

# Search4Code: Code Search Intent Classification Using Weak Supervision

Nikitha Rao\*, Chetan Bansal\*, Joe Guan†

\*Microsoft Research, †Microsoft

{t-nirao, chetanb, zhgua}@microsoft.com

**Abstract**—Developers use search for various tasks such as finding code, documentation, debugging information, etc. In particular, web search is heavily used by developers for finding code examples and snippets during the coding process. Recently, natural language based code search has been an active area of research. However, the lack of real-world large-scale datasets is a significant bottleneck. In this work, we propose a weak supervision based approach for detecting code search intent in search queries for C# and Java programming languages. We evaluate the approach against several baselines on a real-world dataset comprised of over 1 million queries mined from Bing web search engine and show that the CNN based model can achieve an accuracy of 77% and 76% for C# and Java respectively. Furthermore, we are also releasing Search4Code, the first large-scale real-world dataset of code search queries mined from Bing web search engine. We hope that the dataset will aid future research on code search.

## I. INTRODUCTION

Searching for code is a common task that developers perform on a regular basis. There are many sources that developers use to search for code: web search engines, code repositories, documentation, online forums, etc. Code searches typically contain a query composed of natural language and expect a code snippet result. Natural language based code search has been looked at by different approaches such as traditional information retrieval techniques [1], [2], [3], deep learning [4], and hybrid approaches [5] that combine various methodologies. One commonality that exists is the requirement of a sufficiently large dataset composed of code and the corresponding natural language labels. Traditionally, researchers have used different methods to gather data, including using the associated docstring of the code snippet and the question title from coding related forums (e.g. StackOverflow). However, these natural language labels do not accurately represent how developers perform searches for code in a typical search engine. While there exist some datasets that include human-annotated labels for code [6], these are limited in size and quantity.

We present a dataset compiled from query logs comprised of millions of queries from Bing web search engine. This dataset contains aggregated queries, which have been anonymized, and classified as either having a code search intent or not for C# and Java programming languages. The dataset also contains the most frequently clicked URLs and a popularity metric denoting the query frequency. To create a large-scale dataset of code search queries, it is crucial to automatically detect code search intent in search queries. Previous research in the area of search query classification [7], [8] has focused

primarily on the classification of web queries in categories such as Debug, API, and HowTo using heuristics and rule-based methods which tend to overfit.

In this paper, we introduce a novel weak supervision based model to classify code search intent in search queries. We define a query as having code search intent if it can be sufficiently answered with a snippet of code. To the best of our knowledge, this is the first usage of weak supervision in the software engineering domain. In summary, our main contributions are:

- A novel weak supervision based model to detect code search intent in queries.
- A large-scale dataset of queries<sup>1</sup>, mined from Bing web search engine, that can be used for code search research.

## II. BACKGROUND AND MOTIVATION

Our work builds on recent advances in the areas of code search, search query intent classification, and weak supervision. In this section, we provide a brief background of the same.

**Code Search:** Code search is a sub-field in information retrieval that focuses on finding relevant code snippets given a natural language query. Code search is an integral part of the software development process [9], [10], [11] as developers often search for code using search engines, documentation, and online forums. However, a significant bottleneck in this area is the lack of datasets for building and experimenting with new techniques. The most recent work in curating a dataset contains 99 human-annotated queries across multiple languages [6] and 287 question-answer pairs extracted from StackOverflow [12]. We aim to contribute a new method to generate a code search dataset by mining query logs from Bing web search engine. Additionally, we open-source this dataset to aid future research on code search.

