

Smelling Source Code Using Deep Learning

A vertical stack of deep learning layer names and their corresponding symbols, rendered in a dark blue color. The layers from top to bottom are: Dense, sigmoid, LSTM, 1-D Convolution, Flatten, Dropout, 2-D Convolution, Embedding, Max pooling, and ReLU.

Dense

sigmoid

LSTM

1-D Convolution

Flatten

Dropout

2-D Convolution

Embedding

Max pooling

ReLU

Tushar Sharma

<http://www.tusharma.in>

Refactoring
Code Design
Code quality
Tools
C#
Change impact analysis

Java
Python
Tools
C#
Code quality
Software design
Architecture
Training
Books
Public speaking

Code smells
IaC
Java
SE research
Open source



What is a smell?



...certain structures in the code that suggest (sometimes they scream for) the possibility of refactoring.

- Kent Beck

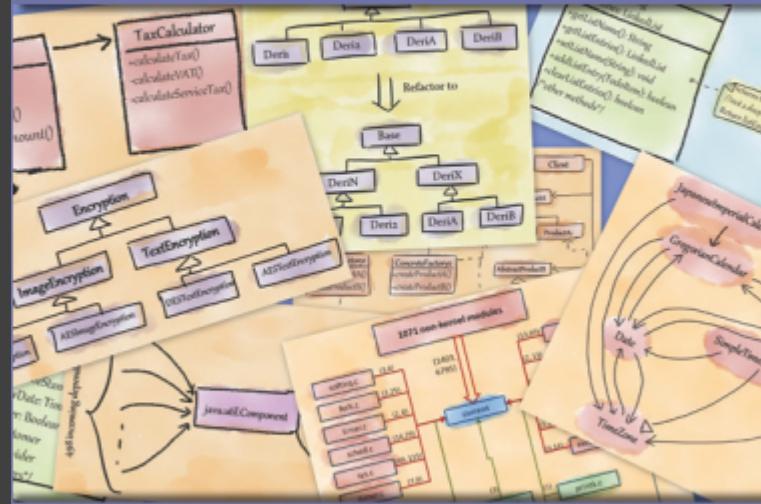
20 Definitions of smells: <http://www.tusharma.in/smells/smellDefs.html>
Smells' catalog: <http://www.tusharma.in/smells/>

Implementation smells

Duplicate Code
Long Parameter List
Complex Conditional
Complex Method
Magic Number
Empty Catch Block
Long Method
Missing Default
Long Identifier
Long Statement

Refactoring for Software Design Smells

Managing Technical Debt



Girish Suryanarayana,
Ganesh Samartham, Tushar Sharma
Forewords by Grady Booch and Stéphane Ducasse



Design Smells

- Wide Hierarchy
- Multipath Hierarchy
- Rebellious Hierarchy
- Speculative Hierarchy
- Duplicate Abstraction
- Unutilized Abstraction
- Unexploited Encapsulation
- Cyclically-dependent Modularization
- Insufficient Modularization
- Hub-like Modularization
- Unnecessary Abstraction
- Deficient Encapsulation
- Incomplete Abstraction
- Broken Modularization
- Unfactored Hierarchy
- Deep Hierarchy
- Missing Hierarchy
- Leaky Encapsulation
- Cyclic Hierarchy
- Missing Abstraction

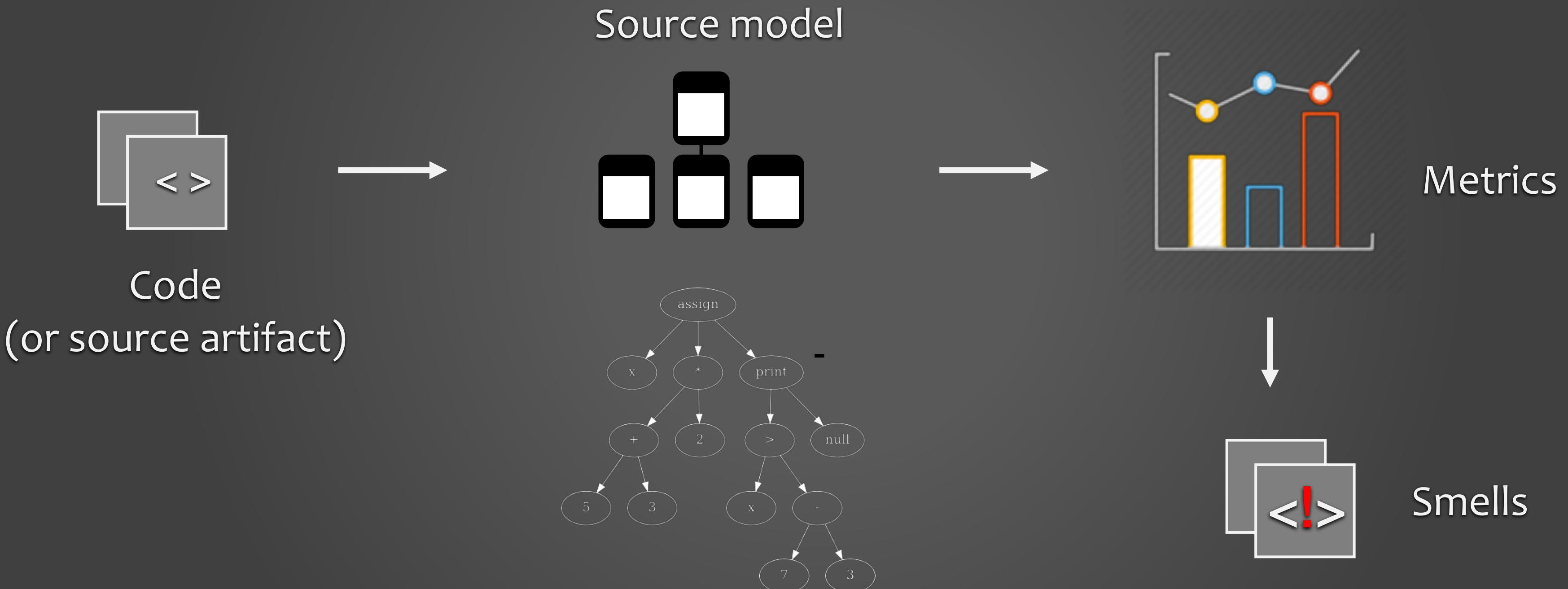
Architecture Smells

Cyclic Dependency
Feature Concentration
God Component
Scattered Functionality
Dense Structure
Ambiguous Interface
Unstable Dependency

