

# Pre-trained Models for SD Activities



**Dr. Antonio Mastropaolo**

*Instructor*

**Mr. Alvi Haque**



*Teaching Assistant*



**WILLIAM & MARY**

CHARTERED 1693

**Spring 2026**



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities

Transformer Model

Vaswani et al.

---

## Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer



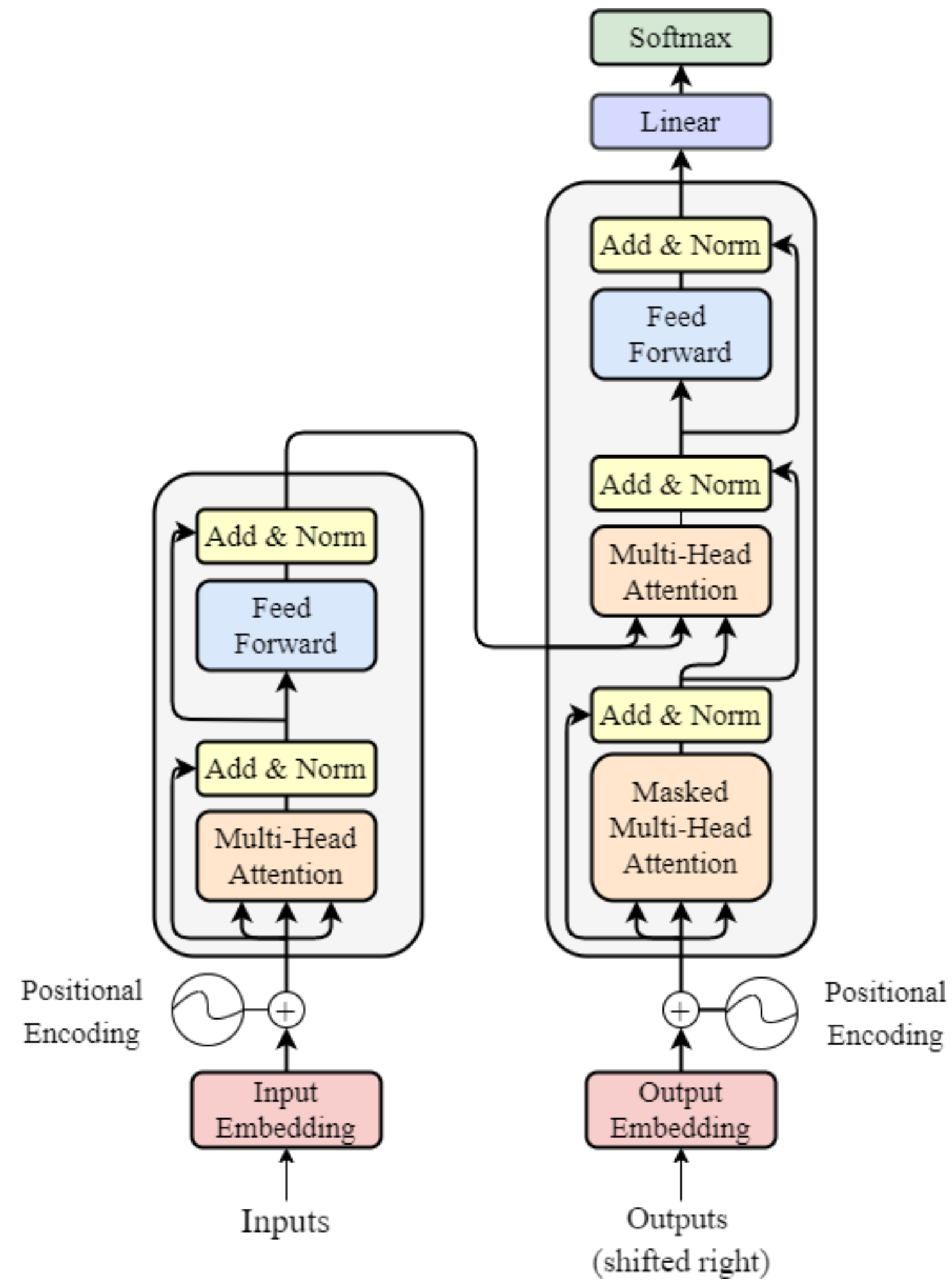
antoniomastropaolo.com



aura-se-lab.github.io

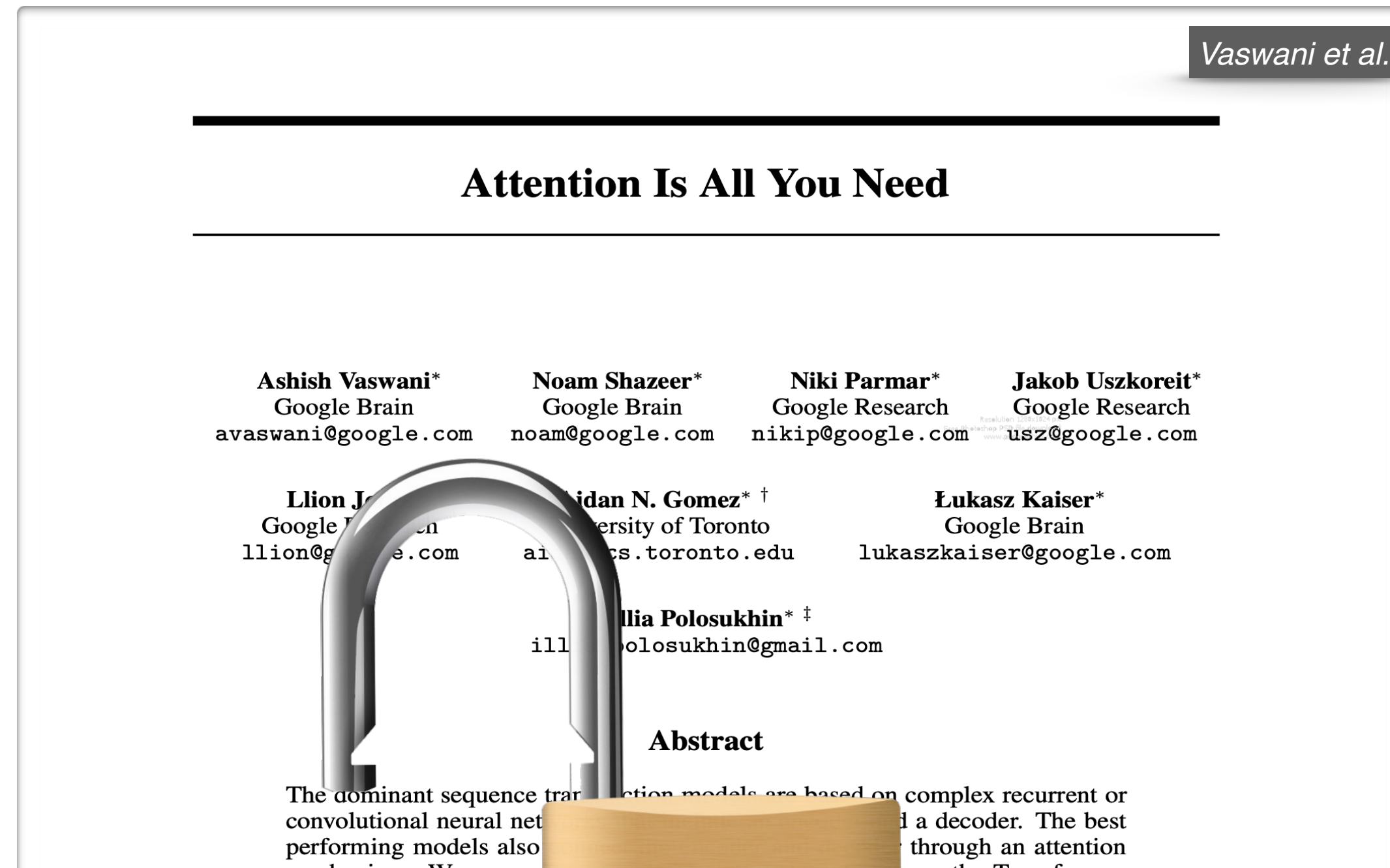


# Pre-trained Models for SD Activities



# Pre-trained Models for SD Activities

Transformer Model



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities

## Long-Range Dependencies

```
public static InetAddress getEmbeddedIPv4ClientAddress
    (InetAddress ip) {
    ip1 = New IP("127.0.0.1")
    if (isCompatIPv4Address(ip)) {
        return getCompatIPv4Address(ip);
    } if (is6to4Address(ip)) {
        return get6to4IPv4Address(ip);
    }
    if (isTeredoAddress(ip)) {
        return getTeredoInfo(ip).getClient();
    }
    return ip1.getDefault();
}
```



# Pre-trained Models for SD Activities

## Long-Range Dependencies

```
public static InetAddress getEmbeddedIPv4ClientAddress
    (Inet6Address ip) {
    ip1 = New IP("127.0.0.1")
    if (isCompatIPv4Address(ip)) {
        return getCompatIPv4Address(ip);
    } if (is6to4Address(ip)) {
        return get6to4IPv4Address(ip);
    }
    if (isTeredoAddress(ip)) {
        return getTeredoInfo(ip).getClient();
    }
    return ip1.getDefault();
}
```



# Pre-trained Models for SD Activities

## Long-Range Dependencies

```
public static InetAddress getEmbeddedIPv4ClientAddress
    (Inet6Address ip) {
    ip1 = New IP("127.0.0.1")
    if (isCompatIPv4Address(ip)) {
        return getCompatIPv4Address(ip);
    } if (is6to4Address(ip)) {
        return get6to4IPv4Address(ip);
    }
    if (isTeredoAddress(ip)) {
        return getTeredoInfo(ip).getClient();
    }
    return ip1.getDefault();
}
```



# Pre-trained Models for SD Activities

Pre-train—then—Fine-tune

**Pre-training:** Instill as much general **knowledge** as possible into the model, such as the **structure** and **patterns** of a language either natural or formal



# Pre-trained Models for SD Activities



Pre-train—then—Fine-tune

**Fine-tuning:** The model is **specialized** for automating a **specific tasks** at hand i.e., sentiment analysis



# Pre-trained Models for SD Activities



A **pre-trained** model is a **deep learning** model that has been trained on a **vast** dataset, allowing the **transfer** of the **knowledge** it has gained to be **re-used** during the **fine-tuning** stage.



# Pre-trained Models for SD Activities

## CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation

Yue Wang<sup>1</sup>, Weishi Wang<sup>12</sup>, Shafiq Joty<sup>12</sup>, and Steven C.H. Hoi<sup>1</sup>  
<sup>1</sup> Salesforce Research Asia  
<sup>2</sup> Nanyang Technological University, Singapore  
{wang.y, weishi.wang, sjoty, shoi}@salesforce.com

## CodeBERT: A Pre-Trained Model for Programming and Natural Languages

Zhangyin Feng<sup>1\*</sup>, Daya Guo<sup>2\*</sup>, Duyu Tang<sup>3</sup>, Nan Duan<sup>3</sup>, Xiaocheng Feng<sup>1</sup>, Ming Gong<sup>4</sup>, Linjun Shou<sup>4</sup>, Bing Qin<sup>1</sup>, Ting Liu<sup>1</sup>, Daxin Jiang<sup>4</sup>, Ming Zhou<sup>3</sup>  
<sup>1</sup> Research Center for Social Computing and Information Retrieval, Harbin Institute of Technology, China  
<sup>2</sup> The School of Data and Computer Science, Sun Yat-sen University, China  
<sup>3</sup> Microsoft Research Asia, Beijing, China  
<sup>4</sup> Microsoft Search Technology Center Asia, Beijing, China  
{yfeng, xcfeng, qinb, tliu}@ir.hit.edu.cn  
{guody5@mail2.sysu.edu.cn, duan, migon, lisho, djiang, mingzhou}@microsoft.com

## GRAPHCODEBERT: PRE-TRAINING CODE REPRESENTATIONS WITH DATA FLOW

Daya Guo<sup>1\*</sup>, Shuo Ren<sup>2\*</sup>, Shuai Lu<sup>3\*</sup>, Zhangyin Feng<sup>4\*</sup>, Duyu Tang<sup>5</sup>, Shujie Liu<sup>5</sup>, Long Zhou<sup>5</sup>, Nan Duan<sup>5</sup>, Alexey Svyatkovskiy<sup>6</sup>, Shengyu Fu<sup>6</sup>, Michele Tufano<sup>6</sup>, Shao Kun Deng<sup>6</sup>, Colin Clement<sup>6</sup>, Dawn Drain<sup>6</sup>, Neel Sundaresan<sup>6</sup>, Jian Yin<sup>1</sup>, Daxin Jiang<sup>7</sup>, and Ming Zhou<sup>5</sup>  
<sup>1</sup>School of Computer Science and Engineering, Sun Yat-sen University.  
<sup>2</sup>Beihang University, <sup>3</sup>Peking University, <sup>4</sup>Harbin Institute of Technology,  
<sup>5</sup>Microsoft Research Asia, <sup>6</sup>Microsoft Devdiv, <sup>7</sup>Microsoft STCA

### ABSTRACT

Pre-trained models for programming language have achieved dramatic empirical improvements on a variety of code-related tasks such as code search, code completion, code summarization, etc. However, existing pre-trained models represent a code snippet as a sequence of tokens, while ignoring the inherent structure. We provide crucial code semantics and would enhance the model's performance. We present GraphCodeBERT, a pre-trained model that considers the inherent structure of code. Instead of representing code as a sequence of tokens, we use a graph-based representation of code like abstract syntax tree (AST), which is a semantic-level structure of code that captures the "value-comes-from" relationship between variables. Such a graph-based representation is more complex and does not bring an unnecessary decrease in efficiency. We use a graph-based attention mechanism of which makes the model more efficient. We use a graph-based attention mechanism on Transformer. In addition to using the task of code completion, we introduce two structure-aware pre-training tasks. One is to align representations between edges, and the other is to align representations between nodes. We implement the model in an efficient way by using a graph-based attention mechanism. We implement the model in an efficient way by using a graph-based attention mechanism. We implement the model in an efficient way by using a graph-based attention mechanism. We implement the model in an efficient way by using a graph-based attention mechanism.

### 1 INTRODUCTION

Pre-trained models such as ELMo (Peters et al., 2018), GPT (Radford et al., 2018) have led to strong improvement on numerous natural language processing tasks. These pre-trained models are first pre-trained on a large unsupervised text

## Multi-task Learning based Pre-trained Language Model for Code Completion

Fang Liu  
Key Lab of High Confidence Software Technology, MoE (Peking University) Beijing, China  
liufang816@pku.edu.cn

Yunfei Zhao  
Key Lab of High Confidence Software Technology, MoE (Peking University) Beijing, China  
zhaoyunfei@pku.edu.cn

Ge Li\*  
Key Lab of High Confidence Software Technology, MoE (Peking University) Beijing, China  
lige@pku.edu.cn

Zhi Jin\*  
Key Lab of High Confidence Software Technology, MoE (Peking University) Beijing, China  
zhijin@pku.edu.cn

### ABSTRACT

Code completion is one of the most useful features in the Integrated Development Environments (IDEs), which can accelerate software development by suggesting the next probable token based on the contextual code in real-time. Recent studies have shown that statistical language modeling techniques can improve the performance of code completion tools through learning from large-scale software repositories. However, these models suffer from two major drawbacks: a) Existing research uses static embeddings, which map a word to the same vector regardless of its context. The differences in the meaning of a token in varying contexts are lost when each token is associated with a single representation; b) Existing language model based code completion models perform poor on completing identifiers, and the type information of the identifiers is ignored.

### CCS CONCEPTS

• Computing methodologies → Artificial intelligence; • Software and its engineering → Software maintenance tools.  
code completion, multi-task learning, pre-trained language model, transformer networks

### ACM Reference Format:

Fang Liu, Ge Li, Yunfei Zhao, and Zhi Jin. 2020. Multi-task Learning based Pre-trained Language Model for Code Completion. In *35th IEEE/ACM International Conference on Automated Software Engineering (ASE '20)*, September 21–25, 2020, Virtual Event, Australia. ACM, New York, NY, USA, 1–11. <https://doi.org/10.1145/3616433.3616439>

### ABSTRACT

CodeBERT, a *bimodal* pre-trained model for programming language (PL) and natural language (NL). CodeBERT learns contextual representations that support NLP applications such as natural language search, code documentation, etc. We develop CodeBERT with a hybrid objective function that supports the pre-training task of code completion, which is to detect masked words from an artificially masked input sequence. The success of pre-trained models in NLP also drives a surge of multi-modal pre-trained models, such as ViLBERT (Lu et al., 2019) for language-image and VideoBERT (Sun et al., 2019) for language-video, which are learned from *bimodal* data such as language-image pairs with *bimodal* self-supervised objectives. We evaluate CodeBERT on various applications by fine-tuning model

and RoBERTa (Liu et al., 2019) have dramatically improved the state-of-the-art on a variety of natural language processing (NLP) tasks. These pre-trained models learn effective contextual representations from massive unlabeled text optimized by self-supervised objectives, such as masked language modeling, which predicts the original masked word from an artificially masked input sequence. The success of pre-trained models in NLP also drives a surge of multi-modal pre-trained models, such as ViLBERT (Lu et al., 2019) for language-image and VideoBERT (Sun et al., 2019) for language-video, which are learned from *bimodal* data such as language-image pairs with *bimodal* self-supervised objectives.

In this work, we present CodeBERT, a *bimodal* pre-trained model for natural language (NL) and programming language (PL) like Python, Java

## CCT5: A Code-Change-Oriented Pre-trained Model

Bo Lin\*  
linbo19@mudt.edu.cn  
College of Computer Science, National University of Defense Technology Changsha, China

Yepang Liu  
liuyep1@sustech.edu.cn  
Department of Computer Science and Engineering, Southern University of Science and Technology Shenzhen, China

Shangwen Wang\*  
wangshangwen13@mudt.edu.cn  
College of Computer Science, National University of Defense Technology Changsha, China

Xin Xia  
xin.xia@acm.org  
Zhejiang University Hangzhou, China

Zhongxin Liu  
liu\_zx@zju.edu.cn  
Zhejiang University Hangzhou, China

Xiaoguang Mao\*  
xgmao@mudt.edu.cn  
College of Computer Science, National University of Defense Technology Changsha, China

code changes and two tasks specific to the code review process. Results show that CCT5 outperforms both conventional deep learning approaches and existing pre-trained models on these tasks.

### CCS CONCEPTS

• Software and its engineering → Software maintenance tools; Maintaining software; Software evolution.

### KEYWORDS

Code Change, Pre-Training, Deep Learning.  
ACM Reference Format:  
Bo Lin, Shangwen Wang, Zhongxin Liu, Yepang Liu, Xin Xia, and Xiaoguang Mao. 2023. CCT5: A Code-Change-Oriented Pre-trained Model. In *Proceedings of the 31st ACM Joint European Software Engineering (ESEC/FSE) Symposium on the Foundations of Software Engineering (FSE/FOSE)*, September 3–9, 2023, San Francisco, CA, USA. ACM, New York, NY, USA, 1–11. <https://doi.org/10.1145/3616433.3616439>

### 1 INTRODUCTION

Software undergoes continuous changes during the phase to fix defects, change execution logic, make more efficient, or introduce new features. Developers are often faced with the task of identifying and fixing these changes. This process is often tedious and error-prone. In this paper, we propose a new pre-trained model, CCT5, which is designed for code review tasks. CCT5 is trained on a large amount of code changes and two tasks specific to the code review process. Results show that CCT5 outperforms both conventional deep learning approaches and existing pre-trained models on these tasks.



