

*Topic :*  
**Deep learning based  
vulnerability detection in source code**

*Researcher :* **Dyna Soumhane Ouchebara**  
Working on CyberExcellence project from Feb 2024 to Dec 2025  
PHD student at UMONS (2024-2027)

*Supervisor :* **Pr. Stéphane Dupont**



## Vulnerability detection in source code

### Vulnerability ?

A **flaw** in the code, which **when exploited**, can cause a **violation of the security policy**

### Vulnerability detection ?

**Detect vulnerabilities** in the code **at early stages of development** in order to fix them (preventive approach)

## Vulnerability detection in source code

### Example : SQL injection

#### - Code :

```
$user_id = $_GET['user_id'];  
$sql_query = "SELECT * FROM users WHERE id = '$user_id'";  
$result = mysqli_query($connexion, $sql_query);
```

→ Vulnerability

#### - Attack :

http://website.com/?user\_id=' OR '1'='1'

#### - Consequence :

SELECT \* FROM users WHERE id = " OR '1'='1' → always "True"  
→ Select all users data from database

## Conventional techniques

Automated, **but...**

Manual inspection



Static analysis



Dynamic analysis



Reliable



Time+Effort  
consuming



-Scalable  
-Early stage



-Can't spot  
runtime  
vulnerabilities



Spot runtime  
vulnerabilities



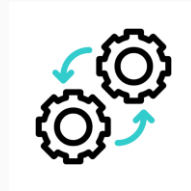
Resource  
intensive

# Deep learning

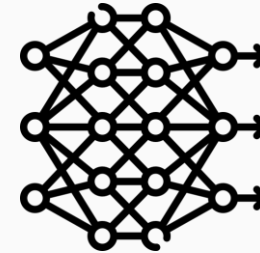
Labeled Data

Code	Label
class Person {...}	Vulnerable
def calculate {...}	Non vulnerable
...	...

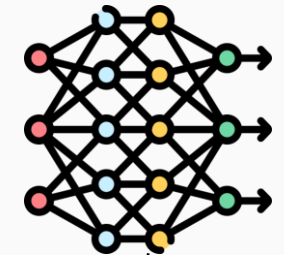
*Training*



Neural network



Trained model



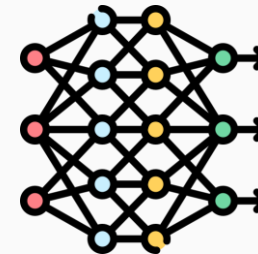
New data

Code	Label
def get_url {...}	?

*Prediction*



Trained model



Label

Label  
Vulnerable

(2016) Deep Learning : Code as Natural Language (*Text*)

# VulDeePecker: A Deep Learning-Based System for Vulnerability Detection

(5 Jan 2018)

Zhen Li<sup>\*†</sup>, Deqing Zou<sup>\*†‡</sup>, Shouhuai Xu<sup>§</sup>, Xinyu Ou<sup>\*</sup>, Hai Jin<sup>\*</sup>,  
Sujuan Wang<sup>\*</sup>, Zhijun Deng<sup>\*</sup> and Yuyi Zhong<sup>\*</sup>

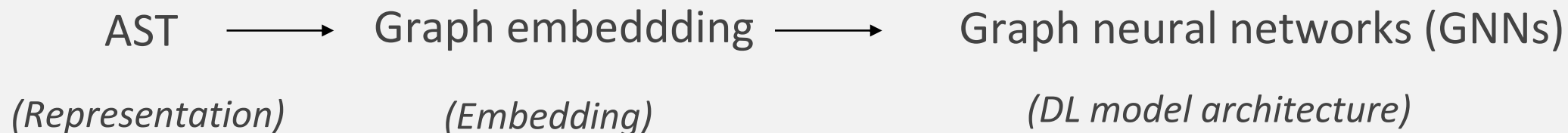
Tokens  $\longrightarrow$  Word2vec  $\longrightarrow$  Recurrent neural networks (RNNs)  
(Representation) (Embedding) (DL model architecture)

(2019) Deep Learning : Code as a Structure (*Graph/Tree*)

***Devign: Effective Vulnerability Identification by  
Learning Comprehensive Program Semantics via  
Graph Neural Networks***