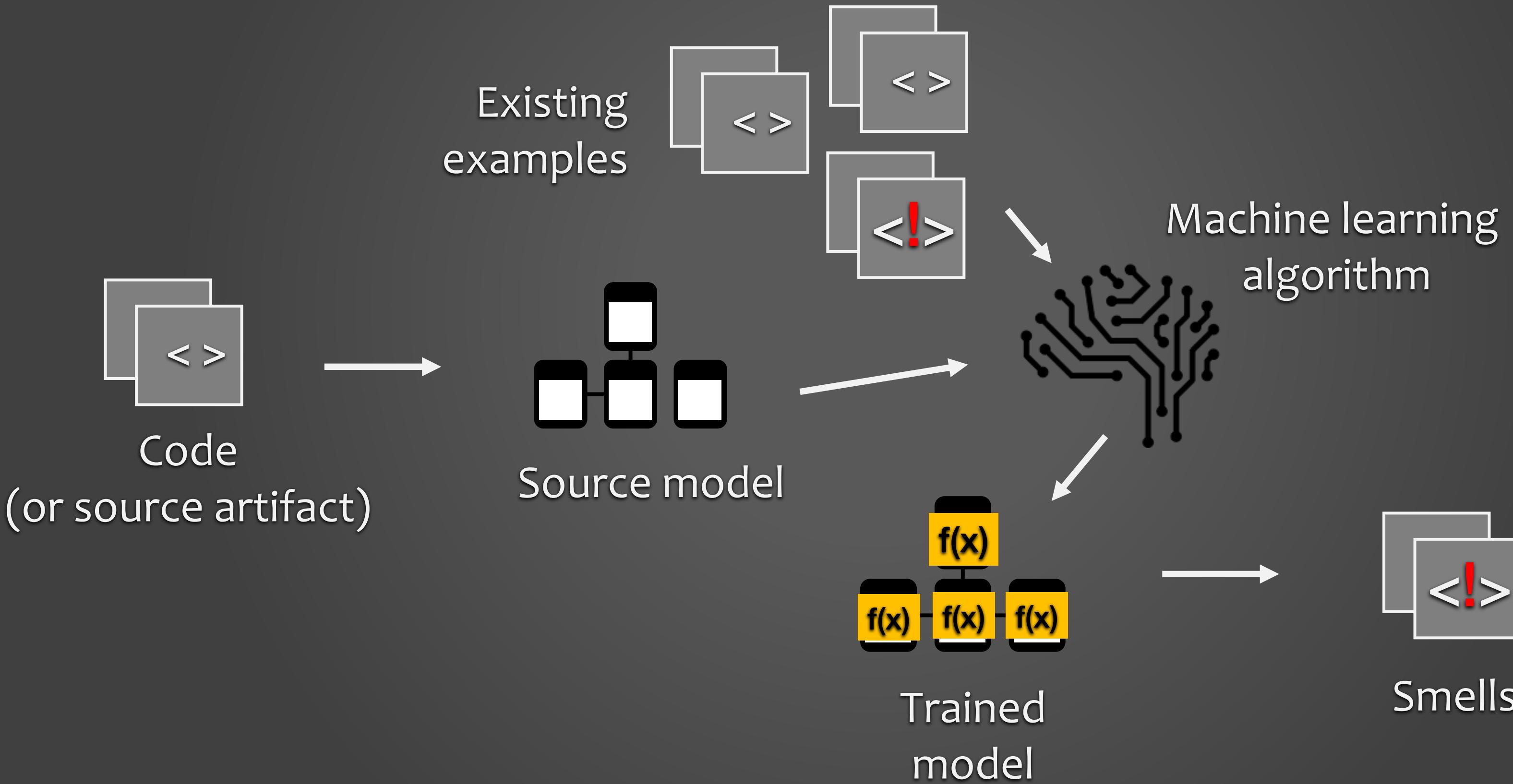
How smells get detected?

History-based
Heuristics-based
Metrics-based
Machine learning-based
Optimization-based

Metrics-based smell detection



Machine learning-based smell detection



Machine learning-based smell detection

Existing academic work:

- Support vector machines
- Bayesian belief network
- Logistic regression
- CNN

Take metrics as the features/input

Validation on balanced samples



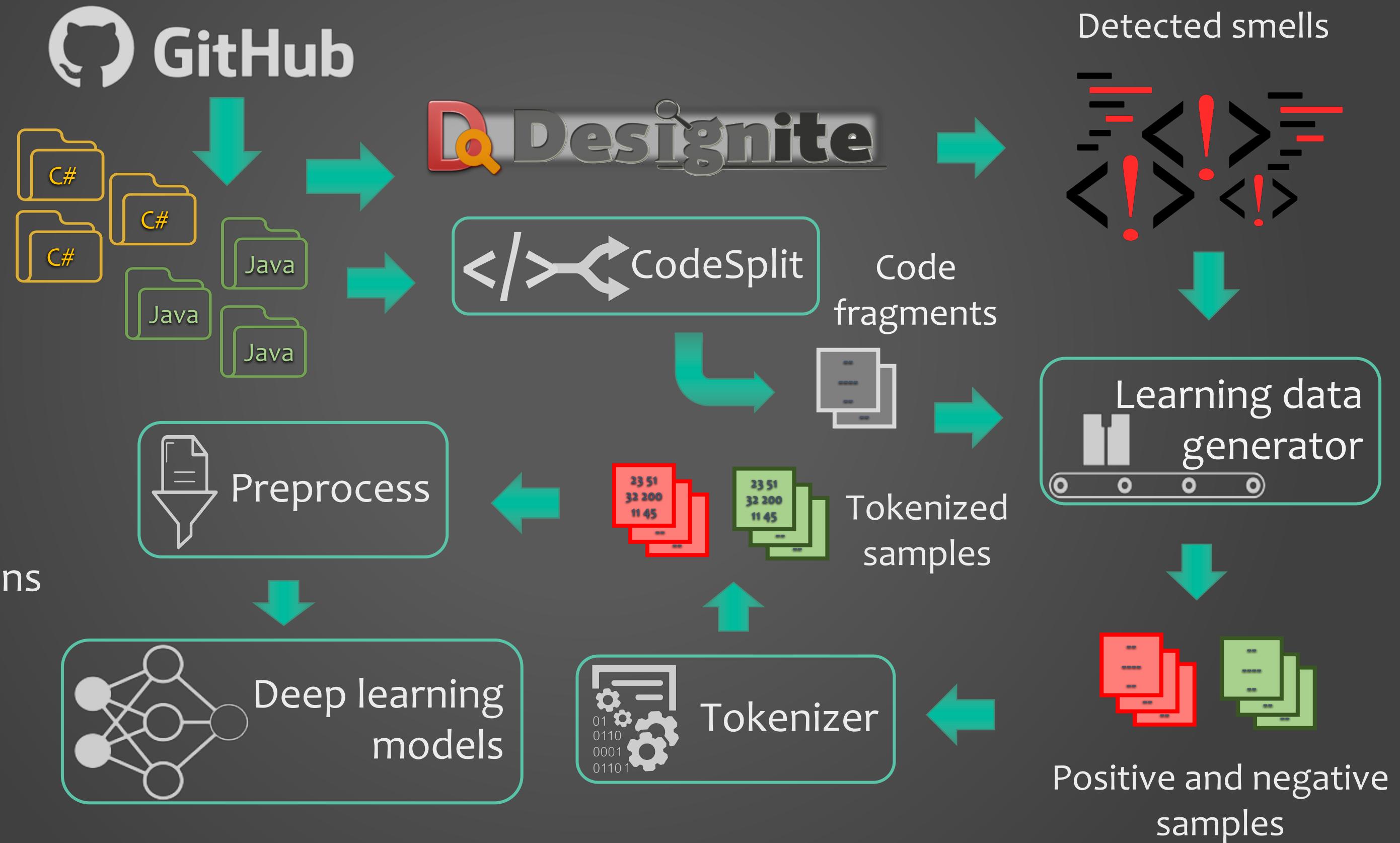
Research questions

RQ1: Would it be possible to use deep learning methods to detect code smells?

RQ2: Is transfer-learning feasible in the context of detecting smells?

Transfer-learning refers to the technique where a learning algorithm exploits the commonalities between different learning tasks to enable knowledge transfer across the tasks

