

# **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 Regelungssystem 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 von Ali El-Amine\*, Mauricio Iturralde\*, Hussein Al Haj Hassan<sup>†</sup> and Loutfi Nuaymi\* im IEEE publizierten Paper, welches den Namen „A\_Distributed Q-Learning Approach for Adaptive Sleep Modes in 5G Networks“ trägt. Im Folgenden schildern wir unsere Erwartungen, Erfahrungen/Schwierigkeiten sowie Erkenntnisse, welche wir im Laufe des Projektes erlangt haben.

## **Ziele welche wir uns für das Institutsprojekt setzten**

Zu Beginn sei zu erwähnen, dass wir uns zusammen mit unserem Betreuer Dr. Michael Reyes 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 als Student immer auch auf die Hilfe und Unterstützung von Kommilitonen angewiesen ist. Diese Erkenntnis lasse 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 zu 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, 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 Skill 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 zu nutzen, 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 modellieren wir mit 3 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 (Delay).

Es stellt sich somit die Frage welchen Modus jede Basisstation wählen soll um bei möglichst guter Performanz 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 möglichst optimale Policy zu erhalten. Hierunter versteht man das gelernte Verhalten des Agenten, welches ihm sagt, welche Aktion er in einem bestimmten Zustand in einer bestimmten Umgebung ausführen soll. Diese Policy wird mittels sogenannter Q-Values (Formel 3) 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 stetig aktualisiert um wie bereits oben erwähnt eine möglichst optimale Lösung eines gegebenen Problems zu erlangen. Die Aktualisierung geschieht wie in Abb.1 zu sehen anhand des vorherigen Q-Values und den neu gemachten Erfahrungen bei der Erkundung der Umgebung des Agenten. Hierbei regelt der Parameter Alpha (Learning Rate), inwieweit bzw. wie schnell der Algorithmus aus der gegebenen Aktion Schlüsse zieht. Der Diskontinuitätsfaktor Gamma 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 zunächst gemäß Formel 4 eine aktuelle Aktion ausgewählt. Dies geschieht zunächst mittels der epsilon-greedy Methode. Nachdem nun die aktuelle Aktion ermittelt wurde, wird diese ausgeführt. Es wird somit anhand Formel 3 einem Tupel  $s_m^t$  und  $a_m^t$  zum Zeitpunkt  $t$  ein zukünftiges Tupel aus  $s_m^{t+1}$  und  $t_m^{t+1}$  ermittelt. Das kleine  $m$  steht hier für die Basisstation und das  $k$  für den Consumer. Schließlich werden die Q-Values anhand Formel 3 aktualisiert.

$$Q(s_m^t, a_m^t) = Q(s_m^t, a_m^t) + \alpha \left[ r_m^t + \gamma \max_a Q(s_m^{t+1}, a) - Q(s_m^t, a_m^t) \right] \quad (3)$$

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, den Trade-Off zwischen dem weiteren Handeln nach der verfolgten Strategie und einer zufälligen Aktion.

$$a_m^t = \begin{cases} \underset{a}{\operatorname{argmax}} Q(s_m^t, a), & \text{if } y > \epsilon \\ \operatorname{rand}(\mathcal{A}), & \text{otherwise} \end{cases} \quad m = 1, \dots, M. \quad (4)$$

Unsere 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 (4). Die rand()-Funktion haben wir als gleichverteilte Zufallsfunktion implementiert.

## 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 und somit unvorteilhaft.

Sleep level	Deactivation duration	Minimum sleep duration	Activation duration
SM 1	35.5 $\mu$ s	71 $\mu$ s	35.5 $\mu$ s
SM 2	0.5 ms	1 ms	0.5 ms
SM 3	5 ms	10 ms	5 ms
SM 4	0.5 s	1 s	0.5 s

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

Consumer werden Log-Normal verteilt, sodass im Durschnitt 1 C/Km<sup>2</sup> erscheint.

Diese werden automatisch der nächsten Basisstation 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 skaliert, sondern auch durch die Übertragungsqualität, welche mit Hilfe die Signal-to-Interference-plus-Noise Ratio (SINR) berechnet wird.

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

Hierbei ist  $P$  die Übertragungsleistung der jeweiligen Basisstation, bei uns sind diese alle gleich,  $h(k)$  ist die Kanalverstärkung, welche den Pfadverlust und „shadowing-effect“ repräsentiert. Zuletzt ist  $\sigma^2$  die additive Leistungsdichte des Gauß'schen weißen Rauschens. Somit kommen wir zur gelieferten Übertragungsrate durch das Shannon-Hartley Theorem, welches lautet wie folgt:

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

dabei repräsentiert  $W$  die Bandbreite und  $\alpha$  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 in jeder Episode 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. (Abb.5)

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

Anhand des Parameters  $\eta$  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 (6).

Parameter	value
Antenna height	30 m
BS Tx Power	45 dBm
Bandwidth	20 MHz
Thermal noise	-174 dBm/Hz
Pathloss	$128.1 + 37.6 \log_{10}(d)$ dB
Shadowing	Log-normal (6 dBm)
User's arrival	Log-normal, $\lambda_a = 1$ , $v = \lambda_a / 10$
Service type	file with mean=4 Mb
Scale parameter	$\lambda = 441.305$
Shape parameter	$k = 0.8$

(6)

**[Hier muss Samuel nochmal checken ob es wirklich so implementiert ist.]**

Um den zeitlichen Ablauf zu simulieren, jedoch nicht diese Zeit warten zu müssen, entschieden wir uns, diese durch Events zu ausrechnen zu lassen. Dies wird durch das triggern von Events realisiert. Ein Event ist dabei das Eintreffen eines Consumers, das Ablaufen seiner Übertragungszeit, aber auch das Aufwachen 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 empfohlene, MatLab, hinauslaufen wird. 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.

**[Hier fehlt noch Text]**

**Hier was irgendwas zu Ergebnissen**

Um den Einfluss unseres Algorithmus bewerten zu können haben wir einerseits unsere Simulation mit dem Q-Learning und andererseits ohne unseren Algorithmus durchgeführt und schließlich

### **Referenzen:**

- (1) Ali El-Amine\*, Mauricio Iturralde\*, Hussein Al Haj Hassan<sup>†</sup> and Loutfi Nuaymi\*  
A\_Distributed\_Q-Learning\_Approach\_for\_Adaptive\_Sleep\_Modes\_in\_5G\_Networks
- (2) Fatma Ezzahra Sale Management of advanced sleep modes for energy-efficient 5G networks