# UNIVERSITY OF GOTHENBURG



# "AuTopEx": Automated Topic Extraction Techniques Applied in the Software Engineering Domain
## The design and evaluation of an approach for Automated Topic Extraction

*Bachelor of Science Thesis Software Engineering & Management*

JONATHAN KLEMETZ
MAGNUS JOHANSSON

# "AuTopEx": Automated Topic Extraction Techniques Applied in the Software Engineering Domain

The design and evaluation of an approach for Automated Topic Extraction

Jonathan Klemetz
Magnus Johansson

Examiner: Christian Berger
Supervisor: Alessia Knauss
Supervisor: Hang Yin

University of Gothenburg
Chalmers University of Technology
Department of Computer Science and Engineering
SE-412 96 Göteborg
Sweden
Telephone + 46 (0)31-772 1000

# "AuTopEx": Automated Topic Extraction Techniques Applied in the Software Engineering Domain

Magnus Johansson
Gothenburg University
Software Engineering & Management
Lindholmsplatsen 1, 412 96 Gothenburg, Sweden
gusjmagn03@student.gu.se

Jonathan Klemetz
Gothenburg University
Software Engineering & Management
gusklejoa@student.gu.se

## Abstract

*Automatically extracting topics from scientific papers can be very beneficial when a researcher needs to classify a large number of such papers.*

*In this thesis we develop and evaluate an approach for Automatic Topic Extraction, AuTopEx. The approach is comprised of four parts:*
*1) Text pre-processing.*
*2) Training a Latent Dirichlet Allocation model on part of a corpus.*
*3) Manually identifying relevant topics from the model.*
*4) Querying the model using the rest of the corpus.*

*We show that it is possible to automatically extract topics by applying AuTopEx on a corpus of scientific papers on autonomous vehicles.*

*According to our evaluation AuTopEx works better on full-text articles than texts consisting of just title, abstract and key-words.*

*Finally we show that this approach is vastly faster than human annotators, although not as accurate.*

The source code used to build AuTopEx can be found at:
(https://github.com/Klemetz/TopicExtraction).

## 1   Introduction

In this thesis we design and evaluate an approach for Automated Topic Extraction. Which is evaluated on papers in the Software Engineering domain, more specifically on autonomous vehicles.

### 1.1   Background

Automated Topic Analysis and Automated Topic Extraction allow researchers to extract the potential topics that are contained in a large text corpus. This has been tried in other scientific domains but (to the best of our knowledge) not in the field of Software Engineering.

### 1.2   Problem Domain & Motivation

In order to find relevant information, researchers often need to read a large number of published articles. This is especially true when conducting work like mapping studies or Systematic Literature Reviews and can be a very time-consuming process.

There is a lack of automated approaches to topic extraction that could support activities such as Systematic Mapping Studies, especially in the Software Engineering domain.

### 1.3   Research Goal & Research Questions

Our research goal is to investigate the automation of Topic Extraction from scientific papers in order to support time-consuming activities such as Systematic Mapping Studies [18]. We investigate extraction from both full-text articles and texts containing only title, abstract and keywords using a topic model called Latent Dirichlet Allocation (LDA).

The research goal has been divided into three research questions.

RQ 1: How can we support Automatic Topic Extraction for scientific papers in the Software Engineering domain?

RQ 2: Which approach is better for Automatic Topic Extraction: a) Extraction from title, abstract and keywords or b) Extraction from full text paper?

RQ 3: How well does the approach of using Latent Dirichlet Allocation (with suitable pre-processing) perform compared to a manual method?

## 1.4 Contributions

In this paper we present an approach (which we call "AuTopEx") for applying Automated Topic Extraction on a large number of scientific papers.

From the extracted data, relevant topics are identified and labeled. Researchers can then also automate the process of finding which papers in the corpus that are most likely to deal with the relevant topics.

Researchers will benefit from AuTopEx as it shows the applicability of Natural Language Processing (NLP) techniques in the Software Engineering Domain.

## 1.5 Scope

We construct and evaluate an approach for Automatic Topic Extraction using Latent Dirichlet Allocation (LDA). This could contribute in making Automatic Topic Extraction a viable approach in Software Engineering research. We do not evaluate other statistical models such as the n-gram model or term frequency-inverse document frequency. However, their potential use in our approach is discussed in the Conclusions and Future Work section.

## 1.6 Structure of the Article

Section 2 presents related work on Automated Topic Extraction. Section 3 covers our Research Strategy. In Section 4 we answer our research questions. First we describe AuTopEx and go into detail about it's implementation. Then we evaluate the results from implementing AuTopEx compared to human performance. Evaluations are made on two corpuses, one corpus containing articles in full-text and the other only title, abstracts and keywords. Section 5 contains analysis and discussion of the results. In this section we also discuss the validity threats to our findings. Section 6 concludes our findings and discusses what implications they may have for future research.

## 2 Related Work

We have not found any articles dealing with tools specifically tailored towards automation of Systematic Mapping Studies. These do however share some similarities with Systematic Literature Reviews (SLR). Hence we can discuss tools that support the latter.

According to Marshall and Brereton [16], the two most popular frameworks for tools that support SLR:s are the Projection Explorer Pex and ReVis which both make use of Visual Text Mining techniques. Projection Explorer Pex [6] can create a visualization from a set of textual documents either by building a vector representation of the text corpus (which is handled as table data to derive similarity information). It can also compute similarities by directly comparing text against text. It could possibly be used in helping with document classification during a mapping study.

ReVis [5] supports primary studies selection during SLRs. Among it's tools is the possibility of visualizing the relationships of potential primary studies. A 2D document map shows content and similarities of different documents. This is based on converting the documents into multi-dimensional vectors which can be reduced using stemming, by eliminating stop words and using projection techniques. ReVis only uses title, abstract and keywords for this document map however, and we ideally want to use full-text articles to discover topics.

CitNetExplorer analyzes citation patterns in scientific literature. The tool collects bibliographical data and constructs a citation network which can then be analyzed and visualized[11]. This could be useful for mapping a specific research topic, since key word searches could miss out on papers that do not contain these key words.

VOSviewer [10] is a tool for creating and visualizing bibliographical networks. It can also use text mining to create term maps from a text corpus. Part of speech-tagging is used to identify noun phrases and a technique for choosing the most relevant noun phrases is applied. Maps and clusters can then be created and visualized.

In [12] the creators of CitNetExplorer and VOSviewer discuss the limitations of both tools. They argue that the loss of information occurring when applying these techniques is very hard to measure and that they should be used as a complement rather than substitute to expert judgment.

One approach for speeding up topic extraction could be automatic summarization of articles. The abstract of a scientific paper is meant to provide a quick overview, but does not necessarily provide enough key information for the researcher. Automatic summarization techniques can capture scientific concepts such as Hypotheses, Method and Background on a sentence level [14] and thus provide more information than just an abstract. However, this work builds on having many domain experts manually annotate a large number of scientific papers used for training the machine learning

classifiers [15]. Such an undertaking is out of scope for this thesis.

The "Latent Dirichlet Allocation" method is widely used and applicable in the discipline of Natural Language Processing (NLP) [2]. As Blei puts it "The simple LDA model provides a powerful tool for discovering and exploiting the hidden thematic structure in large archives of text" [2]. When attempting to extract topics from a large corpus as the purpose is for the AuTopEx approach, a tool like the LDA method is very compelling. Blei has also provided some comparisons with other models which makes the choice of applying LDA an attractive option. He shows that even though LDA is meant to perform "in the spirit of LSI" (Latent semantic indexing) [3], the LDA method outperforms the LSI method regarding perplexity measures [3].

But can we be sure that NLP tools such as the LDA method are fitting for the Software Engineering domain? Studies such as the one performed by Hindle et al [9], shows us that they can be applied but might not always be suitable. Hindle presents in his paper that in the domain of software engineering, neither LDA nor the n-gram analysis approach may be suitable if the intended goal is to extract topics from files containing computer code. However, we do not expect that code snippets will make up anything but a very small portion of the scientific articles that we want to apply Automatic Topic Extraction on.

## 3 Research Strategy

We have chosen Design Science as our research strategy. Hevner et. al [8] present seven guidelines for conducting, evaluating and presenting Design Science research. These address design as an artifact, problem relevance, design evaluation, research contributions, research rigor, design as a search process, and research communication.

The artifact is in our case an approach based on Natural Language Processing techniques. In this approach we apply a number of pre-processing steps on a large corpus of texts. Then we automatically extract topics from the corpus. Finally we automatically classify the papers based on the extracted topics.

The problem relevance is the fact that doing this manually is a very time-consuming process.

Evaluation will be done by comparing the topics that the machine learning algorithm produces with annotation made manually by humans. Both of the authors will first do manual annotation of the same papers separately and then confirm that there exists an inter-annotator agreement. Basically that both authors have identified the same topics in each paper.

As for research contributions we are transferring knowledge from the domain of language technology to the area of Software Engineering. We also believe the resulting approach can be helpful in future research where speeding up topic extraction can be beneficial.

When it comes to research rigor we are actively selecting and applying appropriate theories and methods both when constructing and evaluating the resulting artifact.

Design as a research process has been a must from the start, since this approach has not been applied on scientific articles about Software Engineering before. Investigating existing tools and techniques that are already being used for similar purposes and how to implement them properly, is the whole foundation of constructing our approach. We have worked in iterations during the whole process, constantly improving our approach by trying out different ways of working with existing tools and developing our own tools where needed.

Finally research communication must be taken into account. Since this thesis is concerned with research, we ensure that we explain our methods as thoroughly as possible so that other researchers can evaluate the approach. All of our code is open source and available under the Gnu General Public License Version 2 at (`https://github.com/Klemetz/TopicExtraction`) along with proper documentation so that others can apply the approach themselves.

## 4 Results

### 4.1 AuTopEx

In order to answer **RQ1: "How can we support Automatic Topic Extraction for scientific papers in the Software Engineering domain?"**, we have developed the approach AuTopEx, which supports Automatic Topic Extraction.

AuTopEx can be broken down into four steps:

1. Pre-processing the articles of a large corpus of scientific paper.
2. Training a Latent Dirichlet Allocation (LDA) model using 10 percent of the pre-processed papers.
3. Manual identification and labeling of relevant topics returned by the model.
4. Automatic classification (extracting the topics) of the rest of the corpus by querying the LDA model.

Some of the papers will be annotated manually before automatic classification. We can evaluate the accu-

racy of the model by comparing this manual annotation with the automatic classification of the same papers.

We have chosen existing tools to help answering our research questions and to construct AuTopEx. Calibre (`https://calibre-ebook.com/`) is used for pdf-to -text conversion. The Natural Language ToolKit(NLTK) (`http://www.nltk.org/`) is used for pre-processing individual texts and Gensim (`https://radimrehurek.com/gensim/`) is used for applying the machine learning algorithms. Complementary tools in the form of Python scripts have been developed by the authors when needed.

### 4.1.1 Text Pre-Processing

Pre-processing the scientific papers follows a pipe-line of six consecutive steps:

1. Pdf-to-text conversion
2. Converting all text to lower-case
3. Tokenization
4. Removal of stop words, numbers and punctuation
5. Lemmatization
6. Removal of references section

#### Pdf-to-text conversion

The first step of preparing the individual scientific papers is to convert them from pdf to text format. We use the free and open-source Calibre software. The reasons are two-fold: a) Unlike other tools we tried, Calibre handles ligatures well and b) dehyphenation. If a word is cut off with a hyphen at the end of a column, the program checks if the hyphenated word exists elsewhere in the document without the hyphen and de-hyphenates it if that is the case. This ensures that more accurate words remain in the document.

#### Converting all text into lower-case

All text is transformed into lower-case format, to ensure correct multiplicity of words even if they appear at the start of a sentence. It is also important that the list of stop words (introduced below) are all in lowercase.

#### Tokenization

Tokenization means we break down the stream of text into meaningful elements. In our case this means individual words (all contiguous alphabetic characters become part of a token) which are then separated from symbols such as punctuation. The Natural Language ToolKit(NLTK) has a number of tokenizers, we recommend their "regexptokenizer" for this.

#### Removal of stop words, numbers and punctuation

Latent Dirichlet Allocation uses a Bag-of-words model [3]. This means that neither grammar or word order is important, only the multiplicity of the words. Thus we can now safely remove all punctuation and also stop words (such as "a", "and", "if", "or" etcetera).

Greek letters are often used as mathematical notation in scientific articles. Such a symbol on it's own has little to no semantic value for an annotator examining our results. Neither do we expect numbers from the articles to hold any semantic importance in the topic extraction so these are removed as well. This is easily solved by only allowing alphabetical words in the tokenizer.

Punctuation is also removed using regular expressions.

Stop words are removed by using a stop word list. We use the default stop words list from NLTK, and supplement it with more words that we deem have no semantic value. For example the word "fig." very commonly appears next to images and graphs in research papers. This word pollutes the results rather than give the topics any semantic meaning.

#### Lemmatization

Within a document a word can use several forms (such as "organize", "organizing" or "organizes") while all referring to the same concept, and we are interested in the multiplicity of this concept. Stemming is the process of reducing inflected words to their word stem. A stemmer only operates on the word at hand by cutting of the word stem. "Organize", "organizing" and "organizes" would all be reduced to "organ" using a stemmer, which is not what we want since the word now has an entirely new meaning.

Lemmatizing is closely related to stemming in that it reduces inflected words, however it reduces the word form to linguistically valid lemmas, using algorithms that deal with grammar and a built-in dictionary. For this purpose we use the WordNet Lemmatizer included with NLTK.

Here are a few example sentences. "They walk down the road. She walked by him. The elephant walks on four legs while we are used to walking on two."

Using the most common stemmer (Porter) we get: "They walk down the road . She walk by him . The eleph walk on four leg while we are use to walk on two ."

Using the WordNet lemmatizer we instead get: "They walk down the road . She walked by him . The

Figure 1: The steps involved in text pre-processing.

elephant walk on four leg while we are used to walking on two ."

The differences between stemmers and the basic differences between stemmers and lemmatizers are discussed in [13].

### Removal of References section

Finally, all of the papers have a "references" section at the end. We do not want the words in this section to pollute the article at hand. This section is removed by finding the last occurrence of the word "reference" (remember we have lemmatized all words) and removing all remaining words in the document including "reference". In the rare event that a reference has the word "reference" in it, some references might remain in the document. On the whole however we do not expect this to have a major impact on the results.

### Pre-processing Title, Abstract, Keyword texts

Pre-processing these texts uses the same pipe-line as the full text method outlined at the beginning of this section but with the first and last step removed. This is because we already have the abstracts available in text format and they don't contain any references.

Extracting title, abstracts and keywords from a full-text paper can be difficult, since not all articles are formatted in the same way. Some papers do not even include the keywords in the article document. Therefor we chose to extract this information using meta data stored in the software Endnote used by the researchers who provided us with the data used for the evaluation. Endnote can produce a single text file that contains author names, publishing year, publication, title, abstract and keyword for all articles that you want to perform topic extraction on.

A Python script extracts the relevant meta data (title, abstract and keywords) and saves a separate text file for each article (the model needs an entire corpus of papers to work with) naming them in the format author-publication-year-title for identification purposes when doing the evaluation.

Then we clean each document the same way as we did for the full-text articles (tokenization, removal of stop-words, lemmatization).

#### 4.1.2 Training the model

### Latent Dirichlet Allocation

The intent of using LDA in this study is to get topics from the documents in the supplied data sets. LDA can do this through its probabilistic, generative functionalities. So a trained LDA model will be able to point out topics for documents[19]. There are however quite a few ways of training an LDA model to achieve the queryable functionalities [19].

Most of the ways of training a model boils down to a guessing game. This guessing game begins when for every document every word has been assigned to a random topic. A topic is a list of words and how many instances of them there are. A word can reoccur several times in a topic, this gives it an increased chance to be dominant within this topic, and get more instances of itself within this topic. Then the algorithm, for every word in every document, looks for what topic the current word could fit in as well, then moves it there. Depending on how many instances of the current word there are in that topic already, the chance that the word will be moved there varies. Now when the current word is moved, the current words new topic which has received the current word will have an increased chance of receiving another instance of this word.

In short the guessing game can be described as that the LDA model gets better at guessing as it keeps at it, and a measurement of measuring how well a model guesses is it's perplexity value [3].

### Using the Gensim Framework

Gensim is a framework that is accessible through the programming language Python. The Gensim framework allows the user to build and train their own unique LDA model based on the users own corpora. Gensim also offers other kinds of machine learning algorithms outside the scope of LDA [19].

When creating an LDA model with Gensim, it requires a corpus that has been tokenized. In the case

of AuTopEx, the models trained are handed a number of the pre-processed text files. AuTopEx can then through the Gensim framework train an LDA model given a sample from the whole corpus.

Throughout the training process, the Gensim framework tells the user whether or not the model is improving by printing out what is called a perplexity measure [2]. An indication of whether the model is improving is, if the perplexity measure is decreasing for each iteration [3].

Then, when the perplexity measure is down to a predetermined value, the LDA model can be saved down on the hard drive of the users system and reused on the entire corpus. This is where AuTopEx can return which topics are deemed most relevant for each document.

To find measures that act as good examples when training a model to perform as well as possible one can observe Bleis experiments[3]. When asking for more than a hundred topics in these experiments, the perplexity measure is not improving as much anymore unlike when the number of topics approach a hundred. Blei also presents in his paper that when training models, if training a model with a larger sample of ten percent, of the entire corpus, the gain in accuracy is not significant. However, approaching ten percent of the entire corpus for a training sample the gain in accuracy is certainly appealing [3].

### 4.1.3   Identifying and labeling relevant topics

The trained model provides us with up to 100 topics, each topic consisting of a set number of words. A script exports this data to a spreadsheet for easy access by the annotators.

