

# Multi-task Deep Learning in the Software Development domain

Silvia Severini, Garching, 27.05.19  
Advisor: Ahmed Elnaggar

Chair of Software Engineering for Business Information Systems (sebis)  
Faculty of Informatics  
Technische Universität München  
[www.matthes.in.tum.de](http://www.matthes.in.tum.de)

- Motivation
- Introduction
- Research questions
- Methodology
- Tasks
- Model architecture overview
- Timeline of the thesis
- References

- **Motivation**
- Introduction
- Research questions
- Methodology
- Tasks
- Model architecture overview
- Timeline of the thesis
- References

## A recent growing interest

- A variety of tasks in the software development domain can benefit from the aid of Machine Learning and Deep Learning
- Deep learning has achieved competitive performance against previous algorithms on about 40 SE tasks
- Industrial practitioners are also interested in integrating Deep Learning into their SE solutions

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a dark blue rectangular background.The Google logo, featuring the word "Google" in its characteristic multi-colored font.The IBM logo, consisting of the letters "IBM" in a blue, horizontally-striped font.The DeepMind logo, featuring a blue circular icon with a white swirl and the word "DeepMind" in a grey sans-serif font.The Microsoft logo, consisting of a four-colored square icon (red, green, blue, yellow) followed by the word "Microsoft" in a grey sans-serif font.

# Tasks in the Software development domain



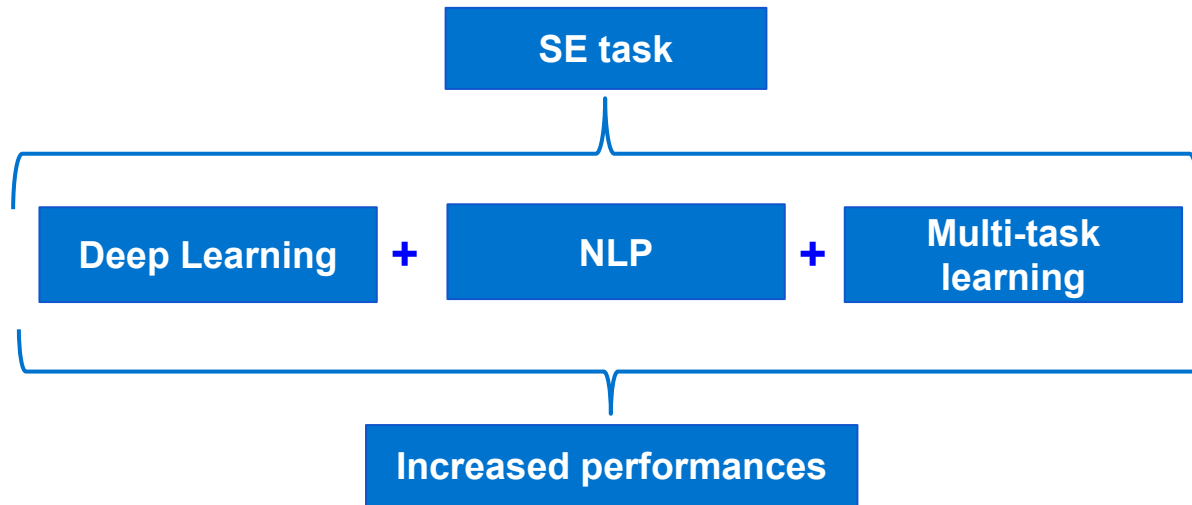
# Tasks in the Software development domain

