

## Sprachaktivitätserkennung mittels eines Mustererkenners für Atemschutzmasken

Michael Brodersen<sup>1,2</sup>, Achim Volmer<sup>1</sup>, Marcus Romba<sup>1</sup>, Gerhard Schmidt<sup>2</sup>

<sup>1</sup> Dräger Safety AG & Co. KGaA, 23560 Lübeck, E-Mail: michael.brodersen/achim.volmer/marcus.romba@draeger.com

<sup>2</sup> Christian-Albrechts-Universität zu Kiel, 24143 Kiel, E-Mail: mibr/gus@tf.uni-kiel.de

### Einleitung

Im Atemschutzeinsatz ist die Kommunikation unter Feuerwehrleuten aufgrund der starken Dämpfung der Atemschutzmaske und der lauten Umgebungsgeräusche sehr erschwert. Um die Kommunikation zu verbessern, werden Kommunikationssysteme für Atemschutzmasken eingesetzt, die das Sprachsignal beispielsweise an ein taktisches Funkgerät weiterleiten. Entscheidend ist, dass die Sprache mit besonders wenig Anstrengung verstanden werden kann. Dies kann durch den Einsatz einer Sprachaktivitätsdetektion unterstützt werden, die Störgeräusche wie beispielsweise Atemgeräusche aus dem Signal entfernt. In diesem Beitrag wird eine Sprachaktivitätsdetektion basierend auf einem Mustererkenner vorgestellt. Die Leistungsfähigkeit wird für unterschiedliche Ansätze wie für ein Neuronales Netz und ein Codebuch evaluiert und die Eignung für eine Echtzeitimplementierung auf einem eingebetteten, batteriebetriebenen System diskutiert.

### Eigenschaften von Atemschutzmasken

Atemschutzmasken schützen das Gesicht sowie die Atemwege vor toxischen Gasen und Rauch (siehe [1]). Die Maske wird durch eine Dichtlinie um das Gesicht abgedichtet, Nase und Mund werden von der Innenmaske bedeckt, welche die Ausatemluft so lenkt, dass das Visier nicht beschlägt (siehe Abb. 1).



Abbildung 1: Atemschutzmaske auf einem Kunstkopf.

Der Raum vor Mund und Nase wird mit frischer Luft aus einer auf dem Rücken getragenen Flasche versorgt. Durch die Abdichtung der Atemschutzmaske am Gesicht, wird die Sprache stark gedämpft. Um diese Dämpfung zu minimieren, ist vor dem Mund eine Sprechmembran angeordnet.

### Kommunikationseinheiten für Atemschutzmasken

Um die Kommunikation weiter zu verbessern, werden Kommunikationseinheiten verwendet, welche die Sprache mit einem Mikrofon vor der Sprechmembran aufzeichnen, verarbeiten und verstärkt auf die Lautsprecher oder ein taktisches Funkgerät ausgeben. Dieser Signalfluss ist in Abb. 2 dargestellt.

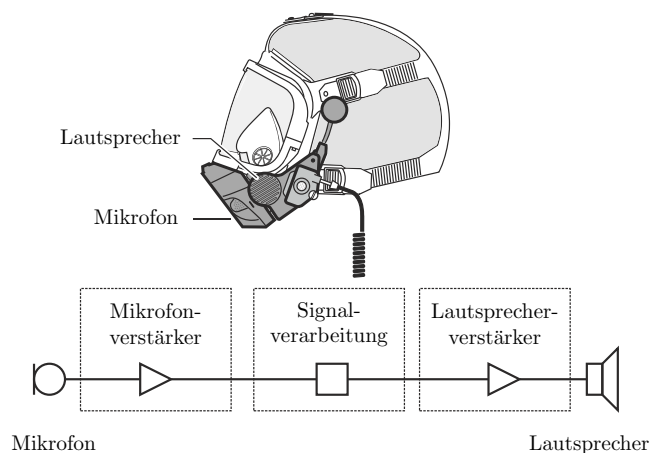


Abbildung 2: Struktureller Überblick der Kommunikationseinheit [1].