An example of a complex topic (10 words) outputted after training could look like this:

(0.005):    0.012*communication  +  0.009*channel + 0.008*packet + 0.008*velocity + 0.007*protocol + 0.006*follower + 0.006*platoon + 0.006*leader + 0.006*transmission + 0.006*controller.

What can be observed from this topic is that the words that follow a number and a star is related to the topic. Inside this topic there are several expressions, for example "0.012*communication". This expression and all other expressions that follows inside this topic will combined provide the interpreter guidance towards labeling the topic.

At a first glance it seems like the topic could be labeled as one of the words that it already contains, "Communication". The way this could be argued is be-

cause the topic also contains "packet" and "protocol". These words are tightly related with communication solutions/properties in software and computers in general. At a closer look there are other options for the label. Since the topic contains "platoon", "follower" and "leader" which all probably refer to a platoon of vehicles (the corpus being related to autonomous vehicles). The label could then arguably be something like "Networked vehicles" or "Cooperating Vehicles". Then again "Transmission" might be related to communication but could also refer to gearbox and we also have "velocity" and "controller" in the topic.

As you can see the labeling phase can prove quite difficult based on the number of words and their semantic relations.

We chose 7 words per topic but we encourage those who want to try this approach to experiment with the number of words per topic. In our experience, with fewer words the topics became more general (e.g. "Network") and with more words the topics became more specific ("Networked vehicles grouped in platoon"). Seven to us seemed like a good compromise because for this specific corpus of papers it gave us a large amount of varied topics.

It's important to note that not all topics from the trained model will be interpret-able by humans. This is due to the generative and probabilistic nature of the LDA model. A model will produce a number of bad topics with low scores. From our experience these topics will never be assigned to papers during the classification, so this is not a problem.

A very large majority of the topics from our corpus were however indeed interpret-able and covered a wide variety of areas (please refer to the Appendix for examples of the topics we got from the model and how we labeled them).

### 4.1.4   Querying the model

When an LDA model is finished and saved onto the hard drive, one can query this model with preprocessed documents in order to get the models opinion of what topics might exist in each specific document.

As an example, when we ask the model to return three possible topics for the paper "A Real-Time Multi-Sensor Fusion Platform for Automated Driving Application Development", the model outputs: "(37, 0.81872937773255205), (55, 0.078783842923631039), (78, 0.034186934006349756)"

This means that according to the model, topic 37 is the most probable topic for this document, followed by topics 55 and 78. These results are exported to a spreadsheet, where the researcher can look up the

A — Corpus of 2000 research papers on autonomous vehicles

F — 50 papers chosen at random for evaluation purposes

G — Annotator #1 labels all 50 papers using only labels from (E)

H — Inter-Annotator Agreement (making sure both annotators agree about which topics are in which papers) The three most dominant topics for each paper are chosen as the result.

B — Each paper is pre-processed (tokenization, lemmatization etc.)

Annotator #2 labels all 50 papers using only labels from (E)

I — After pre-processing the papers, we ask the trained model (D) which topics they are about

C — Machine learning algorithm (Latent Dirichlet Allocation) is trained on 200 of the papers

K — Final Evaluation: How close was the results of the manual annotation to the automatic classification from LDA?

D — Results in trained model containing 100 topics:
1: [car, autonomous, drive]
2: [pattern, architecture, style]
… and so on ….

J — Automatic classification of the 50 random papers:
Paper 1 is about topic 4, 37, 78
Paper 2 is about topic 0, 16, 54
...and so on ….

E — Categories manually labeled:
1 = "Autonomous driving"
2 = "Software architecture"
… and so on ….

Evaluating automatic topic extraction technique AuTopEx

Figure 2: A simplified overview of how we evaluate the AuTopEx approach.

labels corresponding to these numbers.

## 4.2 Evaluating AuTopEx

### 4.2.1 Setting up the evaluation

The data set consists of 425 scientific articles related to autonomous vehicles. These papers had been screened based on certain inclusion and exclusion criteria for an actual Systematic Mapping Study being performed by researchers at Chalmers University of Technology, thus we deemed it an excellent data set for performing our evaluation.

For each of our two evaluations 200 of the 425 scientific papers were selected at random for training the LDA model. This number was chosen because we expect the final mapping study to include at least 2000 articles, and it is considered good practice to use ten percent of the data set for training purposes when implementing LDA.

The 100 topics (containing 7 words each) from the model are now manually labeled by the authors. First each author labels all of the topics on their own and then check whether they disagree on any topic label. Any disagreements are solved by discussing the topic at hand. The labeling phase is arguably the most difficult part of the entire process because it requires the

annotators to have very good language skills as well as domain expertise. More on that in the "Threats to Validity" section of this thesis.

### 4.2.2 Evaluation method

From the remaining 225 papers 50 are chosen at random for evaluation purposes. We use a Python script for random selection as well, in order to eliminate any potential bias where an annotator could choose documents with very clear titles that were similar to the topics we already knew existed in the corpus.

All of the 50 documents are now read and annotated by each of the authors, if a document talks about a topic labeled in the previous step, it gets the same label.

After the human labeling is completed we process the same 50 documents using our trained LDA model. A Python script exports the most probable topics for each paper to a spreadsheet. We chose a two-fold approach for both full-text and title-abstract-keyword evaluations here: First we export the three topics with the highest probability weight according to the algorithm. Then we separately export all probable topics, no matter how low the probability is. This might give us insights into both how the Gensim implementation

of LDA works as well as tell us something about the documents being analyzed (mainly the number of probable topics per paper and their respective probability weight according to the algorithm).

For the purpose of supporting tasks such as document classification in Systematic Mapping Studies we are interested in knowing whether AuTopEx performs better with a data set consisting of full-text articles or a set where the articles only contain titles, abstract and keywords. In order to evaluate this, as well as getting a measure on how well the human annotators and the system agree with each other, we use an evaluation technique called precision and recall [1].

Before one can calculate the values for precision and recall one must first collect the required data. Rather than just presenting this data in tabular form, it helps to produce a confusion matrix, consisting of four fields. See the model below as an example. The four fields are labeled true positive, false positive, true negative and false negative.

In this study, true positives are the topics that are deemed by both the machine and the annotator as relevant for the given articles. False negatives are topics that have not been deemed relevant by the machine but have been deemed relevant by the human annotator. False positives are topics that the machine is returning as relevant topics but have not been deemed relevant by a human annotator. Lastly true negatives, are basically just the rest of the topics that have not been returned by the machine and that should not have been returned according to the human.

Figure 3: Example of the confusion matrix

**Returned by LDA?**

|  | **Yes** | **No** |
|---|---|---|
| **Yes** | True Positive | False Negative |
| **Relevant?** |  |  |
| **No** | False Positive | True Negative |

These boxes would be filled with the values that has been described previously in the respective box. So to show an example of how this would be performed, please refer to

the data supplied in the first appendix. When looking at the first sheet in this spreadsheet, there are four columns, true positive, false negative, false positive and true negative that are of importance. The papers are listed on the left and for each papers corresponding row the values for each of these elements are represented. Since this study is focusing on how the different data sets (full-text vs title, abstract and keywords) perform against each other, one can observe at the bottom part of the sheets, the sums of all the precision and recall values are stored. Here the values from the entire data set are added together and presented. It is these sums of the true positives, false negatives, false positives and true negatives for each data set that are used and later presented inside these confusion matrices that is exemplified above.

When this data has been collected, the following equations can be applied to get the values of precision and recall.

$$Precision = \frac{TP}{TP + FP} \qquad 1$$

$$Recall = \frac{TP}{TP + FN} \qquad 2$$

To bring a bit more clarity to what these values will indicate in the case of this study, lets quickly summarize. **Precision** serves as an indication of how many of the topics that are returned as relevant, are truly relevant. **Recall** represents how many relevant topics were returned by the system.

This study investigates if there is any preference for what type of documents to use when performing Automatic Topic Extraction. Thus, a value called an F-measure, which is a harmonic mean of precision and recall will be used in comparing the different results [1] . The F-measure can be a number between 0 and 1 and measures the accuracy of the test. The closer the result is to 1 the better.

$$FMeasure = 2 * \frac{Precision * Recall}{Precision + Recall} \qquad 3$$

The harmonic mean from precision and recall gives us a good measure of which method is better: Applying LDA on full-text papers or on title, abstract and keywords.

With every query executed in the two LDA models (one for full-text, another for title, abstract and keywords) and all the human annotated data collected, we will now outline what the confusion matrices looks like with the corresponding values.

### 4.2.3 Evaluation 1: All LDA topics, full-text articles vs title, abstract & keywords

Figure 4: Full text, all topics

**Returned by LDA?**

|  | Yes | No |
|---|---|---|
| **Yes** | 102 | 36 |
| **Relevant?** | | |
| **No** | 813 | 3198 |

The values generated by the table above is:

$$Precision = \frac{102}{102 + 813} = 0,111 \qquad 4$$

$$Recall = \frac{102}{102 + 36} = 0,739 \qquad 5$$

$$FMeasure = 2 * \frac{0,111 * 0,739}{0,111 + 0,739} = \mathbf{0{,}193} \qquad 6$$

Figure 5: Title, abstract and keywords, all topics

**Returned by LDA?**

|  | Yes | No |
|---|---|---|
| **Yes** | 83 | 54 |
| **Relevant?** | | |
| **No** | 591 | 2273 |

The values generated by the table above is:

$$Precision = \frac{83}{83 + 591} = 0,123 \qquad 7$$

$$Recall = \frac{83}{83 + 54} = 0,606 \qquad 8$$

$$FMeasure = 2 * \frac{0,123 * 0,606}{0,123 + 0,606} = \mathbf{0{,}204} \qquad 9$$

Regarding "**RQ 2: Which approach is better for Automatic Topic Extraction: a) Extraction from title, abstract and keywords or b) Extraction from full text paper?**" The F-Measure is slightly higher for title, abstract and keywords. However with such a small difference we can't safely say that one type is better than the other.

To answer "**RQ 3: How well does the approach of using Latent Dirichlet Allocation (with suitable pre-processing) perform compared to a manual method?**" We assume that the human performance is perfect, since that is what is accepted and applied today in the Software Engineering domain. So when looking at the amount of false negatives stored (the amount of topics that should have been returned by the machine, but were not) in these two confusion matrices. The full-text gives us 36 and the title, abstract and keyword set 54. So that tells us that full-text data set returns the relevant topics more often than the title, abstract and keywords data set. So the full-text missed 36 topics that the humans had deemed relevant and the title, abstract and keyword missed 54. This is the indication of how much the humans and the algorithm disagree

### 4.2.4 Evaluation 2: Top 3 LDA topics, full-text articles vs title, abstract & keywords

Figure 6: Full text, top three topics

**Returned by LDA?**

|  | Yes | No |
|---|---|---|
| **Yes** | 46 | 103 |
| **Relevant?** | | |
| **No** | 103 | 3897 |

The values generated by the table above is:

$$Precision = \frac{46}{46 + 103} = 0,309 \qquad 10$$

$$Recall = \frac{46}{46 + 103} = 0,309 \qquad 11$$

$$FMeasure = 2 * \frac{0,309 * 0,309}{0,309 + 0,309} = \mathbf{0,309} \qquad 12$$

Figure 7: Title, abstract and keywords, top three topics

**Returned by LDA?**

|  | **Yes** | **No** |
|---|---|---|
| **Yes** | 35 | 115 |
| **Relevant?** | | |
| **No** | 115 | 2735 |

The values generated by the table above is:

$$Precision = \frac{35}{35 + 115} = 0,233 \qquad 13$$

$$Recall = \frac{35}{35 + 115} = 0,233 \qquad 14$$

$$FMeasure = 2 * \frac{0,233 * 0,233}{0,233 + 0,233} = \mathbf{0,233} \qquad 15$$

In regards of "**RQ 2: Which approach is better for Automatic Topic Extraction: a) Extraction from title, abstract and keywords or b) Extraction from full text paper?**" When we ask the model to only return the three most probable topics per paper we get a higher F-measure for the full-text articles and the title, abstract and keywords, than when we asked it to return all topics. This is probably because when the Gensim framework only returns three topics, it returns fewer false positives, thus the value of precision is higher. Though due to probability, there is a smaller chance for the annotators to agree with the machine with only three returned topics. So the recall value is smaller, due to the higher value of false negatives.

Full-text also performs somewhat better than title, abstract and keywords when looking at the top 3 most probable topics.

Regarding "**RQ 3: How well does the approach of using Latent Dirichlet Allocation (with suitable pre-processing) perform compared to a manual method?**" The topics that the machine should have returned. When only using the three most likely topics, there are a lot more topics in the false negative boxes than when returning all topics. There is a bigger chance when returning all topics that the topic the annotator deemed relevant will show up. However, between the two data sets when only returning the three most likely topics, yet again, the full-texts model returns more relevant topics than the title, abstract and keywords. This is since the full-texts confusion matrices only contains 103 false negatives and the other 115.

#### 4.2.5 Evaluation 3: Most probable LDA topic, full-text articles vs title, abstract & keywords

However we are also interested in looking at the most probable topic for each paper (the topic with the highest probability weight according to the algorithm) and comparing this to the human evaluation.

Therefor (for each paper) we also do a simple binary comparison to see if the most probable topic according to the machine is among the three topics identified by the human annotators.

Figure 8 provides a simplified overview of how this evaluation was performed. First we compared the labeling made by human annotators with the machines categorization for the full-text articles and secondly we compared the same results for title, keyword and abstracts.

For **RQ 2: Which approach is better for Automatic Topic Extraction: a) Extraction from title, abstract and keywords or b) Extraction from full text paper?** its a bit difficult to motivate using precision and recall since if the machine would correctly return a relevant topic, there would still be two false negatives left. So a more simple approach is applied for this evaluation. One where if the machine returned a topic that was among the three the humans had deemed relevant it is labeled as a hit. The data

| Most probable topic according to the Model | Documents with a hit | Missed documents | Hit-ratio |
| --- | --- | --- | --- |
| Full-text articles | 17 | 33 | 0,34 |
| Title/abstract/key-words | 13 | 37 | 0,26 |

Figure 8: Only the top favorable topic returned from the queries

This is the result of a comparison of how often the most favorable topic returned from a query was among the three topics assigned from the annotators

sets model with most hits should therefor have returned the most relevant topic as their most probable topic. In the case of this study, please refer to figure 8 to observe that the full-text has a hit rate of 0.34 and title, abstract and keywords only have 0.26. So in this case it seems that the full-text data set has out performed the title, abstract and keywords.

This is our final evaluation in regards to "**RQ 3: How well does the approach of using Latent Dirichlet Allocation (with suitable pre-processing) perform compared to a manual method?**".

We simply check if the most probable topic according to LDA is among the three topics chosen by the human annotators for each article (see figure 8). Here the model also performs slightly better on full-text articles than on title, abstract and keywords. For 17 out of 50 documents, the most probable topic according to the model is also among the topics chosen by the annotators. For title, abstract and keywords. the same number is 14 out of 50. This gives a hit-ratio of 0.34 for full-text and 0.26 for title, abstract and keywords.

# 5 Analysis & Discussion

## 5.1 Analysis

With the result from the human annotators compared to the model, it seems fair to argue that the machine and humans agree more when both are supplied the articles in their entirety.

From the evaluation results using Precision And Recall we can see that the algorithm performs better when evaluating full-text articles rather than title, abstract and keywords and only looking at the top 3 topics.

When comparing the full-text, all topics result with the Abstract and keywords, all topics result, the F-measure of the lastly mentioned is however actually 0.011 higher than the F-measure of the full-text evaluation.

The reason why it still seems fair to argue that the full-text evaluation outperforms the Title, abstract and keywords, is because of when the machine presents its most probable choices of topics. Then the F-measure is much higher in the full-text evaluation. Just to add to this reasoning, another comparison was made with the singular most probable topic according to the machines and the annotators topics, as shown in figure 8. Yet again (with other measurements however) it is clear that when supplying full-text data sets to the machine, it performs better.

Worth mentioning is that when the model for title, abstract and keywords had been trained, it generated far fewer interpret-able topics when the time came to label them. In fact, for the full text model, 83 clear and usable topics were generated as for the abstract and keywords model, only 60 clear and usable topics were generated. So that explains the lower values of the true negatives in the abstract and keywords data sets.

Another reason why we wanted to compare the differences between the results of asking the model for all topics with the model's top 3 topics was to show how the Gensim implementation of LDA produces a lot of topics for some documents with this data set. A lot of these topics get very low probability scores (see appendix) which is why there are a lot less false positives when we just look at the top 3 topics.

## 5.2 Discussion

### Using AuTopEx for Topic Extraction

With all the tools in place a researcher only needs to do the following in order to perform automatic topic extraction:

1. Batch-convert all desired pdf:s.
2. Run the pre-processing script.
3. Train the LDA model using part of the corpus.
4. Query the model with the desired number of

remaining documents from the corpus.

From our experience document conversion and text-cleaning takes the longest time. For a large corpus ($>$ 2000 scientific papers for example) each of these steps can take several hours. The researcher however does not need to be present while the programs are running. Training a model on 200 full-text papers took 40 minutes using a cheap laptop with a Celeron processor clocked at 2.0 GHz (utilizing two of the processor cores). Querying the trained model with 50 papers using the same computer is done in a couple of minutes.

Seeing as how it takes a human reader many hours to read and annotate 50 scientific articles, using an approach such as AuTopEx can greatly speed up topic extraction. Especially during tasks that require a researcher to read a large amount of articles, (such as when doing document classification in a Systematic Mapping Study).

Of course this requires that the model classifies the papers accurately enough, and there is room for improving AuTopEx here.