**Intent Classification:** Applications of intent classification in web search include several domains like healthcare [13], security [14] and e-commerce [15], [16]. Wang et al. have leveraged intent understanding for improving effort estimation in code reviews [17], [18]. Recently, software engineering related search queries have been analyzed and classified into

<sup>1</sup><https://github.com/microsoft/Search4Code/>

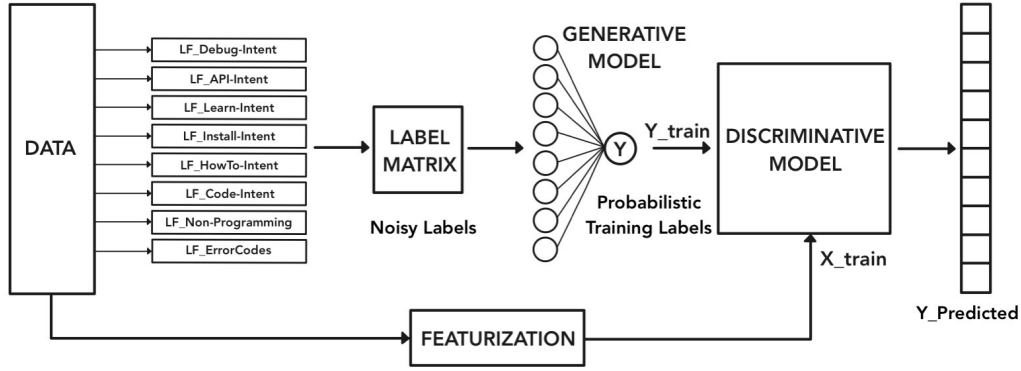


Fig. 1: Overview of the pipeline.

different categories by using distant supervision [7] and token-level intent aggregation [8] to better understand developer behaviour. Our goal is to further improve upon these methods by introducing a weak supervision based approach for code search intent classification.

**Weak Supervision:** One of the primary challenges in supervised learning is to obtain large-scale labelled data. As mentioned above, this obstacle exists in the code search space as well. Weak supervision [19], [20] leverages ‘weak’ or ‘noisy’ learning functions to automatically assign labels to a large amount of unlabeled data.

Formally speaking, given a set of unlabeled data points,  $X$ , the objective of weak supervision is to estimate the ground truth label by using a set of  $n$  learning functions. Each learning function has a probability to abstain and a probability to correctly label a data point as positive or negative. The learning functions are applied over  $m$  unlabeled data points to create a matrix of label outputs,  $\Lambda$ . The generative model then takes  $\Lambda$  as input and returns the probability scores for each class based on the agreements and disagreements between the learning functions. The predicted label distribution output can then be used as probabilistic training labels by a discriminative classifier for a downstream classification task. We use weak supervision to generate the train labels for the code search intent classification task.

### III. APPROACH

In this section, we elaborate on our approach for code intent classification. First, we build the generative model using weak supervision to get the labels for the training data using snorkel, a weak supervision framework by Stanford [21]. We then use this data to train discriminative models to classify queries as having code search intent or not. Figure 1 provides an overview of the entire pipeline.

#### A. Generative Model Pipeline

**Data Collection:** We randomly sample 1 million search queries each for C# and Java, collected from 1<sup>st</sup> September, 2019 to 31<sup>st</sup> August, 2020 from Bing web search engine. We identify queries related to each programming language by doing a simple keyword-based pattern matching (‘c#’, ‘c sharp’ and ‘csharp’ for C# and ‘java’ for Java) [22]. We apply additional filters to ensure that all the queries are in English

locale from the USA region and we eliminate any traffic from bots and other services. Additionally, we exclude queries that have multiple programming languages in them such as ‘c# vs java’, ‘how hard is c# compared to java or c++?’, ‘java to c# converter’ and so on to better isolate queries to an individual programming language.

**Learning Functions (LFs):** As discussed in Section II, we use several ‘weak’ or ‘noisy’ learning functions, described in Table I, that are combined in a weighted manner by the generative model. Weak supervision sources generally include external knowledge bases, patterns, dictionaries and even domain-specific heuristics. In the context of code search intent classification, we leverage the software engineering sub-intent classifiers (such as Debug, HowTo, etc.) proposed by Rao et al. [7]. We also introduce learning functions to identify patterns that indicate code examples, error codes and exceptions. Each learning function acts as a binary classifier that identifies either code search or not code search intent and abstains otherwise. We use the label 1 for code search intent, 0 for not code search intent and  $-1$  for abstain. The label for each learning function is chosen after manually analyzing a sample of queries. Table I provides the target label and description of heuristics used for each of the learning functions used along with a few example queries.