# Pre-trained Models for SD Activities

## TranslationEnglish2Thai

### Fine-tuning Dataset



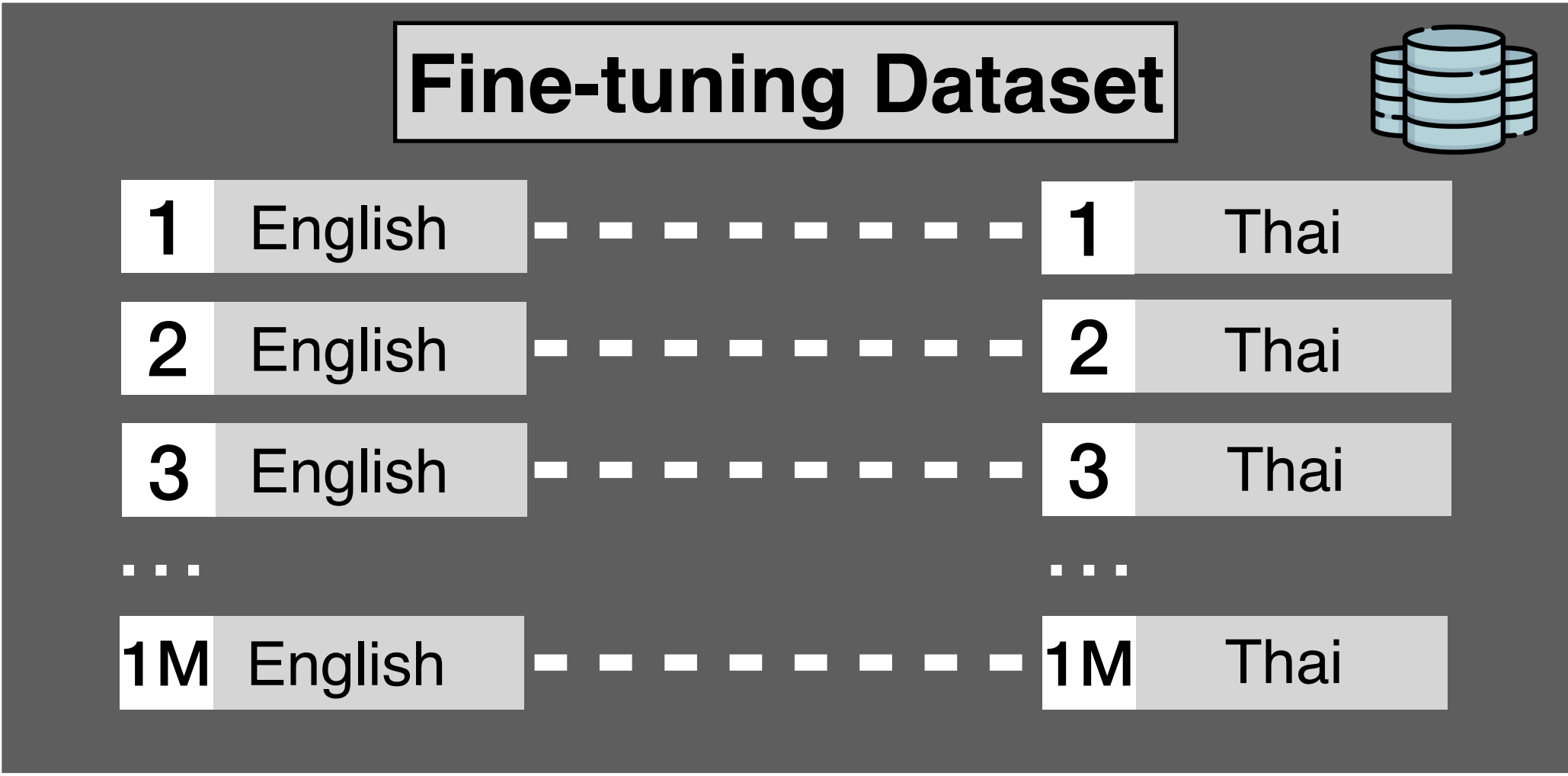
<b>1</b> English	-----	<b>1</b> Thai
<b>2</b> English	-----	<b>2</b> Thai
<b>3</b> English	-----	<b>3</b> Thai
...		...
<b>1M</b> English	-----	<b>1M</b> Thai



# Pre-trained Models for SD Activities

Supervised Dataset

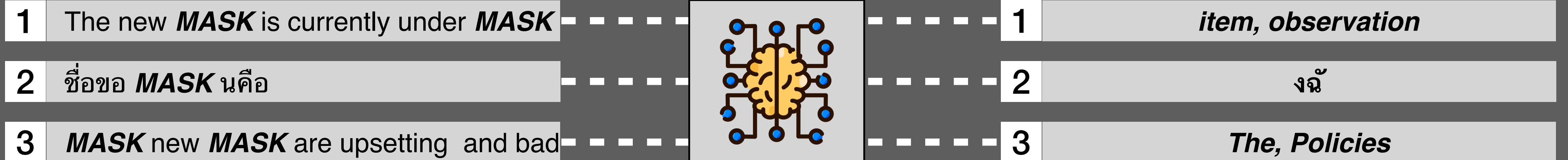
## TranslationEnglish2Thai



# Pre-trained Models for SD Activities

Mask Language Modeling

## Pre-Training Procedure



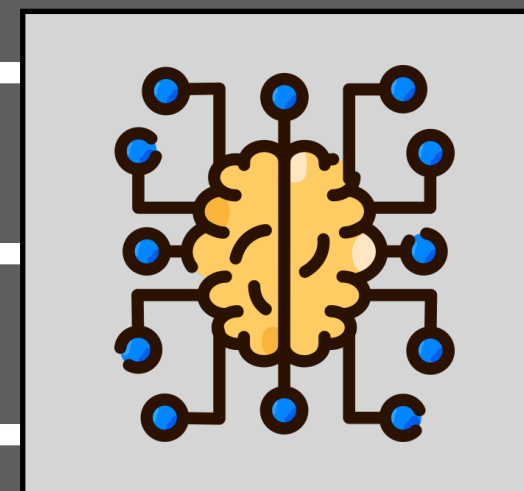
# Pre-trained Models for SD Activities

Mask Language Modeling

## Pre-Training Procedure

Self-supervised Dataset

- 1 The new **MASK** is currently under **MASK**
- 2 ชื่อขอ **MASK** นคือ
- 3 **MASK** new **MASK** are upsetting and bad



- 1 *item, observation*
- 2 ังฉั
- 3 *The, Policies*

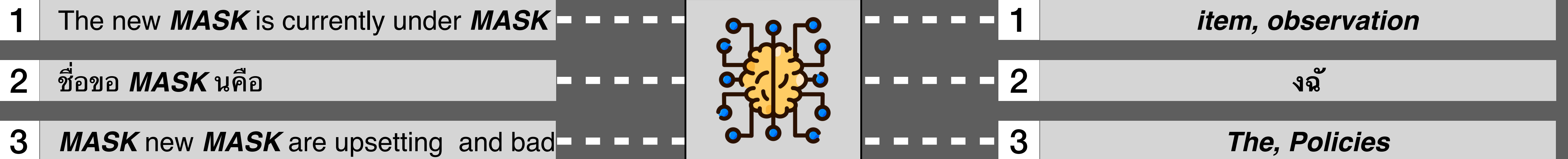


# Pre-trained Models for SD Activities

Mask Language Modeling

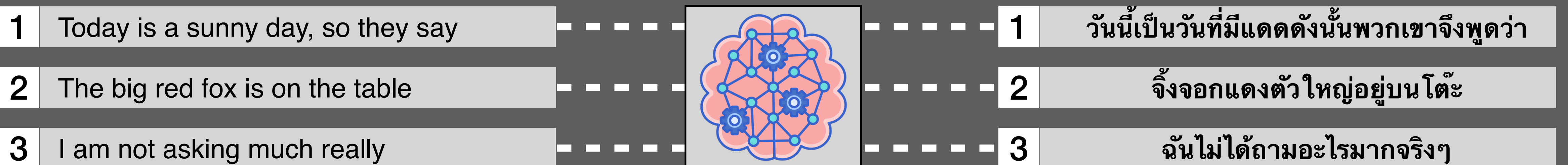
## Pre-Training Procedure

Self-supervised Dataset



## Fine-tuning Procedure

Supervised Dataset



# Pre-trained Models for SD Activities

2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)

## Leveraging Small Software Engineering Data Sets with Pre-trained Neural Networks

Romain Robbes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

Andrea Janes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

**Abstract**—Many software engineering data sets, particularly those that demand manual labelling for classification, are necessarily small. As a consequence, several recent software engineering papers have cast doubt on the effectiveness of deep neural networks for classification tasks, when applied to these data sets. We provide initial evidence that recent advances in Natural Language Processing, that allow neural networks to leverage large amount of unlabelled data in a pre-training phase, can significantly improve performance.

**Index Terms**—Data Sets, Deep Learning, Transfer Learning

### I. INTRODUCTION

Deep neural networks have revolutionized the field of machine learning, leading to significant improvements over the state of the art in a large variety of tasks and domains (including image recognition, speech recognition, machine translation, summarization, and others). In Software Engineering, the results have been more mixed. While many approaches have promising results (e.g., code completion, variable naming, code translation and summarization), several recent papers have cast doubt about the effectiveness of neural networks in some settings [1], [2]. One reason is the performance on small SE data sets, which we examine in this paper.

The use cases where deep neural networks outperform the state of the art are those where large amount of data is available, such as the ImageNet data set [3], which features more than 14 million labelled images in its largest version. However, many data sets are much smaller, particularly in the case of supervised learning, as each data point must be manually labelled. This manual task is often an expensive process, and all the more in Software Engineering: while ImageNet labelling can be crowdsourced, labelling a SE data set requires significant expertise. It cannot be crowdsourced easily, and is thus vastly more expensive. In practice, many SE data sets are manually labelled by researchers, and are thus limited to hundreds or thousands of data points only. Recent work reports disappointing performance when training deep neural networks from scratch on small SE data sets.

In this paper, we investigate whether the small data set issue can be alleviated in Software Engineering, starting with natural language (NL) SE datasets (we are actively exploring the source code case). We show that a form of transfer learning, namely unsupervised pre-training on large data sets, can be effectively used to improve the performance on small NL SE data sets. Several approaches—detailed in Section III—have

shown this strategy to be effective in Computer Vision and (lately) in Natural Language Processing. These approaches are all a variation of a common idea: when training a neural network from scratch, all the model's weights are initialized randomly. Pre-training instead initializes the weights with values obtained by training a model on a related tasks, for which data is easier to gather; in this case, by using large amounts of unlabelled data, no manual labelling is necessary.

### II. SMALL DATA SETS IN SOFTWARE ENGINEERING

Many labelling tasks, such as identifying objects in images, drawing bounding boxes for these objects, or identifying the sentiment of a movie review, do not require specific expertise. They can hence be crowdsourced at a reasonable cost, particularly thanks to crowdsourcing platforms such as Amazon Mechanical Turk.

As manual labelling in Software Engineering requires expertise in Software Engineering, it is usually done by the researchers themselves. This approach is expensive, and due to the amount of effort involved, is often limited to hundreds or thousands of data points. This is all the more true since the effort is often duplicated to verify whether annotators are in agreement with each other. Some examples follow.

Bacchelli et al. manually linked 2,139 emails from 6 open-source projects with the source code entities they referenced [4]. Fakhoury et al. labelled 1,700 Java methods (with comments) according to the linguistic antipatterns they exhibited [2]. Ortu et al. labelled 2,000 JIRA issue reports, and 4,000 sentences with the sentiments expressed in them. Novielli et al. labelled 4,423 stack Overflow posts with sentiment [5]. Lin et al. labelled the sentiment of *each word* of 1,500 Stack Overflow sentences according to their sentiment [1]; they also labelled 341 app reviews for sentiment. Villaroel et al. [6] classified 1,200 app reviews, while Maalej et al. [7] classified 4,400 app reviews, in both cases in categories such as bug reports, and feature requests. Zhou et al. hired students to analyze 1,674 API elements to determine whether they have API directive defects [8]. These data sets are several orders of magnitude smaller than the ones where Deep Learning approaches are the state of the art, such as ImageNet.

In two cases, some of these data sets have been used as input for deep learning techniques, with lackluster results. Fakhoury et al. applied a Convolutional Neural Network (CNN) for linguistic antipattern detection; they find that the CNN matches

978-1-7281-1758-4/19/\$31.00 ©2019 IEEE  
DOI 10.1109/ICSE-NIER.2019.00016

29

ICSE'19 — NIER

Authorized licensed use limited to: Universita della Svizzera Italiana. Downloaded on December 13, 2023 at 05:58:04 UTC from IEEE Xplore. Restrictions apply.



antoniomastropaolo.com



aura-se-lab.github.io



# Pre-trained Models for SD Activities

2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)

## Leveraging Small Software Engineering Data Sets with Pre-trained Neural Networks

Romain Robbes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

Andrea Janes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

**Abstract**—Many software engineering data sets, particularly those that demand manual labelling for classification, are necessarily small. As a consequence, several recent software engineering papers have cast doubt on the effectiveness of deep neural networks for classification tasks, when applied to these data sets. We provide initial evidence that recent advances in Natural Language Processing, that allow neural networks to leverage large amount of unlabelled data in a pre-training phase, can significantly improve performance.

**Index Terms**—Data Sets, Deep Learning, Transfer Learning

### I. INTRODUCTION

Deep neural networks have revolutionized the field of machine learning, leading to significant improvements over the state of the art in a large variety of tasks and domains (including image recognition, speech recognition, machine translation, summarization, and others). In Software Engineering, the results have been more mixed. While many approaches have promising results (e.g., code completion, variable naming, code translation and summarization), several recent papers have cast doubt about the effectiveness of neural networks in some settings [1], [2]. One reason is the performance on small SE data sets, which we examine in this paper.

The use cases where deep neural networks outperform the state of the art are those where large amount of data is available, such as the ImageNet data set [3], which features more than 14 million labelled images in its largest version. However, many data sets are much smaller, particularly in the case of supervised learning, as each data point must be manually labelled. This manual task is often an expensive process, and all the more in Software Engineering: while ImageNet labelling can be crowdsourced, labelling a SE data set requires significant expertise. It cannot be crowdsourced easily, and is thus vastly more expensive. In practice, many SE data sets are manually labelled by researchers, and are thus limited to hundreds or thousands of data points only. Recent work reports disappointing performance when training deep neural networks from scratch on small SE data sets.

In this paper, we investigate whether the small data set issue can be alleviated in Software Engineering, starting with natural language (NL) SE datasets (we are actively exploring the source code case). We show that a form of transfer learning, namely unsupervised pre-training on large data sets, can be effectively used to improve the performance on small NL SE data sets. Several approaches—detailed in Section III—have

shown this strategy to be effective in Computer Vision and (lately) in Natural Language Processing. These approaches are all a variation of a common idea: when training a neural network from scratch, all the model's weights are initialized randomly. Pre-training instead initializes the weights with values obtained by training a model on a related tasks, for which data is easier to gather; in this case, by using large amounts of unlabelled data, no manual labelling is necessary.

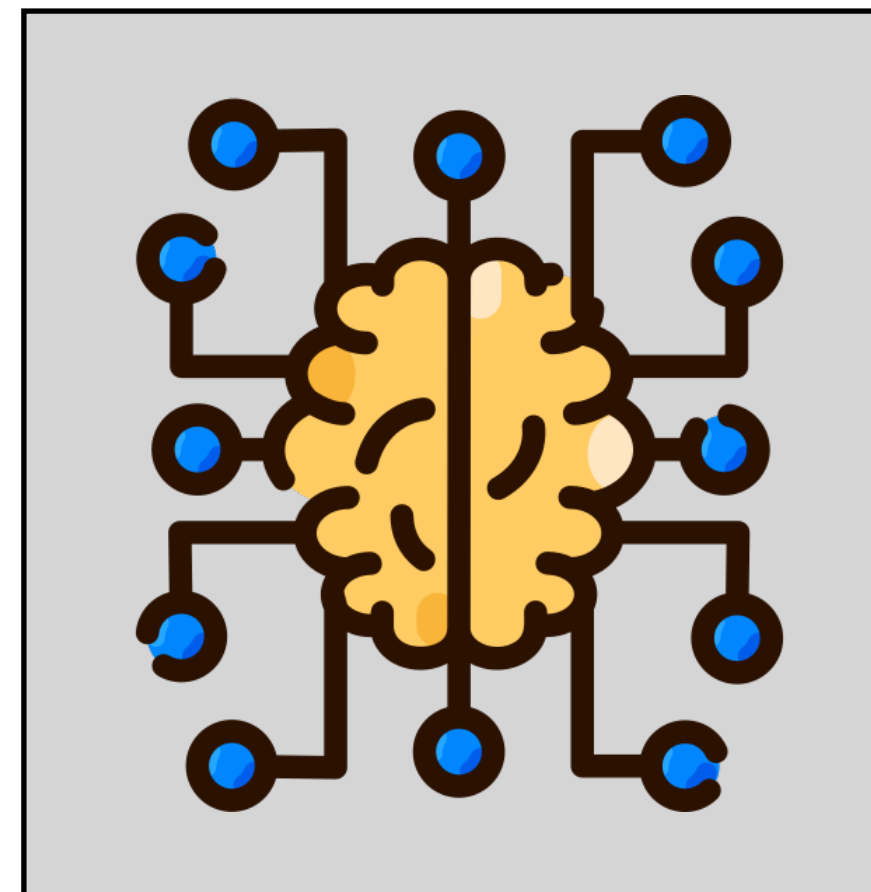
### II. SMALL DATA SETS IN SOFTWARE ENGINEERING

Many labelling tasks, such as identifying objects in images, drawing bounding boxes for these objects, or identifying the sentiment of a movie review, do not require specific expertise. They can hence be crowdsourced at a reasonable cost, particularly thanks to crowdsourcing platforms such as Amazon Mechanical Turk.

As manual labelling in Software Engineering requires expertise in Software Engineering, it is usually done by the researchers themselves. This approach is expensive, and due to the amount of effort involved, is often limited to hundreds or thousands of data points. This is all the more true since the effort is often duplicated to verify whether annotators are in agreement with each other. Some examples follow.

Bacchelli et al. manually linked 2,139 emails from 6 open-source projects with the source code entities they referenced [4]. Fakhoury et al. labelled 1,700 Java methods (with comments) according to the linguistic antipatterns they exhibited [2]. Ortu et al. labelled 2,000 JIRA issue reports, and 4,000 sentences with the sentiments expressed in them. Novielli et al. labelled 4,423 stack Overflow posts with sentiment [5]. Lin et al. labelled the sentiment of *each word* of 1,500 Stack Overflow sentences according to their sentiment [1]; they also labelled 341 app reviews for sentiment. Villaroel et al. [6] classified 1,200 app reviews, while Maalej et al. [7] classified 4,400 app reviews, in both cases in categories such as bug reports, and feature requests. Zhou et al. hired students to analyze 1,674 API elements to determine whether they have API directive defects [8]. These data sets are several orders of magnitude smaller than the ones where Deep Learning approaches are the state of the art, such as ImageNet.

In two cases, some of these data sets have been used as input for deep learning techniques, with lackluster results. Fakhoury et al. applied a Convolutional Neural Network (CNN) for linguistic antipattern detection; they find that the CNN matches



978-1-7281-1758-4/19/\$31.00 ©2019 IEEE  
DOI 10.1109/ICSE-NIER.2019.00016

29

ICSE'19 — NIER

Authorized licensed use limited to: Universita della Svizzera Italiana. Downloaded on December 13, 2023 at 05:58:04 UTC from IEEE Xplore. Restrictions apply.



antoniomastropaolo.com



aura-se-lab.github.io



# Pre-trained Models for SD Activities

2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)

## Leveraging Small Software Engineering Data Sets with Pre-trained Neural Networks

Romain Robbes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

Andrea Janes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

**Abstract**—Many software engineering data sets, particularly those that demand manual labelling for classification, are necessarily small. As a consequence, several recent software engineering papers have cast doubt on the effectiveness of deep neural networks for classification tasks, when applied to these data sets. We provide initial evidence that recent advances in Natural Language Processing, that allow neural networks to leverage large amount of unlabelled data in a pre-training phase, can significantly improve performance.