Für eine gute Kommunikation ist es notwendig, nur das Nutzsignal und nicht das Störsignal zu übertragen, damit der Inhalt verständlicher wird. Bei der Masken-Kommunikationseinheit sind besonders die Atemgeräusche störend. Sprache und Atemgeräusche sind jedoch im Mikrofonsignal vorhanden. Um diese störenden Signale nicht über die Lautsprecher oder über ein taktisches Funkgerät auszugeben, ist eine Sprachaktivitätserkennung und eine anschließende Filterung notwendig, welche diese Störgeräusche herausfiltern kann. Ein möglicher Ansatz dafür sind Mustererkenner, welche auf das Atemgeräusch trainiert werden. Im Folgendem werden ein Neuronales Netz und ein Codebuch in Bezug auf die Leistungsfähigkeit einer Störgeräuschunterdrückung gegenübergestellt.

### Merkmalsextraktion

Ein Mustererkenner nutzt Signaleigenschaften für die Klassifikation aus. Daher wird das Mikrofonsignal zuerst einer Merkmalsextraktion unterzogen. Diese verarbeitet das Eingangssignal so, dass die relevanten Merkmale herausgearbeitet werden. Diese Eigenschaften des Eingangssignal müssen durch möglichst wenige signifi-

kante Merkmale beschrieben werden, um den Rechenaufwand zu minimieren. Bei der Merkmalsextraktion wird das Eingangssignal in den Frequenzbereich transformiert und die Merkmale anhand des Betragsspektrums, einer Melfilterung und einer Logarithmierung ermittelt. Aus 65 komplexen Stützstellen werden nach der Transformation  $N_M = 12$  Merkmale extrahiert und in dem Vektor

$$\mathbf{X}_M(k) = [X_M(0, k), \dots, X_M(N_M - 1, k)]^T \quad (1)$$

zusammengefasst, wobei  $k$  der Rahmenindex ist. Daraufhin werden die Merkmale dem Mustererkenner präsentiert. In Abb. 3 ist der Signalfussgraph der Sprachaktivitätserkennung abgebildet.

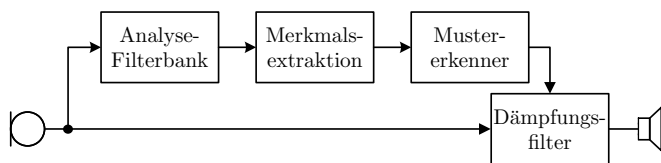


Abbildung 3: Signalfussgraph der Sprachaktivitätserkennung

Dieser Signalfussgraph ist stark vereinfacht und bezieht sich hierbei nur auf die Verarbeitung des Mustererkenners und nicht auf die Problematik des Gesamtsystems. Wird das Signal vom Mustererkenner als Sprache klassifiziert, wird dieses auf dem Lautsprecher ausgegeben.

## Neuronales Netz

Das Neuronale Netz orientiert sich an der Funktion des Gehirns, welches aus Neuronen besteht [2]. Diese Neuronen, auch Knoten genannt, sind miteinander verknüpft und die Übergangspfade sind mit einer Gewichtung  $w_{ij}$  versehen, womit der Übergang von den Knoten  $i$  zu  $j$  beschrieben wird. Das hier verwendete Neuronale Netz besteht aus einer Eingangsschicht, einer verdeckten Schicht und einer Ausgangsschicht. In der Eingangsschicht werden die einzelnen Elemente des Merkmalsvektors  $\tilde{X}(m, k)$  mit

$$\tilde{X}(m, k) = 2 \cdot \frac{X_M(m, k) - X_{M, \min}(m)}{X_{M, \max}(m) - X_{M, \min}(m)} - 1 \quad (2)$$

auf den Bereich  $-1 < \tilde{X}(m, k) < 1$  normiert, wobei  $m$  der Merkmalsindex ist. Dabei sind in  $X_{M, \min}(m)$  die minimalen und in  $X_{M, \max}(m)$  die maximalen Werte des zugehörigen Merkmals der Trainingsdaten, welche zum Training des Neuronalen Netzes benutzt werden, hinterlegt. Daraufhin werden die normierten Eingangsdaten von den zugehörigen Knoten zu den Knoten der verdeckten Schicht verteilt und mit dem zugehörigen Gewicht  $w_{v, in}$  versehen. Zu den eingehenden Pfaden in dem Knoten wird zusätzlich ein Bias  $B_{V_i}$  für die verdeckte Schicht und  $B_{A_j}$  für die Ausgangsschicht an dem Knoten  $i$  hinzu addiert. Der Bias ist als Schwellenwert des jeweiligen Knoten anzusehen. Die eingehenden Signale in jedem Pfad werden aufaddiert und daraufhin mit einer linearen