**Generative Model:** We apply all the individual learning functions to the data and construct a label matrix that is then fed to the generative model. The generative model then uses a weighted average of all learning functions outputs, based on the agreements and disagreements between the learning functions, to return the probability scores for each class. Each datapoint is then assigned a label based on the class having the higher probability score.

#### B. Discriminative Model Pipeline

**Data:** We use the output of the generative model as the train labels ( $Y_{train}$ ) for the data we collected earlier. We then preprocess and featurize the data before passing it to the discriminative model.

**Preprocessing and Featurization:** We first tokenize the queries based on non-alphanumeric characters and remove all stopwords. We then transform the query text into its vector

Learning Function	Label	Description	Examples
API Intent	1	Keywords like ‘api’, ‘function’, ‘method’, ‘call’, etc. that indicate a specific API usage.	‘c# example of restful post api call form url encode’, ‘java immutablelist api’.
Debug Intent	0	Keywords like ‘error’, ‘exception’, ‘fail’, ‘not working’, ‘debug’, etc. that indicate an error or issue.	‘500 internal server error in web api c#’, ‘java createnewfile not working’.
HowTo Intent	1	Keyword ‘how’ is present to indicate the need to accomplish a specific task.	‘c# asp.net how to implement click event for textbox’, ‘how to do quicksort in java’.
Learn Intent	0	Keywords like ‘tutorial’, ‘what’, ‘why’, ‘difference’, ‘versus’, etc. that indicate learning new topics.	‘block body vs lambda method c#’, ‘what is the order of precedence for java math’.
Install Intent	0	Keywords like ‘install’, ‘download’, ‘update’ etc. that indicate installing software packages.	‘c# .net install .msi remotely’, ‘download selenium web driver jars for java’
Code Search Intent	1	Keywords like ‘example’, ‘sample code’, ‘snippet’, ‘implementation’, etc. that indicate code search.	‘proxysocket c# code sample’, ‘java void method no parameters example’
Non-Programming	0	Keywords like ‘interview’, ‘jobs’, etc. that indicate non-programming related queries.	‘c# array questions for interviews’, ‘part time java coding jobs’
Error Codes	0	Regex based patterns to find C# error codes or Java exceptions	‘cs7038 wcf c# failed to emit module’, ‘java.io.eofexception: postman’.

**Table I:** Overview of the learning functions.

representation using Word2Vec [23] to capture any semantic similarities. We retrain the Word2Vec model from scratch on our query data since the pre-trained Word2Vec models don’t generalize well to queries related to programming languages. We compute the word embeddings for each token in a query using the trained Word2Vec model and compute query embedding as the average of all token embeddings. This forms the training data ( $X_{train}$ ) for the discriminative models.

**Discriminative Model:** Using the generated training labels ( $Y_{train}$ ) along with the featurized train data ( $X_{train}$ ) data, we train several supervised machine learning and deep learning models to tackle the problem of code search intent detection in search queries. We further elaborate on the various discriminative models used in Section IV-A.

#### IV. EXPERIMENTS AND RESULTS

In this section, we first describe the experimental setup. We then present the evaluation for the generative model that is used to derive train data labels. Lastly, we evaluate the efficacy of various discriminative models for code search intent classification in search queries. We evaluate the performance of each model against the overall test accuracy along with the precision, recall and F1 scores for both classes. Note that we train and evaluate the models for C# and Java separately.

##### A. Experimental Setup

**Dataset:** The featurized data described in Section III-B along with the generated train labels from Section III-A is used as the training data for the various discriminative models.

For the test data, we uniformly sample a set of 200 queries based on query length for both C# and Java. Three annotators then manually label the data independently. We compute the inter-rater agreement score to be 0.75 using Fleiss’ Kappa [24], which translates to substantial agreement. The final label is obtained by taking a majority vote. We find the distribution of queries with code search intent in the manually

labelled test data to be 62.0% for C# and 34.5% for Java.