## General Discussion

For the full text evaluation, the most probable topic identified by the algorithm was indeed a topic in the paper in 34 % of the cases according to the human annotators. This might not sound as a huge percentage, but seeing as this was the very first evaluation of the AuTopEx approach it seems very promising. Especially when one compares the many hours it takes for a human to read 50 scientific papers compared to the mere minutes it took the algorithm to produce this result.

It can be a good idea to perform word analysis on the corpus using NLTK after text pre-processing, for example checking a lot of the most popular words in the corpus. While time-consuming it can give insights into if some of the pre-processing steps might need adjusted. For example, perhaps there are still words in the corpus that could be considered stop words.

If batch-converting a large number of documents we recommend that the file sizes of the documents are checked afterwards. If any of the text-files have a size of 0 kilobytes the conversion has failed.

## Discussion on Topics and their labeling

Labeling topics manually when performing evaluation can be a very difficult task. It requires both language skills as well as domain knowledge. Sometimes the words in a topic are acronyms or words that have no meaning to those not familiar with the domain. Making sure that you found the correct meaning of the acronym (often an acronym has a number of meanings in a multitude of fields) or finding an explanation of a very niche word can be quite time consuming.

Interestingly enough adjectives were very uncommon in the results from our corpus. Besides "autonomous", which came up in 17 topics, the results were dominated by nouns, followed by verbs. For the full-text experiment only two other adjectives appeared, "intelligent" and "content", and the latter is also a noun. For title/abstracts/keywords the adjectives were more varied: "Intelligent", "dynamic", "generalized", "industrial", "artificial", "automatic" and "natural" appeared. The word "real" appeared three times (and was always accompanied by "time" in the same topic).

The dominance of nouns was quite helpful when labeling the technology-oriented topics often found in the software engineering domain. This is especially true when doing classification that does not take positives and negatives into account (we don't need colorful adjectives criticizing or praising something in a topic). Words like "car", "architecture" or "network" tells us a great deal on their own. Verbs are helpful in a supporting role (such as "driving" appearing in a topic with "autonomous" and "vehicle").

Another interesting note was that even though the entire corpus consisted of scientific papers, none of the topics produced in either of the two evaluation experiments were about scientific methodology. This is useful data to extract when performing tasks like Systematic Mapping Studies.

We found it interesting how the LDA model produced a lot of potential topics with low probabilities on the corpus on autonomous vehicle research. This could be due to how the scientific articles are written, but requires further study before any conclusions can be drawn.

## 5.3 Threats to Validity

It's important to remember that LDA is a probabilistic topic model, thus we are dealing with probabilities. If a human claims that a paper is about a certain topic and the machine claims that this probability is high, we only argue that the likelihood of this to be true is very high.

Properly labeling topics and scientific papers requires a lot from the human annotators. They must have excellent language skills as well as domain expertise in order to interpret each topic supplied by the model. One misunderstanding of a word could result

in an improper label, and this could impact the results of the evaluation.

We mitigated this by reading about concepts we were not familiar with before finishing the topic labeling, and looking up the meaning of any acronyms that appeared in the results. Both authors are software engineering students. Having previously studied concepts such as image processing or lane following for autonomous vehicles meant that we had a good understanding of a large majority of the topics produced by the model.

Then again, labeling 100 topics and reading 50 scientific articles (for two separate evaluations) can be difficult for humans. Stress, fatigue or just having a bad day can impact the accuracy both when performing topic labeling and when manually assigning topics to documents. We tried to mitigate this by taking breaks regularly during the evaluation. However if other researchers would redo our evaluation, using the same articles, we can not say for certain that they would label every single topic or classify the papers in exactly the same way.

Our main mitigation strategy for human error was that there were two of us doing the same work in parallel. We continuously compared the results between ourselves and where there were any disagreements regarding topic labels or which topics belong to a certain paper, we tried to reason with each other until we came to a result we could both agree on.

Another thing to consider when performing this kind of automatic topic extraction is that there is no way of handling positives and negatives. A paper that deals with a certain topic may actually reject the idea behind that topic. We mitigate this by not making any specific claims regarding the documents. We only state that in the results where human and machine agree that a topic exists in a document, that topic is indeed discussed in that specific document.

AuTopEx has only been evaluated on a corpus in the scientific domain of Autonomous Vehicles. We can't say with certainty how different the evaluation results would be if applying the approach on corpora from other domains. However, steps have been taken to make AuTopEx as generally applicable as possible. Especially by only limiting the tokenization to alphabetical words and using a very general stop word list. We recommend that anyone who uses this approach carefully consider if there are any special measures to be taken in the text pre-processing stage (e.g. adding words to the stop words list).

During the testing phase, we noticed that on some occasions several words would appear together as a single token after the texts had been cleaned and we sus-

pect this is due to bad quality of some of the original pdf:s. While it would be far too time-consuming to check the entire corpus manually for this we believe that this should very seldom occur in the data set we used for evaluation. This is because this data has been screened by researchers and only contains pdf:s published in 2005 or later.

## 6    Conclusions and Future Work

In this thesis we presented an approach for Automatic Topic Extraction which we call AuTopEx. This approach uses Natural Language Processing tools and techniques to pre-process the scientific articles of a corpus. Topics from this corpus are then extracted by training and querying a Latent Dirichlet Allocation (LDA) model. This model can be used to automatically classify the documents of the corpus (identifying which topics exist in which articles).

According to our results, Automatic Topic Extraction with Latent Dirichlet Allocation works better on full-text scientific articles than documents that consist of title, abstract and keywords. This is true both when querying the model for the most probable topic per article as well as when asking the model for the three most probable topics per article.

In our evaluation, the model's most probable topic was among the three relevant topics (according to the human annotators) in 34 % of the full-text documents evaluated. While the model is not as accurate than the human annotators it is important to note that this was the first evaluation of AuTopEx and perhaps most important of all: The model does this work in a couple of minutes while it takes humans many hours to perform the same task.

We believe that by refining this approach it will be possible to speed up topic extraction tremendously compared to manually reading and annotating papers.

**Future work**

One possible future experiment could be to allow the use of n-grams in the data set before performing the machine learning algorithm. If for example "autonomous vehicle" was considered a single word it could free up more space for other words to occur together with it in topics, possibly allowing for more meaningful interpretations by human readers. This process could also easily be automated. NLTK for example has the tool Collocations which performs n-gram analysis on documents.

Another idea that could possibly improve the results of our approach is to apply tf-idf on the text corpus

before training the model. Tfidf is the product of two statistics, term frequency and inverse document frequency. Term Frequency is the number of times a term occurs in a document. Inverse Document Frequency is a factor that diminishes the weight of terms that occur very frequently in the document set and increases the weight of terms that occur rarely. Thus a word like "the" will have a very low weight in tf-idf.

A high weight in tfidf is reached by a high term frequency (in the given document) and a low document frequency of the term in the whole collection of documents; the weights hence tend to filter out common terms. This could potentially be used for stop word removal.

The results would depend on how focused the language is in the different articles. An article which uses very broad language (using many synonyms for the same word) will produce different results than an article with very focused language. One idea could also be to duplicate the title of the paper a couple of times in each document before applying tf-idf. Seeing as how the title should reflect what the text is about this would help ensure that the most important words of the papers get a higher weight. Another experiment with tf-idf could be to give all nouns and verbs higher weight since they convey a lot of information about technologically-oriented topics.

A domain-specific lemmatizer for text pre-processing could be useful. This would however require a lot of work by several domain experts for a gold standard to be achieved and might be an unrealistic thing to wish for.

Automatization of the labeling stage could make the threat towards validity smaller while making the entire process quicker and easier to use, since there is less required input from the user. Such tools are already being applied[17].

## Acknowledgements

## References

[1] S. Bird, E. Klein, and L. Edward. *Analyzing Text with the Natural Language Toolkit.* O'Reilly Media, Sebastopol, California, United States, 2009.

[2] M. D. Blei. Introduction to probabilistic topic models. *Prinston University*, pages 1–16, 2011.

[3] M. D. Blei and et.al. Latent dirichlet allocation. *Journal of Machine Learning Research*, (3), 2003.

[4] J. Chang, J. Boyd-Graber, and et.al. Reading tea leaves: How humans interpret topic models. *Neural Information Processing Systems*, pages 1–9, 2009.

[5] R. Felizardo, Katia, N. Salleh, M. Martins, Rafael, E. Mendes, G. MacDonell, Stephen, Maldonado, and C. Jose. Using visual text mining to support the study selection activity in systematic literature reviews. *International Symposium on Empirical Software Engineering and Measurement*, pages 77–86, 2011.

[6] P. Fernando, V, C. Maria, O. F, and M. Rosane. The projection explorer: A flexible tool for projection-based multidimensional visualization. *Analytical and Bioanalytical Chemistry Volume 400, Number 4*, pages 1153 – 1159, 2011.

[7] R. Giuseppe and et.al. Semantic enrichment for recommendation of primary studies in a systematic literature review. *Digital Scholarship in the Humanities Advance Access*, pages 1–14, 2015.

[8] Henver, S. T. March, J. Park, and S. Ram. Design science in information systems research. *MIS Quarterly*, pages 1–14, 2004.

[9] A. Hindle and et.al. On the naturalness of software. *International Conference on Software Engineering (ICSE)*, pages 837–847, 2012.

[10] N. Jan van Eck and L. Waltman. Text mining and visualization using vosviewer. *ISSI newsletter*, pages 50–54, 2011.

[11] N. Jan van Eck and L. Waltman. Citnetexplorer: A new software tool for analyzing and visualizing citation networks. *Journal of Informetrics*, pages 802–823, 2014.

[12] N. Jan van Eck and L. Waltman. *Visualizing bibliometric networks. In Y. Ding, R. Rousseau, & D. Wolfram (Eds.), Measuring scholarly impact: Methods and practice.* Springer Publishing Company, 11 West 42nd Street, 15th Floor New York, NY 10036, 2014.

[13] A. G. Jivani. A comparative study of stemming algorithms. *International Journal of Computer Technology and Applications*, (Vol 2: Issue 6), 2011.

[14] M. Liakata, S. Dobnik, S. Saha, C. Batchelor, and D. Rebholz-Schuhmann. A discourse-driven content model for summarising scientific articles evaluated in a complex question answering task. *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, page 747757, 2013.

[15] M. Liakata, S. Saha, S. Dobnik, B. Colin, and D. Rebholz-Schuhmann. Automatic recognition of conceptualization zones in scientific articles and two life science applications. *Bioinformatics*, pages 991–1000, 2012.

[16] C. Marshall and P. Brereton. Tools to support systematic literature reviews in software engineering: A mapping study. *2013 ACM / IEEE International Symposium on Empirical Software Engineering and Measurement*, pages 296 – 299, 2013.

[17] Q. Mei, X. Shen, and C. Zhai. Automatic labeling of multinomial topic models. pages 1–10, 2007.

[18] K. Petersen, R. Feldt, S. Mujtaba, and M. Mattsson. Systematic mapping studies in software engineering. pages 1–10, 2008.

[19] R. Rehurek and P. Sojka. Software framework for topic modelling with large corpora. *Natural Language Processing Laboratory Masaryk University*, pages 1–5, 2010.

| title | true positive | false negative | false positive | true negative | total number of t | 17 skräp | 83 bra topics | |
|---|---|---|---|---|---|---|---|---|
| Fisheye optics for omnidirectional perception | 2 | 1 | 16 | 64 | 18 | | | |
| Data age based retransmission scheme for reliable control data exchange in platooning applications | 2 | 0 | 19 | 62 | 21 | | | 89, 4, 16, 52, |
| Obstacle Avoidance in Real Time with Nonlinear Model Predictive Control of Autonomous Vehicles | 3 | 0 | 15 | 65 | 18 | | | |
| Intelligent Cruise Control Stop and Go with and without Communication | 3 | 0 | 13 | 67 | 16 | | | |
| Autonomous Navigation: Achievements in Complex Enviro | 1 | 0 | 25 | 57 | 26 | | | |
| Bayesian Network Based Collision Avoidance | 2 | 1 | 17 | 63 | 19 | | | |
| Experience, Results and Lessons Learned from Automated Driving on Germany's Highways | 3 | 0 | 19 | 61 | 22 | | | |
| Multi-Objective Path Planning using Spline Represent | 3 | 0 | 20 | 60 | 23 | | | |
| A Study on Autonomous Vehicle Development Process at University* | 3 | 0 | 13 | 67 | 16 | | | |
| Road Surface Recognition Using Laser Radar for Automatic Platooning | 1 | 2 | 20 | 60 | 21 | | | |
| Building a Prototype for Power-Aware Automatic Parking System | 2 | 0 | 19 | 62 | 21 | | | |
| A Computer Vision System for Detection and Avoidance for Automotive Vehicles | 3 | 0 | 11 | 68 | 15 | | | |
| Path Tracking of Autonomous Ground Vehicle Based on Fractional Order PID Controller Optimized by PSO | 2 | 0 | 12 | 69 | 14 | | | |
| Off-road Path Following using Region Classification and G Constraints* | 3 | 0 | 19 | 61 | 22 | | | |
| Self-Tuning PID Controller for Autonomous Car Tracking in Urban Traffic | 2 | 1 | 18 | 62 | 20 | | | |
| Shared Control of Autonomous Vehicles based on Velocity Space Optimization | 2 | 1 | 19 | 61 | 21 | | | |
| A 13,000 km Intercontinental Trip with Driverless Vehicles: | 2 | 1 | 11 | 69 | 13 | | | |
| Real-Time Coordination of Autonomous Vehicles | 1 | 1 | 17 | 64 | 18 | | | |
| Accurate and Efficient Traffic Sign Detection Using Discrim | 3 | | 16 | 64 | 19 | | | |
| DeepDriving: Learning Affordance for Direct Perception in | 2 | 1 | 16 | 64 | 18 | | | |
| A Robust Algorithm for the Detection of Vehicle Turn Signa | 2 | 0 | 5 | 76 | 7 | | | |
| Constrained Global Path Optimization for Articulated Steeri | 3 | 0 | 19 | 61 | 22 | | | |
| 360° detection and tracking algorithm of both pedestrian an using fisheye images | 3 | 0 | 15 | 64 | 19 | | | |
| State your position | 2 | 0 | 18 | 63 | 20 | | | |
| A robotic platform to evalute autonomous driving systems | 3 | 0 | 17 | 61 | 22 | | | |
| Coordinated control of multiple vehicles with discrete-time periodic communications | 1 | 3 | 17 | 63 | 17 | | | 89, 4, 16, 52, |
| Real-time Implementation of a Novel Safety Function for Pr | 3 | 0 | 12 | 68 | 15 | | | |
| Coordinated Path Following Control for a Group of Car-like | 1 | 2 | 9 | 71 | 10 | | | |
| A Combined Model- and Learning-Based Framework for In | 2 | 1 | 17 | 63 | 19 | | | |
| Towards a Framework for Testing Drivers' Interaction with | 1 | 1 | 16 | 65 | 17 | | | |
| Adopting WirelessHART for In-Vehicle-Networking | 2 | 0 | 18 | 63 | 20 | | | |
| Terrain Mapping for Off-road Autonomous Ground Vehicle | 3 | 0 | 18 | 62 | 21 | | | |
| Incremental Sampling-based Algorithm for Minimum-violation Motion Planning | 1 | 2 | 21 | 58 | 23 | | | |
| Vision-based Nighttime Vehicle Detection and Range Esti | 3 | 0 | 19 | 62 | 21 | | | |
| Design and Comparative Analysis of a Driveless LED light | 0 | 3 | 4 | 76 | 4 | | | |
| Local Path Planning for Off-Road Autonomous Driving With Avoidance of Static Obstacles | 3 | 0 | 24 | 56 | 27 | | | |
| HOG Based Multi-object Detection for Urban Navigation | 2 | 1 | 17 | 63 | 19 | | | |
| Genetic Algorithm Approach for Locating Automatic Vehicl Identification Readers | 0 | 1 | 17 | 65 | 17 | | | |
| Reliable Intersection Protocols Using Vehicular Networks | 3 | 0 | 11 | 69 | 14 | | | |
| INTELLIGENT TRAFFIC WITH CONNECTED VEHICLES | 2 | 1 | 18 | 62 | 20 | | | |
| MCMC Particle Filter for Real-Time Visual Tracking of Vehi | 2 | 1 | 22 | 58 | 24 | | | |
| Globally Asymptotically Stable Filter for Navigation aided b and Depth Measurements | 0 | 3 | 22 | 58 | 22 | | | |
| A Real-Time Multi-Sensor Fusion Platform for Automated | 1 | 2 | 2 | 78 | 3 | | | |
| A full-3D Voxel-based Dynamic Obstacle Detection for Urban Scenario using Stereo Vision | 3 | 0 | 11 | 69 | 14 | | | |
| A Real-Time Trajectory Control of Two Driving Mobile Rob | 1 | 2 | 18 | 64 | 19 | | | |
| Vehicle Automation in Cooperation with V2I and Nomadic Devices Communication | 2 | 1 | 19 | 61 | 21 | | | |
| Automatic vehicle classification and tracking method for ve movements at signalized intersections | 3 | 0 | 19 | 61 | 22 | | | |
| Multi-Target Tracking using a 3D-Lidar Sensor for Autono | 1 | 2 | 19 | 61 | 20 | | | |
| Traffic Sign Representation using Sparse-Representations | 1 | 1 | 20 | 61 | 21 | | | |
| Speed Profile Optimization for Vehicles Crossing an Intersection Under a Safety Constraint | 3 | 0 | 14 | 66 | 17 | | | |
| Sum | 102 | 36 | 813 | 3198 | 918 | | | |
| | true positive | false negative | false positive | true negative | total number of topics | | | |
| | | | | | | | | |
| | | | | | | | | |
| | Precision = | 0,111 | | | | | | |
| | Recall = | 0,739 | | | | | | |
| | F = | 0,1930094118 | | | | | | |