Overview



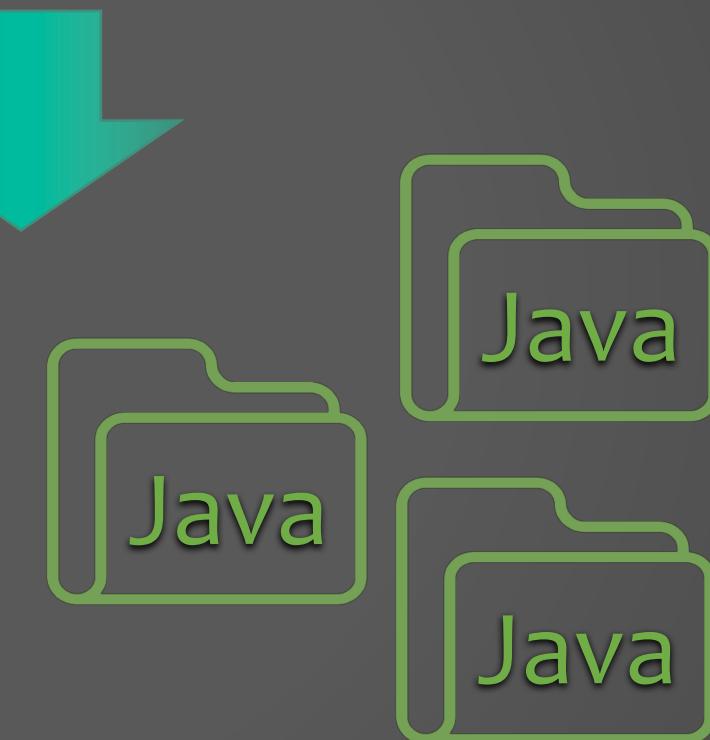
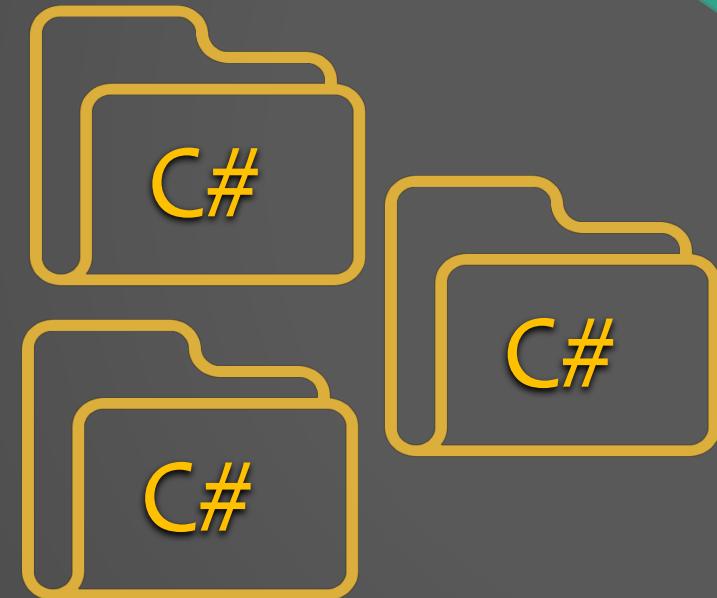
Data Curation



Repositories download



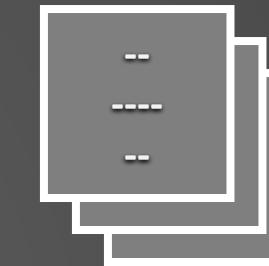
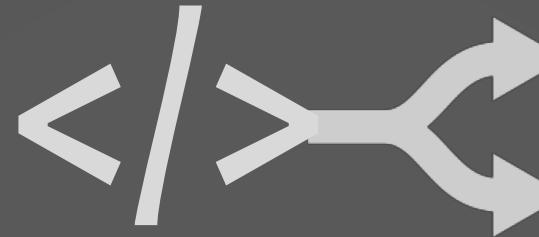
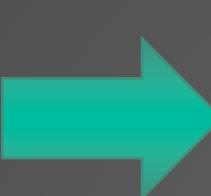
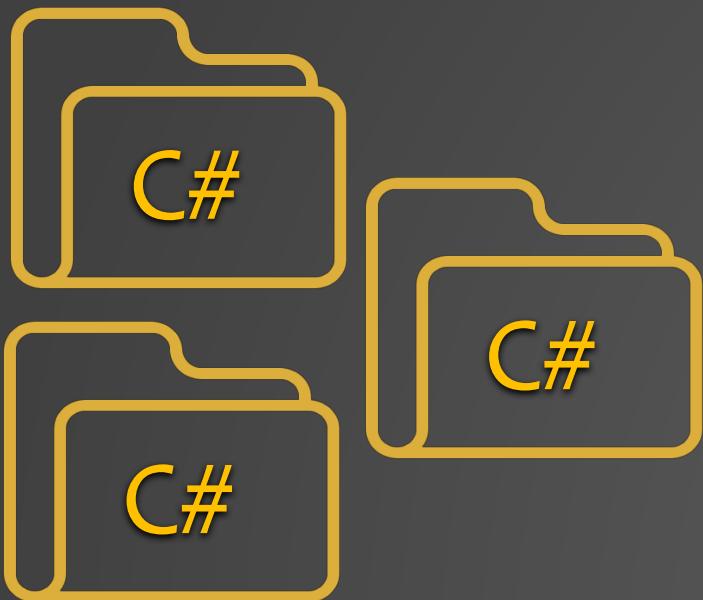
1,072
repositories



2,528 and
selected 100
repositories

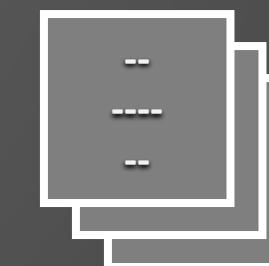
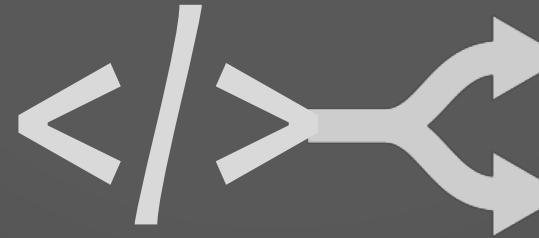
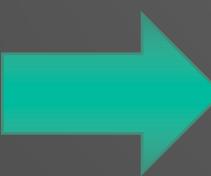
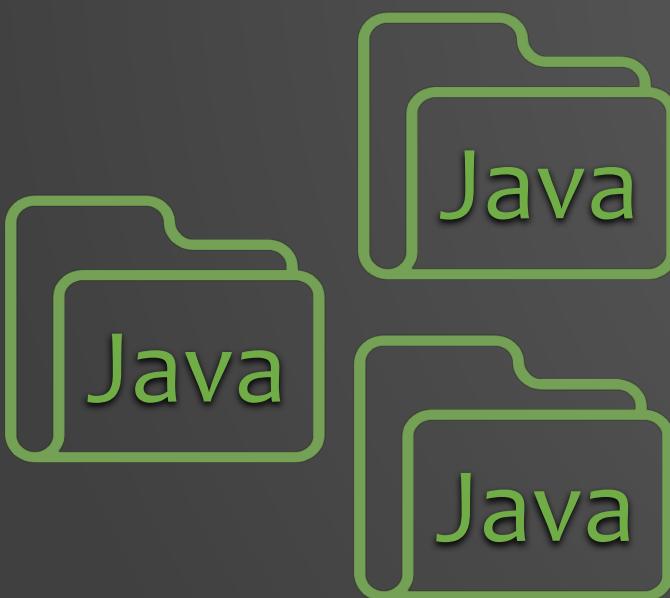
Architecture	Community	CI	Documentation	History	License	Issues	Unit test	Stars
✓	✗	✓	✓	✓	✗	✓	✓	✓

Splitting code fragments

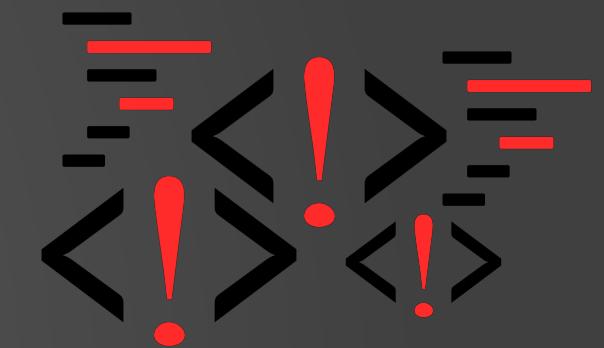
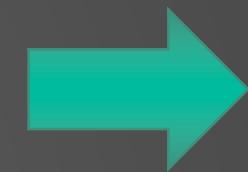
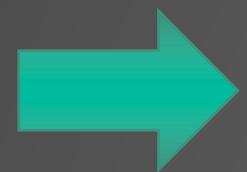
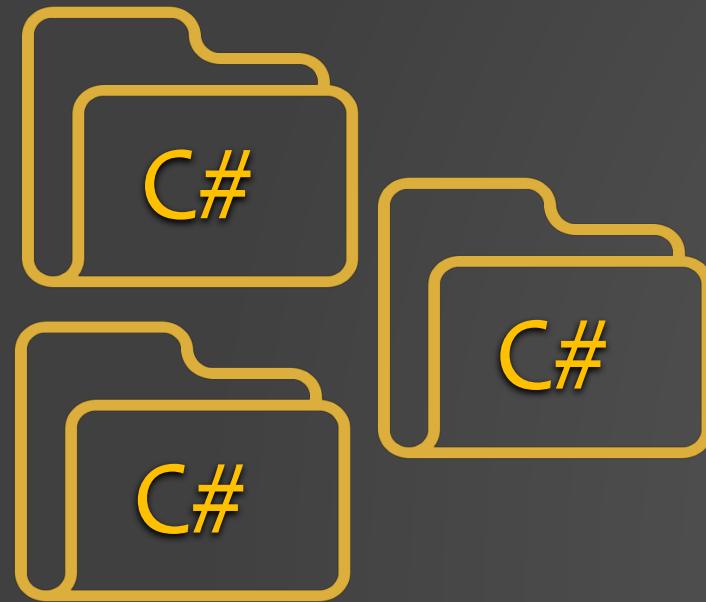