begrenzten Übertragungsfunktion gewichtet, welche das Verhalten

$$f(x) = \begin{cases} 1 & , \text{ wenn } x > 1, \\ -1 & , \text{ wenn } x < -1, \\ x & , \text{ sonst,} \end{cases} \quad (3)$$

besitzt. Dadurch ergibt sich der Zusammenhang der Übertragung von der Eingangsschicht zur verdeckten Schicht inklusive der linearen begrenzten Übertragungsfunktion gemäß Gl. (3) zu

$$X_V(i, k) = f \left( B_{V_i} + \sum_{n=0}^{N_M-1} \tilde{X}(n, k) \cdot w_{V_{in}} \right), \quad (4)$$

für  $0 \leq i < N_V$ .

wobei  $N_M$  die Anzahl der Merkmale,  $w_{V_{in}}$  die Gewichte von der Eingangsschicht zur verdeckten Schicht und  $X_V(i, k)$  der anliegende Wert an dem Knoten  $i$  beschreibt. Der Ergebnisvektor der Ausgangsschicht wird mit  $\mathbf{X}_A(k) = [X_A(0, k), \dots, X_A(N_A - 1, k)]^T$  beschrieben und dessen Berechnung ergibt sich zu:

$$X_A(j, k) = f \left( B_{A_j} + \sum_{i=0}^{N_V-1} X_V(i, k) \cdot w_{A_{ji}} \right), \quad (5)$$

für  $0 \leq j < N_A$ .

Dabei stellt  $N_V$  die Anzahl der Knoten der verdeckten Schicht,  $w_{A_{ji}}$  das Gewicht von der verdeckten Schicht zur Ausgangsschicht und  $X_A(j, k)$  das Ergebnis an dem Knoten  $j$  dar. Der Ergebnisvektor der Ausgangsschicht hat  $N_A = 5$  Elemente, was der Anzahl der zu erkennenden Klassen entspricht. Dabei wird zwischen Geräusch, Pause, Ausatmen, Einatmen und Sprache unterschieden. Diese werden in der Menge  $M$  zusammengefasst, so dass  $M = \{\text{Geräusch, Pause, Ausatmen, Einatmen, Sprache}\}$  ist. Die Klasse Geräusch beinhaltet die auffälligen Geräusche wie beispielsweise Knackgeräusche vom Betätigen der Schalter an der Atemschutzmaske. Zur Detektion, welcher Pfad am Wahrscheinlichsten ist, wird der Index des maximalen Eintrags des Vektors  $\mathbf{X}_A(k)$  bestimmt:

$$d_{\max}(k) = \operatorname{argmax}_{p \in M} \{ \mathbf{X}_A(p, k) \}. \quad (6)$$

Mit dem Index des maximalen Elements  $d_{\max}(k)$  kann nun durch das dabei angegebene  $p$  ermittelt werden, welches Element der Menge  $M$  aktiv ist. Solch ein verwendetes *feed forward*-Netz ist in Abb. 4 zu sehen. Damit die Berechnungen des Neuronalen Netzes vollzogen werden können, muss das Neuronale Netz durch eine Trainingsphase initialisiert werden.