| title | true positive | false negative | false positive | true negative | total number of topics | 17 skräp | 83 bra topics |
|---|---|---|---|---|---|---|---|
| Fisheye optics for omnidirectional perception | 1 | 2 | 2 | 78 | 3 | | |
| Data age based retransmission scheme for reliable control data exchange in platooning applications | 2 | 0 | 0 | 80 | 3 | | |
| Obstacle Avoidance in Real Time with Nonlinear Model Predictive Control of Autonomous Vehicles | 0 | 3 | 3 | 77 | 3 | | |
| Intelligent Cruise Control Stop and Go with and without Communication | 0 | 3 | 3 | 77 | 3 | | |
| Autonomous Navigation: Achievements in Complex Enviro | 0 | 3 | 3 | 77 | 3 | | |
| Bayesian Network Based Collision Avoidance | 1 | 2 | 2 | 78 | 3 | | |
| Experience, Results and Lessons Learned from Automated Driving on Germany's Highways | 1 | 2 | 2 | 78 | 3 | | |
| Multi-Objective Path Planning using Spline Represent | 1 | 2 | 2 | 78 | 3 | | |
| A Study on Autonomous Vehicle Development Process at University* | 0 | 3 | 3 | 77 | 3 | | |
| Road Surface Recognition Using Laser Radar for Automatic Platooning | 1 | 2 | 2 | 78 | 3 | | |
| Building a Prototype for Power-Aware Automatic Parking System | 1 | 2 | 2 | 78 | 3 | | |
| A Computer Vision System for Detection and Avoidance for Automotive Vehicles | 1 | 2 | 2 | 78 | 3 | | |
| Path Tracking of Autonomous Ground Vehicle Based on Fractional Order PID Controller Optimized by PSO | 0 | 3 | 3 | 77 | 3 | | |
| Off-road Path Following using Region Classification and G Constraints* | 1 | 2 | 2 | 78 | 3 | | |
| Self-Tuning PID Controller for Autonomous Car Tracking in Urban Traffic | 1 | 2 | 2 | 78 | 3 | | |
| Shared Control of Autonomous Vehicles based on Velocity Space Optimization | 1 | 2 | 2 | 78 | 3 | | |
| A 13,000 km Intercontinental Trip with Driverless Vehicles: | 0 | 3 | 3 | 77 | 3 | | |
| Real-Time Coordination of Autonomous Vehicles | 0 | 3 | 3 | 77 | 3 | | |
| Accurate and Efficient Traffic Sign Detection Using Discrim | 0 | 3 | 3 | 77 | 3 | | |
| DeepDriving: Learning Affordance for Direct Perception in | 1 | 2 | 2 | 78 | 3 | | |
| A Robust Algorithm for the Detection of Vehicle Turn Signa | 2 | 1 | 1 | 79 | 3 | | |
| Constrained Global Path Optimization for Articulated Steeri | 2 | 1 | 1 | 79 | 3 | | |
| 360° detection and tracking algorithm of both pedestrian an using fisheye images | 1 | 2 | 2 | 78 | 3 | | |
| State your position | 1 | 2 | 2 | 78 | 3 | | |
| A robotic platform to evalute autonomous driving systems | 1 | 2 | 2 | 78 | 3 | | |
| Coordinated control of multiple vehicles with discrete-time periodic communications | 2 | 1 | 1 | 79 | 3 | | |
| Real-time Implementation of a Novel Safety Function for Pr | 1 | 2 | 2 | 78 | 3 | | |
| Coordinated Path Following Control for a Group of Car-like | 1 | 2 | 2 | 78 | 3 | | |
| A Combined Model- and Learning-Based Framework for In | 1 | 2 | 2 | 78 | 3 | | |
| Towards a Framework for Testing Drivers' Interaction with | 0 | 3 | 3 | 77 | 3 | | |
| Adopting WirelessHART for In-Vehicle-Networking | 3 | 0 | 0 | 80 | 3 | | |
| Terrain Mapping for Off-road Autonomous Ground Vehicle | 0 | 3 | 3 | 77 | 3 | | |
| Incremental Sampling-based Algorithm for Minimum-violation Motion Planning | 1 | 2 | 2 | 78 | 3 | | |
| Vision-based Nighttime Vehicle Detection and Range Esti | 2 | 1 | 1 | 79 | 3 | | |
| Design and Comparative Analysis of a Driveless LED light | 0 | 3 | 3 | 77 | 3 | | |
| Local Path Planning for Off-Road Autonomous Driving With Avoidance of Static Obstacles | 2 | 1 | 1 | 79 | 3 | | |
| HOG Based Multi-object Detection for Urban Navigation | 0 | 3 | 3 | 77 | 3 | | |
| Genetic Algorithm Approach for Locating Automatic Vehicl Identification Readers | 0 | 3 | 3 | 77 | 3 | | |
| Reliable Intersection Protocols Using Vehicular Networks | 2 | 1 | 1 | 79 | 3 | | |
| INTELLIGENT TRAFFIC WITH CONNECTED VEHICLES | 1 | 2 | 2 | 78 | 3 | | |
| MCMC Particle Filter for Real-Time Visual Tracking of Vehi | 0 | 3 | 3 | 77 | 3 | | |
| Globally Asymptotically Stable Filter for Navigation aided b and Depth Measurements | 0 | 3 | 3 | 77 | 3 | | |
| A Real-Time Multi-Sensor Fusion Platform for Automated | 1 | 2 | 2 | 78 | 3 | | |
| A full-3D Voxel-based Dynamic Obstacle Detection for Urban Scenario using Stereo Vision | 1 | 2 | 2 | 78 | 3 | | |
| A Real-Time Trajectory Control of Two Driving Mobile Rob | 2 | 1 | 1 | 79 | 3 | | |
| Vehicle Automation in Cooperation with V2I and Nomadic Devices Communication | 2 | 1 | 1 | 79 | 3 | | |
| Automatic vehicle classification and tracking method for ve movements at signalized intersections | 1 | 2 | 2 | 78 | 3 | | |
| Multi-Target Tracking using a 3D-Lidar Sensor for Autono | 1 | 2 | 2 | 78 | 3 | | |
| Traffic Sign Representation using Sparse-Representations | 1 | 2 | 2 | 78 | 3 | | |
| Speed Profile Optimization for Vehicles Crossing an Intersection Under a Safety Constraint | 1 | 2 | 2 | 78 | 3 | | |
| Sum | 46 | 103 | 103 | 3897 | | | |
| | true positive | false negative | false positive | true negative | total number of topics | | |
| | | | | | | | |
| Precision = | | 0,309 | | | | | |
| Recall = | | 0,309 | | | | | |
| F = | | 0,309 | | | | | |

| title | true positive | false negative | false positive | true negative | total number of topics | | 40 skräp | 60 bra topics |
|---|---|---|---|---|---|---|---|---|
| A._B._P._C._S._D._M._C._ | 1 | 2 | 14 | 43 | 15 | | | |
| A._B._S._D._M._P._P._P._ | 2 | 1 | 9 | 48 | 11 | | | |
| A._B.-N._C._Grand_2012_ | 3 | 0 | 11 | 46 | 14 | | | |
| A._C._C._D._Gillet_2014_S | 3 | 0 | 12 | 45 | 15 | | | |
| A._C._L._N._S._M._M._N._ | 0 | 3 | 8 | 49 | 8 | | | |
| B._B._H._Giese_2008_Incr | 1 | 0 | 15 | 44 | 16 | | | |
| B._W._A._K._M._P._T._A._ | 0 | 2 | 16 | 42 | 16 | | | |
| B.-M._S._Chung,_Jin-Woo; | 0 | 3 | 12 | 45 | 12 | | | |
| C._C._J._Liu_2010_A_Rein | 3 | 0 | 13 | 44 | 16 | | | |
| C._C._Y._H._F._G._C._B._ | 2 | 1 | 12 | 45 | 14 | | | |
| C._L._B._N._T._M._C._S._ | 2 | 1 | 14 | 43 | 16 | | | |
| C._W._Axelrod_2015_Enfor | 0 | 1 | 16 | 43 | 16 | | | |
| D._B._W._M._I._Posner_20 | 3 | 0 | 13 | 44 | 16 | | | |
| D._C._S._D._B._P._Stone_ | 3 | 0 | 8 | 49 | 11 | | | |
| G._A._J._I._N._E._M._Neb | 2 | 1 | 13 | 44 | 15 | | | |
| H._x._E._C,_;ne,;T._Sattler; | 2 | 1 | 14 | 43 | 16 | | | |
| J._A._C._S._A._Pascoal_2 | 3 | 0 | 10 | 47 | 13 | | | |
| J._C._S._U._B._L._M._Mau | 3 | 0 | 11 | 46 | 14 | | | |
| J._S._B._P._H._H._Chen_2 | 3 | 0 | 12 | 45 | 15 | | | |
| K._B._H._M._A._Zell_2012 | 1 | 1 | 9 | 49 | 10 | | | |
| K._C._F._J._T._R._S._J._B | 2 | 1 | 14 | 43 | 16 | | | |
| L._C._A._F._L._Pallottino_2 | 1 | 1 | 15 | 43 | 16 | | | |
| L._x._F._J._M._Alvarez;F._ | 2 | 1 | 10 | 47 | 12 | | | |
| M._A._A._R._M._J._M._Ekl | 2 | 1 | 9 | 48 | 11 | | | |
| M._A._J._M._Dolan_2011_ | 1 | 2 | 11 | 46 | 12 | | | |
| M._A._P._F._C._O._J._Sjo | 1 | 2 | 13 | 44 | 14 | | | |
| M._A.-M._W._S._M._Y._W | 1 | 2 | 9 | 48 | 10 | | | |
| M._B._C._H._A._L._M._A._ | 2 | 1 | 15 | 43 | 16 | | | |
| M._B._Z._G._P._Z._M._B._ | 0 | 3 | 11 | 46 | 11 | | | |
| M._C._D._P._M._Pasquier_ | 0 | 3 | 15 | 42 | 15 | | | |
| M._H._Ang_2015_Achievin | 0 | 0 | 14 | 46 | 14 | | | |
| M._J._B._C._M._Veth_201 | 1 | 1 | 11 | 47 | 12 | | | |
| M._J._H._Berg;R._Olsson; | 1 | 2 | 10 | 47 | 11 | | | |
| M._x._E._A,_;yr,;x00E,;M._ | 0 | 1 | 5 | 54 | 5 | | | |
| N._C.-B._A._M._J._R._M._ | 3 | 0 | 13 | 44 | 16 | | | |
| N._T._Atsuhiro,_Yamaguchi | 2 | 1 | 10 | 47 | 12 | | | |
| P._B._D._K._C._B._J._Dick | 3 | 0 | 13 | 44 | 16 | | | |
| P._V._K._B._S._Vidas_201 | 2 | 1 | 14 | 43 | 16 | | | |
| P._V._M._E._O._J._R._d._ | 2 | 1 | 11 | 46 | 13 | | | |
| Q._B._M._P._C._Laugier_2 | 1 | 2 | 5 | 52 | 6 | | | |
| S._A._G._B._R._R._P._Mu | 2 | 1 | 11 | 46 | 13 | | | |
| S._A._S._C._Y._S._Alj_201 | 2 | 1 | 9 | 48 | 11 | | | |
| S._B._M._M.-P._R._M.-P._ | 2 | 1 | 14 | 43 | 16 | | | |
| S._D._B._B._E._A._Speran | 1 | 2 | 15 | 42 | 16 | | | |
| S._J._A._S._B._K._K._I._J. | 1 | 2 | 6 | 51 | 7 | | | |
| S._P._B._R._W._Sadowski | 1 | 2 | 15 | 42 | 16 | | | |
| T._A._M._M._M._Ali_2015_ | 2 | 1 | 14 | 43 | 16 | | | |
| Y._A._P._P._F._P._A._Burr | 2 | 1 | 14 | 43 | 16 | | | |
| Z._B._J._J._N._Y._S._Linc | 3 | 0 | 10 | 47 | 13 | | | |
| Z._K._x._E._T._Akg;x00Fc, | 3 | 0 | 13 | 44 | 16 | | | |
| Sum | 83 | 54 | 591 | 2273 | 673 | | | |
| | true positive | false negative | false positive | true negative | total number of topics | | | |
| | | | | | | | | |
| | | | | | | | | |
| | Precision = | 0,123 | | | | | | |
| | Recall = | 0,606 | | | | | | |
| | F = | 0,2044938272 | | | | | | |

| title | true positive | false negative | false positive | true negative | total number of topics | | 40 skräp | 60 bra topics |
|---|---|---|---|---|---|---|---|---|
| A._B._P._C._S._D._M._C. | 0 | 3 | 3 | 54 | 3 | | | |
| A._B._S._D._M._P._P._P. | 1 | 2 | 2 | 55 | 3 | | | |
| A._B.-N._C._Grand_2012 | 1 | 2 | 2 | 55 | 3 | | | |
| A._C._C._D._Gillet_2014_ | 0 | 3 | 3 | 54 | 3 | | | |
| A._C._L._N._S._M._M._N. | 1 | 2 | 2 | 55 | 3 | | | |
| B._B._H._Giese_2008_Inc | 0 | 3 | 3 | 54 | 3 | | | |
| B._W._A._K._M._P._T._A. | 1 | 2 | 2 | 55 | 3 | | | |
| B.-M._S._Chung,_Jin-Woo | 1 | 2 | 2 | 55 | 3 | | | |
| C._C._J._Liu_2010_A_Rei | 2 | 1 | 1 | 56 | 3 | | | |
| C._C._Y._H._F._G._C._B. | 1 | 2 | 2 | 55 | 3 | | | |
| C._L._B._N._T._M._C._S. | 0 | 3 | 3 | 54 | 3 | | | |
| C._W._Axelrod_2015_Enf | 0 | 3 | 3 | 54 | 3 | | | |
| D._B._W._M._I._Posner_2 | 1 | 2 | 2 | 55 | 3 | | | |
| D._C._S._D._B._P._Stone | 0 | 3 | 3 | 54 | 3 | | | |
| G._A._J._I._N._E._M._Ne | 1 | 2 | 2 | 55 | 3 | | | |
| H._x._E._C,_;ne,;T._Sattle | 2 | 1 | 1 | 56 | 3 | | | |
| J._A._C._S._A._Pascoal_ | 0 | 3 | 3 | 54 | 3 | | | |
| J._C._S._U._B._L._M._Ma | 3 | 0 | 0 | 57 | 3 | | | |
| J._S._B._P._H._H._Chen_ | 0 | 3 | 3 | 54 | 3 | | | |
| K._B._H._M._A._Zell_201 | 1 | 2 | 2 | 55 | 3 | | | |
| K._C._F._J._T._R._S._J._ | 1 | 2 | 2 | 55 | 3 | | | |
| L._C._A._F._L._Pallottino | 1 | 2 | 2 | 55 | 3 | | | |
| L._x._F._J._M._Alvarez;F. | 1 | 2 | 2 | 55 | 3 | | | |
| M._A._A._R._M._J._M._E | 0 | 3 | 3 | 54 | 3 | | | |
| M._A._J._M._Dolan_2011 | 0 | 3 | 3 | 54 | 3 | | | |
| M._A._P._F._C._O._J._Sj | 0 | 3 | 3 | 54 | 3 | | | |
| M._A.-M._W._S._M._Y._ | 1 | 2 | 2 | 55 | 3 | | | |
| M._B._C._H._A._L._M._A. | 0 | 3 | 3 | 54 | 3 | | | |
| M._B._Z._G._P._Z._M._B. | 0 | 3 | 3 | 54 | 3 | | | |
| M._C._D._P._M._Pasquier | 0 | 3 | 3 | 54 | 3 | | | |
| M._H._Ang_2015_Achievi | 1 | 2 | 2 | 55 | 3 | | | |
| M._J._B._C._M._Veth_20 | 1 | 2 | 2 | 55 | 3 | | | |
| M._J._H._Berg;R._Olsson; | 0 | 3 | 3 | 54 | 3 | | | |
| M._x._E._A,_;yr,;x00E;;M. | 2 | 1 | 1 | 56 | 3 | | | |
| N._C.-B._A._M._J._R._M. | 2 | 1 | 1 | 56 | 3 | | | |
| N._T._Atsuhiro,_Yamaguc | 0 | 3 | 3 | 54 | 3 | | | |
| P._B._D._K._C._B._J._Dic | 1 | 2 | 2 | 55 | 3 | | | |
| P._V._K._B._S._Vidas_20 | 1 | 2 | 2 | 55 | 3 | | | |
| P._V._M._E._O._J._R._d. | 0 | 3 | 3 | 54 | 3 | | | |
| Q._B._M._P._C._Laugier_ | 0 | 3 | 3 | 54 | 3 | | | |
| S._A._G._B._R._R._P._M | 1 | 2 | 2 | 55 | 3 | | | |
| S._A._S._C._Y._S._Alj_20 | 1 | 2 | 2 | 55 | 3 | | | |
| S._B._M._M.-P._R._M.-P. | 0 | 3 | 3 | 54 | 3 | | | |
| S._D._B._B._E._A._Spera | 0 | 3 | 3 | 54 | 3 | | | |
| S._J._A._S._B._K._K._I._ | 1 | 2 | 2 | 55 | 3 | | | |
| S._P._B._R._W._Sadowsk | 0 | 3 | 3 | 54 | 3 | | | |
| T._A._M._M._M._Ali_2015 | 1 | 2 | 2 | 55 | 3 | | | |
| Y._A._P._P._F._P._A._Bu | 0 | 3 | 3 | 54 | 3 | | | |
| Z._B._J._J._N._Y._S._Lin | 2 | 1 | 1 | 56 | 3 | | | |
| Z._K._x._E._T._Akg;x00F | 1 | 2 | 2 | 55 | 3 | | | |
| **SUM:** | 35 | 115 | 115 | 2735 | | | | |
| | true positive | false negative | false positive | true negative | total number of topics | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | Precision = | 0,233 | | | | | | |
| | Recall = | 0,233 | | | | | | |
| | F = | 0,233 | | | | | | |