CodeSplit

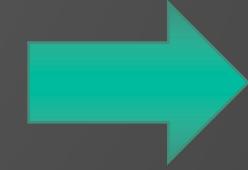
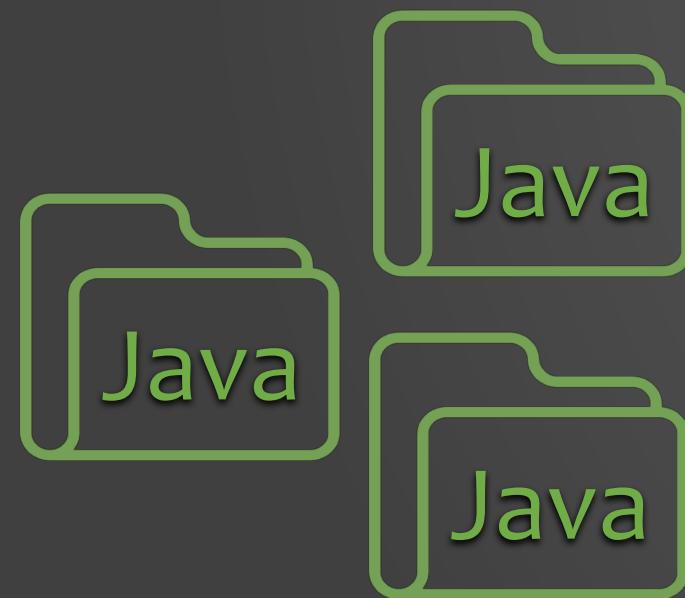
Code fragments
(methods or classes)



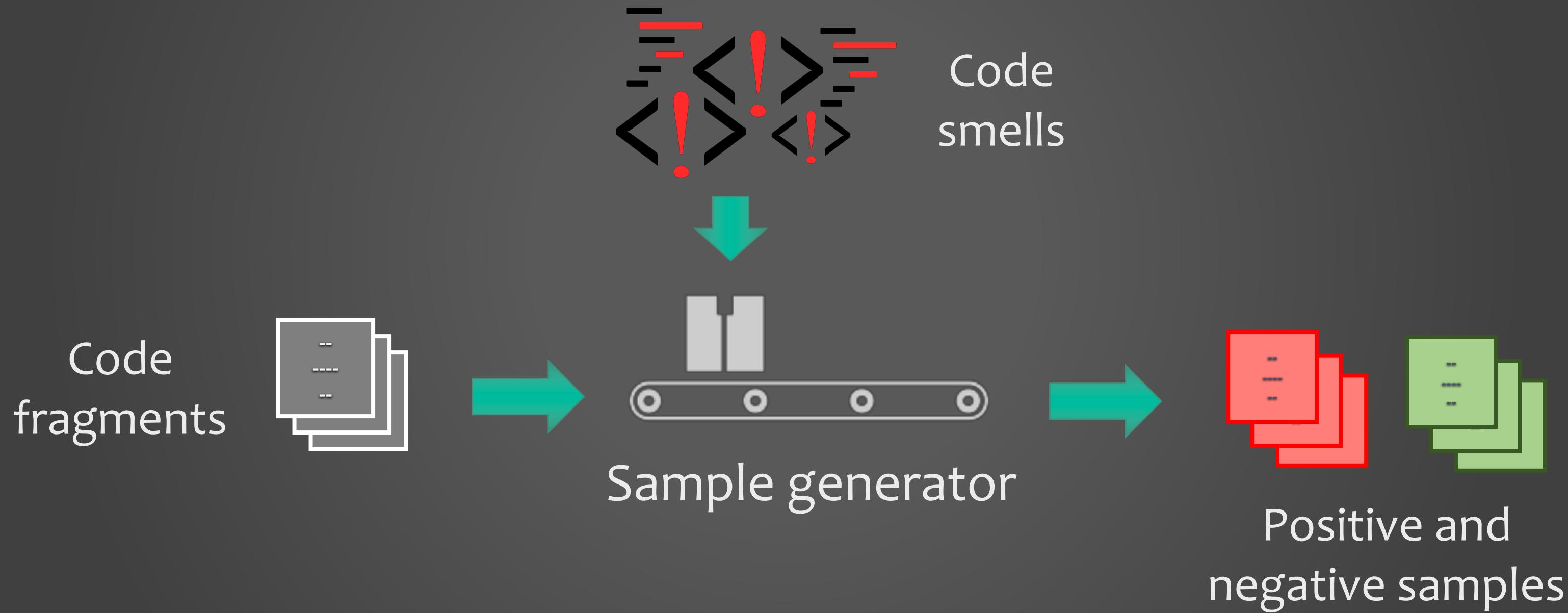
Smell detection



Detected code smells



Generating training and evaluation samples



Tokenizing learning samples



Tokenizing learning samples

```
public void InternalCallback(object state)
{
    Callback(State);
    try
    {
        timer.Change(Period, TimeSpan.Zero);
    }
    catch (ObjectDisposedException)
    {
    }
}
```

