem Heap der Klasse „Map“ gespeichert, sodass diese nacheinander abgearbeitet werden.

Map

Die Observierung, dass Beobachtungsevent für die ganze Karte, ist das einzige Event für Map, wobei alle Consumer und Basisstationen einmal beobachtet werden. Die Datennachfrage und die Datenrate wird neu berechnet und das Übertragungsende jedes Consumers wird hier eingeschätzt, zu welchem Zeitpunkt die Beobachtung dieser Karte wieder getriggert wird. Außerdem werden alle Basisstationen deaktiviert, bei denen es keine Nutzer mehr gibt.

Basisstation

Für die Klasse Basisstation gibt es 5 Events: die Deaktivierung, den Schlafzyklus, die Aktivierung, und den Ruhezustand sowie den aktiven Zustand.

Bei der Deaktivierung wird zunächst mit ε -greedy (Formel (1)) die Anzahl der darauffolgenden Schlafzyklen entschieden. Danach wird der Q-Wert von der letzten Episode mit Formel (2) aktualisiert und dann das Event Schlafzyklus nach der Deaktivierungsdauer vom aktuellen Schlafmodus gestartet.

Beim Schlafzyklus wird geprüft, ob es schon Consumer in der „Bufferlist“ gibt, oder der Schlafzyklus zu Ende ist. Wenn dies der Fall ist, dann wird die Aktivierung für diesen Zeitpunkt ausgelöst. Sonst wird der Zähler für die Anzahl der Schlafzyklen um 1 dekrementiert und ein Event Schlafzyklus nach einer Runde dieses Schlafmodus begonnen.

Bei der Aktivierung wird zunächst entschieden, ob die Basisstation aufwacht, wenn die Bufferlist nicht leer ist, oder im niedrigeren Schlafmodus bleibt. Dies funktioniert anhand einer Treppenfunktion. Diese wurde in unserer zweiten Quelle verwendet und erschien uns auf den ersten Blick sinnvoll. Jedoch fiel uns nach einiger Zeit auf, dass diese die akkurate Berechnung der Q-Werte deutlich komplizierter macht. Da nun jedes Wechseln eines Modus zu einer Entscheidung gemacht und anschließend bewertet werden musste. Danach wird mit ε -greedy (Formel (1)) die Anzahl der nächsten Schlaf- bzw. Idlezyklen wie oben entschieden und der Q-

Wert aktualisiert. Wenn die Bufferlist nicht leer ist, wechselt die Basisstation zum Zustand mit Benutzern, sonst wenn der aktuelle Schlafmodus eins ist, dann wechselt die Basisstation zum Leerlauf, sonst bleibt die Basisstation beim Schlafzyklus.

Beim Ruhezustand wird nichts gemacht, falls die Basisstation schon aktiviert wurde, sonst wird analog zum Schlafzyklus jedes Mal der Zähler um 1 dekrementiert.

Beim aktiven Zustand wird zunächst die Verzögerung und der Energiegewinn berechnet, um den Q-Wert zu aktualisieren. Danach wird die Beobachtung der Karte ausgelöst, wobei alle Basisstationen und Consumer beobachtet und aktualisiert werden.

Consumer

Für Consumer gibt es nur die Ankunft eines Consumers. Zunächst wird dieser Consumer der nächste Basisstation zugeordnet. Wenn die Basisstation noch schläft, wird der Consumer zur Bufferlist, sonst zur Servalist hinzugefügt. Danach wird ein Beobachtungsevent für die Karte gestartet.

Ergebnisse und Auswertung

Um den Einfluss unseres Algorithmus bewerten zu können, haben wir unsere Simulation mit Q-Learning und ohne, also mit zufälligen Entscheidungen, durchgeführt. Die restlichen Parameter und Ergebnisse sind auf der folgenden Abbildung zu sehen.