**Discriminative Models:** We compare the performance of various machine learning and deep learning models to find the best performing code search intent classification model. In particular, we look at the following discriminative models

- First, we look at non-deep learning models like Logistic Regression and Random Forest. We use the default version of the models from scikit-learn to implement them.
- For the deep learning models, we look at Bidirectional LSTM (BiLSTM) with attention and CNN. The BiLSTM is implemented by adding the bidirectional layer on top of the LSTM layer [25]. For the CNN, we use convolution layers with ReLu activation followed by maxpool layers and a dense output layer with sigmoid activation [26], [27]. We implement the models using keras with tensorflow backend.

##### B. Generative Model Evaluation

To evaluate the performance of the generative model for generating the train data labels, we compare the performance of the model with a majority vote model on the test data. The majority vote model assigns the label for each query based on the majority vote of all eight learning functions and ties are settled by assigning a random label. Table II summarizes the evaluation scores for the two models. We find that the generative model outperforms the majority vote model across all metrics with an overall test accuracy of 73% and 72% for C# and Java respectively.

##### C. Discriminative Model Evaluation

To evaluate the efficacy of the various discriminative models for code search intent detection, we first train each model on the train data and compare the performance scores on the test data. Table III summarizes the performance scores of the four models. We find that the CNN model outperforms all the other models across majority of the metrics with an overall test accuracy of 77% and 76% for C# and Java respectively.

Programming Language	Model	Accuracy	Code Intent			Not Code Intent		
			Precision	Recall	F1 Score	Precision	Recall	F1 Score
C#	Majority Vote	66	73	72	72	55	57	56
	Generative	<b>73</b>	<b>80</b>	<b>76</b>	<b>78</b>	<b>63</b>	<b>68</b>	<b>66</b>
Java	Majority Vote	67	52	71	60	81	65	72
	Generative	<b>72</b>	<b>57</b>	<b>80</b>	<b>67</b>	<b>87</b>	<b>68</b>	<b>76</b>

**Table II:** Evaluation of the generative model on the test data.

Programming Language	Model	Accuracy	Code Intent			Not Code Intent		
			Precision	Recall	F1 Score	Precision	Recall	F1 Score
C#	Logistic Regression	71	73	86	79	68	47	56
	Random Forest	73	73	<b>90</b>	80	<b>72</b>	45	55
	CNN	<b>77</b>	<b>79</b>	85	<b>82</b>	<b>72</b>	<b>63</b>	<b>67</b>
	BiLSTM	72	77	78	78	64	62	63
Java	Logistic Regression	74	59	85	<b>70</b>	90	69	78
	Random Forest	73	57	<b>89</b>	<b>70</b>	<b>91</b>	65	76
	CNN	<b>76</b>	<b>63</b>	74	68	85	<b>77</b>	<b>81</b>
	BiLSTM	73	59	76	66	85	72	78

**Table III:** Evaluation of the discriminative models on the test data.

## V. CODE SEARCH QUERY DATASET

In this work, we have built a code search intent classification model based on weak supervision. One of the major impediments for research in this domain is the lack of publicly available large-scale datasets. On this account, we are releasing Search4Code [28], the first large-scale real-world dataset of code search queries for C# and Java mined from Bing web search engine. The dataset is composed of about 4,974 C# queries and 6,596 Java queries. We hope that this dataset will aid future research to not just better code search intent detection but also applications like natural language based code search.

To build the dataset we first collect the anonymized query logs for one year. We then mine the code search queries by following several steps of log mining, processing and aggregation. First, we apply the same filters for locale, bots, etc. and filter out queries that are not related to C# or Java programming languages as described in Section III-A. Next, we apply a  $k$ -anonymity filter [29] with a high value of  $k$ . This filters out queries from the dataset which were entered by less than  $k$  users and could potentially contain sensitive information which was known to less than  $k$  users. Finally, we apply the best performing discriminative model (i.e. CNN) to the queries to identify queries with code search intent.