**Index Terms**—Data Sets, Deep Learning, Transfer Learning

### I. INTRODUCTION

Deep neural networks have revolutionized the field of machine learning, leading to significant improvements over the state of the art in a large variety of tasks and domains (including image recognition, speech recognition, machine translation, summarization, and others). In Software Engineering, the results have been more mixed. While many approaches have promising results (e.g., code completion, variable naming, code translation and summarization), several recent papers have cast doubt about the effectiveness of neural networks in some settings [1], [2]. One reason is the performance on small SE data sets, which we examine in this paper.

The use cases where deep neural networks outperform the state of the art are those where large amount of data is available, such as the ImageNet data set [3], which features more than 14 million labelled images in its largest version. However, many data sets are much smaller, particularly in the case of supervised learning, as each data point must be manually labelled. This manual task is often an expensive process, and all the more in Software Engineering: while ImageNet labelling can be crowdsourced, labelling a SE data set requires significant expertise. It cannot be crowdsourced easily, and is thus vastly more expensive. In practice, many SE data sets are manually labelled by researchers, and are thus limited to hundreds or thousands of data points only. Recent work reports disappointing performance when training deep neural networks from scratch on small SE data sets.

In this paper, we investigate whether the small data set issue can be alleviated in Software Engineering, starting with natural language (NL) SE datasets (we are actively exploring the source code case). We show that a form of transfer learning, namely unsupervised pre-training on large data sets, can be effectively used to improve the performance on small NL SE data sets. Several approaches—detailed in Section III—have

shown this strategy to be effective in Computer Vision and (lately) in Natural Language Processing. These approaches are all a variation of a common idea: when training a neural network from scratch, all the model's weights are initialized randomly. Pre-training instead initializes the weights with values obtained by training a model on a related tasks, for which data is easier to gather; in this case, by using large amounts of unlabelled data, no manual labelling is necessary.

### II. SMALL DATA SETS IN SOFTWARE ENGINEERING

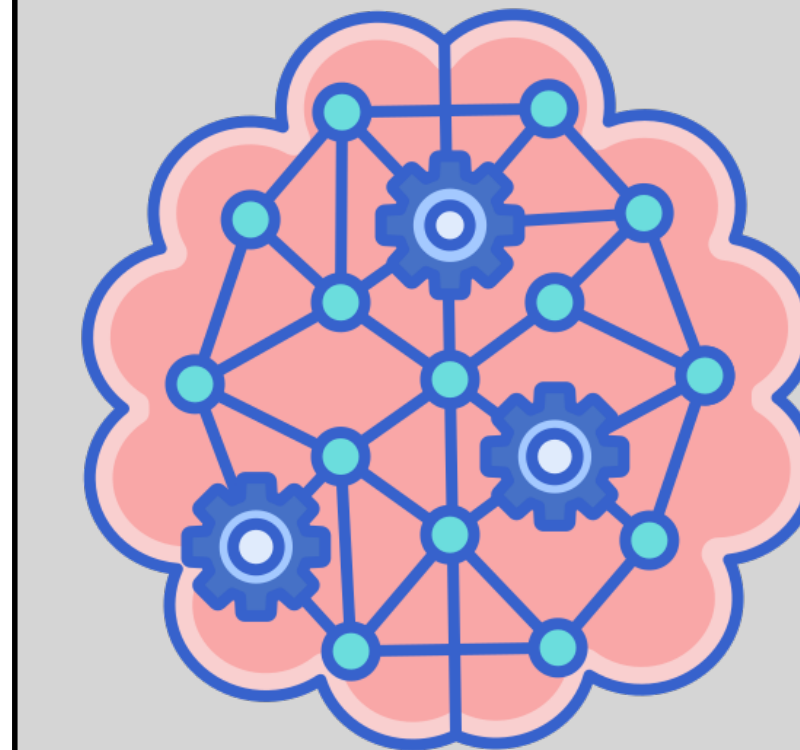
Many labelling tasks, such as identifying objects in images, drawing bounding boxes for these objects, or identifying the sentiment of a movie review, do not require specific expertise. They can hence be crowdsourced at a reasonable cost, particularly thanks to crowdsourcing platforms such as Amazon Mechanical Turk.

As manual labelling in Software Engineering requires expertise in Software Engineering, it is usually done by the researchers themselves. This approach is expensive, and due to the amount of effort involved, is often limited to hundreds or thousands of data points. This is all the more true since the effort is often duplicated to verify whether annotators are in agreement with each other. Some examples follow.

Bacchelli et al. manually linked 2,139 emails from 6 open-source projects with the source code entities they referenced [4]. Fakhoury et al. labelled 1,700 Java methods (with comments) according to the linguistic antipatterns they exhibited [2]. Ortu et al. labelled 2,000 JIRA issue reports, and 4,000 sentences with the sentiments expressed in them. Novielli et al. labelled 4,423 stack Overflow posts with sentiment [5]. Lin et al. labelled the sentiment of *each word* of 1,500 Stack Overflow sentences according to their sentiment [1]; they also labelled 341 app reviews for sentiment. Villaroel et al. [6] classified 1,200 app reviews, while Maalej et al. [7] classified 4,400 app reviews, in both cases in categories such as bug reports, and feature requests. Zhou et al. hired students to analyze 1,674 API elements to determine whether they have API directive defects [8]. These data sets are several orders of magnitude smaller than the ones where Deep Learning approaches are the state of the art, such as ImageNet.

In two cases, some of these data sets have been used as input for deep learning techniques, with lackluster results. Fakhoury et al. applied a Convolutional Neural Network (CNN) for linguistic antipattern detection; they find that the CNN matches

## Pre-trained Model



978-1-7281-1758-4/19/\$31.00 ©2019 IEEE  
DOI 10.1109/ICSE-NIER.2019.00016

29

ICSE'19 — NIER

Authorized licensed use limited to: Universita della Svizzera Italiana. Downloaded on December 13, 2023 at 05:58:04 UTC from IEEE Xplore. Restrictions apply.



antoniomastropaolo.com



aura-se-lab.github.io



# Pre-trained Models for SD Activities

2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)

## Leveraging Small Software Engineering Data Sets with Pre-trained Neural Networks

Romain Robbes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

Andrea Janes  
Faculty of Computer Science  
Free University of Bozen-Bolzano, Italy

**Abstract**—Many software engineering data sets, particularly those that demand manual labelling for classification, are necessarily small. As a consequence, several recent software engineering papers have cast doubt on the effectiveness of deep neural networks for classification tasks, when applied to these data sets. We provide initial evidence that recent advances in Natural Language Processing, that allow neural networks to leverage large amount of unlabelled data in a pre-training phase, can significantly improve performance.

**Index Terms**—Data Sets, Deep Learning, Transfer Learning

### I. INTRODUCTION

Deep neural networks have revolutionized the field of machine learning, leading to significant improvements over the state of the art in a large variety of tasks and domains (including image recognition, speech recognition, machine translation, summarization, and others). In Software Engineering, the results have been more mixed. While many approaches have promising results (e.g., code completion, variable naming, code translation and summarization), several recent papers have cast doubt about the effectiveness of neural networks in some settings [1], [2]. One reason is the performance on small SE data sets, which we examine in this paper.

The use cases where deep neural networks outperform the state of the art are those where large amount of data is available, such as the ImageNet data set [3], which features more than 14 million labelled images in its largest version. However, many data sets are much smaller, particularly in the case of supervised learning, as each data point must be manually labelled. This manual task is often an expensive process, and all the more in Software Engineering: while ImageNet labelling can be crowdsourced, labelling a SE data set requires significant expertise. It cannot be crowdsourced easily, and is thus vastly more expensive. In practice, many SE data sets are manually labelled by researchers, and are thus limited to hundreds or thousands of data points only. Recent work reports disappointing performance when training deep neural networks from scratch on small SE data sets.

In this paper, we investigate whether the small data set issue can be alleviated in Software Engineering, starting with natural language (NL) SE datasets (we are actively exploring the source code case). We show that a form of transfer learning, namely unsupervised pre-training on large data sets, can be effectively used to improve the performance on small NL SE data sets. Several approaches—detailed in Section III—have

shown this strategy to be effective in Computer Vision and (lately) in Natural Language Processing. These approaches are all a variation of a common idea: when training a neural network from scratch, all the model's weights are initialized randomly. Pre-training instead initializes the weights with values obtained by training a model on a related tasks, for which data is easier to gather; in this case, by using large amounts of unlabelled data, no manual labelling is necessary.

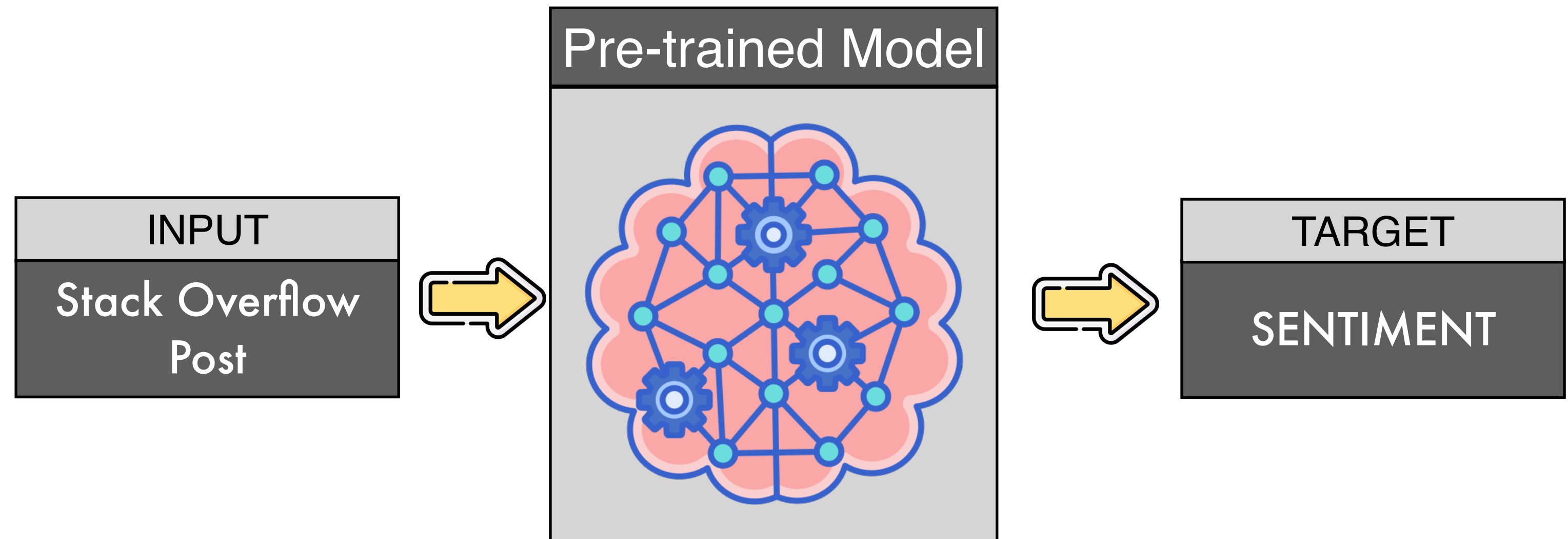
### II. SMALL DATA SETS IN SOFTWARE ENGINEERING

Many labelling tasks, such as identifying objects in images, drawing bounding boxes for these objects, or identifying the sentiment of a movie review, do not require specific expertise. They can hence be crowdsourced at a reasonable cost, particularly thanks to crowdsourcing platforms such as Amazon Mechanical Turk.

As manual labelling in Software Engineering requires expertise in Software Engineering, it is usually done by the researchers themselves. This approach is expensive, and due to the amount of effort involved, is often limited to hundreds or thousands of data points. This is all the more true since the effort is often duplicated to verify whether annotators are in agreement with each other. Some examples follow.

Bacchelli et al. manually linked 2,139 emails from 6 open-source projects with the source code entities they referenced [4]. Fakhoury et al. labelled 1,700 Java methods (with comments) according to the linguistic antipatterns they exhibited [2]. Ortu et al. labelled 2,000 JIRA issue reports, and 4,000 sentences with the sentiments expressed in them. Novielli et al. labelled 4,423 stack Overflow posts with sentiment [5]. Lin et al. labelled the sentiment of *each word* of 1,500 Stack Overflow sentences according to their sentiment [1]; they also labelled 341 app reviews for sentiment. Villaroel et al. [6] classified 1,200 app reviews, while Maalej et al. [7] classified 4,400 app reviews, in both cases in categories such as bug reports, and feature requests. Zhou et al. hired students to analyze 1,674 API elements to determine whether they have API directive defects [8]. These data sets are several orders of magnitude smaller than the ones where Deep Learning approaches are the state of the art, such as ImageNet.

In two cases, some of these data sets have been used as input for deep learning techniques, with lackluster results. Fakhoury et al. applied a Convolutional Neural Network (CNN) for linguistic antipattern detection; they find that the CNN matches



978-1-7281-1758-4/19/\$31.00 ©2019 IEEE  
DOI 10.1109/ICSE-NIER.2019.00016

29

ICSE'19 — NIER

Authorized licensed use limited to: Universita della Svizzera Italiana. Downloaded on December 13, 2023 at 05:58:04 UTC from IEEE Xplore. Restrictions apply.



antoniomastropaolo.com



aura-se-lab.github.io





## Pre-training objectives

*objectives are the specific goals or tasks that a model is designed to solve during the pre-training phase.*

- MLM (**M**asking **L**anguage **M**odeling)
- NSP (**N**ext **S**entence **P**rediction)
- CLM (**C**ausal **L**anguage **M**odeling)
- PLM (**P**ermuted **L**anguage **M**odeling)

# Pre-trained Models for SD Activities



## Pre-training objectives

### MLM (**M**asking **L**anguage **M**odeling)

The model learns the **meaning** and **structure** of words within a sentence by looking at the **full context**.





## Pre-training objectives

### MLM (**M**asking **L**anguage **M**odeling)

The model learns the **meaning** and **structure** of words within a sentence by looking at the **full context**.

The cat is on the table



## Pre-training objectives

### MLM (**M**asking **L**anguage **M**odeling)

The model learns the **meaning** and **structure** of words within a sentence by looking at the **full context**.

The cat **<MASK>** on the table



## Pre-training objectives

### NSP (**N**ext **S**entence **P**rediction)

The model learns how to **connect sentences** and whether there is a **logical progression** between them



## Pre-training objectives

### NSP (**N**ext **S**entence **P**rediction)

The model learns how to **connect sentences** and whether there is a **logical progression** between them

**S1:** The cat is on the table

**S2:** and is asleep



## Pre-training objectives

### NSP (**N**ext **S**entence **P**rediction)

The model learns how to **connect sentences** and whether there is a **logical progression** between them

**S1:** The cat is on the table

**S2:** and is asleep

OUTPUT: 1



## Pre-training objectives

### NSP (**N**ext **S**entence **P**rediction)

The model learns how to **connect sentences** and whether there is a **logical progression** between them

**S1:** and is asleep

**S2:** The cat is on the table

OUTPUT:0

# Pre-trained Models for SD Activities



## Pre-training objectives

### CLM (**C**ausal **L**anguage **M**odeling)

The model is trained to predict the **next word** based only on the words preceding the current one.



# Pre-trained Models for SD Activities



## Pre-training objectives

### CLM (**C**ausal **L**anguage **M**odeling)

The model is trained to predict the **next word** based only on the words preceding the current one.

Does it ring a bell?



# Pre-trained Models for SD Activities



Pre-training objectives

CLM (**C**ausal **L**anguage **M**odeling)

VS

MLM (**M**asking **L**anguage **M**odeling)



# Pre-trained Models for SD Activities

## Pre-training objectives

PLM (**P**ermuted **L**anguage **M**odeling)

The model predicts tokens in a **random, permuted** order that a fixed left-to-right or even unidirectional context.

The cat is on the table and is asleep



# Pre-trained Models for SD Activities

## Pre-training objectives

### PLM (**P**ermuted **L**anguage **M**odeling)

The model predicts tokens in a **random, permuted** order that a fixed left-to-right or even unidirectional context.

The cat is on the table and is asleep

Permuted tokens = [3, 2, 4, 1, 9, 6, 8, 7, 5]

# Pre-trained Models for SD Activities

## Pre-training objectives

### PLM (**P**ermuted **L**anguage **M**odeling)

The model predicts tokens in a **random, permuted** order that a fixed left-to-right or even unidirectional context.

is cat on The asleep is and table the

Permuted tokens = [3, 2, 4, 1, 9, 6, 8, 7, 5]



# Pre-trained Models for SD Activities

UNIDIRECTIONAL

## Pre-training objectives

### PLM (**P**ermuted **L**anguage **M**odeling)

The model predicts tokens in a **random, permuted** order that a fixed left-to-right or even unidirectional context.

is cat on The asleep table is and asleep the

Permuted tokens = [3, 2, 4, 1, 9, 6, 8, 7, 5]



# Pre-trained Models for SD Activities



```
public void bubbleSort(int[] array) {
    boolean sorted = false;
    int temp;
    while (!sorted) {
        sorted = true;
        for (int i = 0; i < array.length - 1; i++) {
            if (a[i] > a[i+1]) {
                temp = a[i];
                a[i] = a[i+1];
                a[i+1] = temp;
                sorted = false;
            }
        }
    }
}
```

## Pre-training



# Pre-trained Models for SD Activities



```
public void bubbleSort(int[] <MASK>) {
    boolean sorted = false;
    int <MASK>;
    while (!sorted) {
        <MASK> = true;
        for (int i = 0; i < <MASK>.length - 1; i++) {
            if (a[i] > a[<MASK> + 1]) {
                <MASK> = a[i];
                a[<MASK>] = a[i <MASK> 1];
                a[i+1] = temp;
                <MASK> = false;
            }
        }
    }
}
<MASK>
```

## Pre-training



# Pre-trained Models for SD Activities



```
public void bubbleSort(int[] array) {
    boolean sorted = false;
    int temp;
    while (!sorted) {
        sorted = true;
        for (int i = 0; i < array.length - 1; i++) {
            if (a[i] > a[i+1]) {
                temp = a[i];
                a[i] = a[i+1];
                a[i+1] = temp;
                sorted = false;
            }
        }
    }
}
```

## Pre-training



# Pre-trained Models for SD Activities



```
public void bubbleSort(int[] array) {  
    boolean sorted = false;  
    int temp;  
    while (!sorted) {  
        sorted = true;  
        for (int i = 1; i < array.length - 1; i++) {  
            if (a[i] > a[i+1]) {  
                temp = a[i];  
                a[i] = a[i+1];  
                a[i+1] = temp;  
                sorted = false;  
            }  
        }  
    }  
}
```

BUG FIXING

## Fine-tuning



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities



```
public void bubbleSort(int[] array) {
    boolean sorted = false;
    int temp;
    while (!sorted) {
        sorted = true;
        for (int i = 0; i < array.length - 1; i++) {
            if (a[i] > a[i+1]) {
                temp = a[i];
                a[i] = a[i+1];
                a[i+1] = temp;
                sorted = false;
            }
        }
    }
}
```

BUG FIXING

## Fine-tuning



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities

## Multi-objective Pre-training

MLM

NSP

### BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova  
Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

#### Abstract

We introduce a new language representation model called **BERT**, which stands for **Bidirectional Encoder Representations from Transformers**. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a re-

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal



# Pre-trained Models for SD Activities

## Transfer Learning

### 1. Pre-train—then—Fine-tune



# Pre-trained Models for SD Activities

## Transfer Learning

1. Pre-train—then—Fine-tune
2. **Multi-task Pre-training**



# Pre-trained Models for SD Activities

## Multi-task Pre-training in the SE domain: A case study

### Using Deep Learning to Generate Complete Log Statements

Antonio Mastropaolo  
antonio.mastropaolo@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

Luca Pascarella  
luca.pascarella@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

Gabriele Bavota  
gabriele.bavota@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

#### ABSTRACT

Logging is a practice widely adopted in several phases of the software lifecycle. For example, during software development log statements allow engineers to verify and debug the system by exposing fine-grained information of the running software. While the benefits of logging are undisputed, taking proper decisions about *where* to inject log statements, *what* information to log, and at which *log level* (e.g., error, warning) is crucial for the logging effectiveness. In this paper, we present LANCE (Log stAtemeNt reCommEnder), the first approach supporting developers in all these decisions. LANCE features a Text-To-Text-Transfer-Transformer (T5) model that has been trained on 6,894,456 Java methods. LANCE takes as input a Java method and injects in it a full log statement, including a human-comprehensible logging message and properly choosing the needed log level and the statement location. Our results show that LANCE is able to (i) properly identify the location in the code where to inject the statement in 65.9% of Java methods requiring it; (ii) selecting the proper log level in 66.2% of cases; and (iii) generate

Although technically possible, logging everything (e.g., every exception) is inefficient and impracticable [53]. On the one hand, a coarse-grained logging risks hiding runtime failures, missing log messages useful for diagnoses [45].