(8 Dec 2019)

Yaqin Zhou<sup>1</sup>, Shangqing Liu<sup>1,\*</sup>, Jingkai Siow<sup>1</sup>, Xiaoning Du<sup>1,\*</sup>, and Yang Liu<sup>1</sup>





(2021) Deep Learning : Code as Text again, but

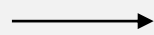
# LineVul: A Transformer-based Line-Level Vulnerability Prediction

(17 Oct 2022)

Michael Fu

Chakkrit Tantithamthavorn\*

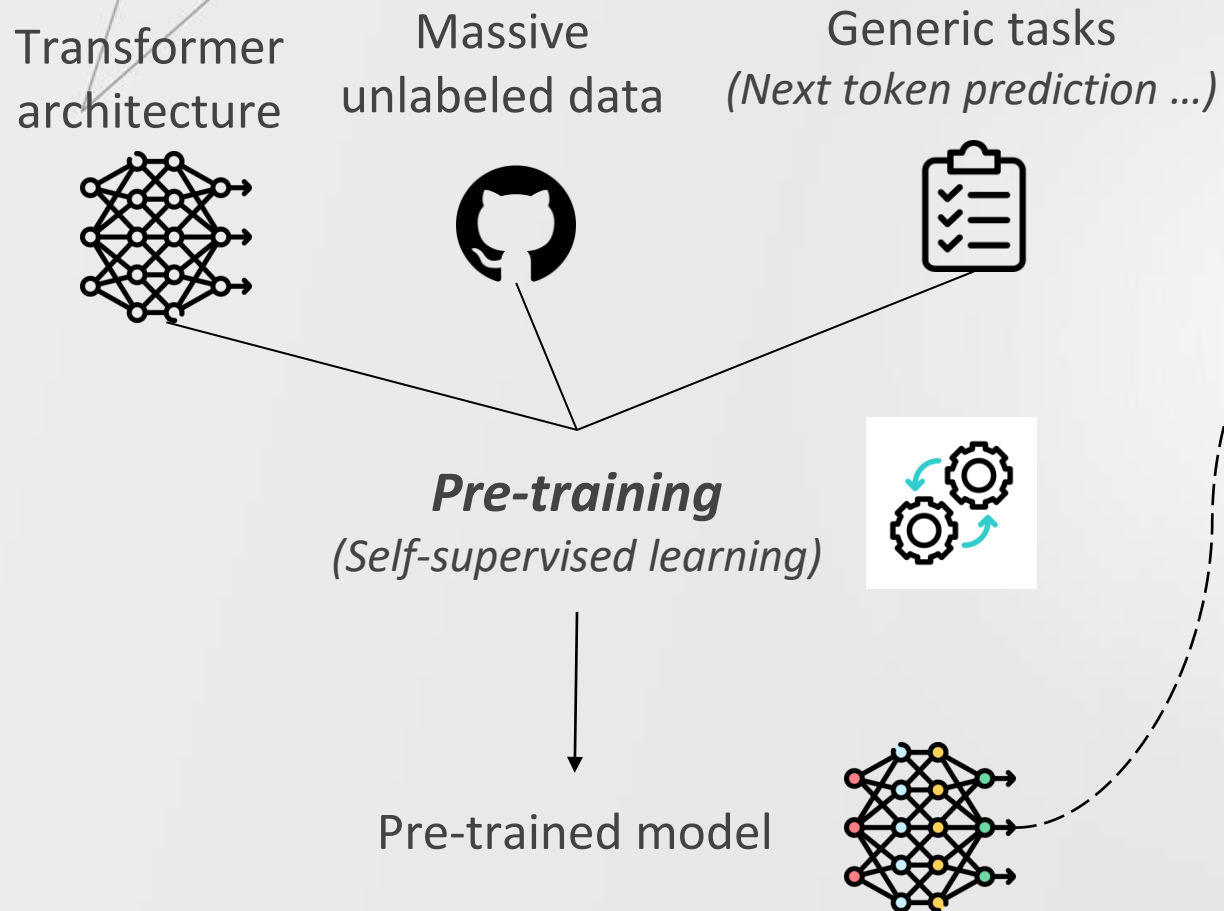
Tokens  
(Representation)