1-D

123 2002 40 2003 41 59 474 123 2004 46 200

2-D

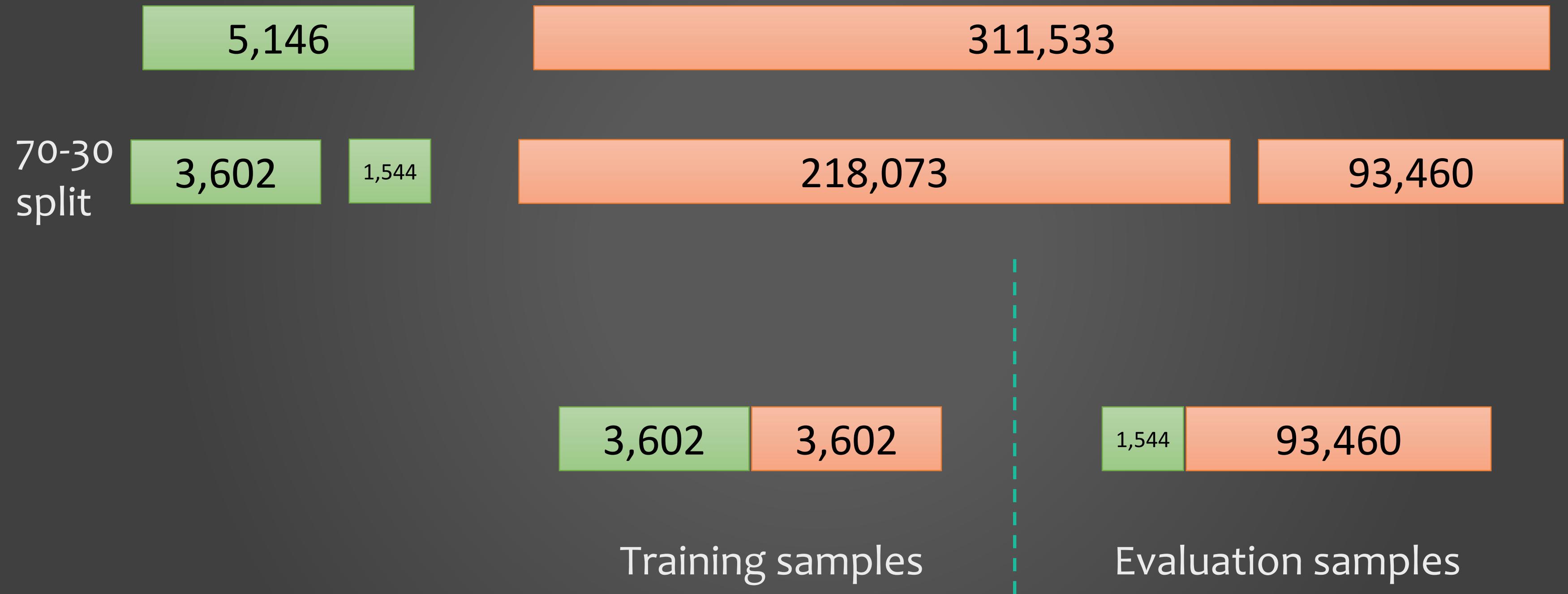
123 2002 40 2003 41 59

474 123 2004 46 2005 40 2006 44 2007 46 200

125 329 40 2009 41 123 125 125



Data preparation



Selection of smells

the method has high cyclomatic complexity

a class has more than one responsibility assigned to it

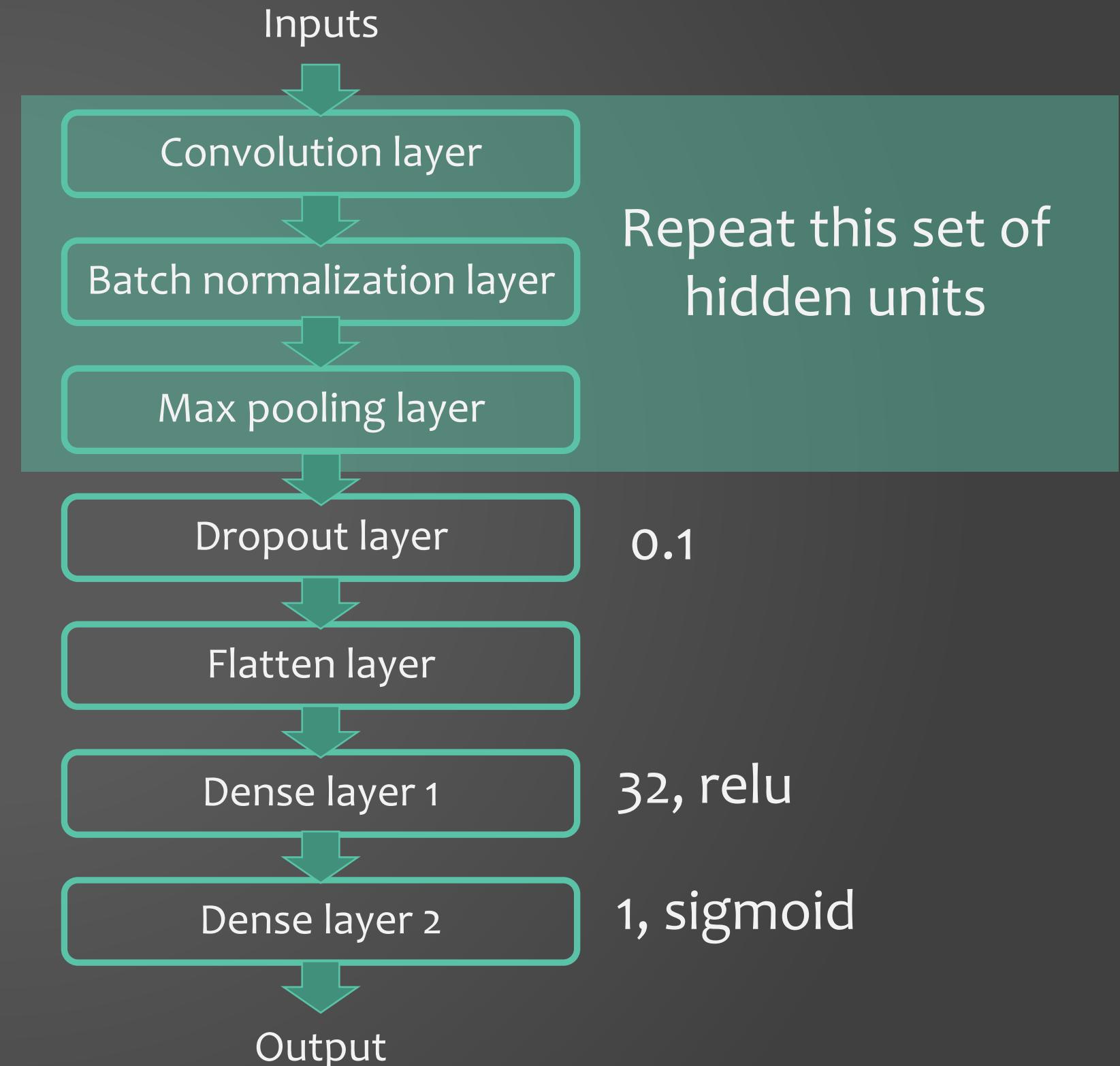
- Complex method
- Magic number
- Empty catch block
- Multifaceted abstraction