# Tasks in requirement (1 paper)	# Tasks in Testing (27 papers)	# Tasks in management (12 papers)
A1 Requirement extraction from natural languages (1)	D1 Defect prediction (9)	<b>F1 Development cost or effort estimation (6)</b>
	<b>D2 Reliability or changeability estimation (8)</b>	F2 Source code classification (4)
	D3 Deep learning testing (3)	F3 Software size estimation (1)
	D4 Energy consumption estimation (1)	F4 Traceable link generation (1)
	<b>D5 Grammar-based fuzzing testing (1)</b>	
	D6 Retesting necessity estimation (1)	
	D7 Reliability model selection (1)	
	D8 Robot testing (1)	
	<b>D9 Test input generation for mobile (1)</b>	
	D10 Testing effort estimation (1)	
	# Tasks in maintenance (27 papers)	
	<b>E1 Malware detection (10)</b>	Industrial practitioners participate in 13 SE tasks (21 papers)
	E2 Bug localization (4)	• C1: <i>DeepMind, Facebook, Google, Microsoft</i> (8 papers)
	E3 Clone detection (3)	• C5: <i>Fiat Chrysler Automobiles</i> (1)
	<b>E4 System anomaly prediction (2)</b>	• C7: <i>Microsoft</i> (1)
	E5 Workload prediction in the cloud (2)	• C11: <i>Clinic Inc.</i> (1)
	E6 Bug report summarization (1)	• C13: <i>Facebook</i> (1)
	E7 Bug triager (1)	• D2: <i>URU Video, Inc.</i> (1)
	<b>E8 Duplicate bug report detection (1)</b>	• D5: <i>Microsoft</i> (1)
	<b>E9 Feature location (1)</b>	• D9: <i>IBM</i> (1)
	E10 Real-time task scheduling (1)	• E1: <i>Baidu, Microsoft</i> (2)
	E11 Test report classification (1)	• E4: <i>Tencent Corporation</i> (1)
		• E8: <i>Accenture Tech.</i> (1)
		• E9: <i>ABB Corporate</i> (1)
		• F1: <i>Motorola Canada Ltd.</i> (1)

Fig. 4. The SE tasks solved by deep learning and participated by industrial practitioners [3]

# Outline

- Motivation
- **Introduction**
- Research questions
- Methodology
- Tasks
- Model architecture overview
- Timeline of the thesis
- References

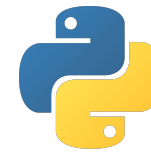




# Natural Language Processing for Source Code

Programming language as a new language like English

- Complexity
- Context awareness
- Unlimited vocabulary
- Dataset scarcity (required GitHub scraper<sup>[5]</sup>)
- Tokenization of each programming language

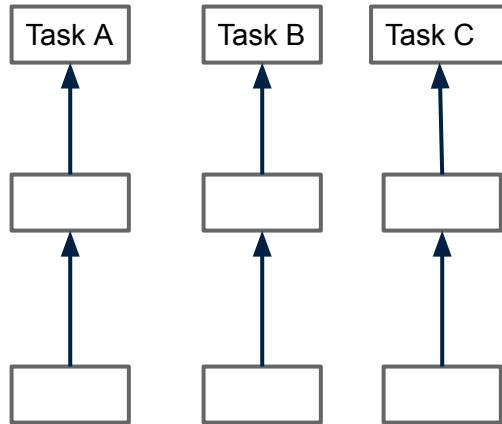


```
public class HelloWorld {  
  
    public static void main(String[] args) {  
        // Prints "Hello, World" to the terminal window.  
        System.out.println("Hello, World");  
    }  
  
}
```

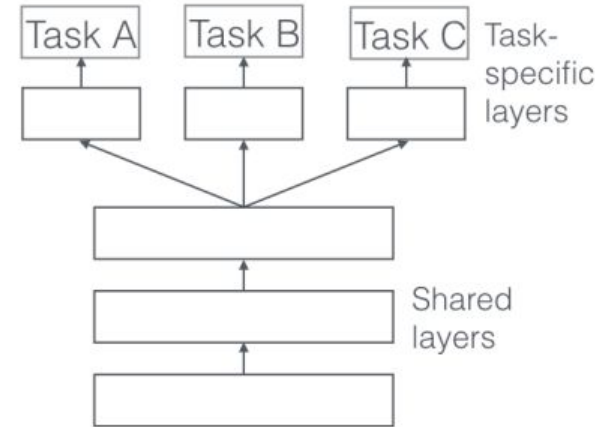


```
public class HelloWorld { public static void main ( String [] args ) { System . out . println ( " Hello, World " ) ; } }
```