On the other hand, a fine-grained logging risks increasing the overhead of log management and analysis [26]. To optimize the quantity and quality of data generated, developers insert log statements in strategic positions, specify appropriate log levels (e.g., error, debug, info), and define compact but comprehensible text messages. Nonetheless, it remains a non-trivial task for developers to decide where, what, and at which level to log [26, 46].

For these reasons, researchers have proposed techniques to support developers in deciding what parts of the system to log [45], the log level for logging statements [26, 27, 36, 46], and the structure of log messages [28]. For example, Jia *et al.* [23] proposed an approach based on association rules to place error logs after code branches such as *if* statements. Li *et al.* [25] studied the use of topic modeling for log placement at method-level. Also, two recent works tackled challenges related to log statements writing based on



# Pre-trained Models for SD Activities

**Logging**, is a widespread activity among developers since **log statements** are used to **record** valuable runtime information about applications.



# Pre-trained Models for SD Activities

**Logging**, is a widespread activity among developers since **log statements** are used to **record** valuable runtime information about applications.

**Log statements** represent the mechanism through which we perform logging. Each logging statement features three different **log components** a **log level**, a **log messages** and the **position** in the underlying code snippet



# Pre-trained Models for SD Activities

```
public void stopContext () {  
    if ( _context != null && _publisher != null ) {  
        try {  
            _publisher.destroy(_context);  
        }  
        catch (Exception e) {  
            logger.err("Cannot destroy the context")  
        }  
    }  
}
```

# Pre-trained Models for SD Activities



## End-to-End Support for Logging Activities

### Using Deep Learning to Generate Complete Log Statements

Antonio Mastropaolo  
antonio.mastropaolo@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

Luca Pascarella  
luca.pascarella@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

Gabriele Bavota  
gabriele.bavota@usi.ch  
SEART @ Software Institute  
Università della Svizzera italiana  
Switzerland

#### ABSTRACT

Logging is a practice widely adopted in several phases of the software lifecycle. For example, during software development log statements allow engineers to verify and debug the system by exposing fine-grained information of the running software. While the benefits of logging are undisputed, taking proper decisions about *where* to inject log statements, *what* information to log, and at which *log level* (e.g., error, warning) is crucial for the logging effectiveness. In this paper, we present LANCE (Log stAtemeNt reCommEnder), the first approach supporting developers in all these decisions. LANCE features a Text-To-Text-Transfer-Transformer (T5) model that has been trained on 6,894,456 Java methods. LANCE takes as input a Java method and injects in it a full log statement, including a human-comprehensible logging message and properly choosing the needed log level and the statement location. Our results show that LANCE is able to (i) properly identify the location in the code where to inject the statement in 65.9% of Java methods requiring it; (ii) selecting the proper log level in 66.2% of cases; and (iii) generate

Although technically possible, logging everything (e.g., every exception) is inefficient and impracticable [53]. On the one hand, a coarse-grained logging risks hiding runtime failures, missing log messages useful for diagnoses [45].

On the other hand, a fine-grained logging risks increasing the overhead of log management and analysis [26]. To optimize the quantity and quality of data generated, developers insert log statements in strategic positions, specify appropriate log levels (e.g., error, debug, info), and define compact but comprehensible text messages. Nonetheless, it remains a non-trivial task for developers to decide where, what, and at which level to log [26, 46].

For these reasons, researchers have proposed techniques to support developers in deciding what parts of the system to log [45], the log level for logging statements [26, 27, 36, 46], and the structure of log messages [28]. For example, Jia *et al.* [23] proposed an approach based on association rules to place error logs after code branches such as *if* statements. Li *et al.* [25] studied the use of topic modeling for log placement at method-level. Also, two recent works



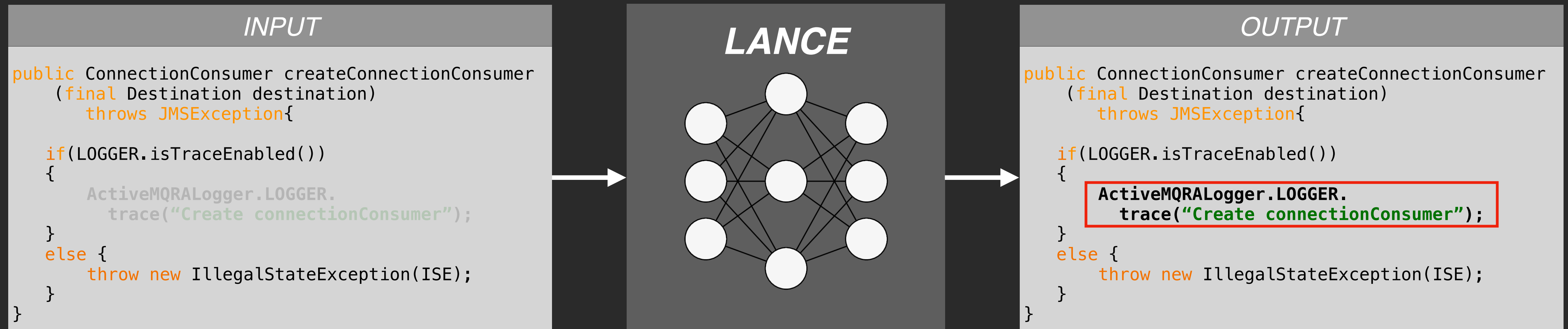
antoniomastropaolo.com



aura-se-lab.github.io



# Pre-trained Models for SD Activities





## Training of LANCE

1. Multiple Pre-training Tasks/Objective
2. Injection of Log Statements

# Pre-trained Models for SD Activities



## Input

1

```
protected boolean isBoundaryEvent(Object obj) {  
    if (obj instanceof <MASK>) {  
        mxCell cell = (mxCell) <MASK>;  
        <MASK> cell.getId().  
            startsWith("boundary-event-");  
    }  
    return false;  
}
```

2

```
protected Polyline(<MASK> dx,  
                  double dy, Polyline next) {  
    this.dx = dx;  
    this.<MASK> = dy;  
    <MASK>.next = next;  
}
```

...

## Target

1

- 1) mxCell
- 2) obj
- 3) return

2

- 1) double
- 2) dy
- 3) this

...

~6.7M

7



# Pre-trained Models for SD Activities



## Input

1

```
@Override public void
actionPerformed(ActionEvent ev){
    boolean ok =
        toolsDialog.showDialog(config, true);
    if(!ok){
        return;
    }
}
```

2

```
public void stopContext (){
    if ( _context != null
        && _publisher != null ){
        try {
            _publisher.destroy(_context);
        }
        catch (Exception e) {}
    }
}
```

...

## Target

1

```
@Override public void
actionPerformed(ActionEvent ev){
    <LOG_STMT>
    boolean ok =
        toolsDialog.showDialog(config, true);
    if(!ok){
        return;
    }
}
```

2

```
public void stopContext (){
    if ( _context != null
        && _publisher != null ){
        try {
            _publisher.destroy(_context);
        }
        catch (Exception e) { <LOG_STMT> }
    }
}
```

...

~133K

2



# Pre-trained Models for SD Activities



## Input

1

```
public LocalDateTimeToStringConverter  
    <MASK>() {  
    return new LocalDateTimeToStringConverter <MASK>  
}  
-----  
private void infolog(String msg){  
    processingLog += msg + ".\n";  
}
```

Masking

Log-Pos

2

```
public String getTableName() <MASK>  
    return <MASK>;  
}  
-----  
public void stopCounting(){  
    runCounting = false;  
    stopLogReport();  
}
```

Masking

Log-Pos

...

## Target

1

```
1) LocalDateTimeToStringConverter  
2) ();  
-----  
private void infolog(String msg){  
    <LOG_STMT>  
    processingLog += msg + ".\n";  
}
```

Masking

Log-Pos

2

```
1) {  
2) tableName  
-----  
public void stopCounting(){  
    runCounting = false;  
    <LOG_STMT>  
    stopLogReport();  
}
```

Masking

Log-Pos

...

~6.8M

3



# Pre-trained Models for SD Activities



## Input

## Fine-tuning

## Target

1

```
public ConnectionConsumer createConnectionConsumer(final
Destination destination) throws JMSEException {
    if (ActiveMQRALogger.LOGGER.isTraceEnabled()) {
        ActiveMQRALogger.LOGGER.trace(
            "Create connectionConsumer");
    } else {
        logger.info("Throwing Exception")
        throw new IllegalStateException(ISE);
    }
}
```

1

```
public ConnectionConsumer createConnectionConsumer(final
Destination destination) throws JMSEException {
    if (ActiveMQRALogger.LOGGER.isTraceEnabled()) {
        ActiveMQRALogger.LOGGER.trace(
            "Create connectionConsumer");
    } else {
        logger.info("Throwing Exception")
        throw new IllegalStateException(ISE);
    }
}
```

2

```
public ConnectionConsumer createConnectionConsumer(final
Destination destination) throws JMSEException {
    if (ActiveMQRALogger.LOGGER.isTraceEnabled()) {
        ActiveMQRALogger.LOGGER.trace(
            "Create connectionConsumer");
    } else {
        logger.info("Throwing Exception")
        throw new IllegalStateException(ISE);
    }
}
```

2

```
public ConnectionConsumer createConnectionConsumer(final
Destination destination) throws JMSEException {
    if (ActiveMQRALogger.LOGGER.isTraceEnabled()) {
        ActiveMQRALogger.LOGGER.trace(
            "Create connectionConsumer");
    } else {
        logger.info("Throwing Exception")
        throw new IllegalStateException(ISE);
    }
}
```

...

...

~100K



# Pre-trained Models for SD Activities

```
public HandlerResult handle(ProcessState state, ProcessInstance process) {
    Secret secret = (Secret) state.getResource();
    String secretValue = secret.getValue();
    if (StringUtils.isNotBlank(secretValue)) {
        try {
            secretsService.delete(secret.getAccountId(), secret.getValue());
        } catch (IOException e) {
            log.error("Error deleting secret {}: {}", secret.getId(), e.getMessage());
            // LANCE log

            log.error("Failed to delete secret from storage [{}]", secret.getId(), e);
            // Target log

            [...]
        }
    }
    return null;
}
```



# Pre-trained Models for SD Activities

```
private synchronized CswSubscription deleteCswSubscription (String subscriptionId) throws CswException {  
    [...]  
    try {  
        ServiceRegistration sr = registeredSubscriptions.remove(subscriptionId);  
        [...]  
    } catch (Exception e) {  
        LOGGER.debug("Exception trying to delete subscription's configuration for subscription ID {}",  
                    logSanitizedId, e); LANCE log  
        LOGGER.debug("Could not delete subscription for {}", logSanitizedId, e); Target log  
    }  
}
```

# Pre-trained Models for SD Activities

## Transfer Learning

1. Pre-train—then—Fine-tune
2. Multi-task Pre-training



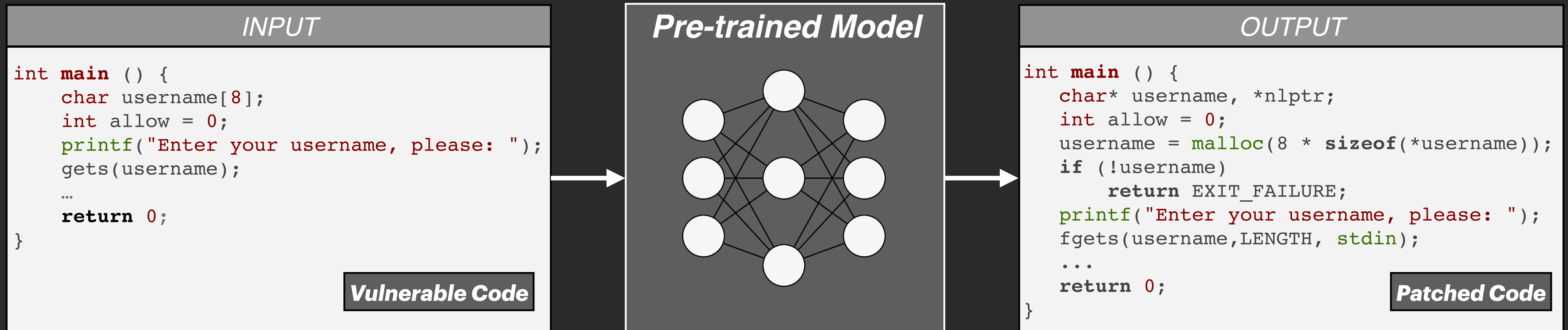
# Pre-trained Models for SD Activities

## Transfer Learning

1. Pre-train—then—Fine-tune
2. Multi-task Pre-training
3. Chaining Tasks



# Pre-trained Models for SD Activities



## DL-based Vulnerability Patching

# Pre-trained Models for SD Activities



Data Scarcity



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities

**Bugs** are a “superclass” of **vulnerabilities**

We have plenty of data in the form of **bug-fixing** commits



# Pre-trained Models for SD Activities

**Bugs** are a “superclass” of **vulnerabilities**

We have plenty of data in the form of **bug-fixing** commits

SUPERVISED PRE-TRAINING



# Pre-trained Models for SD Activities

## Input

1

```
int main () {  
  char username[8];  
  int allow = 0;  
  printf("Enter your username, please: ");  
  gets(username);  
  ...  
  return 0;  
}
```

2

```
int main(int argc, char *argv[]) {  
  char buffer[10];  
  strcpy(buffer, argv[1]);  
  printf("Buffer contents: %s\n", buffer);  
  return 0;  
}
```

...

## Target

1

```
int main () {  
  char* username, *nlptr;  
  int allow = 0;  
  username = malloc(8 * sizeof(*username));  
  fgets(username, LENGTH, stdin);  
  ...  
  return 0;  
}
```

2

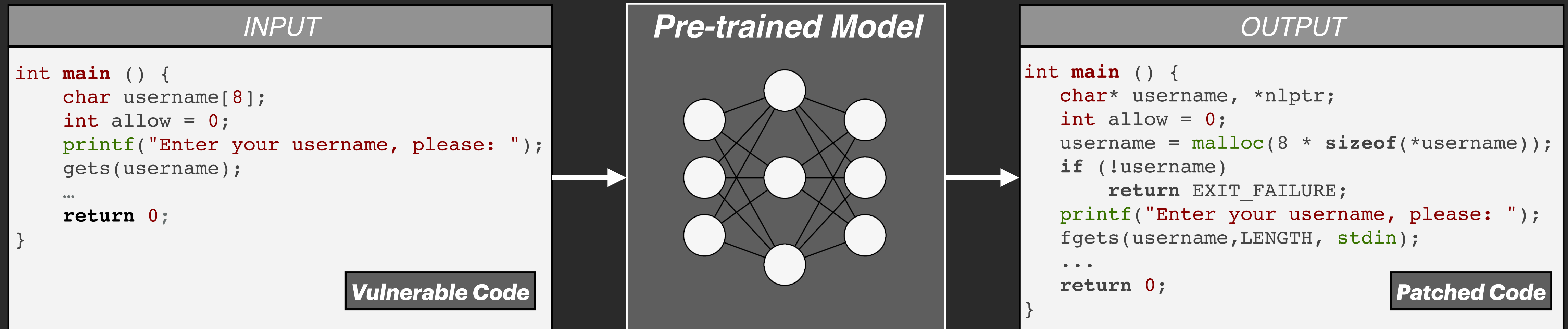
```
int main(int argc, char *argv[]) {  
  char buffer[10];  
  strncpy(buffer, argv[1], sizeof(buffer) - 1);  
  buffer[sizeof(buffer) - 1] = '\0';  
  printf("Buffer contents: %s\n", buffer);  
  return 0;  
}
```

...

~100K



# Pre-trained Models for SD Activities



## DL-based Vulnerability Patching

# Pre-trained Models for SD Activities

## Transfer Learning

1. Pre-train—then—Fine-tune
2. Multi-task Pre-training
3. Chaining Tasks



# Pre-trained Models for SD Activities

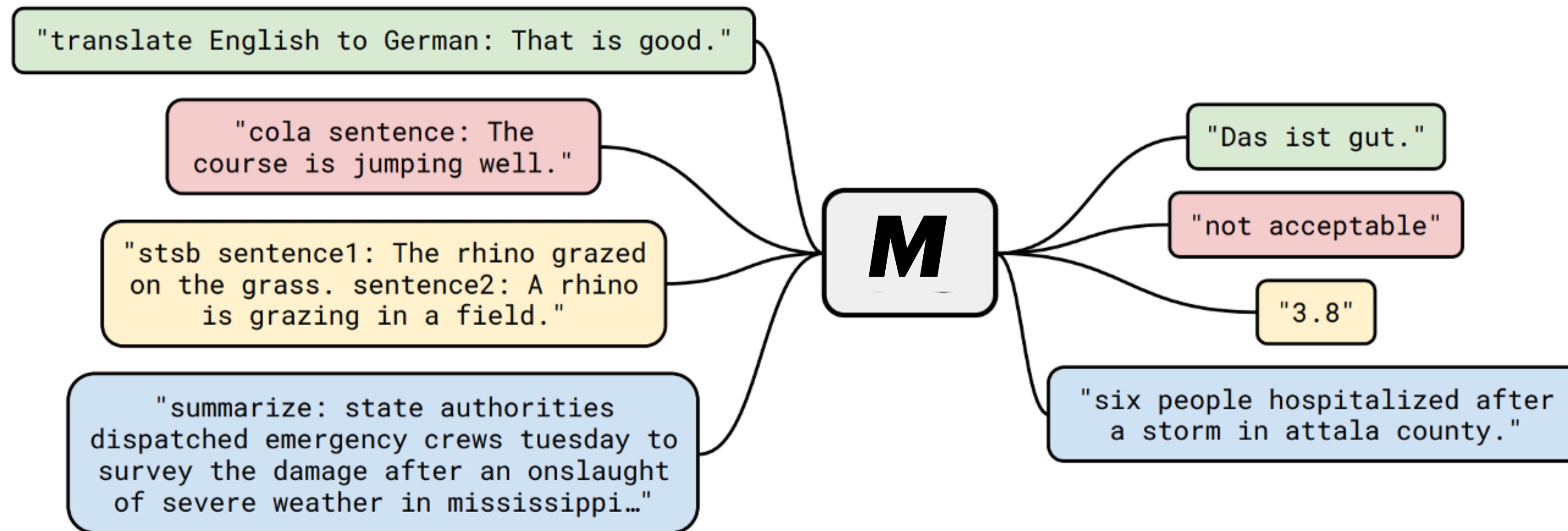
## Transfer Learning

1. Pre-train—then—Fine-tune
2. Multi-task Pre-training
3. Chaining Tasks
4. **Multi-task Fine-tuning**





## Multi-task Fine-tuning



# Pre-trained Models for SD Activities

## Multi-task Fine-tuning

### Using Transfer Learning for Code-Related Tasks

Antonio Mastropaolo<sup>1</sup>, Nathan Cooper<sup>1</sup>, David Nader Palacio, *Student Member, IEEE*,  
Simone Scalabrino<sup>1</sup>, Denys Poshyvanyk<sup>1</sup>, *Senior Member, IEEE*, Rocco Oliveto<sup>1</sup>, and Gabriele Bavota<sup>1</sup>

**Abstract**—Deep learning (DL) techniques have been used to support several code-related tasks such as code summarization and bug-fixing. In particular, pre-trained transformer models are on the rise, also thanks to the excellent results they achieved in Natural Language Processing (NLP) tasks. The basic idea behind these models is to first pre-train them on a generic dataset using a self-supervised task (e.g., filling masked words in sentences). Then, these models are fine-tuned to support specific tasks of interest (e.g., language translation). A single model can be fine-tuned to support multiple tasks, possibly exploiting the benefits of *transfer learning*. This means that knowledge acquired to solve a specific task (e.g., language translation) can be useful to boost performance on another task (e.g., sentiment classification). While the benefits of transfer learning have been widely studied in NLP, limited empirical evidence is available when it comes to code-related tasks. In this paper, we assess the performance of the Text-To-Text Transfer Transformer (T5) model in supporting four different code-related tasks: (i) automatic bug-fixing, (ii) injection of code mutants, (iii) generation of assert statements, and (iv) code summarization. We pay particular attention in studying the role played by pre-training and multi-task fine-tuning on the model's performance. We show that (i) the T5 can achieve better performance as compared to state-of-the-art baselines; and (ii) while pre-training helps the model, not all tasks benefit from a multi-task fine-tuning.