an unexplained numeric literal is used in an expression

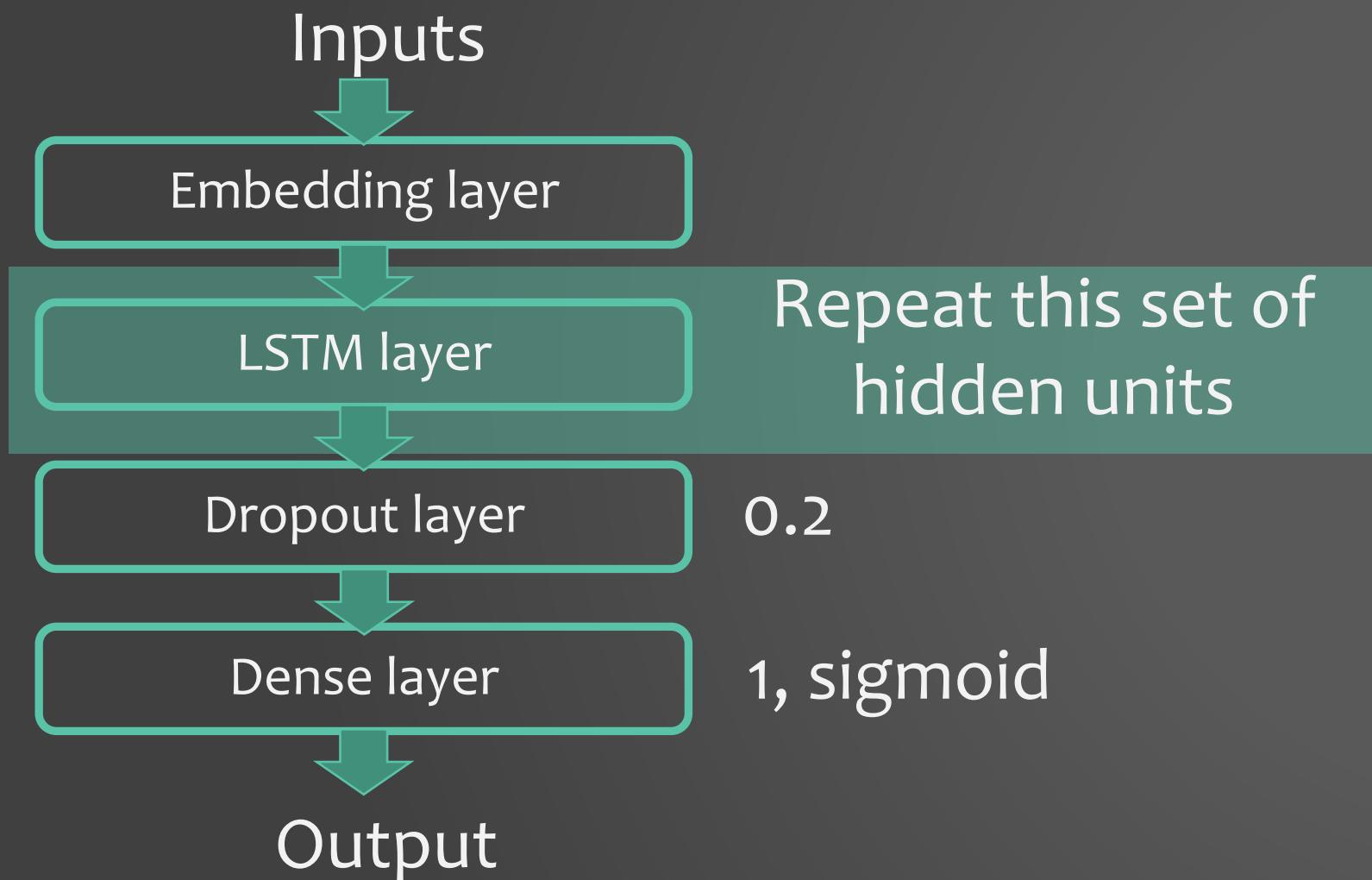
a catch block of an exception is empty

Architecture - CNN

- Filters = {8, 16, 32, 64}
- Kernel size = {5, 7, 11}
- Pooling window = {2, 3, 4, 5}
- Dynamic Batch size = {32, 64, 128, 256}
- Callbacks
 - Early stopping (patience = 5)
 - Model check point



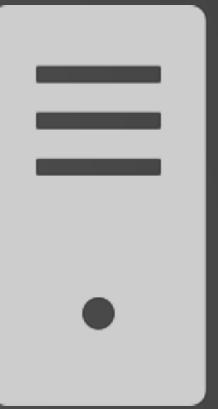
Architecture - RNN



- Dimensionality of embedding layer = $\{16, 32\}$
- LSTM units = $\{32, 64, 128\}$
- Dynamic Batch size = $\{32, 64, 128, 256\}$
- Callbacks
 - Early stopping (patience = 2)
 - Model check point

Running experiments

- Phase 1 – Grid search for optimal hyper-parameters
 - Validation set – 20%
 - Number of configurations
 - CNN = 144
 - RNN = 18
- Phase 2 – experiment with the optimal hyper-parameters



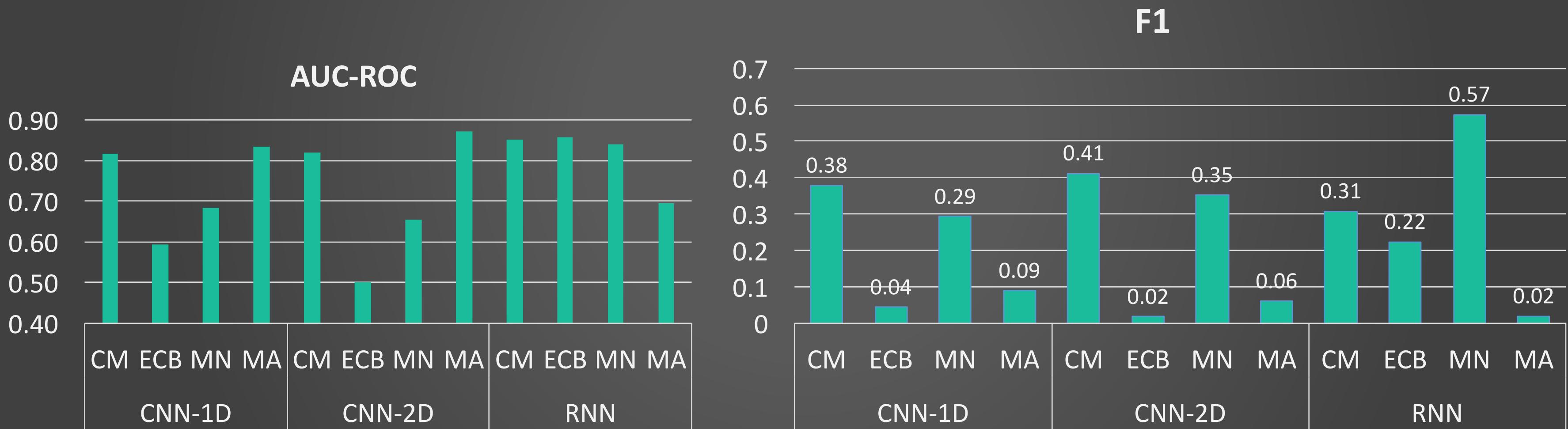
GRNET Super
computing facility

Each experiment using
1 GPU with 64 GB
memory

Results

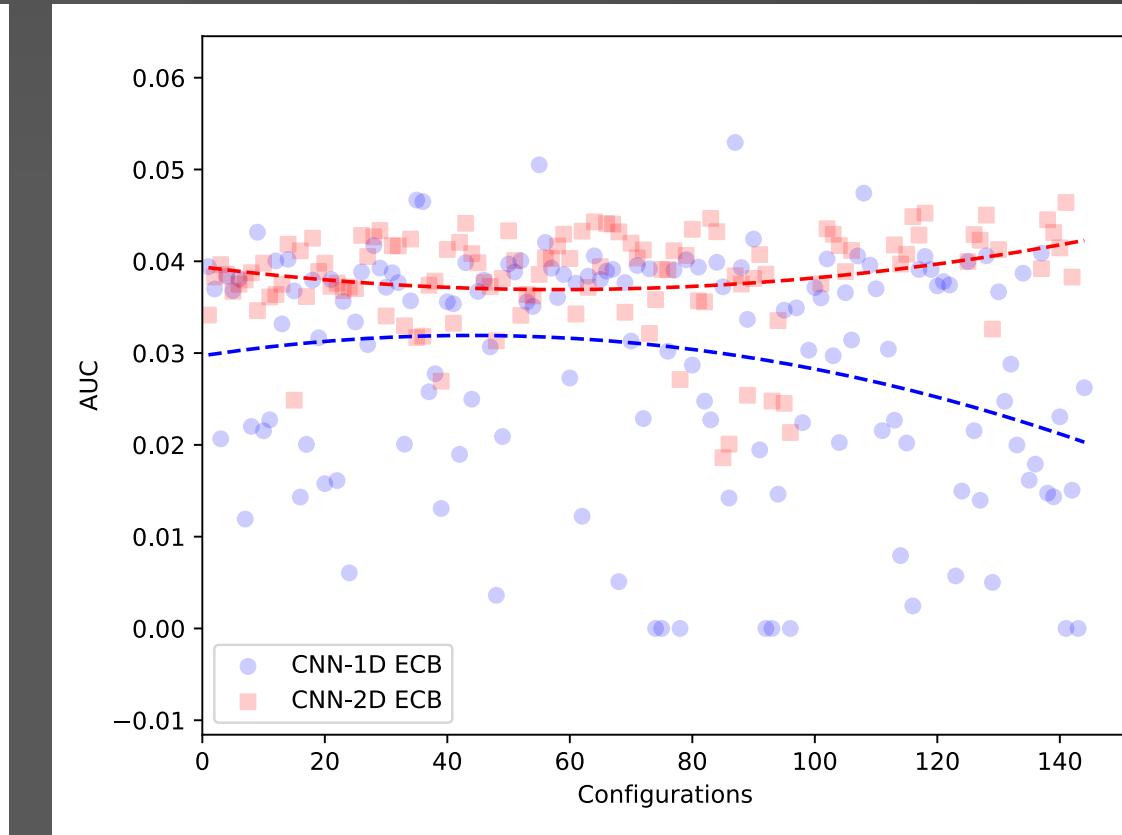
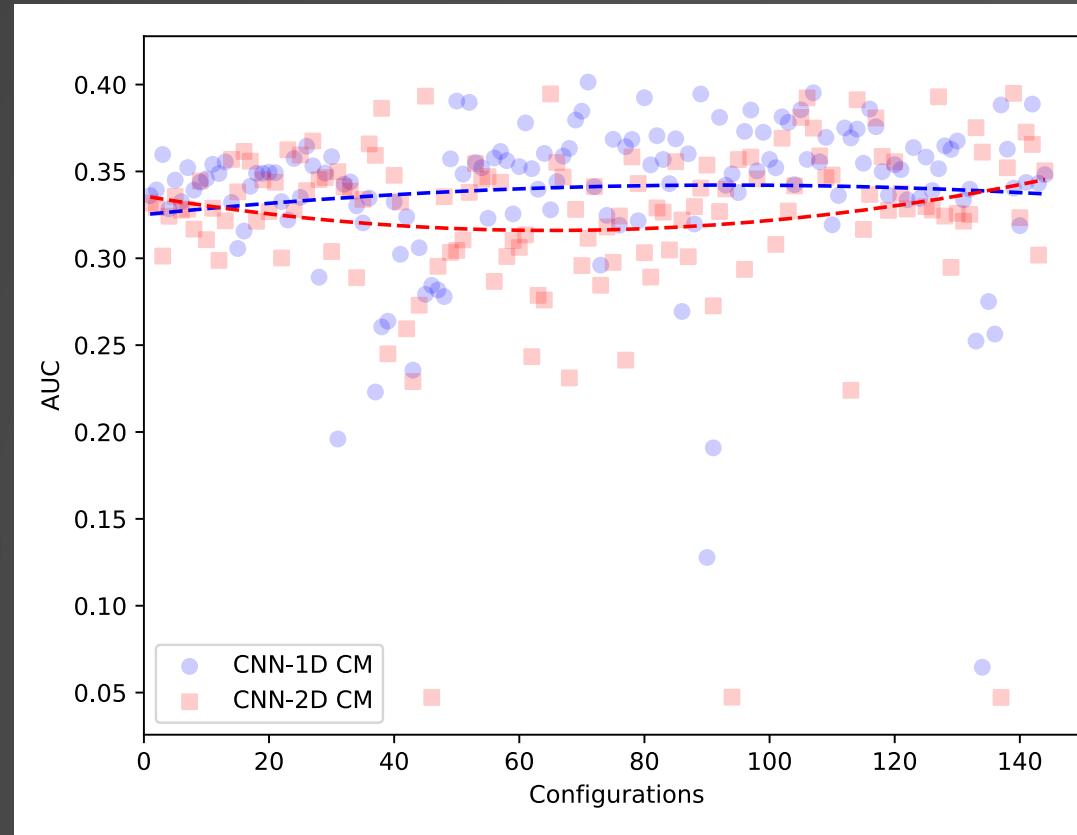