| Title | Topics | Vector |
|---|---|---|
| Fisheye optics for omnidirectional p | 96, 70, 99 | [(4, 0.03645642270614621б), (11, 0.01390029347610463Э), (12, 0.0298957050521015414), (27, 0.0104091008787905174), (32, 0.0084519107906727814), (45, 0.0168179326225496693), (55, … |
| Data age based retransmission sch … control data exchange in platooning | 86, 44 | [(4, 0.028874553229211141), (17, 0.03670069354163217б), (11, 0.0051178409894453223), (27, 0.0332234513454645253), (30, 0.0266026255782015622), (32, 0.0290140848625946934), (34, 0.0131155101716334), (38, … |
| Obstacle Avoidance in Real Time w Model Predictive Control of Autono | 78, 82, 57 | [(5, 0.038388696477847086), (9, 0.01065301332004176), (11, 0.0150753712772000043), (12, 0.0140088796449002024), (18, 0.023323205910297839964), (29, 0.03657937174521405212), (33, … |
| Intelligent Cruise Control Stop and Go with and without Com | 52, 55, 44 | [(5, 0.02763974534367649Ч), (12, 0.0191883622636389938), (25, 0.0223011498623795883), (34, 0.0236530003333320826), (38, 0.0626475742286854771), (42, 0.1047802285152695), (44, 0.243234815937234542), (45, … |
| Autonomous Navigation: Achievem | 87 | [(4, 0.0484607697513722696), (4, 0.001054335170088361Ч), (5, 0.037379429421918356), (6, 0.0152809799872924478), (11, 0.023002253031501745118), (23, 0.0474112651635080896), (57, … |
| Bayesian Network Based Collision | 95, number 4, 5 | [(2, 0.059117322211207163), (11, 0.0347108202025940303), (18, 0.0564250600112055531), (25, 0.11909768008701004449), (32, 0.01781772205520089899), (34, 0.0259069282103189501), (37, 0.0333728695102766944), (44, … |
| Experience, Results and Lessons Learned from Automated Driving on Germany's Highways | 37, 55, 62 | [(5, 0.105235153436964221), (10, 0.04720292987155119932), (11, 0.0897123581897818995193), (32, 0.033636294545462086), (41, … |
| Building a Prototype for Power-Awa Parking System | 47, 62 | [(0.00381732670683068), (9, 0.01565608533126467), (10, 0.01156689981897808), (11, 0.07215838305893408), (52, 0.012460472459398640), (83, 0.02803409020162020408), (84, 0.0951775130301042555), (87, … |
| A Computer Vision System for Dete Avoidance for Automotive Vehicles | 5, 2, 99 | [(2, 0.01617804455552536666), (5, 0.0380000871454568Ч), (5, 0.013530408714260407103), (11, 0.07215833032408457), (24, 0.034671280070427977), (31, 0.034447534279895359), (52, … |
| Multi-Objective Path Planning us | 62, 5, 57 | [(4, 0.072830703289600651), (24, 0.0171872607289007628), (25, 0.0154959587866114924), (34, 0.08326165573420067Э), (36, 0.02626575104342212), (38, 0.01364182625663210), (55, … |
| A Study on Autonomous Vehicle De University" | No relevant topic | [(4, 0.0173398131733975), (06, 0.02015467125371570Б), (61, 0.03387830537040617844), (66, 0.0784455055404029002), (74, 0.1271243522428508Ч), (79, 0.025137330500028009), (82, … |
| Path Tracking of Autonomous Grou Fractional Order PID Controller Opti | 62, vehicle control | [(0.031175422315106771Ч), (9, 0.04105266311468445907), (10, 0.014228096227839072), (18, 0.029844196763067090), (23, 0.0973985886794989009), (25, 0.0291135524708788241), (52, … |
| Off-road Path Following using Regi Constraints" | 48, 28, 63 | [(2, 0.0189666418648630094), (10, 0.01866307988372102), (11, 0.0003984366648324959), (12, 0.02041674212901099625), (15, 0.0165759826561807103), (18, 0.133811027080044132), (48, … |
| Self-Tuning PID Controller for Autonomous Car Tracking in Urban | 37, 95, Vehicle control | [(11, 0.01454683040979200497), (12, 0.0149952876603034576), (24, 0.115717074775080825), (29, 0.044392252640027999), (34, 0.0505121907616541127), (35, 0.0257210779098804Ч6), (42, … |
| Shared Control of Autonomous Veh based on Velocity Space Optimizati | 45, vehicle control, 78 | [(0.00440453192220097701), (62, 0.235653390500849109), (78, 0.0243284040248595109), (82, 0.01390402060825484417), (35, 0.018434660100698936), (38, 0.01098118487520086), (44, … |
| A 13,000 km Intercontinental Trip wi | 95, vehicle control, 82 | [(0.010106961582830335), (4, 0.0020935266188488066), (11, 0.031573857037845321), (34, 0.013940020600205254847), (35, 0.018843465001009393), (38, 0.010981118475a08086), (44, … |
| Real-Time Coordination of Autonom | 35, 86 | [(11, 0.01622743380524038Ч), (17, 0.0098733158028272603), (30, 0.065444383473319300026), (32, 0.0783232534350500471), (35, 0.0444224233944877163), (62, 0.0643518337900727779), (42, … |
| Accurate and Efficient Traffic Sign | 96, 85, 2 | [(2, 0.079236222886240097), (11, 0.0270706803130624084), (12, 0.01159437529804167), (15, 0.01026158800123819Э), (18, 0.108215843732727547), (28, 0.035099324478959924), (29, 0.07821729508144Э5812), (33, … |
| DeepDriving: Learning Affordance f | 77, 88, 2 | [(11, 0.02936725777807299Э2), (18, 0.12161960091533071), (35, 0.01080030607866004181), (37, 0.025123040971068743Z), (42, 0.0170102006029953), (55, 0.0531221066865129б), (62, … |
| A Robust Algorithm for the Detectio | 33, 88 | [(11, 0.020833932677210169), (12, 0.02616204272313882), (33, 0.814367748568701942), (57, 0.0122436495682600Ч), (94, 0.0124074989885117357), (95, 0.0208850438656307Ч2), … |
| Constrained Global Path Optimizati | 62, 78, 10 | [(11, 0.01701709338441009996), (12, 0.055209955333392901), (18, 0.045772637775049592), (32, 0.01222829661128847091), (33, 0.03395694503580011), (45, 0.026130647818693301), (47, 0.01943289873873511), (54, … |
| 360- detection and tracking algorith using fisheye images | 47, 96, 99 | [(11, 0.07164898288085124), (12, 0.0552009553339001), (79, 0.01035911709440092031Ч), (88, 0.05880123482667109657), (94, 0.0195683958629038081), (95, … |
| State your position | 1, 48, | [(1, 0.0259720902606675884), (12, 0.010128326163672211), (15, 0.0175194801117256651), (27, 0.01739219560536906024), (29, 0.0344452934002902049), (35, … |
| A robotic platform to evaluate autono | Vehicle control, 55, 60 | [(4, 0.01895733037892329935), (9, 0.024475029498239301), (47, 0.047207643686773б5), (60, 0.017659064568741523), (64, 0.00814680556874071266), (65, … |
| Coordinated control of multiple vehi discrete-time periodic communicatio | 77, 82, 86 | [(5, 0.0253020530948755579), (8, 0.0310310784837535131), (17, 0.07729467239258813188), (24, 0.01088584807733753013), (29, 0.040887867101693), (34, 0.030573093384064132), (35, … |

| Title | | |
|---|---|---|
| Real-time Implementation of a Nove | Vehicle control, 24, 25 | (4, 0.042323093953479394), (11, 0.0176272974445858600094), (12, 0.0198997138048633221), (17, 0.04361809940319596586), (24, 0.0411054211175794985), (25, 0.1102223186654111181), (27, ... |
| Coordinated Path Following Control | 82, 86, vehicle control | (5, 0.0112219451688871641), (8, 0.527630028354877761), (14, 0.015528531338256996), (16, 0.091083476359100135), (25, 0.0496121121566256619), (39, 0.01531000922606716), (55, 0.0108281763809087245), (62, ... |
| A Combined Model- and Learning-B | 35, 92, vehicle control | (6, 0.0180278050742172338), (12, 0.0208406074347469447), (17, 0.0387640759451490340), (24, 0.0290975750149034005), (24, 0.0193993000348488), (29, 0.010052986783699604816), (35, 0.0605185072188441), (47, ... |
| Towards a Framework for Testing D | 55, 45 | (5, 0.01632372262635422544), (12, 0.01275070022680222396), (18, 0.0101003611657351701), (28, 0.0039861611811722352), (33, 0.030049352181170478), (35, 0.02419829331655016?), (66, ... |
| Adopting WirelessHART for In-Vehi | 86, 61 | (4, 0.01753673346775707), (12, 0.02572142983851503), (14, 0.0117937327164685821), (18, 0.011556642698780133), (28, 0.003986161189117223521), (33, 0.0303049532181170478), (35, 0.02419829331655016), (61, ... |
| Terrain Mapping for Off-road Auton | 96, 99, 18, | (5, 0.0254001433947249309), (6, 0.01112761938387477), (8, 0.001664614512000123157), (11, 0.0231443032372463361), (11, 0.0256750952923920116), (41, 0.0164877925069111039), (59, ... |
| Incremental Sampling-based Algorit | 60, Vehicle control, 78 | (2, 0.00130369816466699), (8, 0.0276434428959587754), (11, 0.0238600301615118), (12, 0.04426784984416832563), (18, 0.03259804032661772286), (15, 0.012729151331878793903), (16, 0.01225116062633530561), (18, ... |
| Minimum-violation Motion Planning | | |
| Vision-based Nighttime Vehicle Det | 51, 96, 55 | (29, 0.03460609857430117178), (37, 0.69863093038102941453), (39, 0.0149464752319068571), (35, 0.1948959878785017951 ... |
| Design and Comparative Analysis o | 0, 33, 9 | (4, 0.0107820176907800890), (8, 0.01069137305569681), (10, 0.0552967724486137021), (11, 0.0249717117605139189247), (12, 0.0442678498416832563), (18, 0.0329804032266411772288), (15, ... |
| Local Path Planning for Off-Road A Driving With Avoidance of Static Ob | 62, 78, 99, | (2, 0.006132274238333372), (12, 0.065492966474895553), (12, 0.01543995258733788699), (13, 0.1317220151592422281), (28, 0.01546771711487384898), (32, 0.010455594380878792), (35, 0.0257754433534383196), (91, ... |
| HOG Based Multi-object Detection f | 98, 51, 87 | (2, 0.0606312727423833372), (55, 0.02031577778002027), (66, 0.02247617713937113), (73, 0.020245244144545974), (82, 0.025475940221899), (87, 0.0143379492777320616), (91, ... |
| Genetic Algorithm Approach for Loc Identification Readers | 36 | (15, 0.144017237457092492), (18, 0.195375235409477748), (27, 0.108873060087309488), (28, 0.045724574910157211), (35, 0.063876388662141468098), (41, 0.067533600008884361), (45, ... |
| Reliable Intersection Protocols Usin | 38, 44, 42 | (13, 0.18385128694123), (82, 0.0612966879881194773), (89, 0.033884413577770451111 ... |
| INTELLIGENT TRAFFIC WITH CO VEHICLES | 44, 49, 20, | (27, 0.014014615273312073), (28, 0.0712075457566209), (32, 0.0194858577917258), (37, 0.0143607469686245486), (42, 0.15184692653754589), (44, 0.06320017294480257), (59, ... |
| MCMC Particle Filter for Real-Time | 33, 88, 57 | (6, 0.034896871638292294), (8, 0.0464381286640256881), (11, 0.011112215673904270) ... |
| Globally Asymptotically Stable Filter and Depth Measurements | 91, 87, 61 | (14, 0.1416066704090924), (70, 0.06206132784724458), (74, 0.0665326233774602868) ... |
| A Real-Time Multi-Sensor Fusion Pl | 37, 71, 70 | (37, 0.81872937773255205), (55, 0.07878384429263103) ... |
| A full-3D Voxel-based Dynamic Obs Urban Scenario using Stereo Vision | 99, 96, 28 | (11, 0.1127627010000428584), (12, 0.212959122167004411), (18, 0.007386671787783663315), (28, 0.03296850148435222221), (33, 0.01839990845403183387), (35, 0.05014794602580639397) ... |
| A Real-Time Trajectory Control of T | vehicle control, 78, 99 | (5, 0.0417442035141044462), (8, 0.0273327110438853115), (11, 0.11313252786801186), (17, 0.022540918693733783), (29, 0.0360262067516344585), (35, 0.0144902192826010019), (37, 0.03200940638885363342), (41, ... |
| Vehicle Automation in Cooperation Nomadic Devices Communication | 55,99,76 | (4, 0.017501547572969963), (6, 0.01287370968319741), (8, 0.0724902882818116334), (12, 0.01060077094638752), (24, 0.01253983607240187), (33, 0.0608007092180063638), (34, 0.0110801534685350), (38, ... |
| Automatic vehicle classification and movements at signalized intersectio | 8, 42, 54 | (6, 0.0958979368012623286), (11, 0.0347851232941994965), (12, 0.0707692514343933932), (18, 0.007151806140576442a), (23, 0.046230642433257107), (35, 0.0104936479235473322), (47, ... |
| Multi-Target Tracking using a 3D-Li | 15, 77, 98 | (2, 0.109646451178223582), (11, 0.00182140346255068448), (25, 0.0412791802087192256), (35, 0.024229053366161248), (37, 0.01007222751730475 71), (38, 0.0100932797001388), (42, ... |
| Traffic Sign Representation using S | 2, number 20 | (5, 0.0185658412652368874), (11, 0.00182140346255068448), (46, 0.066644508450203584488), (52, 0.0102031100623586624), (57, 0.0319077513448185121), (59, 0.0989905145375710611), (62, ... |
| Speed Profile Optimization for Vehi Intersection Under a Safety Constra | 38, 42, vehicle control | (5, 0.05739846727255389598), (94, 0.024910022699195553), (95, 0.02097971202136645418 ... |

| | | | | |
|---|---|---|---|---|
| Fisheye optics f | 96, 70, 99 | (11, 0.1390023 | , (32, 0.084519 | (99, 0.10121410958494333) |
| Data age based control data exc | 86, 44 | (38, 0.0680274 | (44, 0.2029934 | (86, 0.0732155059085260407), |
| Obstacle Avoid Model Predictiv | 78, 82, 57 | (12, 0.0740088 | (52, 0.1608838 | (62, 0.1970401075172302) |
| Intelligent Cruis Stop and Go wit | 52, 55, 44 | (42, 0.1047802 | (44, 0.2432348 | (82, 0.1107893594279438) |
| Autonomous Na | 87 | (48, 0.0769147 | (82, 0.0935574 | (98, 0.1050395954312905 6) |
| Bayesian Netwo | 95, number 4, 5 | (38, 0.1212474 | (82, 0.1893622 | (82, 0.2274413299260052 7) |
| Experience, Re Lessons Learne Automated Drivi Germany's High | 37, 55, 62 | (55, 0.0847935 | (98, 0.1063484170061548 4) | |
| Multi-Objective | 62, 5, 57 | (5, 0.1052351534396422 1), (62, 0.1230114569017088), (84, 0.0951751230301 42555) | | |
| A Study on Auto University* | No relevant topic | (34, 0.0832615 | (74, 0.1271243 | (82, 0.2285297626658 70359), |
| Road Surface R Automatic Plato | 35, 18, 54 | (25, 0.1109076 | (82, 0.0807853 | , (94, 0.1440117082060 0007), |
| Building a Proto Parking System | 47, 62 | (47, 0.1143939 | (66, 0.1301513 | (82, 0.09654512807601 8782), |
| A Computer Visi Avoidance for A | 5, 2, 99 | (82, 0.2103425 | (94, 0.1302130 | (99, 0.1686118185386 129) |
| Path Tracking o Fractional Order | 62, vehicle control | (16, 0.1431923 | (23, 0.0979858 | (89, 0.0640922319101 71715), |
| Off-road Path F Constraints* | 48, 28, 63 | (18, 0.1338110 | (82, 0.1087812 | (94, 0.1106823302525 5438), |
| Self-Tuning PID Autonomous Ca | 37, 95, Vehicle control | (24, 0.1157130 | (52, 0.1390555 | (60, 0.1080876068906 5738), |
| Shared Control based on Veloci | 45, vehicle control, 78 | (17, 0.0792136 | (52, 0.0732116 | (62, 0.2356533950809 1404), |
| A 13,000 km Int | 95, vehicle control, 82 | (66, 0.0681289 | (74, 0.4081724 | (82, 0.1870438346103 6819) |
| Real-Time Coor | 35, 86 | (32, 0.0783232 | (44, 0.1282214 | (55, 0.0854262938093 04457) |
| Accurate and Ef | 96, 85, 2 | (18, 0.1082158 | (55, 0.0854262 | (93, 0.0829513506857 58013) |
| DeepDriving: Le | 77, 88, 2 | (18, 0.1216196 | (66, 0.1038493 | (82, 0.1896091585670 9309) |
| A Robust Algorit | 33, 88 | (12, 0.0261620 | (33, 0.8143677 | (95, 0.0208850438653 60702) |
| Constrained Glo | 62, 78, 10 | (10, 0.1476177 | (57, 0.0813935 | (62, 0.1643083600445 0927) |