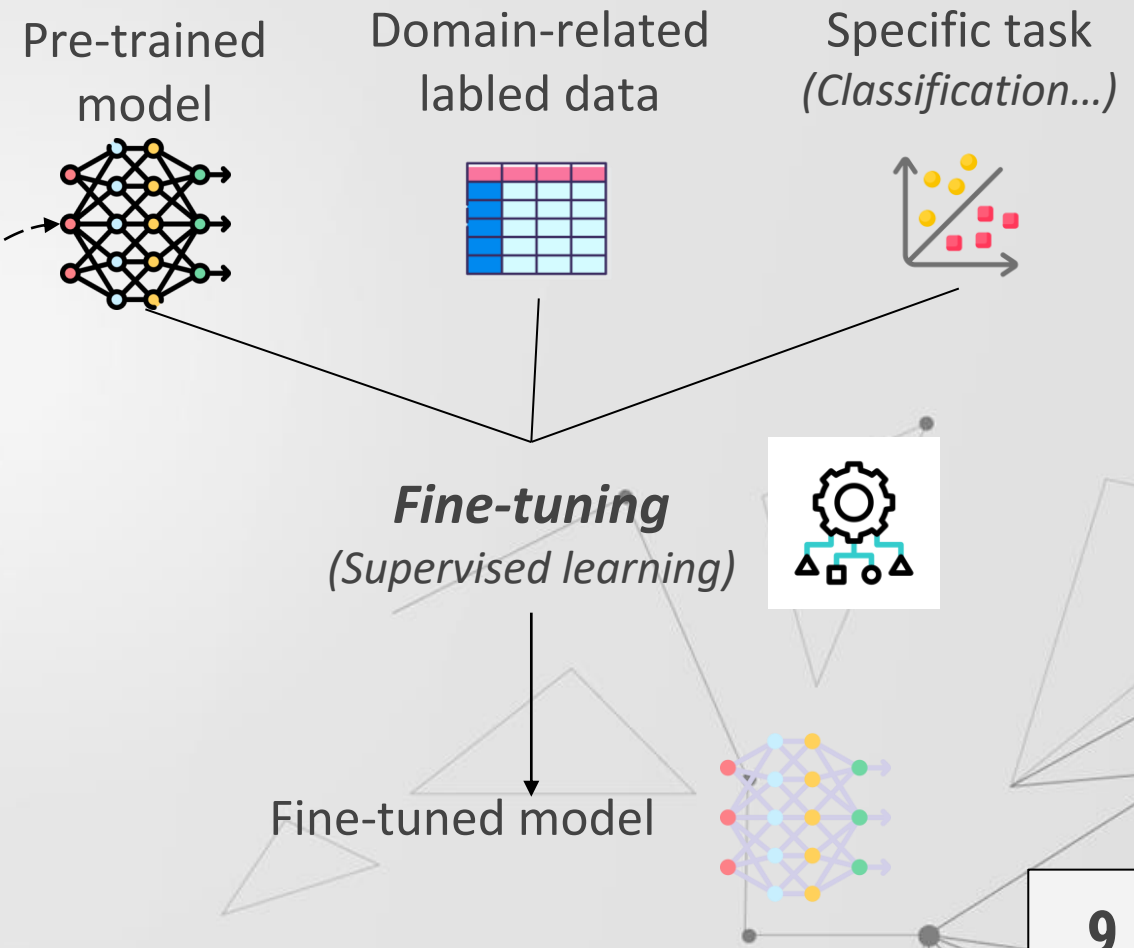
Transformer-based Pre-trained model  
(Embedding + DL model architecture)

# Transformer based pre-trained models

## Phase 1 : Pre-training

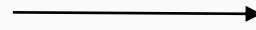


## Phase 2 : Fine-tuning



## Challenges

### 1. Erronous datasets



Dataset	Correct Label
<b>CodeXGLUE [8, 19]</b>	24.0 %
<b>VulnPatchPairs [20]</b>	36.0 %
<b>BigVul [9]</b>	25.0* %
<b>CrossVul [11]</b>	47.8* %
<b>CVEFixes [10]</b>	51.7* %