# Single-task vs Multi-task learning



Single-task learning



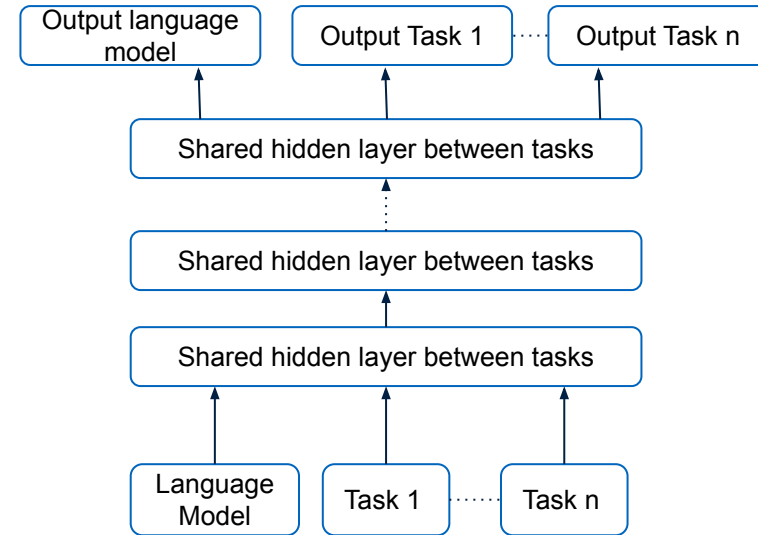
Hard parameter sharing for multi-task learning in deep neural networks [\[1\]](#)

## Why Multi-task learning ?

*“Given  $m$  learning tasks  $\{T_i\}_{i=1}^m$  where all the tasks or a subset of them are related, multi-task learning aims to help improve the learning of a model for  $T_i$  by using the knowledge contained in all or some of the  $m$  tasks.”*

- Implicit data augmentation
- Regularization
- Attention focusing
- Representation bias

=> Augment of the generalization capabilities



- Motivation
- Introduction
- **Research questions**
- Methodology
- Tasks
- Model architecture overview
- Timeline of the thesis
- References

**1** Can multi-task deep learning be beneficial for tasks in the software development domain?

**2** How far is multi-task deep learning from state-of-the-art solutions in the software development domain?

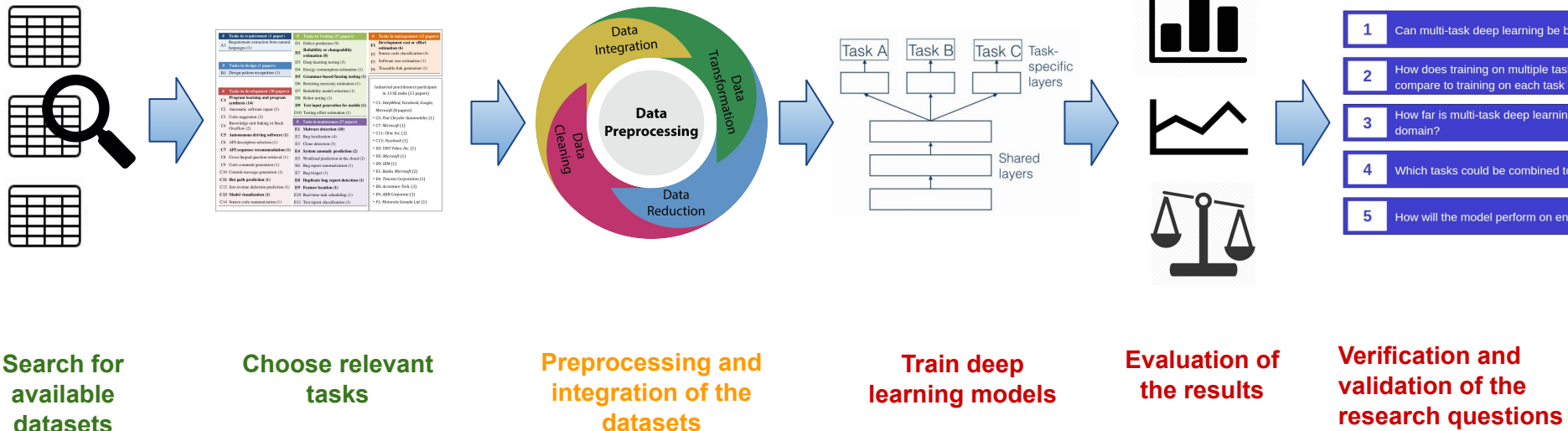
**3** Could the model be trained with English language and programming language together?

**4** How does training on multiple tasks of the software development domain simultaneously compare to training on each task separately?