Die benötigten minimalen und maximalen Eingangswerte, die Übergangsgewichte und die Bias werden durch ein Training generiert. Dazu ist es notwendig, Merkmalsvektoren aus Audiosignalen der jeweiligen Klasse zu erzeugen. Das Training des Neuronalen Netzes wird mittels des *back propagation*-Algorithmus durchgeführt. Das entworfene Neuronale Netz ist in Abb. 4 dargestellt,

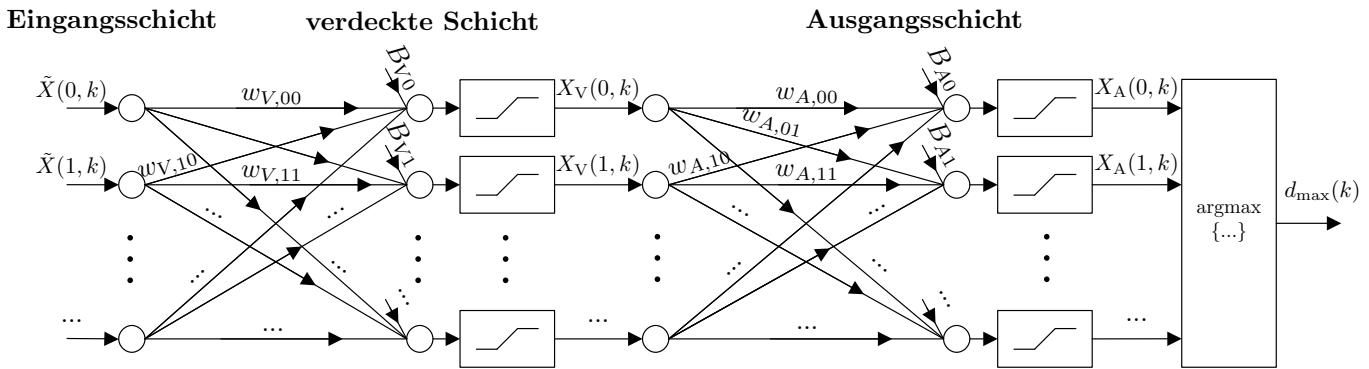


Abbildung 4: Neuronales Netz.

wobei die Übertragungsfunktion aus Gl. (3), die Anzahl der Merkmale und die Anzahl der Schichten der verdeckten Schicht und die Anzahl der Ausgangsmerkmale verwendet werden.

### Codebuch

Das Codebuch ist ein Mustererkenner, der basierend auf einer trainierten Datenbank einen Vergleich zwischen den Datenbankeinträgen und dem aktuellem Merkmalsvektor vollzieht [3]. Der Datenbankeintrag mit minimaler, quadratischer, euklidischer Distanz wird als am wahrscheinlichsten interpretiert, so dass die Klasse dieses Datenbankeintrags als aktuelle Zustandsschätzung verwendet wird. Somit kann zwischen den Klassen Geräusch, Pause, Einatmen, Ausatmen und Sprache unterschieden werden.

Die Einträge der Datenbank werden durch ein Training erzeugt, für welches Merkmalsvektoren der jeweiligen Klasse erzeugt werden. Diese Merkmalsvektoren werden mittels des *K-means Clustering*-Algorithmus [4] so zusammengefasst, dass sich ein kompaktes Codebuch ergibt. Das Codebuch für Sprache besitzt mehr Einträge als die übrigen Klassen, da Sprache eine größere Varianz im Merkmalsraum aufweist. Die Klassen Geräusch, Pause, Ausatmen und Einatmen werden auf 16 Merkmalsvektoren und Sprache auf 64 Merkmalsvektoren für die entsprechenden Codebücher reduziert.

Mit dem trainiertem Codebuch kann die Aktivität der Klassen bestimmt werden. Zu dieser Bestimmung wird die minimale quadratische euklidische Distanz

$$d_{E,p}(\tilde{\mathbf{X}}(k), \mathbf{c}_{s,p}, k) = \min_{s=0 \dots N_c-1} \left\{ \sqrt{\sum_{m=0}^{N_M-1} |\tilde{X}(m,k) - c_{s,p,m}|^2} \right\}, \quad \text{für } p \in M \quad (7)$$

zwischen dem Merkmalsvektor  $\tilde{\mathbf{X}}(k)$  und den Codebucheinträgen  $\mathbf{c}_{s,p}$  berechnet [5]. Der Index  $s$  steht dabei für den Codebuchvektor  $s$  aus dem Set  $p$ , wobei es fünf Sets (Geräusch, Pause, Einatmen, Ausatmen und Sprache) gibt. Die fünf berechneten minimalen quadratischen

euklidischen Distanzen von den Sets werden miteinander verglichen und das Argument des Minimums wird bestimmt:

$$d_{\min}(k) = \operatorname{argmin}_{p \in M} \left\{ d_{E,p}(\tilde{\mathbf{X}}(k), \mathbf{c}_{s,p}, k) \right\}. \quad (8)$$

Anhand der minimalen Distanz kann nun ermittelt werden, welches Element der Menge  $M$  aktiv ist. Die Sprachaktivitätserkennung mittels Codebuch ist in Abb. 5 abgebildet.