RQ1. Would it be possible to use deep learning methods to detect code smells?

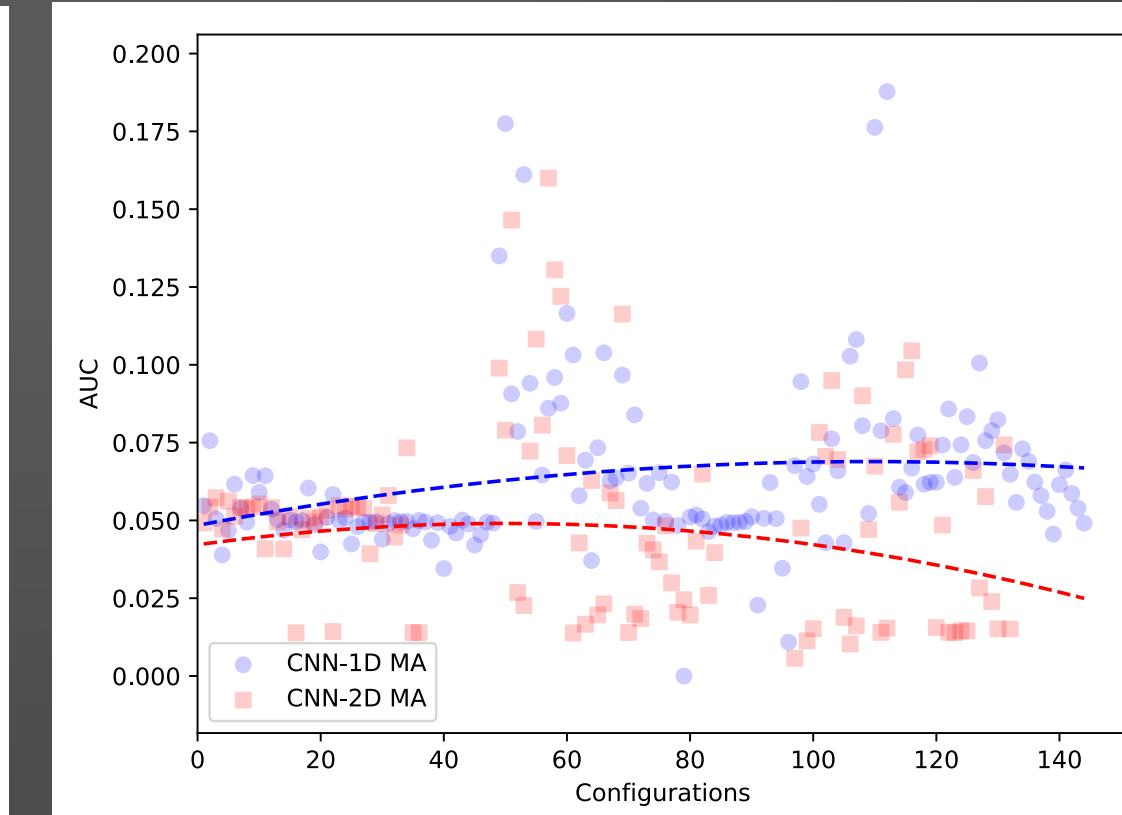
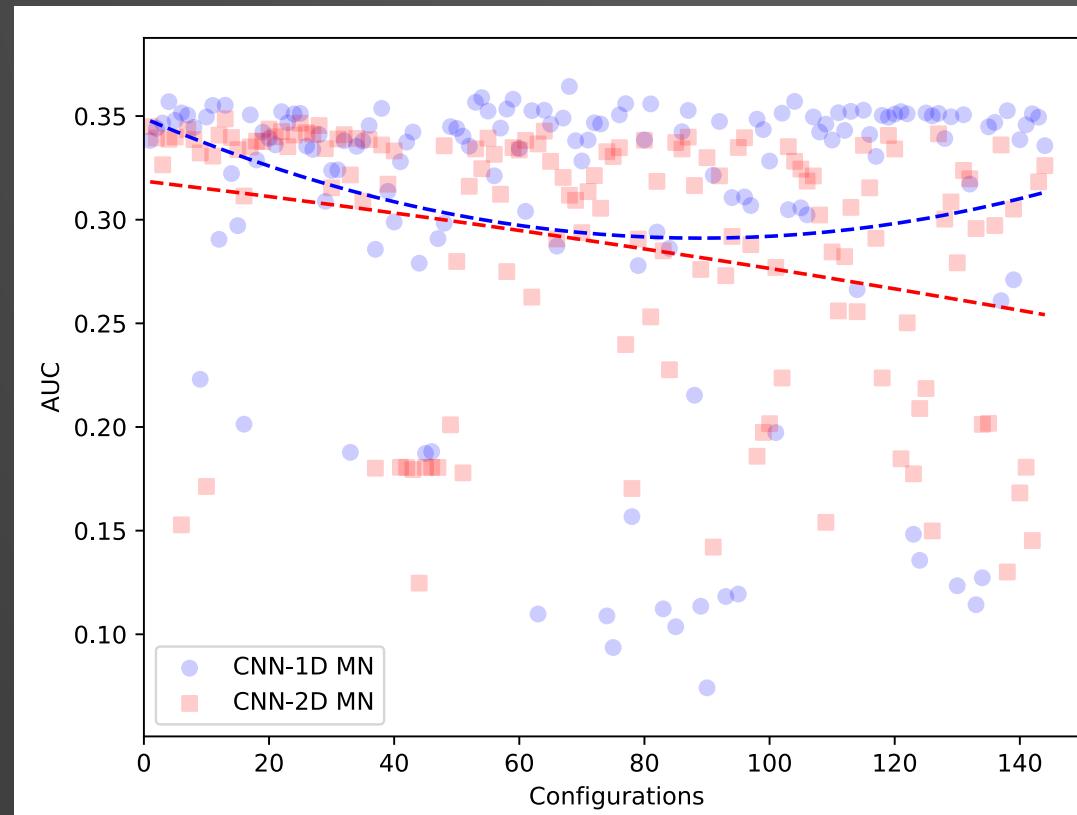