**5** Which tasks could be combined together in order to achieve better performances?

- Motivation
- Introduction
- Research questions
- **Methodology**
- Tasks
- Model architecture overview
- Timeline of the thesis
- References

# Methodology: overview



Search for available datasets

Choose relevant tasks

Preprocessing and integration of the datasets

Train deep learning models

Evaluation of the results

Verification and validation of the research questions

- 1 Can multi-task deep learning be better than training on each task?
- 2 How does training on multiple tasks compare to training on each task?
- 3 How far is multi-task deep learning from the domain?
- 4 Which tasks could be combined together?
- 5 How will the model perform on new tasks?

# Outline

- Motivation
- Introduction
- Research questions
- Methodology
- **Tasks**
- Model architecture overview
- Timeline of the thesis
- References



# Tasks in the Software development domain

# Tasks in requirement (1 paper)	# Tasks in Testing (27 papers)	# Tasks in management (12 papers)
A1 Requirement extraction from natural languages (1)	D1 Defect prediction (9)	<b>F1 Development cost or effort estimation (6)</b>
	<b>D2 Reliability or changeability estimation (8)</b>	F2 Source code classification (4)
	D3 Deep learning testing (3)	F3 Software size estimation (1)
<b># Tasks in design (1 papers)</b>	D4 Energy consumption estimation (1)	F4 Traceable link generation (1)
B1 Design pattern recognition (1)	<b>D5 Grammar-based fuzzing testing (1)</b>	
	D6 Retesting necessity estimation (1)	
<b># Tasks in development (30 papers)</b>	D7 Reliability model selection (1)	Industrial practitioners participate in 13 SE tasks (21 papers)
<b>C1 Program learning and program synthesis (14)</b>	D8 Robot testing (1)	• C1: <i>DeepMind, Facebook, Google, Microsoft</i> (8 papers)
C2 Automatic software repair (2)	<b>D9 Test input generation for mobile (1)</b>	• C5: <i>Fiat Chrysler Automobiles</i> (1)
C3 Code suggestion (2)	D10 Testing effort estimation (1)	• C7: <i>Microsoft</i> (1)
C4 Knowledge unit linking in Stack Overflow (2)	<b># Tasks in maintenance (27 papers)</b>	• C11: <i>Clinic Inc.</i> (1)
<b>C5 Autonomous driving software (1)</b>	<b>E1 Malware detection (10)</b>	• C13: <i>Facebook</i> (1)
C6 API description selection (1)	E2 Bug localization (4)	• D2: <i>URU Video, Inc.</i> (1)
<b>C7 API sequence recommendation (1)</b>	E3 Clone detection (3)	• D5: <i>Microsoft</i> (1)
C8 Cross-lingual question retrieval (1)	<b>E4 System anomaly prediction (2)</b>	• D9: <i>IBM</i> (1)
C9 Code comment generation (1)	E5 Workload prediction in the cloud (2)	• E1: <i>Baidu, Microsoft</i> (2)
C10 Commit message generation (1)	E6 Bug report summarization (1)	• E4: <i>Tencent Corporation</i> (1)
<b>C11 Hot path prediction (1)</b>	E7 Bug triager (1)	• E8: <i>Accenture Tech.</i> (1)
C12 Just-in-time deflection prediction (1)	<b>E8 Duplicate bug report detection (1)</b>	• E9: <i>ABB Corporate</i> (1)
<b>C13 Model visualization (1)</b>	<b>E9 Feature location (1)</b>	• F1: <i>Motorola Canada Ltd.</i> (1)
C14 Source code summarization (1)	E10 Real-time task scheduling (1)	
	E11 Test report classification (1)	

Fig. 4. The SE tasks solved by deep learning and participated by industrial practitioners [3]

Tasks	Description	Number of samples
Program learning and synthesis	Generate programs from natural language descriptions	100.000
API sequence recommendation	Generates relevant API usage sequences given a natural language query	7.500.000
Code comment generation	Automatic generation of code comments	400.000
Commit message generation	Automatically “translate” diffs into commit messages	30.000
Source code summarization	Summarization of source code snippets	80.000