We have defined the schema for the dataset in Table IV. It contains not only the raw queries but also other useful attributes such as top click URLs and rank based on popularity. Here are some of the unique features of the dataset:

- **Real queries:** The queries are sampled from anonymized Bing search logs. We believe this provides a rich dataset indicative of real-world user behavior.
- **Click URLs:** Each query has a list of the three most

Attribute	Description
Id	Identifier for the query.
Query	The raw query issued by the users.
Is Code Search Query	Whether the query has code search intent.
Top Clicked URLs	Top 3 document URLs (comma delimited) by click frequency.
Popularity Rank	Rank based on the query frequency. Most popular query is ranked 1.

**Table IV:** Schema of the code search queries dataset

frequently clicked URLs from the query logs based on user interactions.

- **Popularity score:** Each query is assigned a popularity rank based on the frequency of occurrence.
- **Large scale:** The dataset contains thousands of queries, hence, enabling large scale analysis of search for software engineering.
- **Other applications:** Since the dataset contains both code-search and non code-search queries, it could also be used to analyze other user intents, as described in our prior work [7].

**Limitations:** While we provide the top clicked URLs for each query, the code samples themselves are not provided and will have to be mined from the URLs present. Also, we predict the code search intent solely based on the search queries since we don't have access to the content from the clicked web pages. So, it is possible that not all the Clicked Urls contain code samples. Additionally, since the code search intent labels are generated by the CNN model, they will not be 100% accurate.

## VI. CONCLUSION AND FUTURE WORK

Search is heavily used by developers for various tasks during the software development process. Given the lack of labelled data, we use weak supervision for code search intent classification. We develop a CNN based model for code search intent classification for C# and Java search queries mined from Bing web search engine. We also evaluate it against various baselines which demonstrates the efficacy of the weak supervision based approach. Furthermore, we are releasing the first large-scale real-world code search query dataset comprising more than 11,000 search queries. Our code search intent model can be integrated with several applications such as IDEs, Search Engines and even developer forums like StackOverflow for improving the code search experience.

Future work on code search can leverage the dataset for building and improving natural language based code search techniques. Additionally, to the best of our knowledge, this is the first work to explore the usage of weak supervision in the software engineering domain. Weak supervision can also be leveraged in other tasks such as bug detection and program repair where a limited amount of labelled data is available. Lastly, we plan to experiment with other advanced transformer-based neural model architectures such as BERT [30] to improve the discriminative model performance for code search intent classification.

## VII. ACKNOWLEDGEMENTS

We would like to acknowledge the invaluable contributions of Mark Wilson-Thomas, Shengyu Fu, Nachi Nagappan, Tom Zimmermann and B. Ashok.