CNN-1D vs CNN-2D

CNN-1D (max) - 0.40
CNN-2D (max) - 0.39



CNN-1D (max) - 0.36
CNN-2D (max) - 0.35



CNN-1D (max) - 0.05
CNN-2D (max) - 0.04

CNN-1D (max) - 0.18
CNN-2D (max) - 0.16

CNN vs RNN

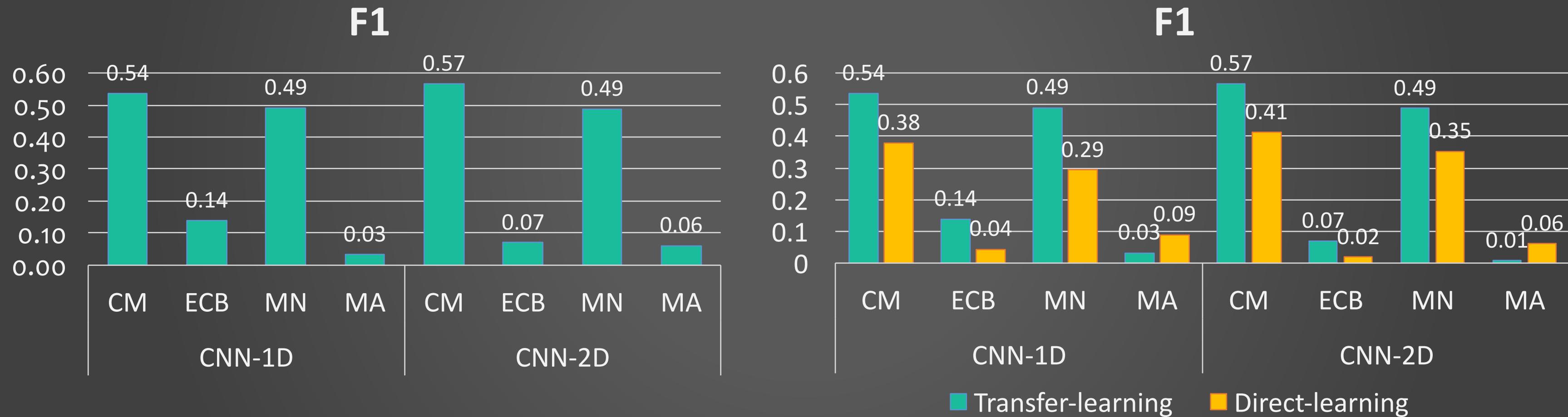
	RNN and CNN-1D	RNN and CNN-2D
CM	-22.94	-33.81
ECB	80.23	91.94
MN	48.96	38.58
MA	-349.12	-205.26

Difference in percentage; comparing max F1

Are more deep layers always good?

	Layers	CM	ECB	MN	MA
CNN- 1D	1	0.36	0.05	0.36	0.08
	2	0.40	0.05	0.36	0.18
	3	0.40	0.05	0.36	0.19
CNN- 2D	1	0.39	0.04	0.35	0.07
	2	0.39	0.04	0.34	0.16
	3	0.39	0.05	0.34	0.10
RNN	1	0.34	0.21	0.48	0.28
	2	0.36	0.24	0.48	0.22
	3	0.37	0.23	0.48	0.20

RQ2: Is transfer-learning feasible in the context of detecting smells?



It is feasible to make the deep learning model learn to detect smells

Transfer-learning is feasible.

Conclusions

Improvements – many possibilities

- Performance
- Add more smells – different kinds

Relevant links



Source code and data

<https://github.com/tushartushar/DeepLearningSmells>



Smell detection tool

Java - <https://github.com/tushartushar/DesigniteJava>

C# - <http://www.designite-tools.com>



CodeSplit

Java - <https://github.com/tushartushar/CodeSplitJava>

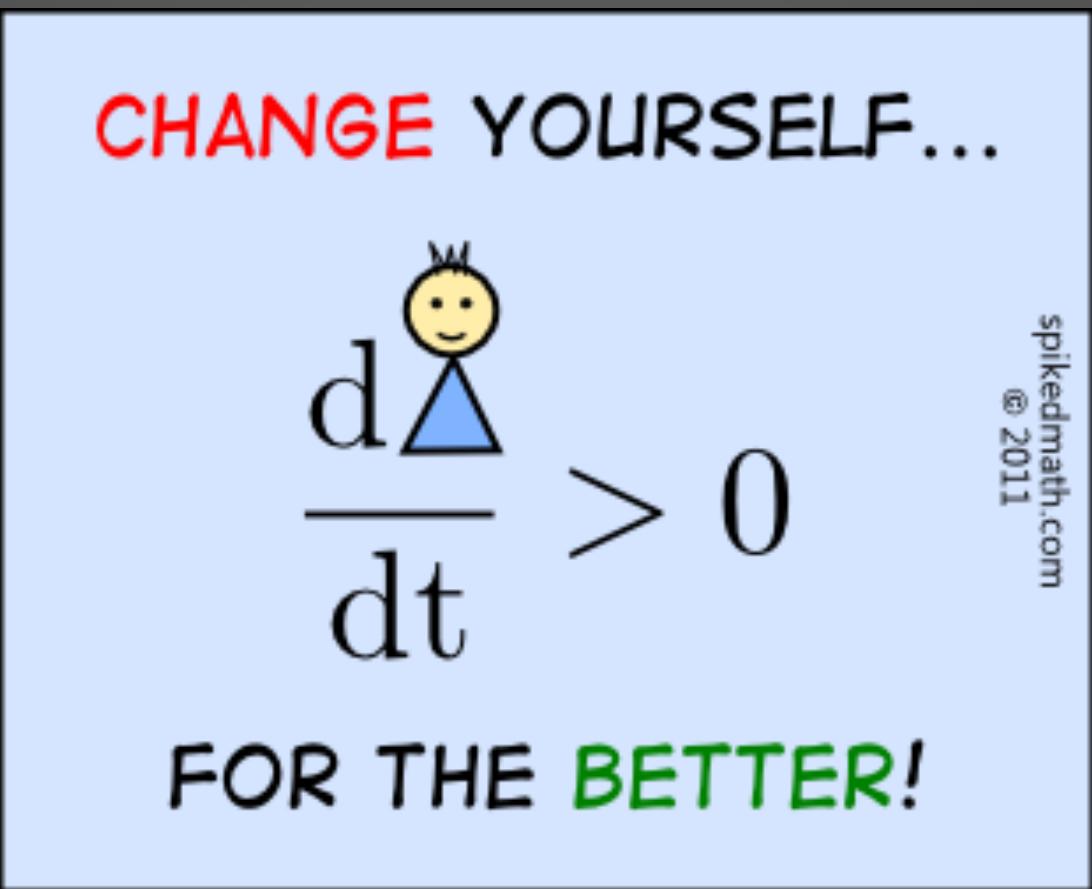
C# - <https://github.com/tushartushar/DeepLearningSmells/tree/master/CodeSplit>



Tokenizer

<https://github.com/dspinellis/tokenizer>

Thank you!!



Courtesy: spikedmath.com