Unsupervised Language model	English: 1 Billion world corpus <a href="#">[4]</a>	300.000.000
	Java from PGA <a href="#">[5]</a>	500.000
	SQL corpus <a href="#">[6]</a>	135.000
	150K Python Dataset <a href="#">[7]</a>	150.000
	C# from PGA <a href="#">[5]</a>	500.000

# Example of input-output pairs

## Source code summarization:

- ```
from pygithub3 import Github\n\nusername = raw_input("Please enter a Github username: ")\npassword = raw_input("Please enter the account password: ")\n\ngh = Github(login=username, password = password)\n\nget_user = gh.users.get()\n\nuser_repos = gh.repos.list().all()\n\nfor repo in user_repos:\n    print repo.language
```
- Getting repository information using pygithub3 for Python

## Code comment generation:

- ```
public void handleEntryExpiredSA(EntryExpiredBusPacket packet) throws Exception {\n    handleEntryExpiredCoreSA(packet.getEntryHolder(),packet.getTransaction(),packet.isFromReplication());\n}
```
- Handles EntryExpired packets.

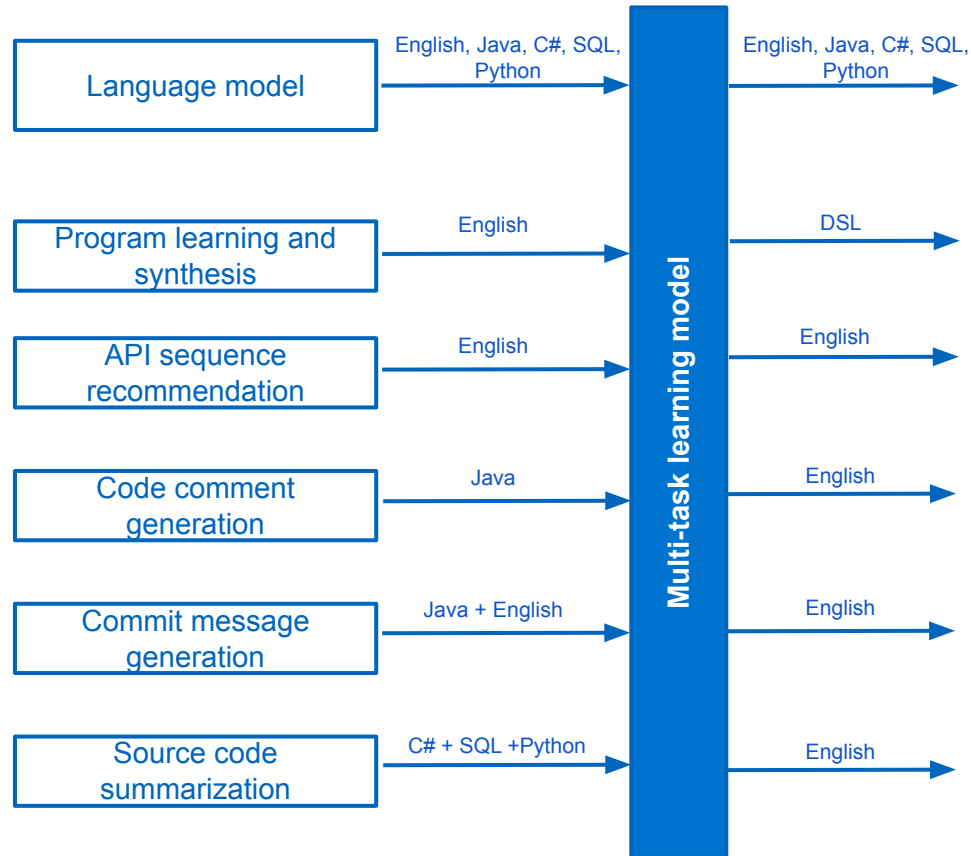
## API sequence recommendation:

- return a printable representation of this exception for debugging purposes
- `StringBuffer . <init> StringBuffer . append StringBuffer . toString`