### 2. Generalization capability

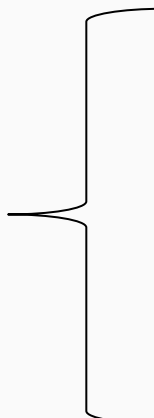
### 3. Detection granularity

### 4. Explainability



## Research directions

### Face the challenges



- Better datasets
- Better generalization capability
- Best detection granularity
- Explainability of the models

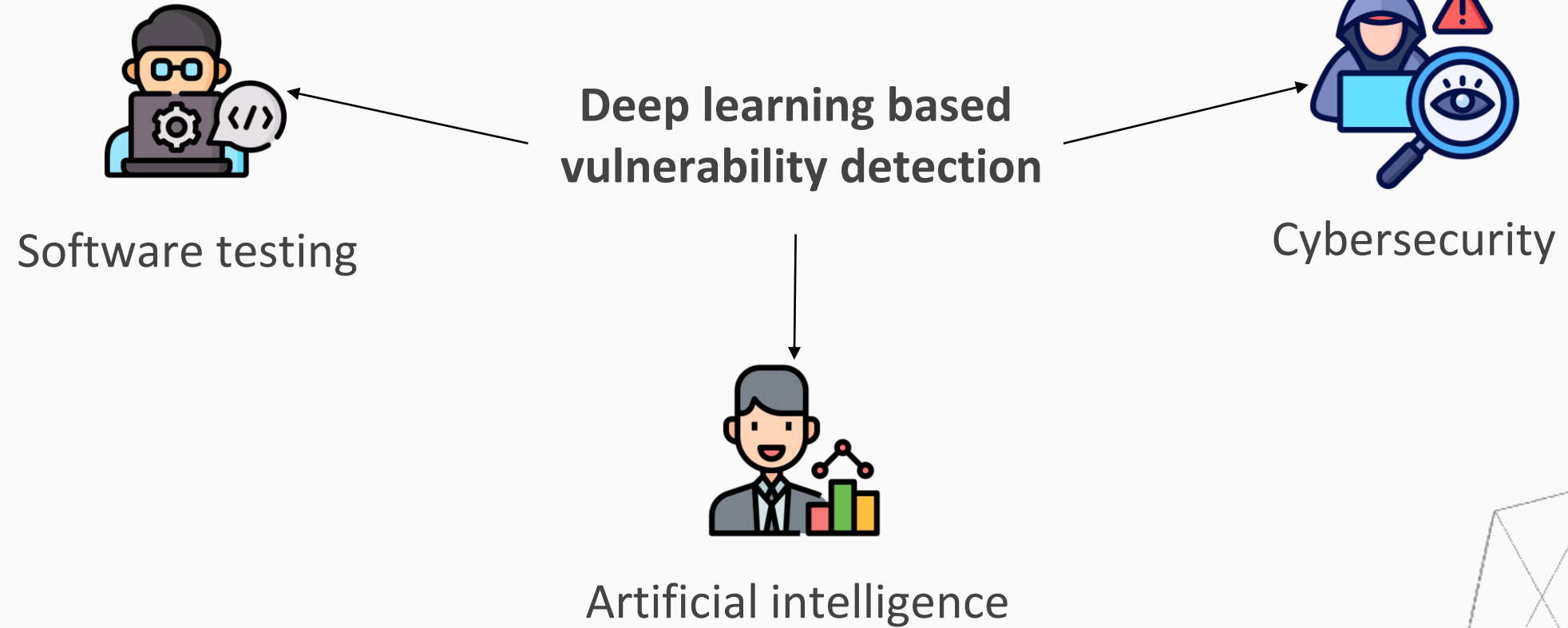
Better representation and/or embedding of code

Better DL architectures

More suited pre-training tasks

Better fine-tuning strategies

## Final note





# Thank you for your attention !

## References :

1. Li Z, Zou D, Xu S, Ou X, Jin H, Wang S, Deng Z, Zhong Y (2018) VulDeePecker: a deep learning-based system for vulnerability detection
2. Zhou Y, Liu S, Siow J, Du X, Liu Y (2019) Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks
3. Michael Fu and Chakkrit Tantithamthavorn (2022) Linevul: A transformer-based line-level vulnerability prediction