**Index Terms**—Deep learning, empirical software engineering

#### 1 INTRODUCTION

SEVERAL code-related tasks have been recently automated by researchers exploiting Deep Learning (DL) techniques [81]. Several of these works customize DL models proposed in the Natural Language Processing (NLP) field to support code-related tasks, and most of them share one common characteristic: *They shape the problem at hand as a text-to-text transformation, in which the input and the output of the model are text strings.* For instance, Tufano *et al.* [78] used

automating bug fixing [15], [25], [48], [75], learning generic code changes [73], supporting code migration [52], [53], code summarization [24], [32], [39], [42], code reviews [77], [78], pseudo-code generation [55], code deobfuscation [31], [79], injection of code mutants [76], generation of assert statements [82], clone detection [74], [83], traceability [49] and code completion [5], [11], [17], [17], [34], [35], [70], [84].

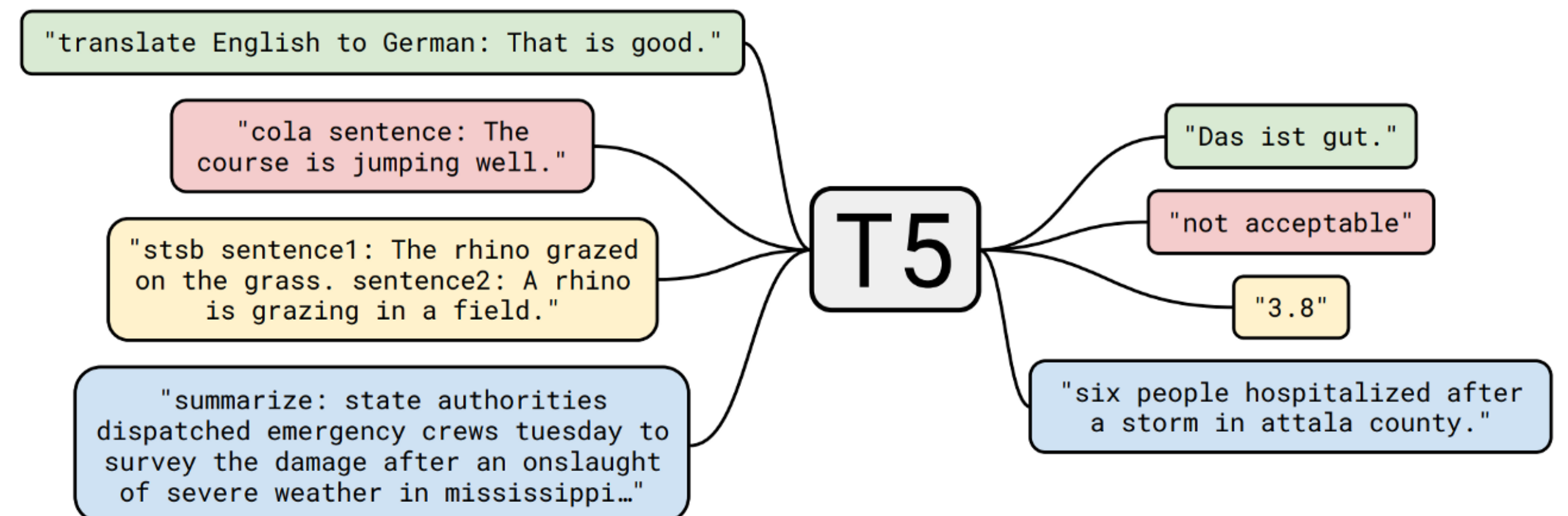
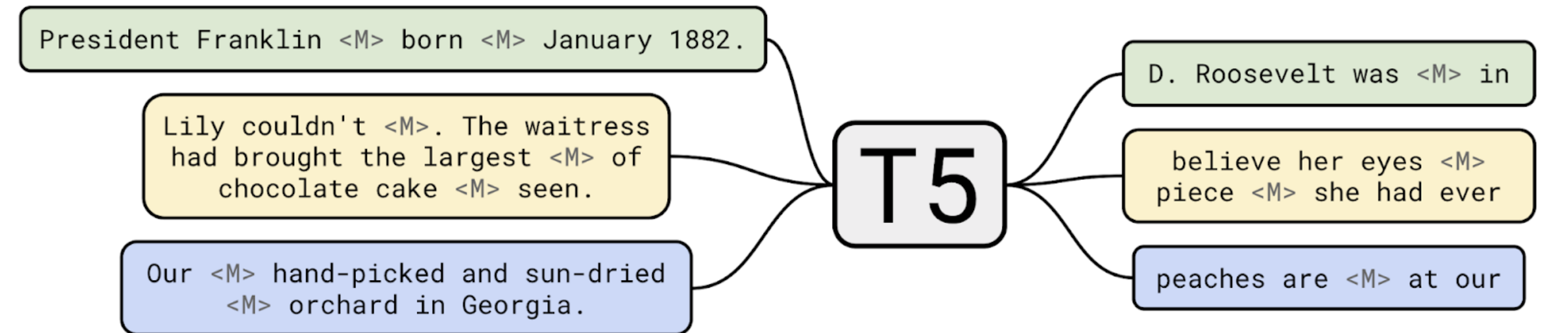
Recent years have seen the rise of *transfer learning* in the field of natural language processing. The basic idea is to first



# Pre-trained Models for SD Activities

## Multi-task Fine-tuning

Text as **Input** and  
Text as **Output**



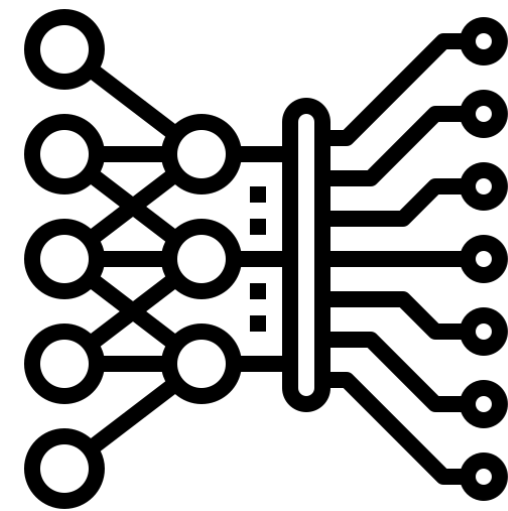
# Pre-trained Models for SD Activities

1



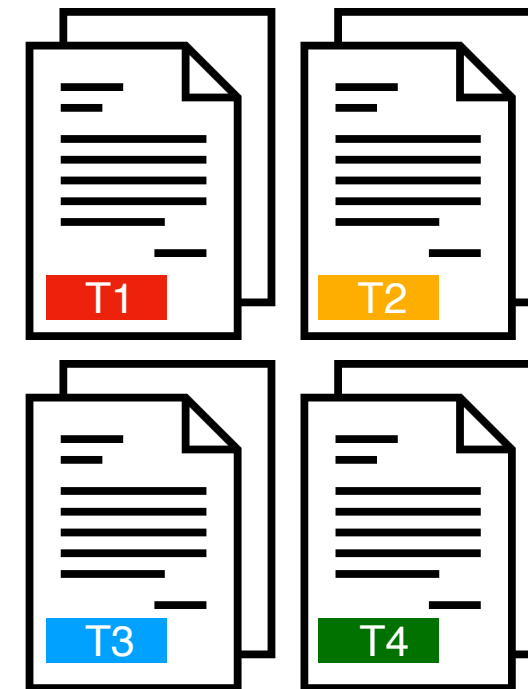
Pre-training dataset

2



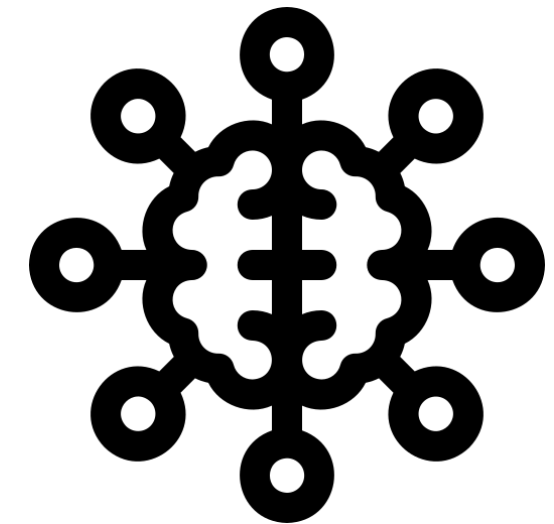
T5 Pre-training

3



Text-to-text task specific dataset

4



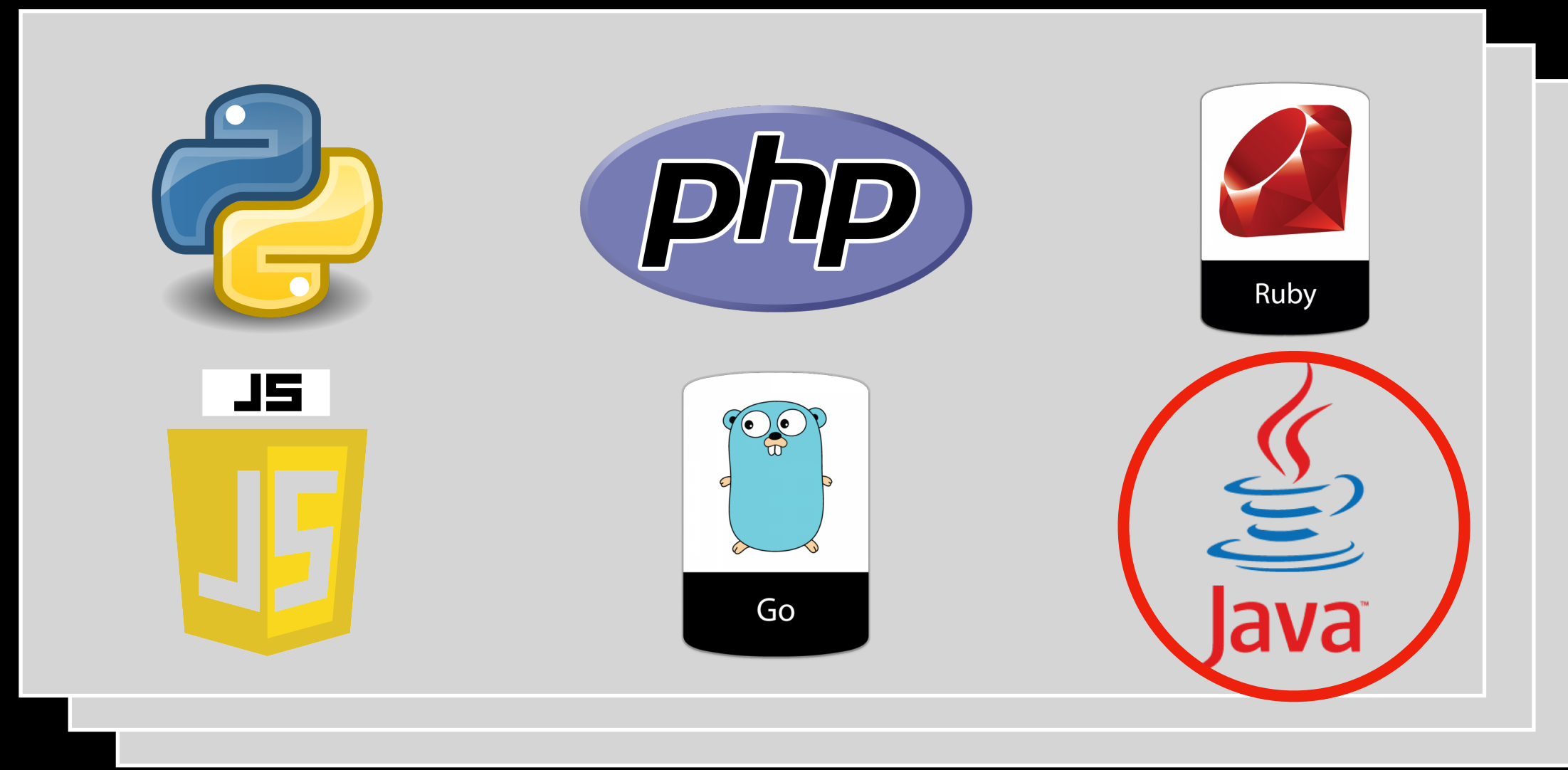
T5 Fine-tuning on tasks-mixture

# Pre-trained Models for SD Activities

1



Pre-training dataset



**CodeSearchNet**  
*Husain et al.*



[antoniomastropaolo.com](http://antoniomastropaolo.com)

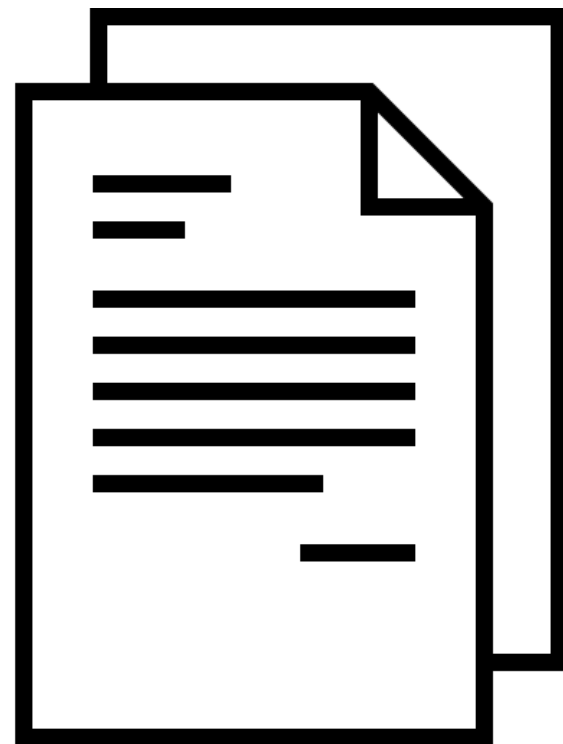


[aura-se-lab.github.io](https://github.com/aura-se-lab)



# Pre-trained Models for SD Activities

1



Pre-training  
dataset

*Creates an UnicastProcessor with the given internal buffer capacity hint*



```
@CheckReturnValue @NonNull
public static <T> UnicastProcessor <T>
create (int capacityHint ) {
    return new UnicastProcessor
    <T>(capacityHint);
}
```

*Creates a text encryptor that uses stronger password - based encryption. Encrypted text is hex - encoded*



```
public static TextEncryptor
delux(CharSequence password,
CharSequence salt){
    return new HexEncodingTextEncryptor
    (stronger(password, salt));
}
```

*Creates a connection to the consumer*



```
public ConnectionConsumer createConnectionConsumer
(final Destination destination) throws JMSEException
    if (ActiveMQRALogger.LOGGER.isTraceEnabled()) {
        ActiveMQRALogger.LOGGER.trace(
            "Create connectionConsumer");
    }
    else {
        logger.info("Throwing Exception")
        throw new IllegalStateException(ISE);
    }
}
```

*Assembles the Distinguished Name that should be used the given username*



```
public DistinguishedName buildDn(
    String username){
    DistinguishedName dn = new
    DistinguishedName (userDnBase);
    dn.add(usernameAttribute, username);
    return dn;
}
```



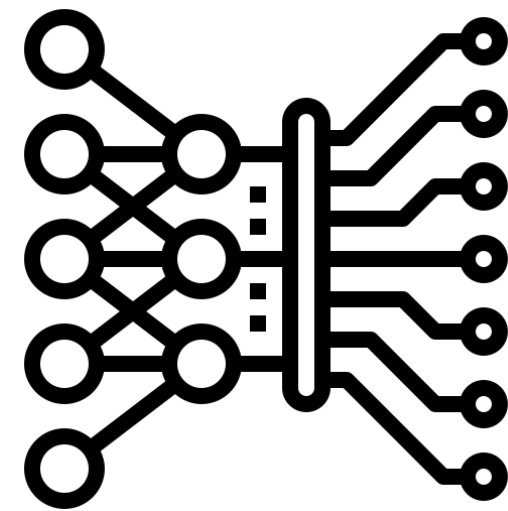
# Pre-trained Models for SD Activities

1



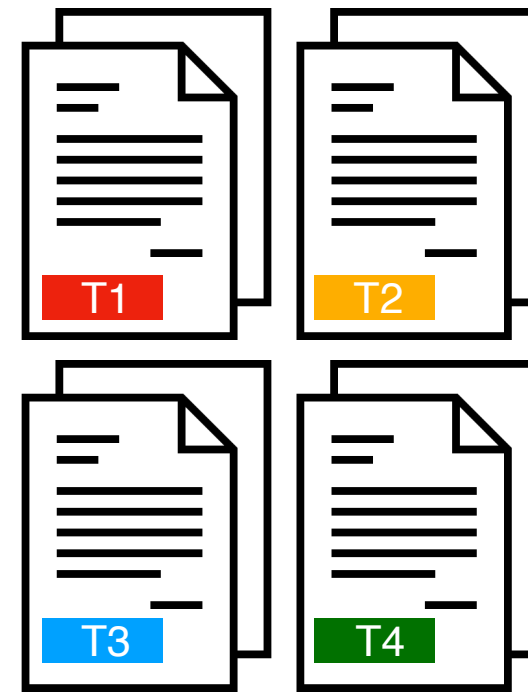
Pre-training  
dataset

2



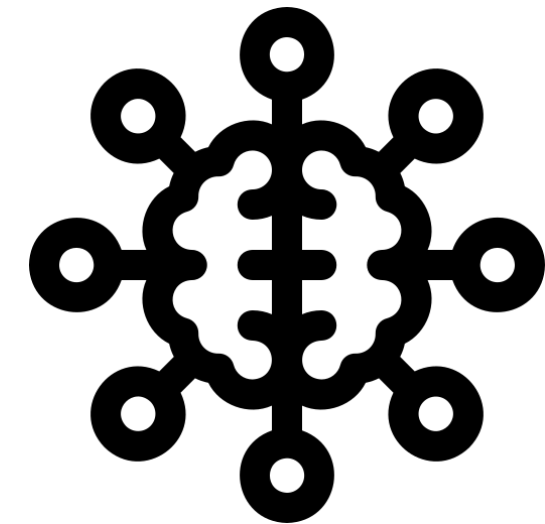
T5 Pre-training

3



Text-to-text task  
specific dataset

4



T5 Fine-tuning on  
tasks-mixture

# Pre-trained Models for SD Activities

## Input

1

```
protected boolean isBoundaryEvent(Object obj) {  
    if (obj instanceof <MASK>) {  
        mxCell cell = (mxCell) <MASK>;  
        <MASK> cell.getId().  
            startsWith("boundary-event-");  
    }  
    return false;  
}
```

2

*Creates an <MASK> with the given internal <MASK> capacity hint.*

...