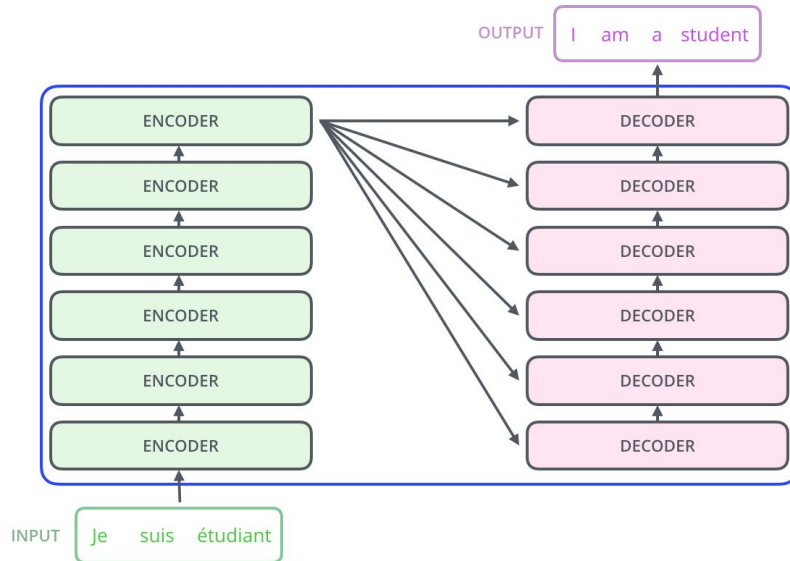
# Outline

- Motivation
- Introduction
- Research questions
- Methodology
- Tasks
- **Model architecture overview**
- Timeline of the thesis
- References