| Title | Values | Col A | Col B | Col C |
|---|---|---|---|---|
| 360◦ detection a using fisheye im | 47, 96, 99 | (82, 0.1103107) | (96, 0.0795476) | (99, 0.16454368464617050091) |
| State your positi | 1, 48, | (1, 0.25972098) | (15, 0.1332502) | (52, 0.077168595128774123) |
| A robotic platfor | Vehicle control, 55, 60 | (11, 0.1056438) | (55, 0.06188856) | (66, 0.081648855621690927) |
| Coordinated co discrete-time pe | 77, 82, 86 | (16, 0.2589560) | (17, 0.0772982) | (62, 0.075266478492064609) |
| Real-time Imple | Vehicle control, 24, 25 | (11, 0.0762797) | (25, 0.1102223) | (52, 0.27489935471687632) |
| Coordinated Pat | 82, 86, vehicle control | (8, 0.52763602) | (16, 0.0910834) | (62, 0.065248220945653648) |
| A Combined Mo | 35, 92, vehicle control | (35, 0.0605185) | (55, 0.5331159) | (76, 0.168404332761006981) |
| Towards a Fra | 55, 45 | (49, 0.0351974) | (62, 0.1552549) | (82, 0.110881759298887201) |
| Adopting Wirele | 86, 61 | (66, 0.1851605) | (74, 0.0803295) | (82, 0.086750851051705574) |
| Terrain Mappin | 96, 99, 18, | (18, 0.1373032) | (96, 0.0928738) | (99, 0.097444633362412547) |
| Incremental Sa Minimum-violati | 60, Vehicle control, 78 | (15, 0.1016502) | (56, 0.11132531) | (62, 0.122461699493354441) |
| Vision-based Ni | 51, 96, 55 | (12, 0.0442678) | (33, 0.2048777) | (82, 0.187238110974212242) |
| Design and Co | 0, 33, 9 | (29, 0.0346098) | (37, 0.6986309) | (55, 0.19485987878501795) |
| Reliable Interse | 38, 44, 42 | (38, 0.3940938) | (42, 0.1518469) | (76, 0.076593369408331127) |
| Genetic Algorith Identification Re | 36 | (15, 0.1440173) | (18, 0.1953752) | (27, 0.108873060087309488) |
| HOG Based Mu | 98, 51, 87 | (18, 0.1317220) | (82, 0.1025475) | (94, 0.135155027552201253), |
| Local Path Plan Driving With Av | 62, 78, 99, | (10, 0.0552967) | (35, 0.0468949) | (78, 0.054223110303390045) |
| INTELLIGENT VEHICLES | 44, 49, 20, | (44, 0.1111465) | (66, 0.0837658) | (82, 0.208644933799386656) |
| MCMC Particle | 33, 88, 57 | (12, 0.0842157) | (18, 0.1213754) | (98, 0.083557085210943058) |
| Globally Asympt and Depth Mea | 91, 87, 61 | (15, 0.0640450) | (16, 0.1830467) | (70, 0.062061327847244458) |
| A Real-Time Mu | 37, 71, 70 | (37, 0.8187293) | (55, 0.0787838) | (78, 0.034186934006349756) |
| A full-3D Voxel- Urban Scenario | 99, 98, 28 | (11, 0.1127621) | (12, 0.2129591) | (96, 0.105005429946000785) |
| A Real-Time Tr | vehicle control, 78, 99 | (11, 0.11313250) | (52, 0.1852250) | (87, 0.1906831953200861) |
| Vehicle Automa Nomadic Devic | 55,99,76 | (11, 0.0868970) | (55, 0.0991782) | (82, 0.210237918696611239) |
| Automatic vehic movements at s | 8, 42, 54 | (8, 0.07249082) | (82, 0.17718510) | (94, 0.071540369671543758) |

| | | | | |
|---|---|---|---|---|
| Multi-Target Tra | 15, 77, 98 | (6, 0.09598793) | (52, 0.0773815) | (98, 0.27334997661984639) |
| Traffic Sign Rep | 2, number 20 | (2, 0.10946451) | (18, 0.1408045) | (63, 0.10563674501161169) |
| Speed Profile O Intersection Un | 38, 42, vehicle control | (38, 0.1000932) | (59, 0.0869065) | (62, 0.16169786391830407) |

| Title | Values | Result |
|---|---|---|
| Fisheye optics f | 96, 70, 99 | (11, 0.13900239476104639) |
| Data age based control data exc | 86, 44 | (44, 0.20299343908608297) |
| Obstacle Avoid Model Predictiv | 78, 82, 57 | (62, 0.1970401075172302) |
| Intelligent Cruis Stop and Go wit | 52, 55, 44 | (44, 0.24323481593723542), (98, 0.10503959543129056) |
| Autonomous Na | 87 | (82, 0.22744132992600527) |
| Bayesian Netwo | 95, number 4, 5 | (82, 0.18936225101487586) |
| Multi-Objective | 62 ,5, 57 | (62, 0.1230114569017088) |
| A Study on Auto University* | No relevant topi | (82, 0.22852976265870359), |
| Road Surface R Automatic Plato | 35, 18, 54 | , (94, 0.14401170820600007), |
| Building a Proto Parking System | 47, 62 | (66, 0.13015139933796271) |
| A Computer Visi Avoidance for A | 5, 2, 99 | (82, 0.21034250761433265), |
| Path Tracking o Fractional Order | 62, vehicle cont | (16, 0.14319230305794198), |
| Off-road Path F Constraints* | 48, 28, 63 | (18, 0.13381102708044132) |
| Self-Tuning PID Autonomous Ca | 37, 95, Vehicle | (52, 0.13905556918544504) |
| Shared Control based on Veloci | 45 , vehicle cont | (62, 0.23565339508091404), |
| A 13,000 km Int | 95, vehicle cont | (74, 0.40817247415538177) |
| Real-Time Coor | 35, 86 | (44, 0.12822145200088794) |
| Accurate and Ef | 96 ,85, 2 | (18, 0.10821584373327547) |
| DeepDriving: Le | 77, 88, 2 | (82, 0.18960915856709309) |
| A Robust Algorit | 33, 88 | (33, 0.81436774856870142) |

| Title | Values | Probability |
|---|---|---|
| Constrained Glo | 62, 78, 10 | (62, 0.164308360044550927) |
| 360° detection a using fisheye im | 47, 96, 99 | (99, 0.164543684617050091) |
| State your positi | 1, 48, | (1, 0.259720980266667584) |
| A robotic platfor | Vehicle control, | (11, 0.105643872598399503) |
| Coordinated co discrete-time pe | 77, 82, 86 | (16, 0.258956074653347326) |
| Real-time Imple | Vehicle control, | (52, 0.274899354716876632) |
| Coordinated Pat | 82, 86, vehicle | (8, 0.527636026354877761) |
| A Combined Mo | 35, 92, vehicle | (76, 0.168404332761069981) |
| Towards a Fra | 55, 45 | (55, 0.533115973677673769) |
| Adopting Wirele | 86, 61 | (66, 0.185160542210530744) |
| Terrain Mappin | 96, 99, 18, | (18, 0.137303287831085759) |
| Incremental Sa Minimum-violati | 60, Vehicle cont | (62, 0.122461699493355441) |
| Vision-based Ni | 51, 96, 55 | (33, 0.204877781945499938) |
| Design and Co | 0, 33, 9 | (37, 0.698630938142941539) |
| Local Path Plan Driving With Av | 62, 78, 99, | (10, 0.055296772486137021) |
| HOG Based Mu | 98, 51, 87 | (94, 0.135155027552012539), |
| Genetic Algorith Identification Re | 36 | (18, 0.195375235405977489), |
| Reliable Interse | 38, 44, 42 | (38, 0.394093817534969259) |
| INTELLIGENT VEHICLES | 44, 49, 20, | (82, 0.208644933799386569) |
| MCMC Particle | 33, 88, 57 | (18, 0.121375476489438839) |
| Globally Asympt and Depth Mea | 91, 87, 61 | (16, 0.183046707575858459) |
| A Real-Time Mu | 37, 71, 70 | (37, 0.818729377732552059) |
| A full-3D Voxel- Urban Scenario | 99, 98, 28 | (12, 0.212959122167044139) |
| A Real-Time Tr | vehicle control, | (87, 0.190683195320086139) |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Vehicle Automa Nomadic Devic | 55,99,76 | | | | | | | (82, 0.210237918696611239) | |
| Automatic vehic movements at s | 8, 42, 54 | | | | | | | (82, 0.17185107072267314) | |
| Multi-Target Tra | 15, 77, 98 | | | | | | | **(98, 0.27334997661984639)** | |
| Traffic Sign Rep | 2, number 20 | | | | | | | (18, 0.140804500332334806), | |
| Speed Profile O Intersection Un | 38, 42, vehicle c | | | | | 17 st | | (62, 0.161697863918304407) | |
| | | | | | | 33 st | | | |

| Label | | | | | | | | | | Same topics (= green if present) |
|---|---|---|---|---|---|---|---|---|---|---|
| clean_A__B__P. | (36, 0.1848332) | (29, 0.1298672) | (76, 0.1187326) | (37, 0.0812135) | (42, 0.0350045) | (34, 0.0335877) | (63, 0.03264458) | (98, 0.0269327) | (94, 0.0170646) | (16, 0.01190954099902100056) | 27 | 24 | 46 | autonomous vehicle control: | 0,40,76 |
| clean_A__B__S. | (72, 0.2558892) | (76, 0.1131356) | (76, 0.1072234) | (15, 0.0879327) | (44, 0.0666556) | (67, 0.0490105) | (42, 0.0439286) | | | | 3 | 37 | | vehicle path planning: | 16,78 |
| clean_A__B__N. | (98, 0.1333643) | (98, 0.1151860) | (6, 0.099980107) | (15, 0.10315660) | (87, 0.0684457) | (39, 0.015528555433854917) | | | | | 8 | 18 | 8 | vehicle detection/tracking: | 54,58,66 |
| clean_A__B__N. | (72, 0.2273818) | (6, 0.08954729) | (59, 0.0605728) | (90, 0.0491177) | (18, 0.0484480) | (15, 0.0469500) | | | | | 81 | 10 | 18 | automated parking: | 35,38 |
| clean_A__C__L. | (8, 0.18097194) | (36, 0.1108507) | (3, 0.0833808) | (6, 0.06049974) | (10, 0.0482495) | (34, 0.0379784) | (51, 0.0325365) | | | | 8 | 24 | 10 | networked vehicles: | 31,34 |
| clean_B__B__H. | (8, 0.52548262) | (83, 0.1011615) | (98, 0.08180036) | (6, 0.07244673) | (78, 0.0384183) | (33, 0.0312181644160492261) | | | | | 14 | | 24 | vehicle communication: | 62,72,99 |
| clean_B__W__A. | (35, 0.1040605) | (36, 0.1027102) | (63, 0.0932942) | (83, 0.0906500) | (0, 0.05276510) | (41, 0.0485640) | (34, 0.0362071) | (27, 0.0349161) | (78, 0.03161168) | (31, 0.0209966) | 95 | 62 | | road detection: | 25,64 |
| clean_C__L. | (2, 0.11201532) | (44, 0.1065425) | (78, 0.1098718) | (15, 0.0538467) | (14, 0.0590098) | (72, 0.0550037) | (95, 0.0532252) | (41, 0.0338320) | (86, 0.0277617) | (16, 0.0183011) | 2 | 31 | 37 | | |
| clean_C__Y. | (2, 0.13368008) | (95, 0.1256067) | (74, 0.0368196) | (74, 0.0398162) | (68, 0.0743063) | (6, 0.0739692) | (22, 0.0907557) | (90, 0.0199901855969662173) | (39, 0.0166852) | (23, 0.0165662) | 49 | 56 | | | |
| clean_C__L__B. | (97, 0.1988303) | (62, 0.1416671) | (63, 0.111173812) | (63, 0.0994831) | (94, 0.0362724) | (55, 0.0507047) | (59, 0.0150594) | (34, 0.0313654542990161580) | (63, 0.024296) | (46, 0.0208169) | 3 | 2 | 97 | | |
| clean_C__W__M. | (31, 0.3009491) | (62, 0.0918091) | (46, 0.0996277) | (37, 0.0381002) | (35, 0.0338037) | (54, 0.0338037) | (78, 0.0262566) | (94, 0.0199999) | (55, 0.0302549) | (94, 0.0241462) | 46 | 90 | | | |
| clean_C__W. | (64, 0.1244015) | (32, 0.0757782) | (46, 0.0670448) | (16, 0.0689471) | (24, 0.0649523) | (23, 0.037974) | (46, 0.0172676) | (76, 0.0146945) | (94, 0.0278960) | (23, 0.0191992) | 32 | 91 | | | |
| clean_D__C__S. | (0, 0.28794799) | (78, 0.0974409) | (31, 0.0491384) | (54, 0.047931) | (72, 0.0286192) | (35, 0.02535958) | (27, 0.013165997448584592) | (36, 0.03205249) | (3, 0.018019870) | | 94 | 27 | | | |
| clean_D__B__W. | (64, 0.23491233) | (35, 0.0941132) | (32, 0.0705237) | (29, 0.0697600) | (56, 0.0627237) | (56, 0.0544620) | (36, 0.0305349) | (36, 0.012665799) | | | 32 | 46 | 91 | | |
| clean_G__A__J. | (0, 0.23981233) | (94, 0.0919301) | (37, 0.1187409) | (31, 0.0491384) | (63, 0.0330859) | (62, 0.0325047) | (25, 0.0286637) | (42, 0.024293219989350316) | | | 85 | 3 | 27 | ABOUT INTERNET OF THINGS | |
| clean_H__x__E. | (64, 0.16220410) | (62, 0.1032673) | (55, 0.011347) | (16, 0.0703622) | (34, 0.0398617) | (63, 0.0336817) | (6, 0.02396606) | (55, 0.0198637) | (38, 0.0174714) | | 16 | 2 | 23 | ABOUT ROBOT | |
| clean_M__A__M. | (67, 0.3247214) | (14, 0.0884570) | (85, 0.0463860) | (99, 0.0405668) | (20, 0.0379839) | (34, 0.0310356) | (67, 0.0287670) | (55, 0.0198050) | | | 29 | 51 | 94 | | |
| clean_M__B__C. | (78, 0.2229050) | (2, 0.05743418) | (22, 0.107247126) | (95, 0.0781148) | (78, 0.0558509) | (71, 0.03053839) | (94, 0.019156179644451415) | | | | 31 | 62 | | | |
| clean_M__C__D. | (37, 0.1695113) | (37, 0.1110608) | (94, 0.0892026) | (95, 0.0885008) | (66, 0.0373135) | (52, 0.0256741) | (94, 0.0215120) | (85, 0.015296817977437) | | | 32 | 58 | | | |
| clean_M__B__Z. | (86, 0.1321063) | (37, 0.1316656) | (14, 0.0892026) | (0, 0.06544911) | (2, 0.0412817) | (71, 0.0448988) | (41, 0.0428818) | (94, 0.03000069) | (35, 0.0280322) | | 68 | 31 | | | |
| clean_M__H__A. | (81, 0.1612689) | (67, 0.1042430) | (78, 0.0767087) | (72, 0.0687707) | (62, 0.03879315) | (63, 0.0593547) | (86, 0.0423050) | (16, 0.0231577) | (50, 0.01919941850410000111) | | 51 | 24 | | | |
| clean_M__J__A. | (54, 0.14282397) | (91, 0.1122323) | (15, 0.1010960) | (27, 0.0908212) | (97, 0.03983686) | (78, 0.03793315) | (3, 0.02005423) | (38, 0.0178973232968601131) | | | 81 | 3 | | | |
| clean_M__J__H. | (36, 0.17363999) | (76, 0.1122080) | (15, 0.0773080) | (16, 0.0929904) | (86, 0.0334008) | (49, 0.02790023) | (16, 0.0236904) | | | | 14 | 14 | | | |
| clean_M__J__H. | (8, 0.3558818) | (3, 0.10158367) | (95, 0.0892093) | (83, 0.03838837) | (44, 0.03602285) | (59, 0.0330463) | (70, 0.03096609) | (78, 0.02102298) | (68, 0.011486738998817859) | | 10 | 3 | | | |
| clean_M__x__E. | (76, 0.65144179) | (64, 0.1343393) | (23, 0.0262232) | (15, 0.032693075487806525) | | | | | | | 16 | 14 | | | |
| clean_N__C__B. | (32, 0.1649800) | (37, 0.1284980) | (72, 0.1182114) | (3, 0.06393005) | (24, 0.0484068) | (64, 0.0462743) | (68, 0.0398392) | (35, 0.0283078) | (67, 0.02806645) | (86, 0.0196094) | 81 | 10 | | | |
| clean_N__T__A. | (20, 0.5435616) | (38, 0.0723725) | (59, 0.0462966) | (3, 0.00400195) | (42, 0.0621802) | (51, 0.0172130) | (98, 0.01144545) | (23, 0.01939) | (20, 0.0193233) | (72, 0.0172744) | 94 | 81 | | | |
| clean_P__B__D. | (64, 0.1508629) | (3, 0.12069143) | (16, 0.1189783) | (78, 0.0633776) | (99, 0.0333977) | (52, 0.0303196) | (42, 0.02013362) | (58, 0.01019998229621947) | (25, 0.0127775126254119499) | (52, 0.013814) | 98 | 94 | | | |
| clean_P__V__K. | (46, 0.1440402) | (31, 0.0866366) | (6, 0.00900087) | (17, 0.06168736) | (32, 0.0596022) | (22, 0.0418362) | (40, 0.0264197) | (2, 0.023680329) | (99, 0.0140082) | (52, 0.0138141) | 42 | 70 | | | |
| clean_P__V__M. | (8, 0.32049936) | (16, 0.08665306) | (78, 0.0584076) | (64, 0.0369604) | (78, 0.0544074) | (22, 0.049245) | (32, 0.0481847) | (42, 0.0379891) | (52, 0.0362766) | (23, 0.02251361) | 8 | 14 | | | |
| clean_Q__B__M. | (64, 0.30902760) | (14, 0.2418771) | (63, 0.0721793) | (37, 0.06950597) | (55, 0.0392101) | (31, 0.0221150) | (61, 0.01207759279241510) | | | | 98 | 2 | | | |
| clean_S__A__G. | (67, 0.4891512) | (37, 0.0882767) | (94, 0.05875663) | (72, 0.0530567) | (94, 0.0207420) | (35, 0.01727252) | (95, 0.0145961789962371168) | | | | 31 | 94 | 72 | | |
| clean_S__A__S. | (38, 0.2684228) | (46, 0.17110749) | (8, 0.00622350) | (29, 0.0559041) | (81, 0.05122411) | (23, 0.01614520192283057) | | | | | 38 | 16 | 16 | | |
| clean_S__B__M. | (35, 0.1151783) | (14, 0.0855717) | (31, 0.0816080) | (63, 0.0628652) | (37, 0.0597022) | (2, 0.0385763397) | (94, 0.0226686) | (39, 0.03318285) | | | 31 | 34 | 34 | | |
| clean_S__D__B. | (10, 0.1621257) | (86, 0.1059432) | (6, 0.07362194) | (23, 0.0728841) | (42, 0.0069907) | (14, 0.0401288) | (22, 0.0386812) | (94, 0.0222010) | (55, 0.0222010) | (37, 0.0019815) | 97 | 42 | | | |
| clean_S__D__A. | (78, 0.8132457) | (2, 0.04427687) | (33, 0.032975572) | (66, 0.0101091799814445656) | | | (29, 0.0439753) | (35, 0.02698612) | | | 81 | 78 | autonomous vehicle control | |
| clean_S__P__B. | (61, 0.1657975) | (42, 0.1143952) | (22, 0.0861378) | (6, 0.05092820) | (41, 0.0507432) | (35, 0.0330778) | (71, 0.0178680) | (2, 0.01660323) | (15, 0.0148087) | (67, 0.01144630) | 46 | 70 | autonomous vehicle control | |
| clean_T__A__M. | (86, 0.1014169) | (67, 0.0907007) | (15, 0.0841750) | (70, 0.0576244) | (37, 0.0545590) | (3, 0.0329805) | (94, 0.0272723) | (72, 0.0253321) | (62, 0.0223435) | (71, 0.0225984) | 86 | 62 | autonomous vehicle control | |
| clean_Y__A__P. | (64, 0.1485641) | (61, 0.1357398) | (61, 0.0505877) | (23, 0.0577777) | (23, 0.0386581) | (22, 0.0298268) | (70, 0.0173934) | (44, 0.0164889) | | | 91 | 32 | autonomous ve | |
| clean_Z__B__J. | (54, 0.1497370) | (37, 0.1132989) | (64, 0.0951638) | (95, 0.0463537) | (76, 0.04001288) | (37, 0.0381580) | (22, 0.0388638) | | | | 46 | 54 | autonomous ve | |
| clean_Z__K__x. | (14, 0.1791914) | (76, 0.1175761) | (37, 0.09686642) | (3, 0.07762989) | (32, 0.03380805) | (14, 0.0323105) | (46, 0.0294377) | (58, 0.0194477) | (22, 0.01220313280269705) | | 3 | 37 | autonomous vehicle control | |