## Target

1

- 1) mxCell
- 2) obj
- 3) return

2

- 1) UnicastProcessor
- 2) buffer

...

~2.6M

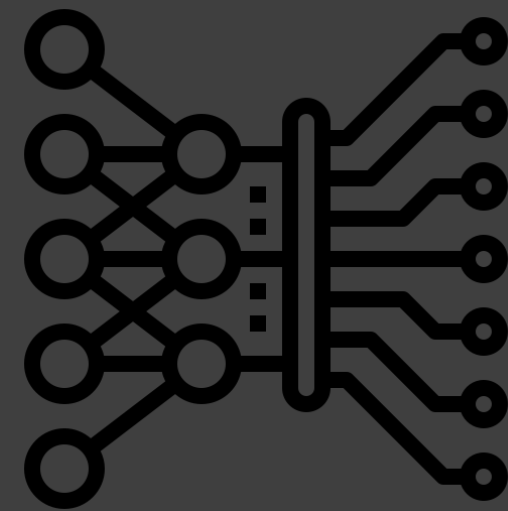
# Pre-trained Models for SD Activities

1



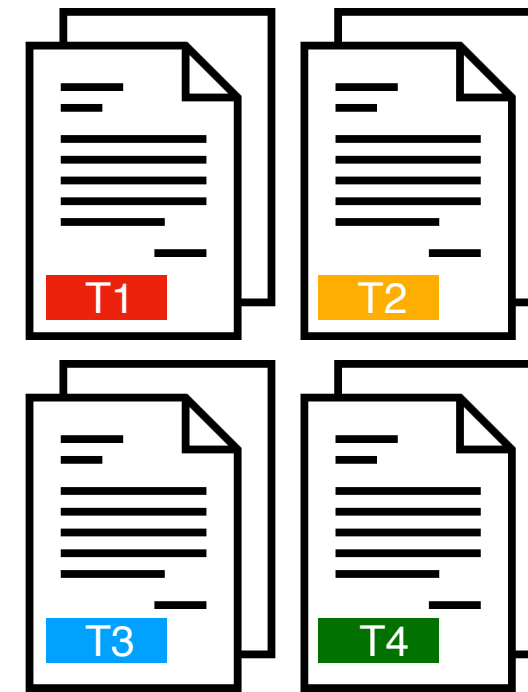
Pre-training dataset

2



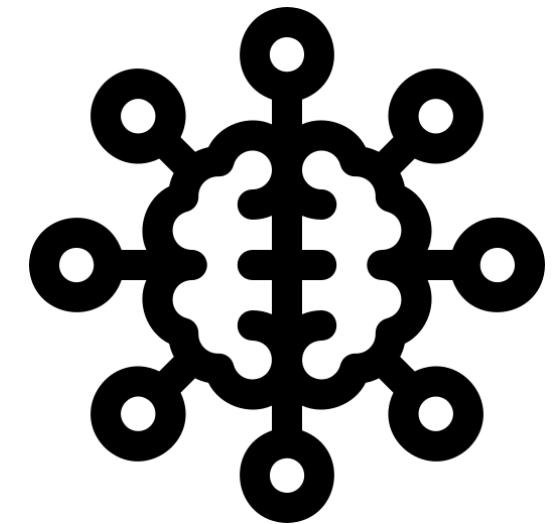
T5 Pre-training

3



Text-to-text task specific dataset

4



T5 Fine-tuning on tasks-mixture



# Pre-trained Models for SD Activities

3



Text-to-text task  
specific dataset

## *Code-Related Tasks Selection*

*Automatic  
Bug-Fixing*

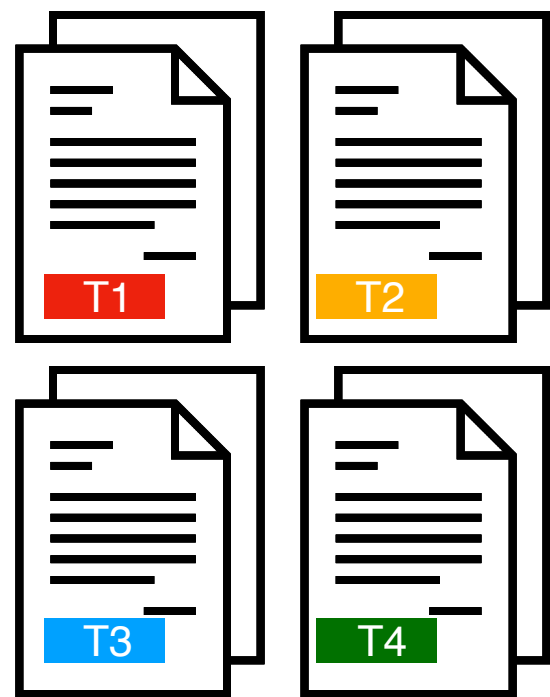
*Code  
Summarization*

*Injection of  
Code Mutants*

*Assert Statements  
Generation*

# Pre-trained Models for SD Activities

3



Text-to-text task  
specific dataset

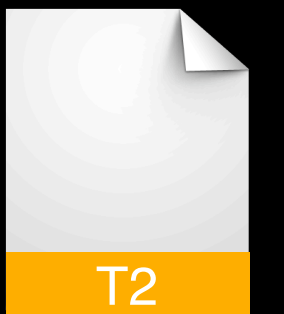
*An empirical study on learning bug-fixes patches in  
the wild via Neural Machine Translation*

TOSEM, —Tufano *et al.*—



*Improved automatic summarisation  
of code routines via Attention to File Context*

MSR, —Haque *et al.*—



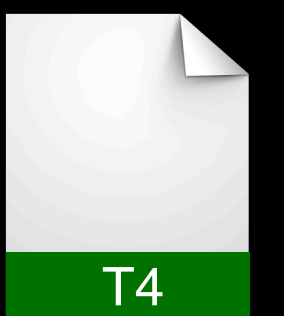
*Learn how to mutate source code from bug-fixes*

ICSME, —Tufano *et al.*—



*On learning meaningful assert statements for  
unit test cases*

ICSE, —Watson *et al.*—



[antoniomastropaolo.com](http://antoniomastropaolo.com)



[aura-se-lab.github.io](https://github.com/aura-se-lab)



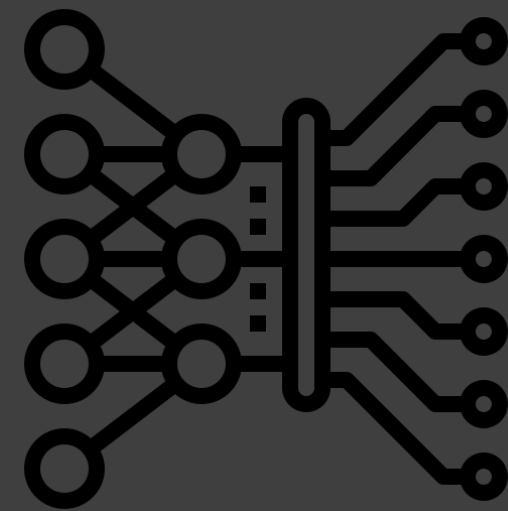
# Pre-trained Models for SD Activities

1



Pre-training  
dataset

2



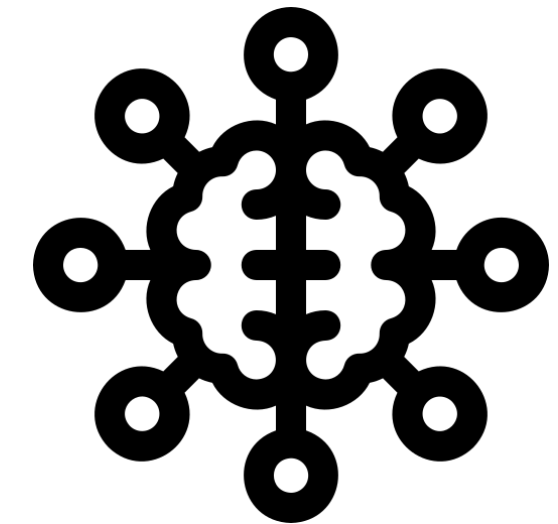
T5 Pre-training

3



Text-to-text task  
specific dataset

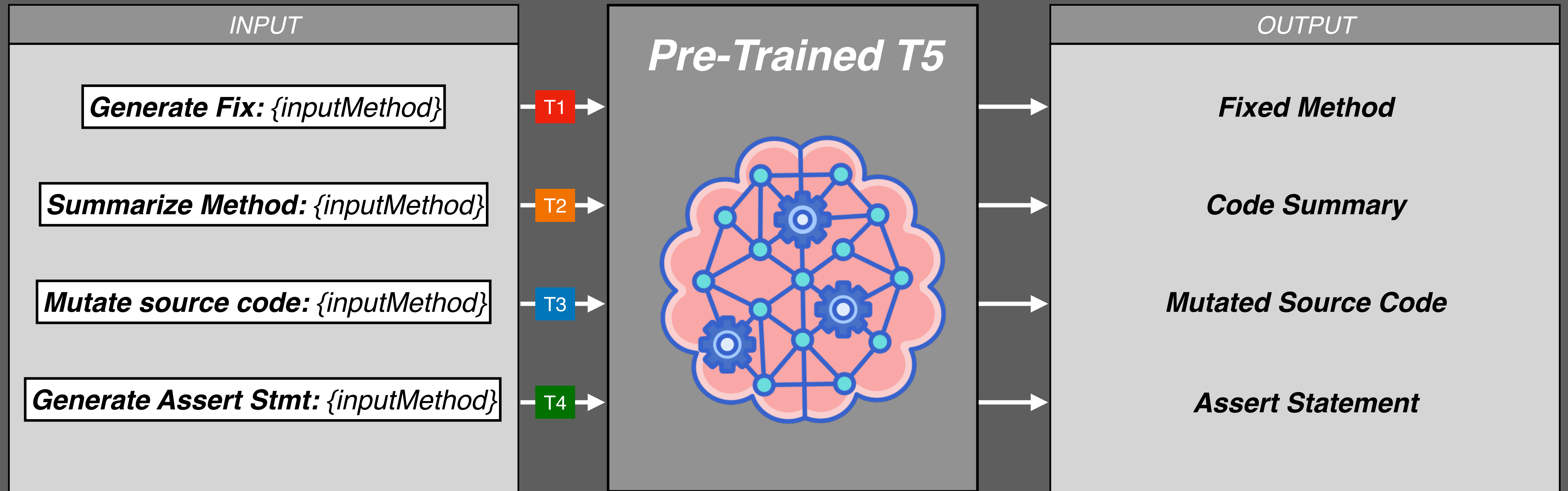
4



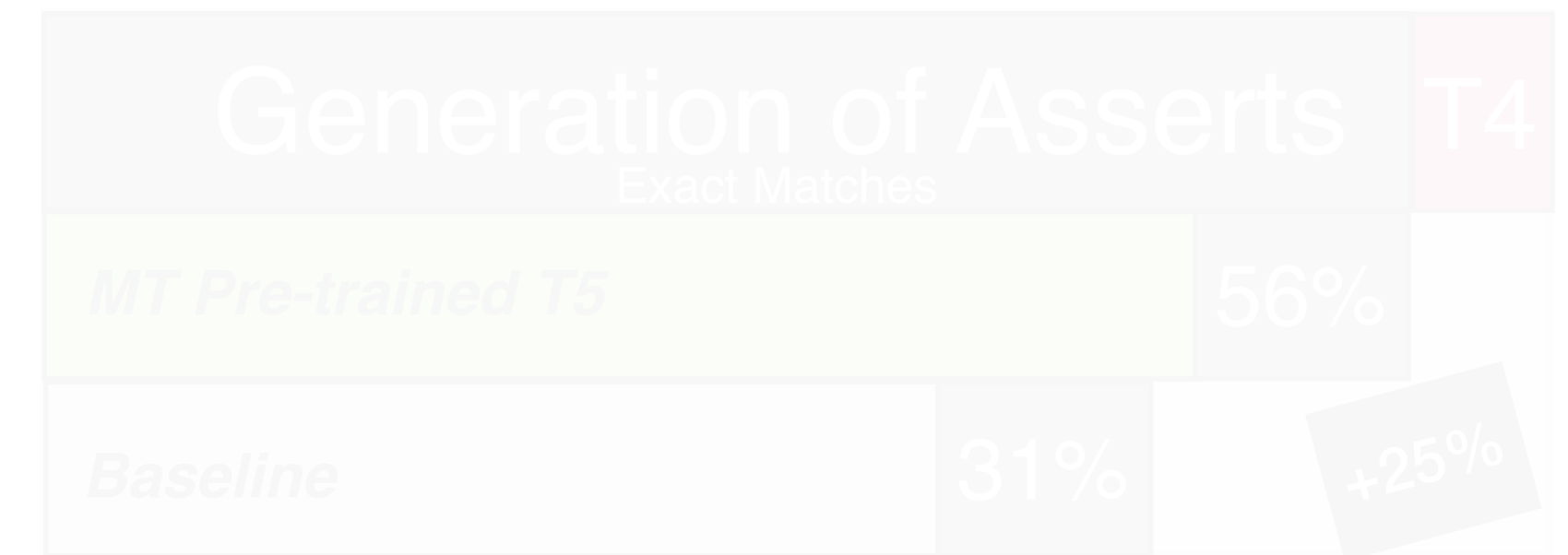
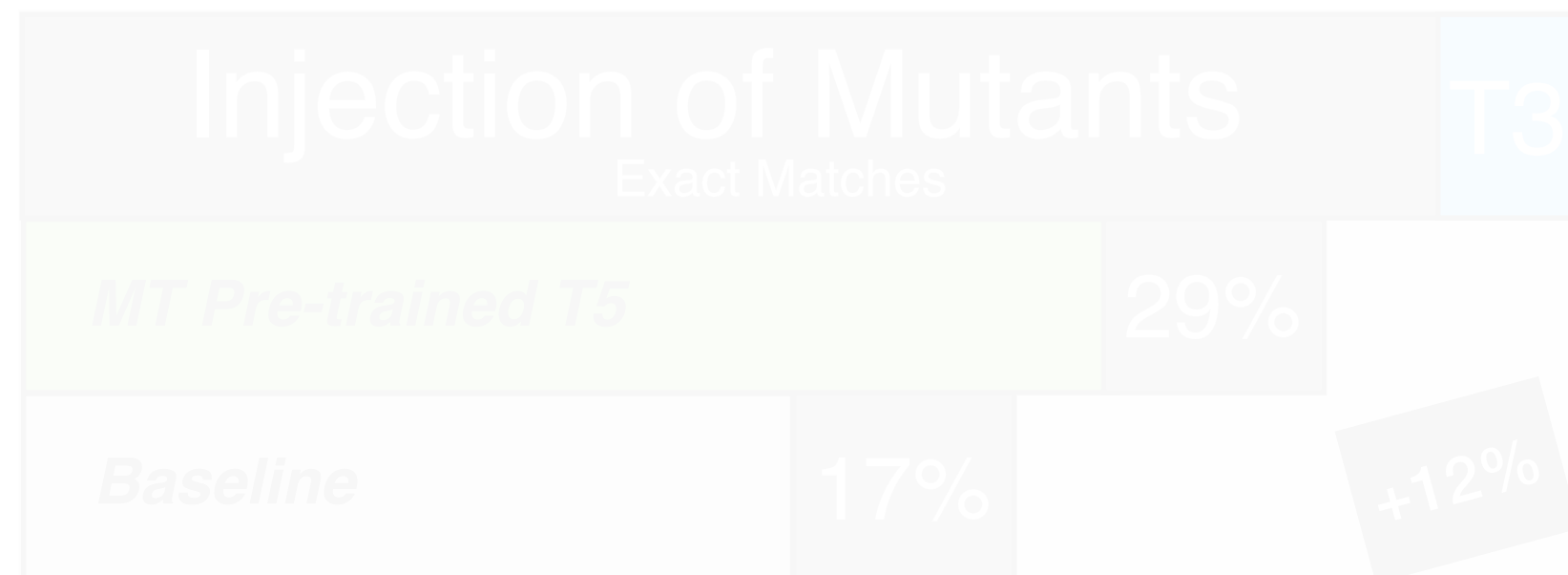
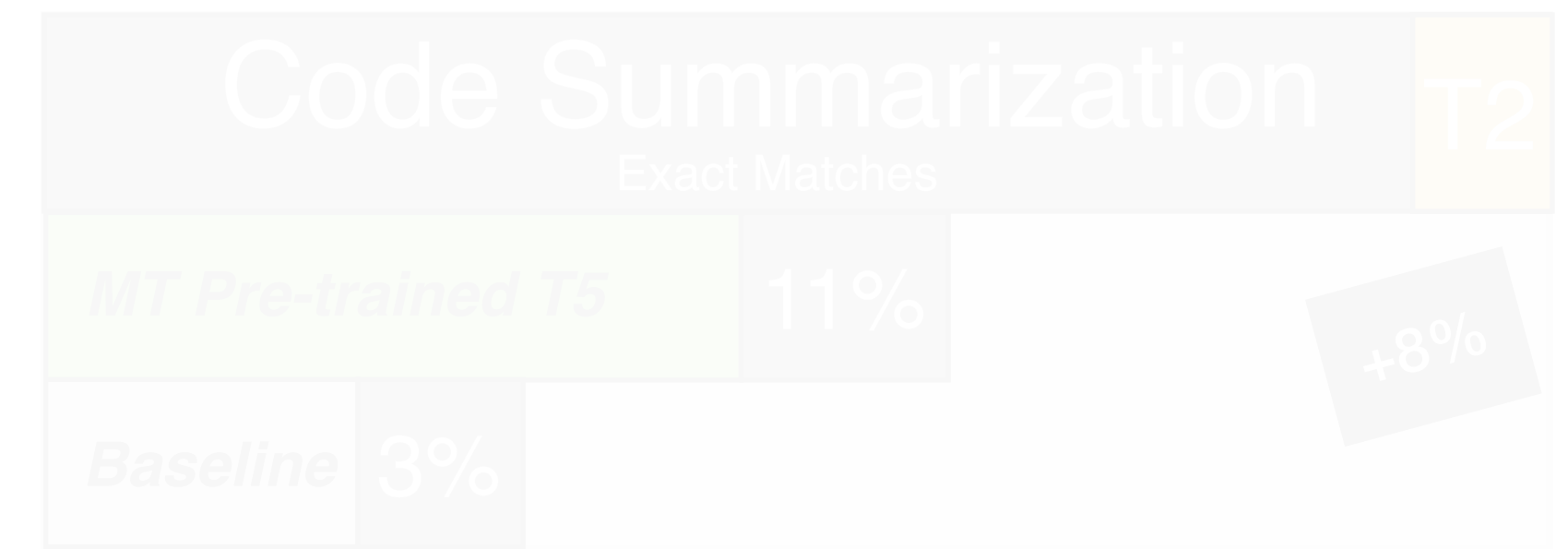
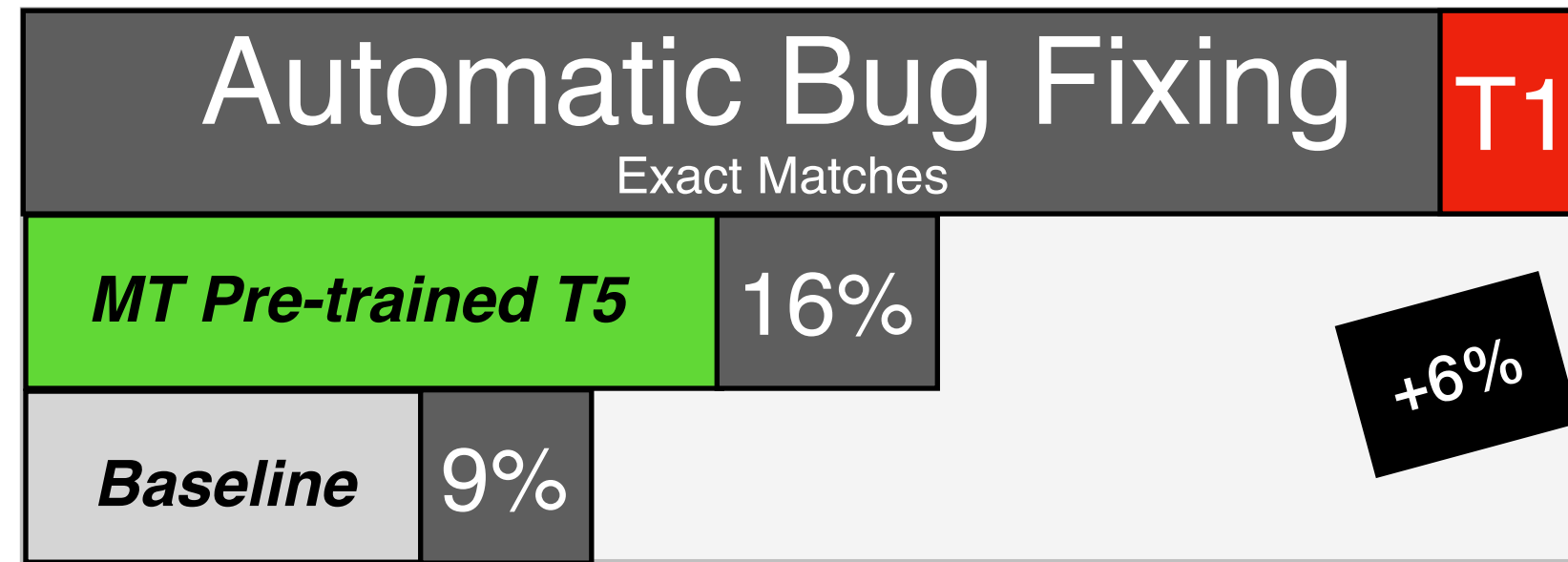
T5 Fine-tuning on  
tasks-mixture