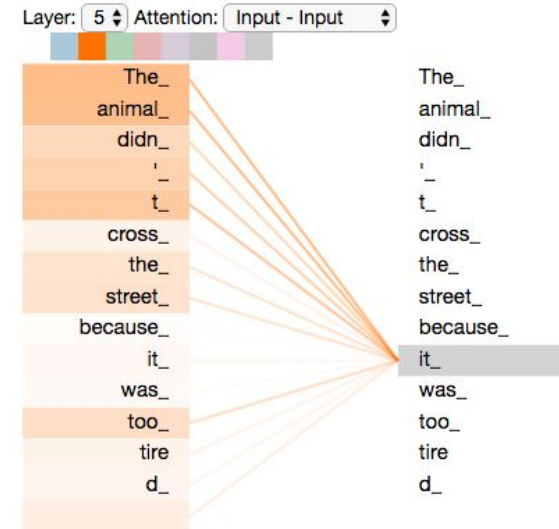
# Model architecture overview



## Encoder - Decoder model

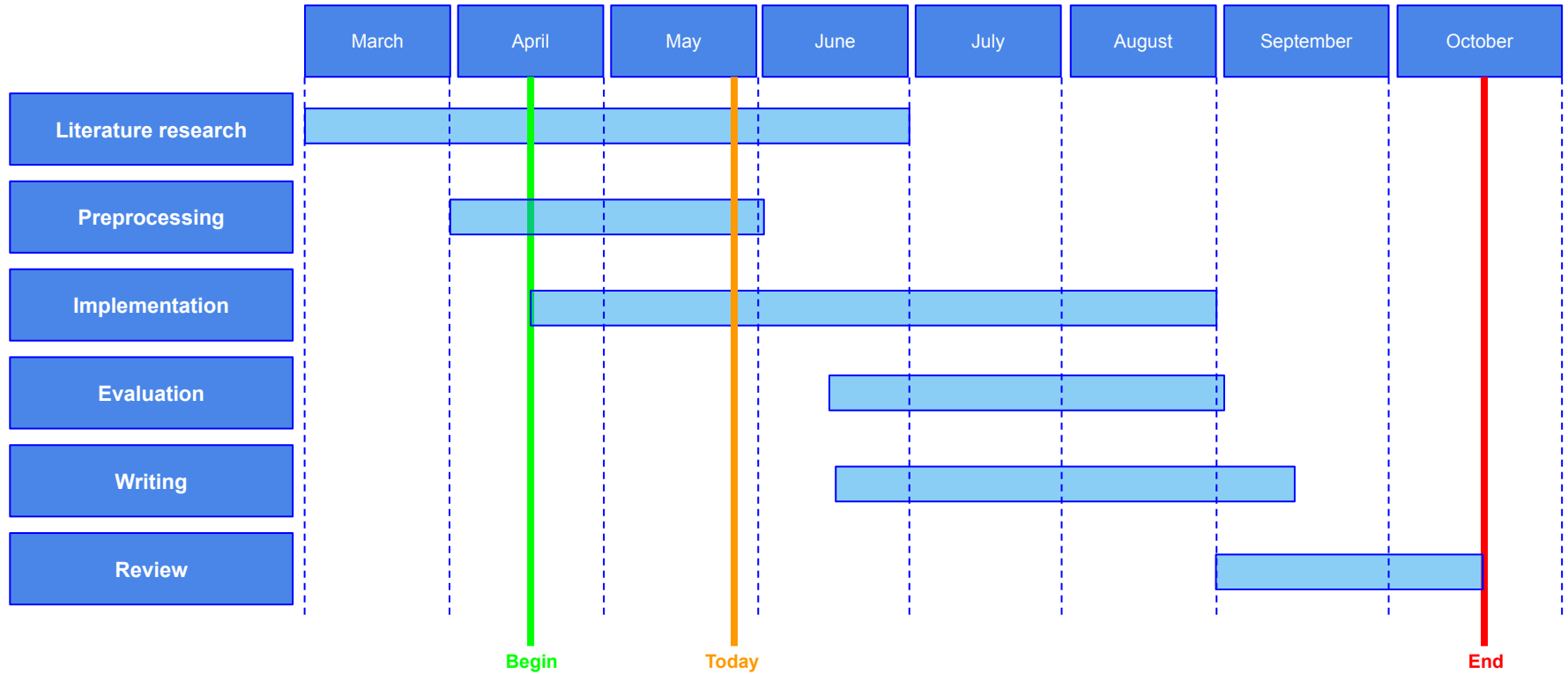


## Attention mechanism



- Motivation
- Introduction
- Research questions
- Methodology
- Tasks
- Model architecture overview
- **Timeline of the thesis**
- References

# Timeline of the thesis





# Outline

- Motivation
- Introduction
- Research questions
- Methodology
- Tasks
- Model architecture overview
- Timeline of the thesis
- **References**

- [1] Ruder, Sebastian. "An overview of multi-task learning in deep neural networks." arXiv preprint arXiv:1706.05098 (2017)
- [2] Zhang, Yu, and Qiang Yang. "A survey on multi-task learning." arXiv preprint arXiv:1707.08114 (2017).
- [3] Li, Xiaochen, et al. "Deep Learning in Software Engineering." arXiv preprint arXiv:1805.04825 (2018).
- [4] <http://www.statmt.org/lm-benchmark/>
- [5] <https://github.com/src-d/datasets/tree/master/PublicGitArchive>
- [6] [https://github.com/LittleYUYU/StackOverflow-Question-Code-Dataset/blob/master/annotation\\_tool/data/code\\_solution\\_labeled\\_data/source/sql\\_how\\_to\\_do\\_it\\_by\\_classifier\\_multiple\\_iid\\_to\\_code.pickle](https://github.com/LittleYUYU/StackOverflow-Question-Code-Dataset/blob/master/annotation_tool/data/code_solution_labeled_data/source/sql_how_to_do_it_by_classifier_multiple_iid_to_code.pickle)
- [7] <https://www.sri.inf.ethz.ch/py150>
- [8] Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.
- [9] <http://jalammar.github.io/illustrated-transformer/>

## Tasks related papers:

- Polosukhin, Illia, and Alexander Skidanov. "Neural program search: Solving programming tasks from description and examples." *arXiv preprint arXiv:1802.04335* (2018).
- Gu, Xiaodong, et al. "Deep API learning." *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*. ACM, 2016.
- Hu, Xing, et al. "Deep code comment generation." *Proceedings of the 26th Conference on Program Comprehension*. ACM, 2018.
- Iyer, Srinivasan, et al. "Summarizing source code using a neural attention model." Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Vol. 1. 2016.
- Jiang, Siyuan, Ameer Armaly, and Collin McMillan. "Automatically generating commit messages from diffs using neural machine translation." Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering. IEEE Press, 2017.



## Silvia Severini

Technische Universität München  
Faculty of Informatics  
Chair of Software Engineering for  
Business Information Systems

Boltzmannstraße 3  
85748 Garching bei München

Tel +49.89.289. 17132  
Fax +49.89.289.17136

matthes@in.tum.de  
[www.matthes.in.tum.de](http://www.matthes.in.tum.de)