## REFERENCES

- [1] W.-K. Chan, H. Cheng, and D. Lo, "Searching connected api subgraph via text phrases," in *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, ser. FSE '12. New York, NY, USA: Association for Computing Machinery, 2012.
- [2] E. Hill, L. Pollock, and K. Vijay-Shanker, "Improving source code search with natural language phrasal representations of method signatures," in *2011 26th IEEE/ACM International Conference on Automated Software Engineering*, 2011, pp. 524–527.
- [3] Meili Lu, X. Sun, S. Wang, D. Lo, and Yucong Duan, "Query expansion via wordnet for effective code search," in *2015 IEEE 22nd SANER*, 2015, pp. 545–549.
- [4] X. Gu, H. Zhang, and S. Kim, "Deep code search," in *Proceedings of the 40th International Conference on Software Engineering*, ser. ICSE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 933–944.
- [5] S. Sachdev, H. Li, S. Luan, S. Kim, K. Sen, and S. Chandra, "Retrieval on source code: A neural code search," in *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, ser. MAPL 2018. New York, NY, USA: Association for Computing Machinery, 2018, p. 31–41.
- [6] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "Codesearchnet challenge: Evaluating the state of semantic code search," in *arXiv*, 2020.
- [7] N. Rao, C. Bansal, T. Zimmermann, A. H. Awadallah, and N. Nagappan, "Analyzing web search behavior for software engineering tasks," *arXiv preprint arXiv:1912.09519*, 2019.
- [8] M. M. Rahman, J. Barson, S. Paul, J. Kayani, F. A. Lois, S. F. Quezada, C. Parmin, K. T. Stolee, and B. Ray, "Evaluating how developers use general-purpose web-search for code retrieval," in *Proceedings of the 15th International Conference on Mining Software Repositories*, ser. MSR '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 465–475.
- [9] V. Bauer, J. Eckhardt, B. Hauptmann, and M. Klimek, "An exploratory study on reuse at google," 06 2014.
- [10] C. Sadowski, K. T. Stolee, and S. Elbaum, "How developers search for code: A case study," in *Joint Meeting of the European Software Engineering Conference and the Symposium on the Foundations of Software Engineering (ESEC/FSE)*, 1600 Amphitheatre Parkway, 2015.
- [11] M. Umarji, S. E. Sim, and C. Lopes, "Archetypal internet-scale source code searching," in *Open Source Development, Communities and Quality*, B. Russo, E. Damiani, S. Hissam, B. Lundell, and G. Succi, Eds. Boston, MA: Springer US, 2008, pp. 257–263.
- [12] H. Li, S. Kim, and S. Chandra, "Neural code search evaluation dataset," in *arXiv*, 2019.
- [13] M. Paul, R. White, and E. Horvitz, "Diagnoses, decisions, and outcomes: Web search as decision support for cancer," 05 2015, pp. 831–841.
- [14] C. Bansal, P. Deligiannis, C. Maddila, and N. Rao, "Studying ransomware attacks using web search logs," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR '20, 2020, p. 1517–1520.
- [15] A. Ahuja, N. Rao, S. Katariya, K. Subbian, and C. K. Reddy, "Language-agnostic representation learning for product search on e-commerce platforms," in *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 7–15.
- [16] N. Rao, C. Bansal, S. Mukherjee, and C. Maddila, "Product insights: Analyzing product intents in web search," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, ser. CIKM '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 2189–2192. [Online]. Available: <https://doi.org/10.1145/3340531.3412090>
- [17] S. Wang, C. Bansal, N. Nagappan, and A. A. Philip, "Leveraging change intents for characterizing and identifying large-review-effort changes," in *Proceedings of the Fifteenth International Conference on Predictive Models and Data Analytics in Software Engineering*, 2019, pp. 46–55.
- [18] S. Wang, C. Bansal, and N. Nagappan, "Large-scale intent analysis for identifying large-review-effort code changes," *Information and Software Technology*, p. 106408, 2020.
- [19] J. Robinson, S. Jegelka, and S. Sra, "Strength from weakness: Fast learning using weak supervision," 2020.
- [20] Z.-H. Zhou, "A brief introduction to weakly supervised learning," *National Science Review*, vol. 5, no. 1, pp. 44–53, 08 2017.
- [21] A. Ratner, S. H. Bach, H. R. Ehrenberg, J. A. Fries, S. Wu, and C. Ré, "Snorkel: Rapid training data creation with weak supervision," *CoRR*, vol. abs/1711.10160, 2017.
- [22] F. Hassan, C. Bansal, N. Nagappan, T. Zimmermann, and A. H. Awadallah, "An empirical study of software exceptions in the field using search logs," in *Proceedings of the 14th ACM / IEEE International Symposium on Empirical Software Engineering and Measurement*, ser. ESEM '20. New York, NY, USA: Association for Computing Machinery, 2020.
- [23] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *NIPS'13 - Proceedings(Volume 2)*, USA, 2013.
- [24] J. L. Fleiss, "Measuring nominal scale agreement among many raters," *Psychological bulletin*, vol. 76, no. 5, p. 378, 1971.
- [25] G. Liu and J. Guo, "Bidirectional lstm with attention mechanism and convolutional layer for text classification," *Neurocomputing*, vol. 337, 02 2019.
- [26] V. Choubey. (2020) Text classification using cnn. [Online]. Available: <https://medium.com/voice-tech-podcast/text-classification-using-cnn-9ade8155dfb9>
- [27] H. T. Le, C. Cerisara, and A. Denis, "Do convolutional networks need to be deep for text classification?" 2017.
- [28] "Search4Code: Web queries dataset for code search," 2020. [Online]. Available: <https://github.com/microsoft/Search4Code/>
- [29] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 557–570, 2002.
- [30] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.