ZUFALL:

1. **Learning Rate: 0**
 2. Epsilon Greedy: 0.5
 3. Reward: 0.5
 4. Discount: 0.5
- Energiegewinn \approx **17.52 kW**
- Gesamtverzögerung \approx **62.4 ms**
- Energiegewinn/Delay \approx **0.28**

Q-LEARNING:

1. **Learning Rate: 1**
 2. Epsilon Greedy: 0.5
 3. Reward: 0.5
 4. Discount: 0.5
- Energiegewinn \approx **9.875 kW**
- Gesamtverzögerung \approx **10.41 ms**
- Energiegewinn/Delay \approx **0.95**

Da der Reward hier so gewählt wurde, gleich auf Energiegewinn und Delay zu achten, kann man klar erkennen, dass der Algorithmus ein sehr gutes Verhältnis von 95 Prozent erreicht hat. Es wurde zwar nur die Hälfte der Energie gespart, jedoch der Delay um das Sechsfache reduziert und dabei ein Gleichgewicht gefunden.

Anschließend haben wir den Einfluss von ungenauen und genaueren Q-Werten gegenübergestellt. Dafür haben wir einmal eine sehr niedrige und einmal eine mittlere Learning Rate gewählt. Die Parameter sowie Ergebnisse sind wie folgt.

SIMULATIONSERGEBNISSE 2

ZUFALL:

1. **Learning Rate: 0**
 2. Epsilon Greedy: 0.6
 3. Reward: 0.9
 4. Discount: 0.2
- Energiegewinn \approx **9.47 kW**
- Gesamtverzögerung \approx **57.99 ms**
- Energiegewinn/Delay \approx **0.16**

Q-LEARNING:

1. **Learning Rate: 0.5**
 2. Epsilon Greedy: 0.6
 3. Reward: 0.9
 4. Discount: 0.2
- Energiegewinn \approx **15.21 kW**
- Gesamtverzögerung \approx **34.24 ms**
- Energiegewinn/Delay \approx **0.44**

Nun wurde der Reward so justiert, dass er Energiegewinn stark präferiert wird. Dies hat der Algorithmus mit einem Faktor von $\sim 1,6$ umgesetzt und dabei sogar noch den Delay reduziert.

Insgesamt lässt sich aus den Ergebnissen schließen, dass das Q-Learning ein gutes Tool ist, um unbekannte, überschaubar zufällige Umgebungen zu analysieren und Entscheidungen zu optimieren.

Ausblick auf weitere Möglichkeiten

Im Rahmen dieses Projekts war es uns zeitlich nicht möglich, alle unsere zuvor gesetzten Ergebnisse zu erzeugen und zu evaluieren. Daher wollen wir hier noch einen Ausblick darauf geben, welche Verbesserungsmöglichkeiten bestehen bzw. welche Ergebnisse wir noch gerne erzeugt hätten.

Um noch aussagekräftigere Ergebnis zu erhalten, würden wir den Epsilon Wert, also die Exploration, über die Zeit niedriger werden lassen. Außerdem bestände hier die Möglichkeit, durch eine sehr große Menge an Simulationen die Werte auf Konvergenz zu prüfen. Da jeder Zufall in unserer Simulation einer Gleichverteilung unterliegt, sollte diese vorliegen. Anschließend würden wir eine Simulation durchführen, bei der von Anfang an gelernte Q-Werte genutzt werden. Die Ergebnisse würden wir dann bei gleicher Simulation jenen gegenüberstellen, welche wir beim Erlernen erhalten haben. Da von Anfang an passende Werte vorliegen sollten, würden wir also auch bessere Ergebnisse erwarten.

Eine weitere Idee zur Verbesserung wären mehr Schlafmodi, um dynamischer zwischen den Modi wechseln zu können. Jedoch ist die praktische Anwendbarkeit fraglich, da sich hier ebenfalls die Frage stellt, ob dies überhaupt durch die Hardware umsetzbar wäre.

Referenzen:

- (1) Ali El-Amine*, Mauricio Iturralde*, Hussein Al Haj Hassan[†] and Loutfi Nuaymi* A Distributed Q-Learning Approach for Adaptive Sleep Modes in 5G Networks, 2019 IEEE Wireless Communications and Networking Conference (WCNC)
- (2) Fatma Ezzahra Sale Management of advanced sleep modes for energy-efficient 5G networks, Institut Polytechnique de Paris, 2019