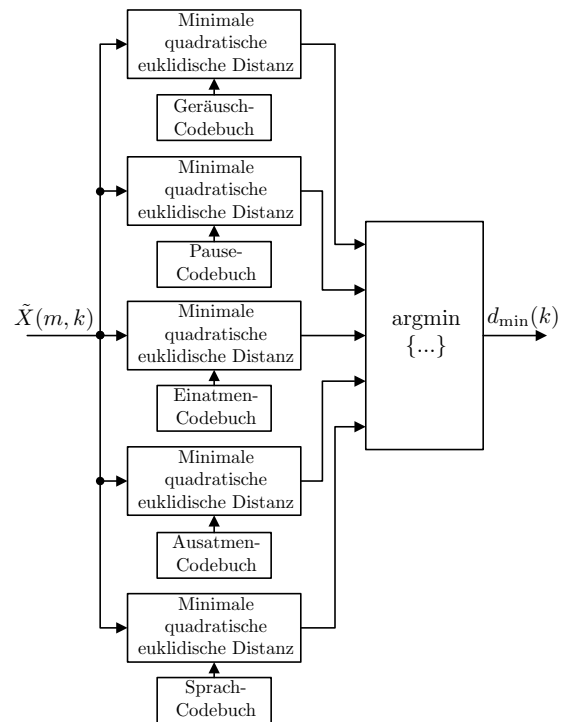


Abbildung 5: Blockschaltbild des Codebuchs

### Performancevergleich der Mustererkenner

Die Kommunikationseinheit der Atemschutzmaske ist ein batteriebetriebenes Gerät, bei welchem die Batterielaufzeit sehr entscheidend ist. Um diese Zeit zu maximieren werden Prozessoren mit sehr geringer Leistung verwendet, wodurch die Anzahl der Rechenoperationen pro Zeit begrenzt sind. Aus diesem Grund werden die verschiede-

nen Mustererkennungsansätze hinsichtlich ihrer Erkennungsraten, benötigten Rechenoperationen und auf den benötigten Speicherbedarf verglichen.

Die Erkennungsraten des beschriebenen Neuronalen Netzes sind in Tab. 1 und die des beschriebenen Codebuches in Tab. 2 dargestellt, wobei die Klassen Geräusch, Pause, Ausatmen, Einatmen und Sprache durch G, P, A, E und S abgekürzt werden.

		Ziel					Erkennung	
		G	P	A	E	S		
Eingang	G	53%	23%	14%	2%	8%	53%	47%
	P	14%	70%	8%	1%	7%	70%	30%
	A	2%	2%	83%	1%	12%	83%	17%
	E	<1%	<1%	<1%	99%	<1%	99%	1%
	S	2%	<1%	2%	<1%	95%	95%	5%
Durchschnittliche Erkennungsrate							79%	21%

**Tabelle 1:** Erkennungsmatrix des Neuronalen Netzes.

		Ziel					Erkennung	
		G	P	A	E	S		
Eingang	G	38%	40%	9%	3%	10%	38%	62%
	P	4%	91%	3%	<1%	2%	91%	9%
	A	6%	10%	68%	1%	14%	68%	32%
	E	2%	<1%	<1%	96%	<1%	96%	4%
	S	3%	3%	5%	<1%	88%	88%	12%
Durchschnittliche Erkennungsrate							76%	24%

**Tabelle 2:** Erkennungsmatrix des Codebuches.

In den Tabellen sind links die Eingangsklassen und oben die Zielklassen abgebildet. Für die Zielklassen sind die Erkennungsraten in % für die jeweilige Eingangsklasse abgebildet; rechts daneben ist die Gesamterkennungsrate für die Eingangsklasse zusammengefasst. Die Erkennungsraten der fälschlich zugewiesenen Merkmale sind in rot und die richtigen in grün dargestellt. In der untersten Zeile ist die durchschnittliche Erkennungsrate für alle Klassen dargestellt. Bei den Erkennungsraten sollte die Verwechslung zwischen Ausatmen und Sprache möglichst klein sein, da beispielsweise Zischlaute ansonsten als Ausatmen klassifiziert werden.