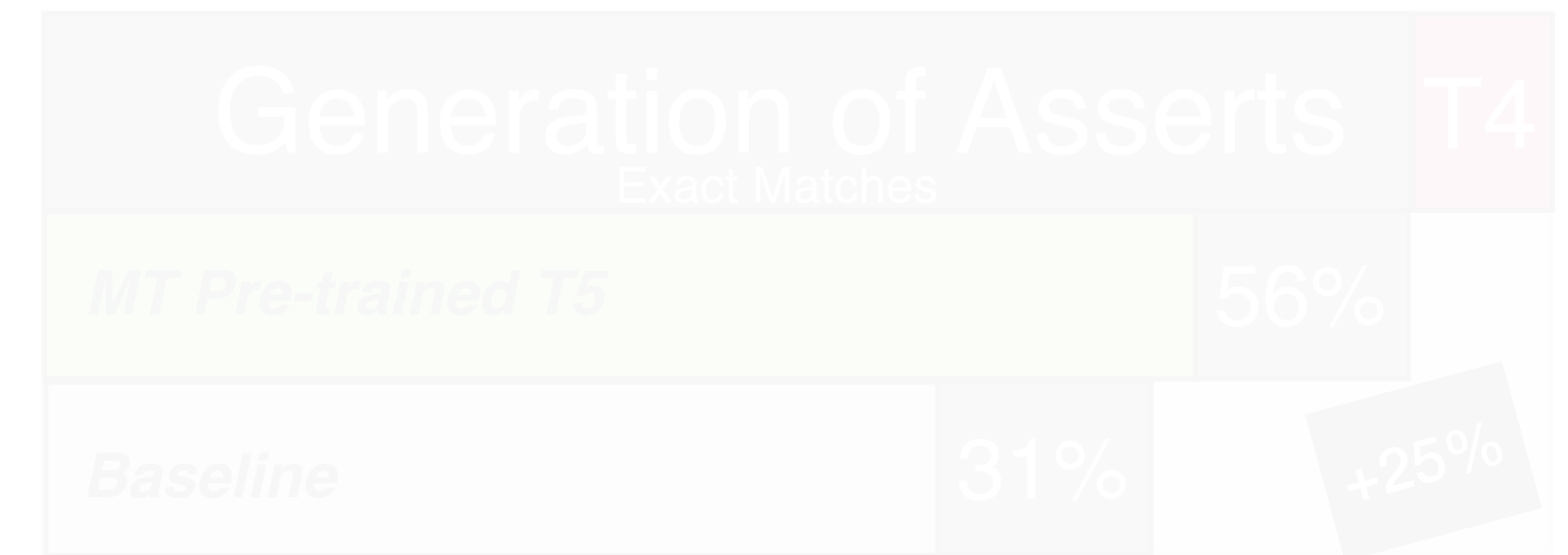
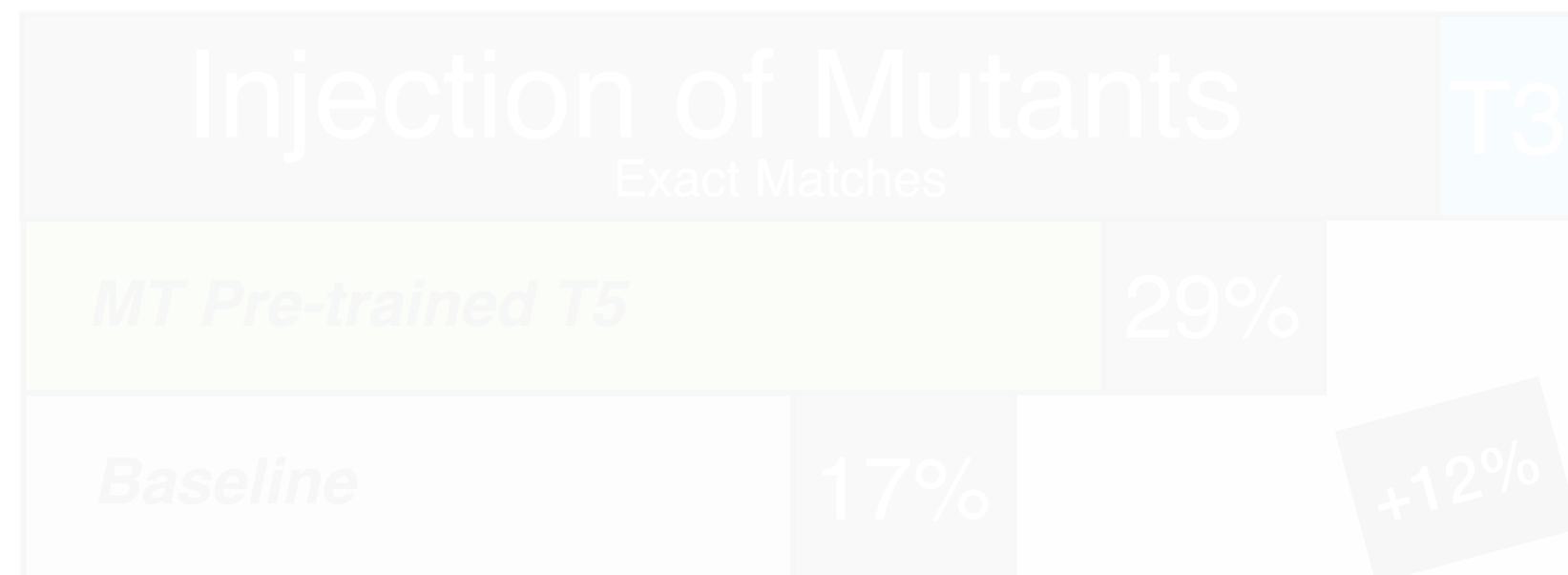
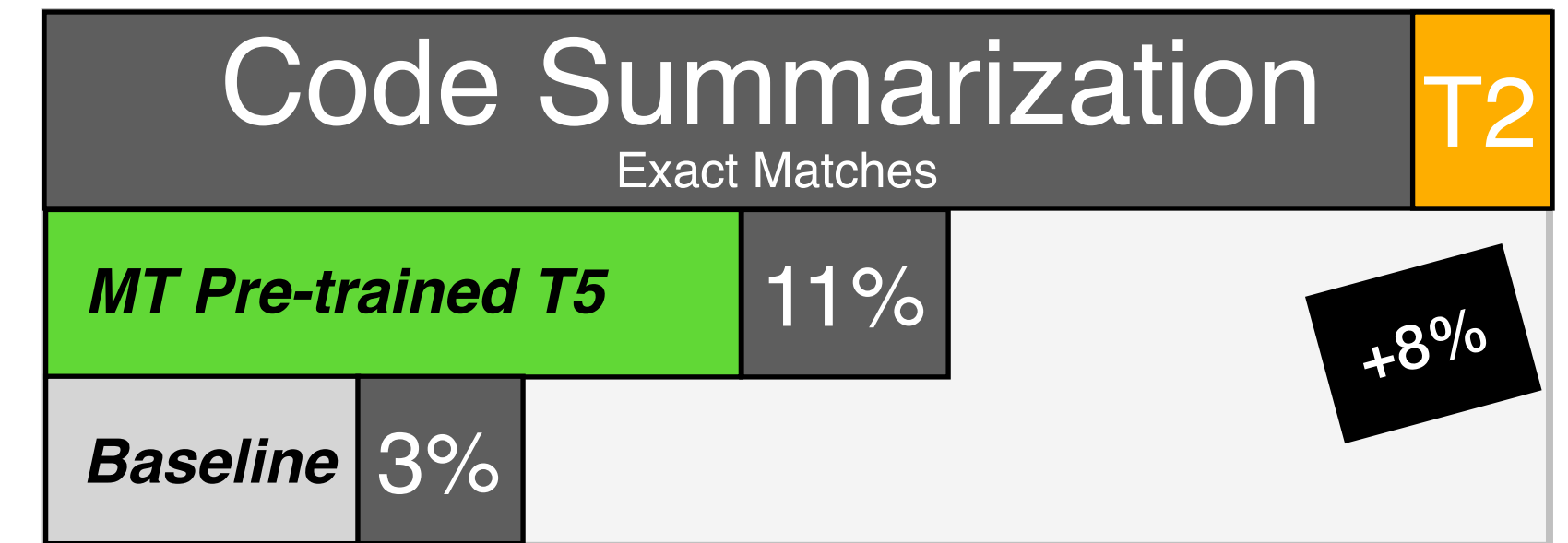
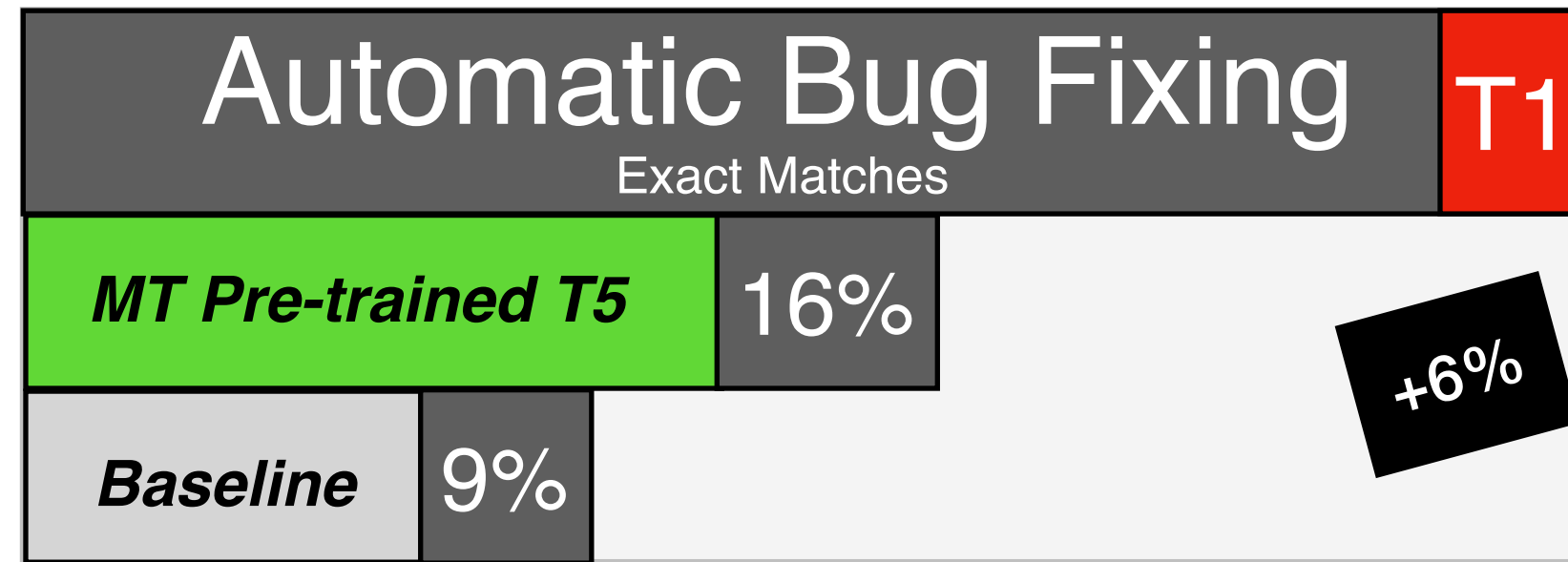
# Pre-trained Models for SD Activities



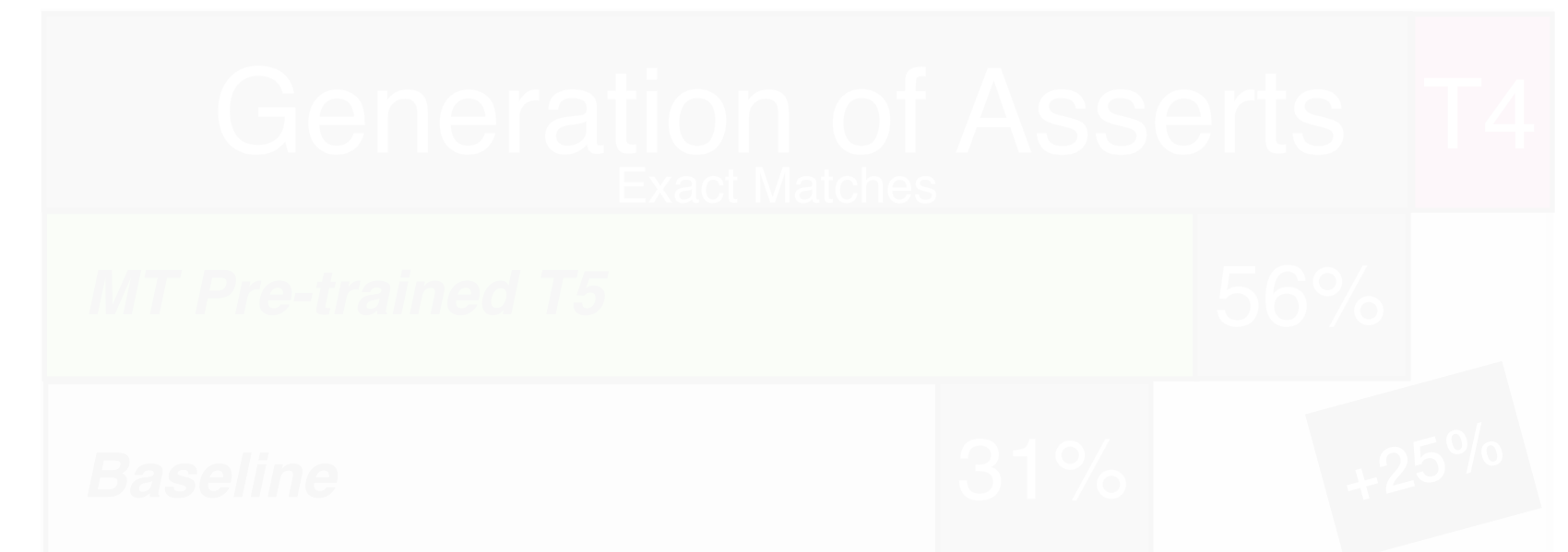
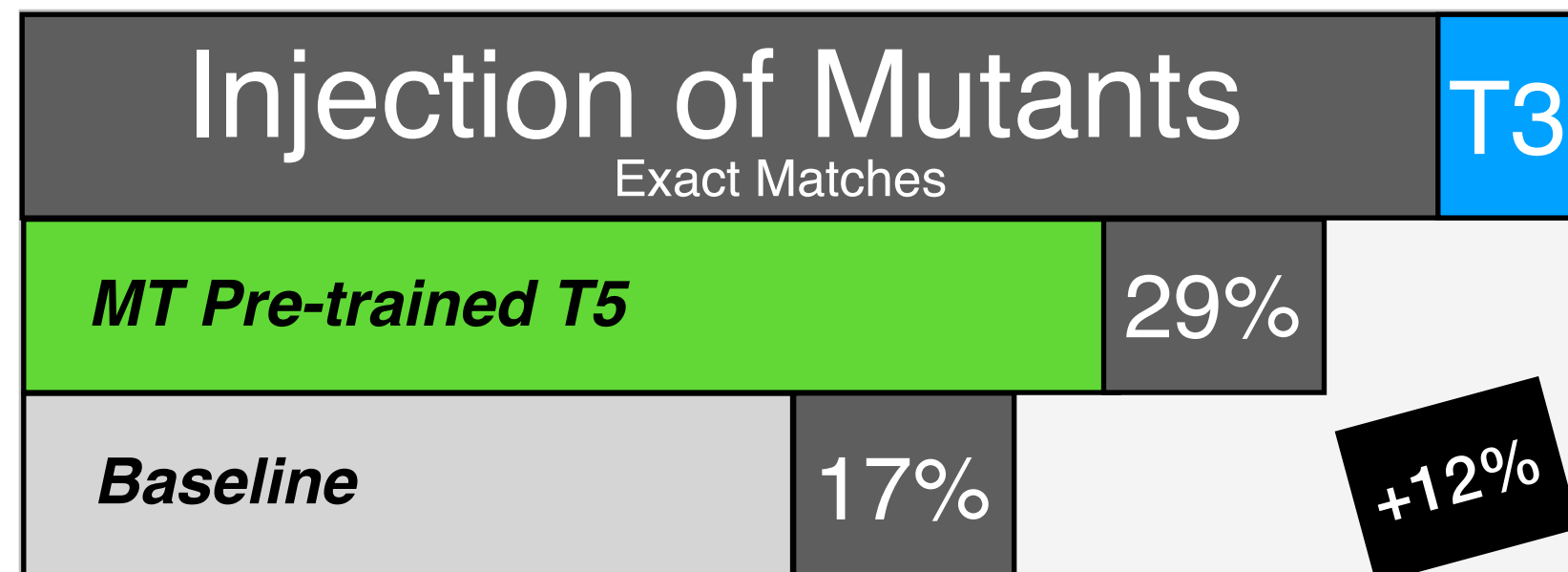
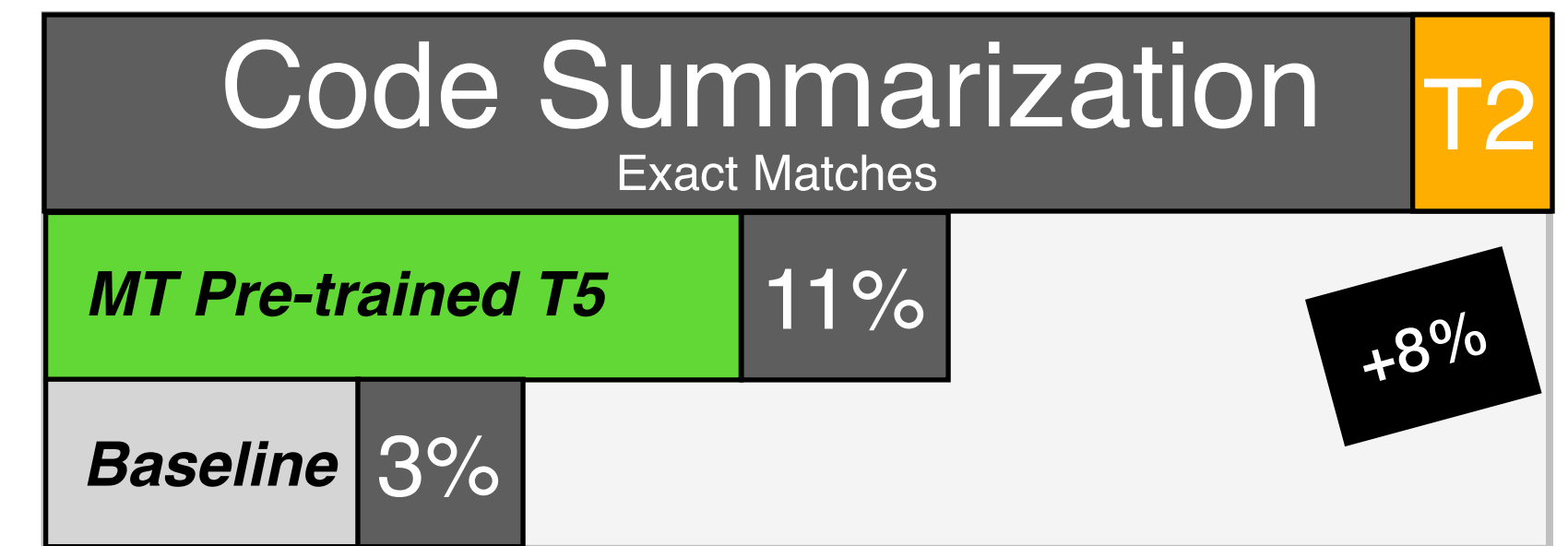
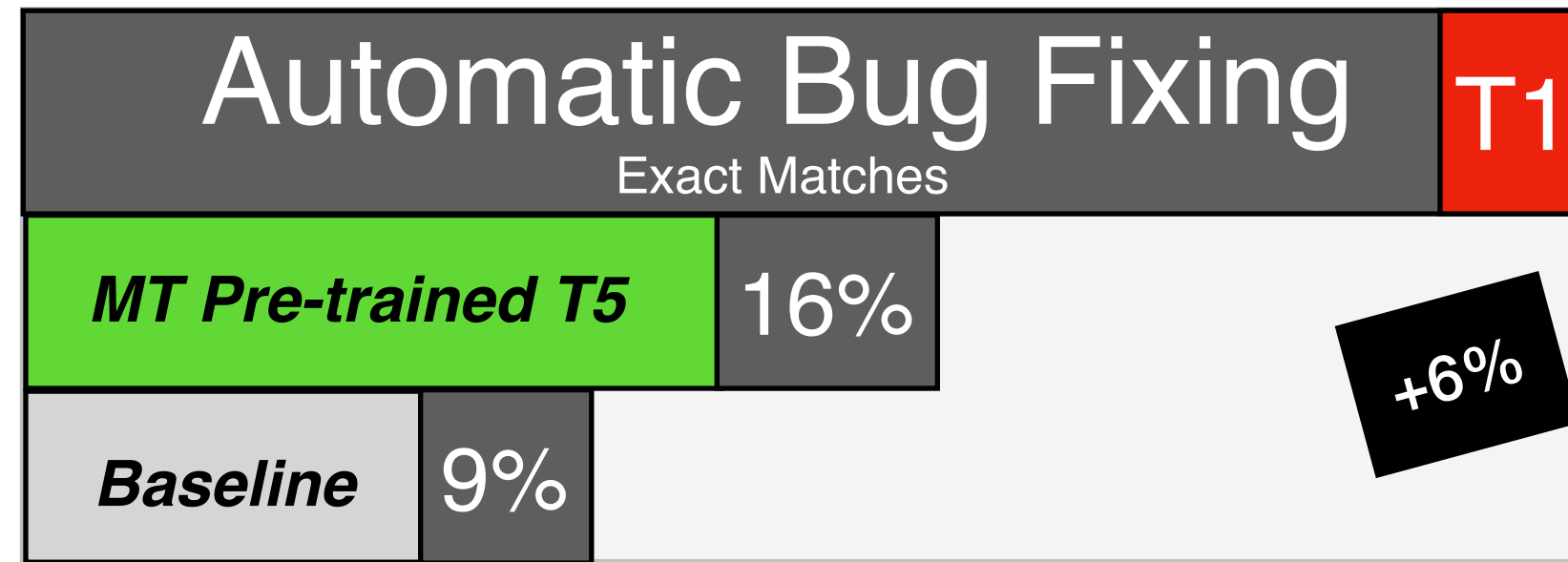
# Pre-trained Models for SD Activities