**Header / topic labels (right side of sheet):**

| Label | Values |
|---|---|
| *Same topics (= green if present)* | 0,40,76 |
| autonomous vehicle control: | 16,78 |
| vehicle path planning: | |
| vehicle detection/tracking | 54,58,66 |
| automated parking: | 35,38 |
| networked vehicles: | 31,34 |
| vehicle communication: | 62,72,99 |
| road detection: | 25,64 |

**Main data table** (green-highlighted cells shown in **bold**):

| C1 | C2 | C3 | C4 | C5 | C6 |
|---|---|---|---|---|---|
| 27 | 24 | 46 | (36, 0.1848564) | (29, 0.1298889) | (76, 0.1187243754951616186) |
| autonomous ve | 3 | 37 | (72, 0.2559601) | (61, 0.1129838) | **(76, 0.1072838672442939928)** |
| autonomous ve | 8 | 18 | (8, 0.18093794) | **(8, 0.13359666)** | (98, 0.1148220405427511112) |
| autonomous ve | 8 | 10 | (8, 0.22751369) | (36, 0.1405603) | (16, 0.110954660679559121) |
| autonomous ve | 81 | 24 | (64, 0.52582454) | (35, 0.1053866) | (83, 0.099992708227210195) |
| 95 | | 14 | **(95, 0.2345570)** | (85, 0.1075159) | (36, 0.1019715825805026) |
| 62 | | 31 | (2, 0.11209784) | (0, 0.1098274) | (78, 0.109746674438585584) |
| 3 | | 2 | **(37, 0.1776671)** | (98, 0.1415837) | (95, 0.114392564117774463) |
| 25 | | 90 | (46, 0.3030772) | (61, 0.2103519) | (23, 0.091408064252341062) |
| 49 | | 2 | (36, 0.1507248) | (64, 0.1845176) | (78, 0.097066065517260055) |
| 2 | | 97 | **(56, 0.1892263)** | (2, 0.13043038) | **(56, 0.1271625380114527)** |
| ABOUT INTERNET OF THINGS | 46 | 39 | (31, 0.3009422) | (62, 0.0918108) | (0, 0.080914136703779838) |
| 32 | | 91 | (64, 0.1208752) | (10, 0.0981494) | (35, 0.085743075471516106) |
| 94 | autonomous ve | 27 | (0, 0.28780847) | (14, 0.1804275) | (67, 0.160989080890118968) |
| 3 | 85 | autonomous ve | (36, 0.2382778) | (37, 0.1099536) | (94, 0.091980057358372974) |
| 16 | | 2 | (64, 0.1820686) | (62, 0.1043574) | (62, 0.1032817582908135) |
| 31 | autonomous ve | 62 | (8, 0.16329881) | (23, 0.1043574) | **(34, 0.1192125247919195)** |
| 32 | | 58 | (14, 0.1694132) | (16, 0.1556911) | (61, 0.13212397282221669) |
| 86 | | 31 | **(86, 0.1676208)** | **(34, 0.1203641)** | **(62, 0.098210500234343259)** |
| autonomous ve | ABOUT ROBOT | 39 | (14, 0.1678848) | (16, 0.1519493) | (2, 0.12432099890414264) |
| 54 | 70 | | (14, 0.1592115) | **(54, 0.1507850)** | (32, 0.074856650682069967) |
| 78 | autonomous ve | 34 | (3, 0.20696931) | **(78, 0.1675713)** | (52, 0.10858727116328266) |
| 64 | 46 | | (91, 0.1982367) | **(64, 0.1259936)** | (61, 0.125673402342225596) |
| 10 | autonomous ve | 14 | (8, 0.18235012) | (78, 0.1422352) | **(14, 0.104092801963043 62)** |
| 16 | 81 | | (36, 0.2004328) | (85, 0.1446369) | (3, 0.13290219339100076) |
| 81 | autonomous ve | 10 | (22, 0.2103915) | (37, 0.2065201) | (14, 0.088431382122653 89) |
| 29 | 51 | | (67, 0.3241753) | (86, 0.2248693) | (31, 0.12382616448151035) |
| 36 | autonomous ve | 3 | (78, 0.2228182) | **(36, 0.1871256)** | (37, 0.10862600579458387) |
| 3 | 24 | | (37, 0.1686570) | (67, 0.1659653) | (86, 0.13212604834057004) |
| 51 | 81 | 3 | (37, 0.1987332) | (67, 0.1021661) | (61, 0.092271868290538586) |
| Very lightweight paper, should produce low proba | | 46 | (81, 0.1769540) | (2, 0.16601571) | (91, 0.11885317918222637) |
| 2 | 46 | | (56, 0.1906900) | (6, 0.17118230) | **(2, 0.14453992095615442)** |
| 46 | autonomous ve | 14 | (61, 0.3349081) | (8, 0.24586743) | **(3, 0.10115871457693873)** |
| 46 | | 3 | (76, 0.6515046) | (64, 0.1347841) | (80, 0.11052799511617664) |
| 97 | 68 | | (32, 0.1648620) | **(97, 0.1284760)** | **(46, 0.1179388406423926)** |
| 20 | 38 | 24 | (81, 0.5448869) | **(20, 0.1467503)** | **(38, 0.066431417501460852)** |
| 17 | 98 | 42 | (64, 0.1508724) | (3, 0.12064299) | (16, 0.11892074962130053) |
| 46 | 42 | 70 | **(46, 0.1449212)** | (31, 0.0869137) | (6, 0.081995218719289387) |
| autonomous ve | 8 | 14 | **(8, 0.32051546)** | (6, 0.14205713) | (3, 0.13477592829968937) |
| 68 | 98 | 2 | (64, 0.3092776) | (14, 0.2421471) | (94, 0.16943680061267874) |
| 31 | 94 | 72 | (67, 0.4859769) | (37, 0.0879113) | (72, 0.12607761826882365) |
| 38 | autonomous ve | 16 | **(38, 0.2684245)** | (46, 0.1710797) | **(94, 0.052500121435727472)** |
| 31 | autonomous ve | 34 | (35, 0.1151748) | (14, 0.0856002) | (6, 0.081638863744275481) |

| | | | | | |
|---|---|---|---|---|---|
| 16 | 97 | 42 | (10, 0.1479399) | (86, 0.1079217) | (14, 0.0859171868892826... |
| autonomous ve | 81 | 78 | (78, 0.8130703) | (2, 0.04279595) | (8, 0.0415220561639050016) |
| 46 | 70 | autonomous ve | (61, 0.1677124) | (16, 0.1457492) | (42, 0.0891276513216518...) |
| 86 | autonomous ve | 62 | (86, 0.1014034) | (67, 0.0901782) | (15, 0.08420761338234857) |
| 91 | 32 | | (64, 0.1485523) | (78, 0.1357089) | (98, 0.075125971005098865) |
| 46 | 54 | 91 | (54, 0.1524063) | (91, 0.1294819) | (37, 0.1093812209030285... |
| 3 | 37 | autonomous ve | (14, 0.1791856) | (76, 0.1176699) | (2, 0.110825369100374...) |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| clean_C__C__J. | (64, 0.3030772) | (61, 0.2103519) | (23, 0.09140806425234101662) | (36, 0.1848566480625556622) | 27 | 24 | 46 |
| clean_G__A__J. | (36, 0.2382778) | (37, 0.1099536) | (94, 0.09198007356837297974) | (72, 0.25596012169338933) | autonomous ve | 3 | 37 |
| clean_Z__B__J. | (54, 0.1524063) | (91, 0.1294819) | (37, 0.10938122093028534) | (2, 0.2275136900047077) | autonomous ve | 8 | 18 |
| clean_M__A__P. | (22, 0.2103915) | (37, 0.2065201) | (14, 0.08843138212265389) | (8, 0.18093794262797613) | 8 | 81 | 10 |
| clean_S__B__M. | (35, 0.1151748) | (14, 0.0856002) | (6, 0.08163886374427548481) | (8, 0.52582454211157947) | autonomous ve | 14 | 24 |
| clean_D__C__S. | (0, 0.28780847) | (14, 0.1804275) | (67, 0.16098908089011896968) | (95, 0.23475704115663953) | 95 | | ABOUT INTERNET OF THINGS |
| clean_A__B-N. | (8, 0.22751369) | (98, 0.13359666) | (98, 0.11482204054275112) | (2, 0.11209784927366824) | 62 | 31 | |
| clean_D__B__W. | (64, 0.1208752) | (10, 0.0981494) | (35, 0.08574307547151610) | (37, 0.17788710077087728) | 3 | 2 | 37 |
| clean_T__A__M. | (86, 0.10140034) | (67, 0.0901782) | (15, 0.08420776133823485857) | (64, 0.30307722290019357) | 2 | 90 | 46 |
| clean_K__C__F. | (14, 0.1592115) | (54, 0.1507850) | (32, 0.07485665068206996967) | (36, 0.15072481588128858) | 25 | 2 | 56 |
| clean_J__C__S. | (14, 0.1694132) | (16, 0.1556911) | (61, 0.13212397282221669) | (97, 0.18982639068292873) | 49 | 56 | |
| clean_M__x__E. | (76, 0.6515046) | (64, 0.1347841) | (80, 0.11052799511617664) | (86, 0.16762208122568538) | 2 | 97 | 62 |
| clean_S__A__S. | (38, 0.2684245) | (46, 0.1710797) | (72, 0.12607618268822365) | (14, 0.16798481714197125) | 86 | 46 | 31 |
| clean_M__J__B. | (56, 0.1906900) | (6, 0.17118230) | (2, 0.14453992095615442) | (31, 0.30094220805326038) | autonomous ve | ABOUT ROBOT | 39 |
| clean_L__x__F_ | (91, 0.1982367) | (64, 0.1295936) | (61, 0.12567340234222596) | (64, 0.12087527306715326) | 32 | 91 | |
| clean_Z__K__x. | (14, 0.1791856) | (76, 0.1176699) | (2, 0.11082569610037747) | (0, 0.28780847536848242) | 94 | autonomous ve | 27 |
| clean_J__A__C. | (8, 0.16329881) | (0, 0.15599647) | (34, 0.11921252479198195) | (36, 0.23827786532158918) | 3 | 85 | autonomous vehicle control |
| clean_C__C__Y. | (36, 0.1507248) | (2, 0.13043038) | (56, 0.12716253801145527) | (76, 0.65150463233416078) | 46 | 46 | |
| clean_S__A__G. | (67, 0.4859769) | (37, 0.0879113) | (94, 0.05250012143572472) | (32, 0.16482082528314469) | 97 | 68 | |
| clean_N__T__At | (81, 0.5448869) | (20, 0.1467503) | (38, 0.06643141750146085) | (61, 0.33490816986681493) | autonomous ve | 14 | 3 |
| clean_S__D__B. | (10, 0.1479399) | (86, 0.1079207) | (14, 0.08591718688928264) | (56, 0.19069006158388141) | 2 | 46 | |
| clean_P__V__M. | (8, 0.32051546) | (6, 0.14205713) | (3, 0.13477592829068937) | (81, 0.17695400107396794) | Very lightweight paper, should produce low probabilities | 46 | |
| clean_M__H__A | (81, 0.1769540) | (2, 0.16601571) | (91, 0.11885317918232637) | (37, 0.19873329096901887) | 51 | 81 | 81 |
| clean_S__P__B. | (61, 0.1677124) | (16, 0.1457492) | (42, 0.08912765132165181818) | (37, 0.16865701426790394) | 3 | 81 | 3 |
| clean_A__C__L. | (8, 0.52582454) | (35, 0.1053866) | (83, 0.09999270822721019...) | (78, 0.22281829261209953) | 36 | autonomous ve | 3 |
| clean_M__A__J. | (36, 0.2004328) | (85, 0.1446369) | (3, 0.132902193391000076) | (67, 0.32417530572754122) | 29 | 51 | 94 |
| clean_P__B__D. | (64, 0.1508724) | (3, 0.12064299) | (16, 0.11892074962130053) | (22, 0.21039159482527384) | 81 | autonomous ve | 10 |
| clean_B__B__H. | (95, 0.2347570) | (85, 0.1075159) | (36, 0.1019715825850502626) | (36, 0.20043282690429736) | 16 | 81 | autonomous vehicle control |
| clean_P__V__K. | (46, 0.1449212) | (31, 0.0869137) | (6, 0.08199652187192893877) | (8, 0.18235012515513513) | 10 | autonomous ve | 14 |
| clean_K__B__H. | (14, 0.1679848) | (16, 0.1519493) | (2, 0.124320998904142664) | (91, 0.19823679216460691) | 64 | 46 | 54 |
| clean_J__S__B. | (72, 0.2559601) | (34, 0.1203641) | (76, 0.1072838672442929928) | (64, 0.18206868421320669) | 78 | 34 | autonomous vehicle control |
| clean_A__B__S. | (61, 0.1129838) | (62, 0.0982105002344325...) | (8, 0.16329881214835552) | (31, autonomous ve) | 2 | 23 | |
| clean_M__B__Z. | (37, 0.1686570) | (67, 0.1659653) | (86, 0.1321260483405700...) | (14, 0.16941322421587424) | 58 | 68 | |
| clean_C__L__B. | (97, 0.1898263) | (64, 0.1845176) | (78, 0.0970660605651726055) | (64, 0.15087242592216636) | 17 | 98 | 42 |
| clean_A__C__C. | (8, 0.18093794) | (36, 0.1405603) | (16, 0.11109546667959121) | (81, 0.54488699580868405) | 20 | 38 | 24 |
| clean_C__C__Y. | (36, 0.1507248) | (2, 0.13043038) | (56, 0.12716253801145527) | (32, 0.16482082528314469) | 97 | 68 | 46 |
| clean_M__B__C. | (78, 0.2228182) | (36, 0.1871256) | (37, 0.1086260057945838...) | (46, 0.14492128773661103) | 46 | 42 | 70 |
| clean_B-M__S. | (37, 0.1778671) | (98, 0.1415837) | (95, 0.114392564117777463) | (8, 0.32051546792194785) | autonomous ve | 8 | 14 |

| | (col2) | (col3) | (col4) | (col5) | | | |
|---|---|---|---|---|---|---|---|
| clean_H_x_E. | (64, 0.1820686) | (23, 0.1043574) | (62, 0.1032817589290815) | (64, 0.3092776500292902) | 68 | 98 | 2 |
| clean_B_W_A. | (2, 0.11209784) | (0, 0.10982743) | (78, 0.109746674438585684) | (67, 0.4859769401259374) | 31 | 94 | 72 |
| clean_M_C_D. | (37, 0.1987332) | (67, 0.1021661) | (61, 0.092271868290538586) | (38, 0.268424539469421) | 38 | 16 | |
| clean_L_C_A. | (3, 0.20696931) | (78, 0.1675713) | (52, 0.108587271663328266) | (35, 0.1151748814137828281) | 31 | autonomous ve | 34 |
| clean_Q_B_M. | (64, 0.3092776) | (14, 0.2421471) | (94, 0.16943680061267874) | (10, 0.147939983163076888) | 16 | 97 | 42 |
| clean_A_B_P. | (36, 0.1848564) | (29, 0.1298889) | (76, 0.11872437549516186) | (78, 0.813070370477201554) | 81 | 78 | |
| clean_M_A.-M. | (67, 0.3241753) | (86, 0.2248693) | (31, 0.123826164481151035) | (61, 0.167771243527230448) | 46 | autonomous vehicle control | 70 |
| clean_M_A_A. | (8, 0.18235012) | (78, 0.1422352) | (14, 0.104092801963404362) | (86, 0.101403477238754423) | 86 | autonomous ve | 62 |
| clean_C_W_A | (31, 0.3009422) | (62, 0.0918108) | (0, 0.080914136703779838) | (64, 0.148552321922200494) | 91 | 32 | |
| clean_N_C.-B. | (32, 0.1648620) | (97, 0.1284760) | (46, 0.117938840642239226) | (54, 0.15240063995367580808) | 46 | 54 | 91 |
| clean_Y_A_P. | (64, 0.1485523) | (78, 0.1357089) | (98, 0.075125971005098865) | (14, 0.179185860023085695) | 3 | autonomous vehicle control | 37 |