Der Vergleich der Erkennungsraten beider Mustererkenner zeigt, dass die Erkennungsraten mit dem Neuronalen Netz für die Klasse Sprache um 7%, für Einatmen um 3%, für Ausatmen um 15% und für Geräusch um 15% höher sind. Lediglich die Erkennungsrate für die Klasse Pause ist beim Codebuch um 21% besser, wobei die Fehlererkennung von Pause beim Neuronalen Netz größtenteils in der Klasse Geräusch wieder zu finden ist. Somit werden diese Fehlerkennungen von Pause nicht fälschlicherweise als Sprache detektiert. Die Verwechslung von Ausatmen und Sprache ist beim Neuronalen Netz geringer und die durchschnittliche Erkennungsrate um 3% höher. Somit ist das Neuronale Netz in den entscheidenden Fällen dem Codebuch vorzuziehen.

Die Erkennungsraten müssen dabei in Bezug auf die benötigten Rechenoperationen und den Speicherbedarf

gesehen werden, welche für das Neuronale Netz (NN) und für das Codebuch (CB) in Tab. 3 dargestellt sind, wobei der Speicherbedarf in Bezug auf die Anzahl der benötigten Parameter im 16 Bit-Format für einen 4 ms Rahmen angegeben ist.

	Additionen	Multiplikationen	Speicherbedarf
NN	188	172	209
CB	3328	1792	1536

**Tabelle 3:** Benötigte Ressourcen der Mustererkenner.

Das Neuronale Netz ist gegenüber dem Codebuch sowohl hinsichtlich der benötigten Rechenleistung als auch des Speicherbedarfs sehr kompakt. Das Codebuch benötigt 17-mal so viele Additionen wie das Neuronale Netz, 10-mal so viele Multiplikationen und den 7-fachen Speicherplatz.

Im Gesamtvergleich der Leistungsfähigkeit ist das Neuronale Netz gegenüber dem Codebuch merklich im Vorteil, da die Erkennungsraten höher sind, allgemein sowie in den entscheidenden Fällen, und da es deutlich weniger Ressourcen benötigt.

## Zusammenfassung

Mit Mustererkennern kann eine robuste Sprachaktivitätserkennung umgesetzt werden, die zwischen Geräuschen beim Bedienen der Atemschutzmaske, Pausen, Ausatmen, Einatmen und Sprache unterscheiden kann. Als Ansätze zur Mustererkennung wurden ein Neuronales Netz und ein Codebuch untersucht und hinsichtlich ihrer Leistungsfähigkeit verglichen. Die Erkennungsraten des Neuronalen Netzes und des Codebuches unterscheiden sich nur geringfügig, wobei das Codebuch für eine ähnlich gute Erkennung deutlich mehr Rechenoperation und Speicherplatz benötigt. Daher ist als Mustererkenner das Neuronale Netz im Vergleich zum Codebuch für die batteriebetriebene Anwendung in der Atemschutzmaske zu bevorzugen.

## Literatur

- [1] Dr.-Ing. A. Volmer, Dr.-Ing. M. Romba, C. Schmidt und M. Houssein Harbi: Optimization of Speech Intelligence for Fire Fighters' Full Face Masks, DAGA, 2013
- [2] C. M. Bishop: Pattern Recognition and Machine Learning, Springer, 2006
- [3] G. A. Fink: Markov Models for Pattern Recognition: From Theory to Applications, Springer, London, 2014
- [4] B. Pfister und T. Kaufmann: Sprachverarbeitung: Grundlagen und Methoden der Sprachsynthese und Spracherkennung, Springer, 2008
- [5] M. Bossert: Kanalcodierung, Oldenbourg Wissenschaftsverlag, 2013