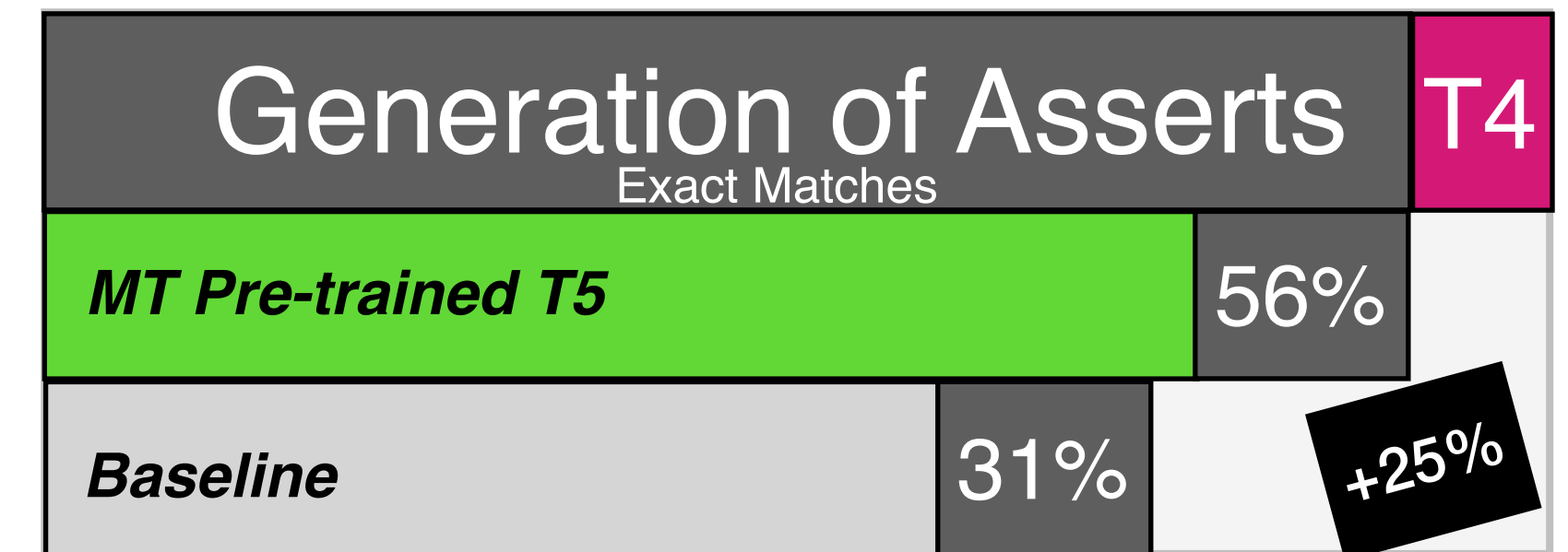
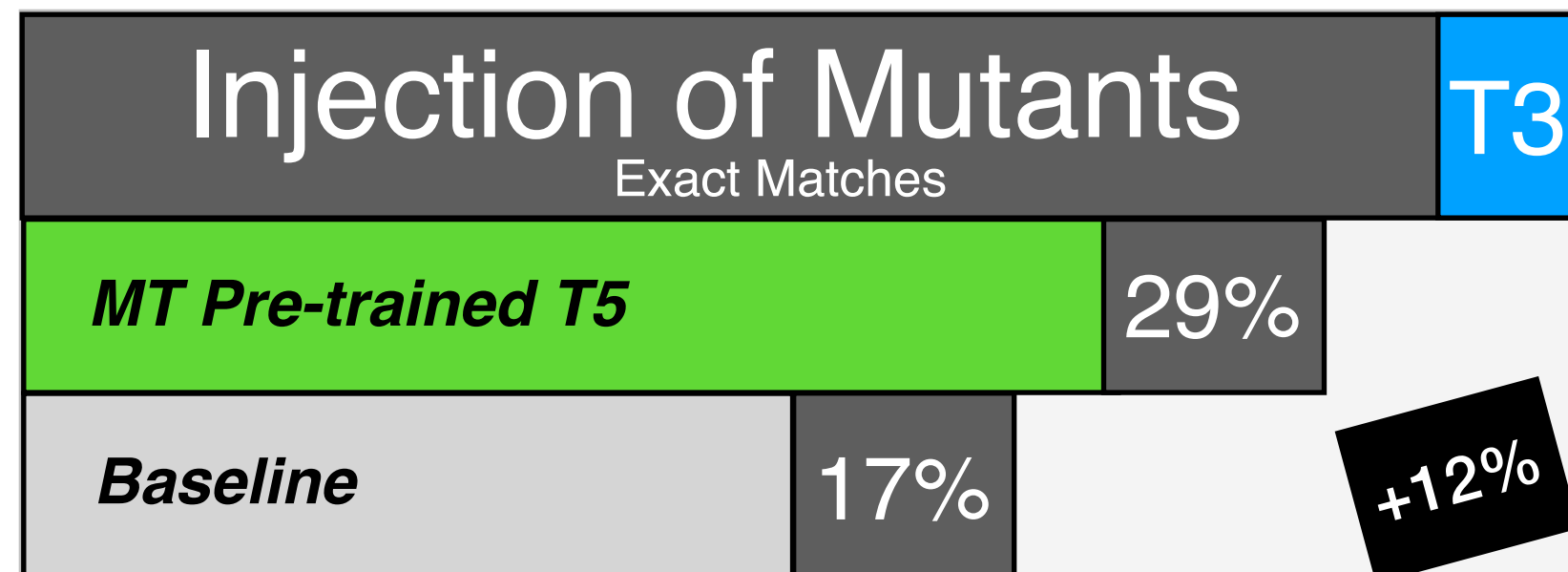
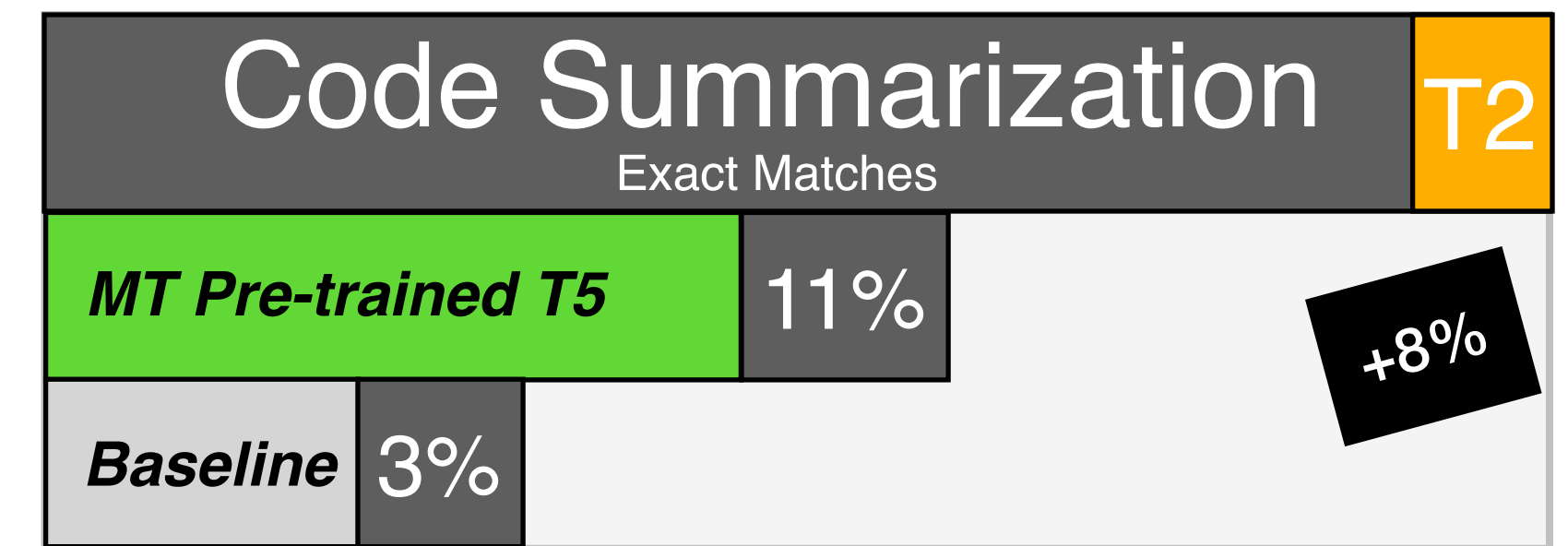
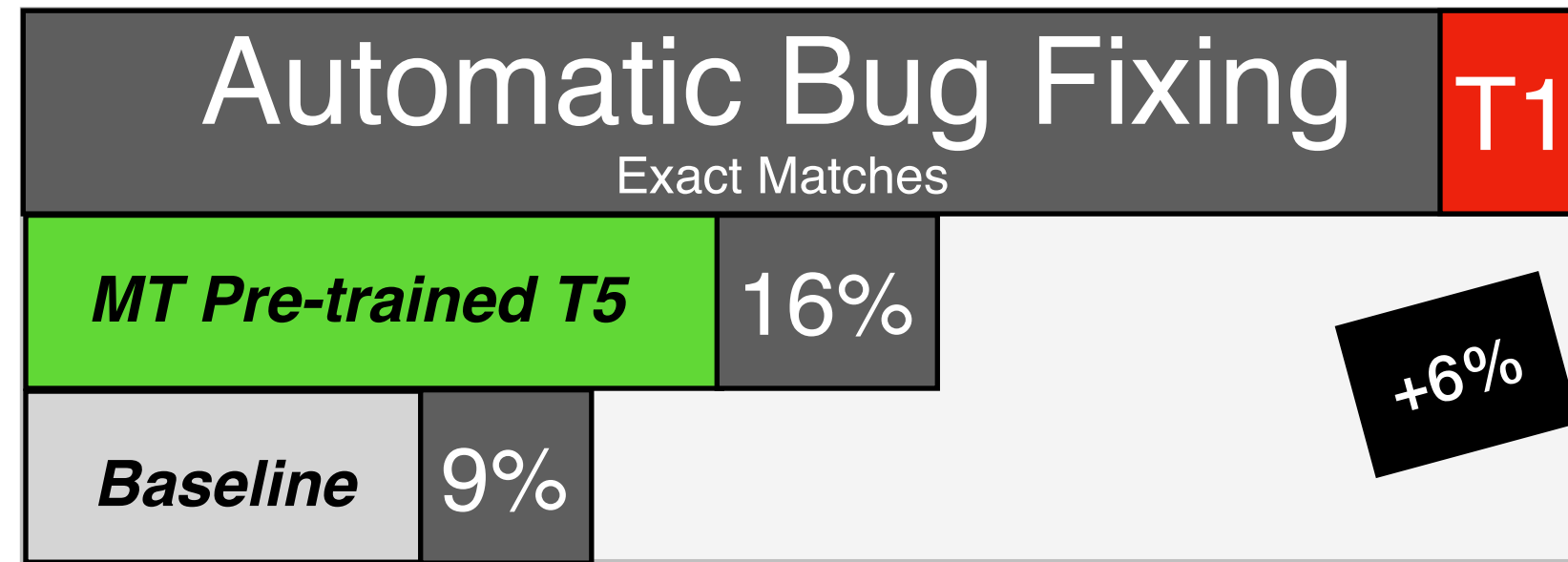
# Pre-trained Models for SD Activities



# Pre-trained Models for SD Activities



# Pre-trained Models for SD Activities



# Pre-trained Models for SD Activities



**Still alive guys?**



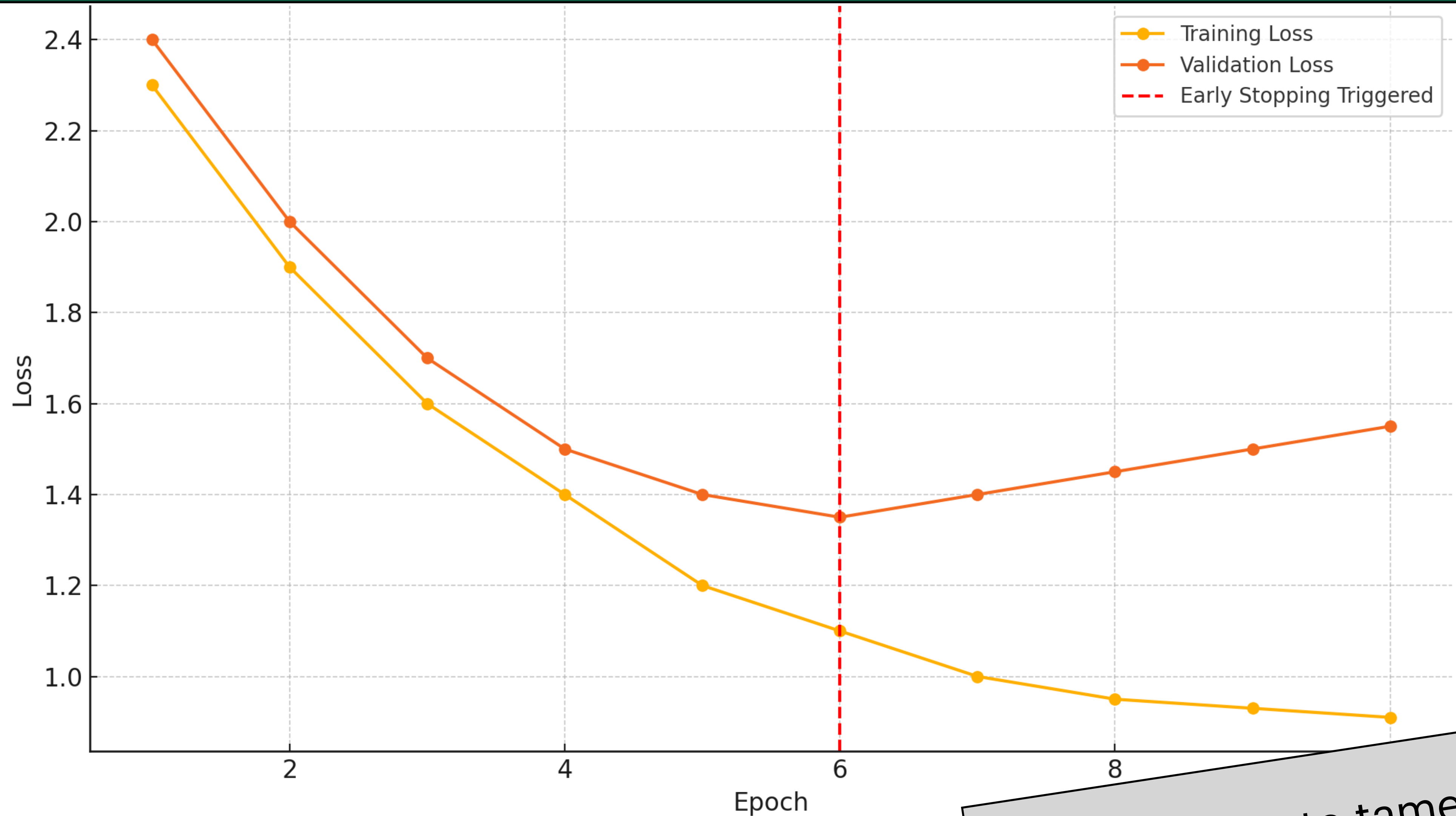
# Pre-trained Models for SD Activities

## Selecting The **Best-Performing** Model

For how many epochs we let the model (un)learn?



# Pre-trained Models for SD Activities



Early Stopping to tame overfitting