13 st

37 st

| image | Word #1 | Word #2 | Word #3 | Word #4 | Word #5 | Word #6 | Word #7 | Label | Acronyms/ word explanations |
|---|---|---|---|---|---|---|---|---|---|
| 0 | agent | coil | power | circuit | policy | aorta | turn | Power management | |
| 1 | uncertainty | parameter | estimation | covariance | method | matrix | system | Probability estimation | |
| 2 | image | recognition | sample | vehicle | model | mask | logo | Image recognition | |
| 3 | oscillation | natural | system | vo | velocity | locomotion | value | | |
| 4 | vehicle | speed | control | lateral | reference | profile | strategy | Vehicle Control | |
| 5 | obstacle | vehicle | task | avoidance | control | path | velocity | Obstacle avoidance managemen | |
| 6 | vehicle | utility | target | feature | eye | assignment | based | Vehicle Utility | |
| 7 | ve | olarak | filtre | için | bir | ekbho | bu | | |
| 8 | vehicle | task | method | control | background | detection | video | Vehicle detection (video feed?) | |
| 9 | sensor | market | imu | system | cost | fusion | data | IMU market | imu = inertial measurement unit, measures crafts velocity orientation and gravitational forces |
| 10 | path | curve | vehicle | bézier | point | | | Path prediction | |
| 11 | system | vehicle | image | obstacle | used | decision | detection | Obstacle detection | |
| 12 | vehicle | set | time | approach | algorithm | system | | Vehicle software | |
| 13 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 14 | rate | ber | vocoder | channel | amr | ecall | mode | | |
| 15 | problem | bound | path | constraint | solution | state | approach | Pathfinding | |
| 16 | control | vehicle | tracking | trajectory | controller | following | time | Vehicle Control | |
| 17 | vehicle | set | trajectory | reachable | constraint | tk | system | Vehicle limitations | |
| 18 | road | model | color | algorithm | point | method | image | Road modelling | |
| 19 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 20 | image | character | plate | network | neural | value | recognition | Image input to machine learning | |
| 21 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 22 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 23 | particle | measurement | fastslam | filter | problem | position | based | position measur | fastslam is an algoritm for localising a robot and mapping it's surroundings |
| 24 | vehicle | threat | figure | position | landmark | avg | time | Vehicle threat assesment | |
| 25 | vehicle | control | wheel | tire | friction | road | longitudinal | Wheel control | |
| 26 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 27 | attack | vehicle | risk | severity | problem | car | number | Security aspects | |
| 28 | cost | disparity | cloud | time | census | pixel | weight | Image processi | disparity could be related to image recognition |
| 29 | agent | ri | source | signal | ci | formation | algorithm | | |
| 30 | process | group | request | pi | node | quorum | message | Obstacle management | |
| 31 | image | obstacle | top | used | system | view | path | Obstacle management | |
| 32 | system | data | obstacle | tentacle | mobility | information | vehicle | Obstacle management | |
| 33 | light | vehicle | detection | tracking | algorithm | frame | signal | Vehicle detection | |
| 34 | control | brake | system | vehicle | design | stopping | model | Brake control/management | |
| 35 | vehicle | state | time | model | fleet | based | road | Model for vehicle fleets | |
| 36 | research | technology | curve | patent | stage | vehicle | | Vehicle research | |
| 37 | vehicle | component | driving | architecture | platform | system | control | System architecture | |
| 38 | vehicle | intersection | figure | type | lane | system | output | Intersection handling | |
| 39 | av | ri | leader | time | follower | one | | | |
| 40 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 41 | tag | reader | antenna | function | rra | position | positioning | | tag might be related to gps technology |
| 42 | intersection | vehicle | traffic | time | group | car | light | Intersection handling | |
| 43 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 44 | vehicle | communication | cooperative | attack | cacc | stream | | Vehicle coopera | cacc could be correlated active clause coverage, A Logic Coverage Criterion from Software Testing |
| 45 | trajectory | passenger | based | test | boarding | alighting | tracking | System test with | alighting might refer to getting off a vehicle |
| 46 | demand | policy | vehicle | condition | tmhp | time | stability | | |
| 47 | parking | vehicle | space | pedestrian | task | state | system | Parking | |
| 48 | image | based | figure | time | navigation | algorithm | environment | Image-based navigation | |
| 49 | car | vehicle | algorithm | caravan | time | aid | three | Vehicle cooperation | |
| 50 | information | platform | service | self | driving | tourist | content | Information service for tourists? | |
| 51 | vehicle | detection | cluster | algorithm | ve | time | FALSE | Vehicle detection? | |
| 52 | vehicle | control | angle | model | lateral | dynamic | system | Vehicle Control | |
| 53 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 54 | task | vehicle | tracking | target | model | platoon | measurement | Vehicle tracking (platoon implies several vehicles) | |
| 55 | driving | driver | automated | vehicle | wa | speed | distance | Driving | |
| 56 | graph | pattern | rule | system | time | set | model | Data representation | |
| 57 | trajectory | vehicle | trim | point | using | time | maneuver | Maneuvering | |
| 58 | task | cbba | bundle | assignment | agent | bid | ij | | |
| 59 | agent | task | node | algorithm | cbba | assignment | time | | |
| 60 | system | driving | change | vehicle | development | software | simulation | Software development for vehicles in simulation | |
| 61 | vehicle | prt | system | wheel | consumption | fuel | campus | Fuel Consumption | |
| 62 | vehicle | path | planning | constraint | trajectory | problem | time | Path planning (Planera körning?) | |
| 63 | road | image | network | feature | learning | pixel | segmentation | Image recognition for roads | |
| 64 | oscillation | vo | system | natural | velocity | locomotion | control | ?? | |
| 65 | ve | için | filtre | olarak | ekbho | bir | bu | | |
| 66 | car | system | model | control | driven | ha | vehicle | System for vehicle control | |
| 67 | fault | bg | system | element | set | model | | | |
| 68 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 69 | frequency | stimulus | signal | autonomous | car | ss | | Signal interpretention from sensors? | |
| 70 | camera | state | imu | estimate | vehicle | model | behavior | Hardware beha | (IMU is a chip that works along with cameras) |
| 71 | system | sensor | vehicle | control | team | rascal | figure | Hardware relation | |
| 72 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 73 | signal | turbulence | sequence | time | set | system | coherence | Signal data behavior | |
| 74 | vehicle | system | data | wa | sensor | test | gps | Vehicle localization/vehicle navigation | |
| 75 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 76 | vehicle | road | time | future | position | forwarding | driving | Maneuver planning (Planera körning?) | |
| 77 | road | lane | semantic | node | map | based | image | Lane following | |
| 78 | vehicle | trajectory | obstacle | car | driving | self | road | Obstacle avoidance | |
| 79 | driving | event | feature | data | driver | speed | classification | Driving data handling | |
| 80 | plate | image | character | license | correction | value | neural | Image recognition of license plates | |
| 81 | arduous | narikiyo | centred | parametrize | hiriko | emplified | pean | | |
| 82 | vehicle | system | lane | line | detection | sensor | wa | Lane following | |
| 83 | potential | control | vehicle | field | sin | co | | ? | |
| 84 | segment | point | track | tolerance | curvature | curve | arc | Curve assessment | |
| 85 | video | tracking | frame | system | tdt | sequence | proposed | Video tracking | |
| 86 | uxv | node | flow | routing | network | packet | path | Network | |
| 87 | robot | mobile | control | navigation | environment | obstacle | motion | Navigation | |
| 88 | camera | image | vehicle | feature | map | track | method | Vehicle tracking through image processing | |
| 89 | car | driver | control | system | track | program | model | Vehicle Control | |
| 90 | prt | vehicle | fuel | system | consumption | campus | amd | Fuel Consumption | |
| 91 | observer | fx | fy | longitudinal | sensor | sliding | fault | Navigation failure management | |
| 92 | probability | input | traffic | markov | chain | state | participant | Traffic predictio | markov = probability mathematician |
| 93 | traffic | light | prior | detection | image | location | score | Traffic Light detection | |
| 94 | road | image | system | edge | method | detection | algorithm | Lane detection | |
| 95 | follower | sensor | velocity | leader | vehicle | test | heading | Follow other vehicle | |
| 96 | image | camera | point | frame | estimation | motion | method | Image processing | |
| 97 | image | visual | database | localization | feature | solution | infrared | Image processing | |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 98 | object | grid | moving | detection | laser | motion | data | Object detection | | | | |
| 99 | obstacle | camera | depth | image | car | map | plane | Obstacle detection | | | | |

| # | | | | | | | Label | Note 1 | Note 2 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | system | vehicle | control | driver | controller | traffic | autonomous | autonomous vehicle c | 40, 76 | |
| 1 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 2 | vehicle | system | autonomous | control | intelligent | path | navigation | autonomous vehicle navigation | | |
| 3 | vehicle | model | driving | road | autonomous | driver | based | autonomous driving | | |
| 4 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 5 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 6 | measurement | vehicle | inertial | noise | path | controller | ground | | | |
| 7 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 8 | control | vehicle | lateral | speed | autonomous | cruise | based | autonomous vehicle speed control | cruise lateral och speed refererar till speed | |
| 9 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 10 | vehicle | obstacle | uncertainty | probabilistic | planning | electric | linear | vehicle obstacle avoidance | | |
| 11 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 12 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 13 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 14 | vehicle | autonomous | system | simulation | sensor | real | time | autonomous vehicle simulation | | |
| 15 | vehicle | dc | lane | battery | system | converter | control | vehicle power management | | |
| 16 | vehicle | trajectory | path | optimization | control | planning | method | Vehicle path planning | Samma som 78 | |
| 17 | radar | automotive | sige | bicmos | technology | packaging | ghz | Radar technology | | |
| 18 | sliding | mode | vehicle | skid | control | observer | force | Sliding mode | | |
| 19 | vision | robot | road | vehicle | obstacle | dynamic | static | | | |
| 20 | vehicle | interface | mobile | human | user | automotive | interaction | Human-vehicle interface/interaction | | |
| 21 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 22 | vehicle | system | autonomous | data | tracking | tire | signal | | | |
| 23 | vehicle | detection | obstacle | algorithm | system | camera | autonomous | obstacle detection | | |
| 24 | pedestrian | feature | behavior | traffic | road | estimation | relevance | traffic behavior | | |
| 25 | road | cue | detection | level | method | low | vision | road detection | samma som 64 | |
| 26 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 27 | road | software | challenge | robot | urban | traffic | system | urban traffic software | | |
| 28 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 29 | system | lighting | autonomous | vehicle | car | intelligent | led | vehicle lights/lighting | | |
| 30 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 31 | vehicle | sensor | system | fleet | network | automated | management | Networked vehicles | | |
| 32 | object | detection | data | vehicle | sensor | fusion | classification | Object detection & classification | | |
| 33 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 34 | vehicle | tracking | distributed | network | system | topology | dynamic | Networked vehicle mapping/tracking | | |
| 35 | vehicle | automated | parking | human | autonomous | system | detection | Automated parking | | |
| 36 | road | maneuver | autonomous | vehicle | driving | prediction | model | Maneuver prediction | | |
| 37 | control | system | driving | vehicle | road | gps | driver | Driving | | |
| 38 | car | system | parking | automated | android | existing | human | 35 | | |
| 39 | cell | sram | write | tfet | characteristic | circuit | noise | Hardware | | |
| 40 | control | system | vehicle | autonomous | integrator | robot | dynamic | 0 | | |
| 41 | software | component | algorithm | robotic | advanced | robot | system | robot software | | |
| 42 | sensor | system | market | calibration | autonomous | fusion | parameter | Sensor system | | |
| 43 | control | warehouse | system | vehicle | net | petri | generalized | vehicle-warehouse system | petri net used for graphical moddeling of formal system | |
| 44 | transport | autonomous | vehicle | model | agent | based | future | Autonomous transport | | |
| 45 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 46 | camera | vehicle | image | road | vision | system | autonomous | Image processing for autonomous vehicles | | |
| 47 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 48 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 49 | system | vehicle | autonomous | forest | control | terrain | ground | Off road vehicle control | | |
| 50 | tracking | driving | autonomous | activity | driver | recognition | classification | | | |
| 51 | neural | network | system | artificial | braking | car | labview | artificial intelligence | | |
| 52 | robot | vehicle | autonomous | state | position | control | industrial | ? | | |
| 53 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 54 | vehicle | video | detection | tracking | traffic | automatic | method | Vehicle detection/track | lik 58, 66 | |
| 55 | vehicle | skew | plate | system | recognition | number | correction | Vehicle number plate | | |
| 56 | semantic | autonomous | mapping | bridge | scale | large | map | Map/mapping | | |
| 57 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 58 | vehicle | obstacle | lane | tracking | system | camera | vision | vehicle/obstacle tracki | lik 54, 66 | |
| 59 | locomotion | system | natural | oscillation | matrix | damping | body | | | |
| 60 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 61 | system | road | lane | detection | vision | vehicle | based | lane detection | | |
| 62 | driving | network | self | communication | vehicular | car | vehicle | vehicle communication | obs samma som 72 och 99 | |
| 63 | fuzzy | control | decision | logic | system | set | paper | | | |
| 64 | image | detection | estimation | vision | stereo | road | based | Road detection | samma som 25 | |
| 65 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 66 | vehicle | tracking | particle | method | time | real | filter | vehicle tracking | lik 54,58 | |
| 67 | intersection | traffic | autonomous | vehicle | transportation | intelligent | road | ? | | |
| 68 | fusion | sensor | track | filter | information | data | environment | environment information | | |
| 69 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 70 | vehicle | vision | road | algorithm | hough | lane | ransac | image analysis | | hough = smart computer guy who does computer vision stuff |
| 71 | graph | control | vehicle | theory | system | dimensional | rigid | ? | | |
| 72 | vehicle | driving | automated | communication | sensor | technology | infrastructure | Vehicles communicati | samma som 62 och 99 | |
| 73 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 74 | learning | terrain | classification | mechanical | visual | supervision | automatic | Terrain classification | | |
| 75 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 76 | robot | vehicle | control | sensor | image | autonomous | mobile | 0 | | |
| 77 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 78 | vehicle | planning | path | control | autonomous | approach | based | path planning | samma som 16 | |
| 79 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 80 | control | vehicle | sensor | tool | nist | robot | guided | guided vehicle | | |
| 81 | vehicle | system | collision | safety | service | avoidance | sensor | safety management | | |
| 82 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 83 | vehicle | driving | personalization | autopilot | safety | automated | control | automated driving | | |
| 84 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 85 | traffic | safety | vehicle | dynamic | participant | verification | markov | safe decision making | | |
| 86 | autonomous | vehicle | system | car | decision | communication | algorithm | decision making | | |
| 87 | controller | hierarchical | intelligent | vehicle | control | model | level | intelligent vehicle controller | | |
| 88 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 89 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 90 | trajectory | profile | generation | velocity | curvature | curve | vehicle | Vehicle curve handling | | |
| 91 | road | detection | representation | pedestrian | image | geometry | classification | image classification in traffic | | |
| 92 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 93 | control | vehicle | field | potential | aircraft | avoidance | spline | POTENTIAL AIRCRAFT AVOIDANCE | | |
| 94 | vehicle | information | behavior | tracking | intersection | position | system | intersection handling | | |
| 95 | vehicle | system | autonomous | time | real | utility | coordination | autonomous vehicle coordination | | |
| 96 | re | maximum | cubic | interpolation | manoeuvre | phase | parametrized | | | |
| 97 | vehicle | localization | visual | method | feature | route | robot | localization/navigation | | |
| 98 | control | vehicle | motion | autonomous | system | dynamic | tracking | autonomous vehicle motion tracking | | |
| 99 | vehicle | communication | system | wireless | intelligent | highway | roadside | vehicle communication | obs samma som 62 och 72 | |