# Exploring the Potential of LLMs for Code Deobfuscation

**David Beste**, *Grégoire Menguy, Mario Fritz, Antonio Emanuele Cinà, Thorsten Holz, Thorsten Eisenhofer and Lea Schönherr*

DIMVA '25 | 10th July 2025

# Motivation

- Obfuscation used by malware authors
- Need for deobfuscation
- Significant success of LLMs in code-related tasks
- Can LLMs aid deobfuscation in a universal way?

https://blogs.vmware.com/security/wp-content/uploads/sites/26/2020/05/fig1_fn_blowfish_init_before_trim.png
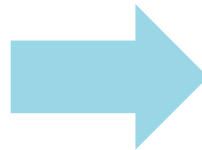
# Research Questions

1. Can LLMs deobfuscate state-of-the-art code obfusca-tion transformations?

2. Can LLMs deobfuscate code in a real-world scenario where multiple transformations are chained?

3. How much is memorization affecting the performance?

# Methodology: Dataset

- Used **Exebench dataset**

  - Dataset of millions of C functions crawled from GitHub

  - According to software complexity metrics representative of real-world code

  - Includes test I/O pairs for correctness checks

- New **deobfuscation dataset** with around 30000 samples

```
 1   __inline static void strtoupper(char *s) {
 2     char *c;
 3     c = s;
 4     while (*c) {
 5       if ((int )*c >= 97) {
 6         if ((int )*c <= 122) {
 7           *c = (char )(((int )*c - 97) + 65);
 8         }
 9       }
10       c ++;
11     }
12     return;
13   }
```

```
 1   void _xa(char *_k0, long _k1) {
 2     char *_k2 ;
 3     unsigned long _k3 ;
 4     int _k4 ;
 5     _k3 = 1UL;
 6     while (1) {
 7       switch (_k3) {
 8       case 4UL: ;
 9       if (97 <= (int )*_k2) {
10         _k3 = 0UL;
11       } else {
12         _k3 = 3UL;
13       }
14       break;
15       case 0UL: ;
16       if (((unsigned int )(((int )*_k2 | -123) & (((int
    ↪ )*_k2 ^ 122) | ~ (122 - (int )*_k2))) >> 31U) & 1U) {
17         _k3 = 7UL;
18       [...]
```

# Methodology: Obfuscation

- **Tigress** C obfuscator
  - State-of-the-art C obfuscator
  - Chose five transformations

- Alter different aspects of the code
- Transformations
  - Encode Arithmetic
  - Encode Branches
  - Flatten
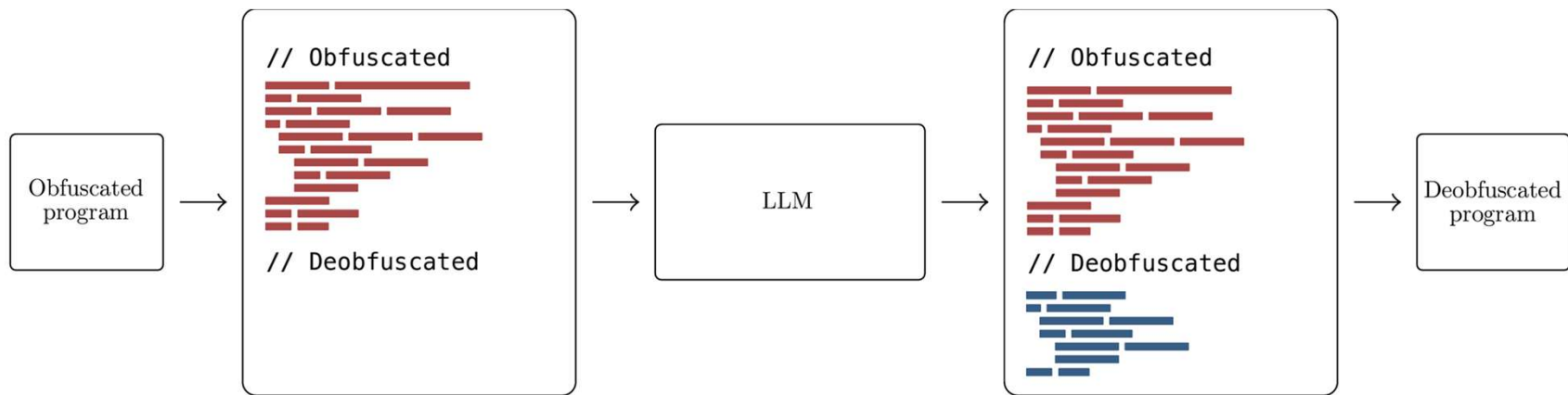  - Opaque Predicates
  - Randomize Arguments

# Methodology: Models and Baselines

- Fine-tuned two local **open-source LLMs** on these samples
- Performed a memorization test on hand-selected samples
- Evaluated on a test set and compared to **GPT-4** in a zero-shot setting
- Also used **Clang** as a sanity check

| Name | Size | Open Access | Instruction Tuned | Coding Specialist |
|---|---|---|---|---|
| DeepSeek Coder | 6.7B | ✓ | ✓ | ✓ |
| Code Llama | 7B | ✓ | ✗ | ✓ |
| GPT-4 | n/a | ✗ | ✓ | ✗ |

# Methodology: Pipeline

# Methodology: Deobfuscation Performance Formula

- Comparison of **original** ($C_{Orig}$), **obfuscated** ($C_{Obf}$) and LLM **deobfuscated** ($C_{Deobf}$) versions' complexity
- Formula computes the "point" at which the LLM returned sample lays between original and obfuscated
  - 0 -> Failure
  - 1 -> Complete Success

$$P_{Deobf} = 1 - \frac{C_{Deobf} - C_{Orig}}{C_{Obf} - C_{Orig}}$$
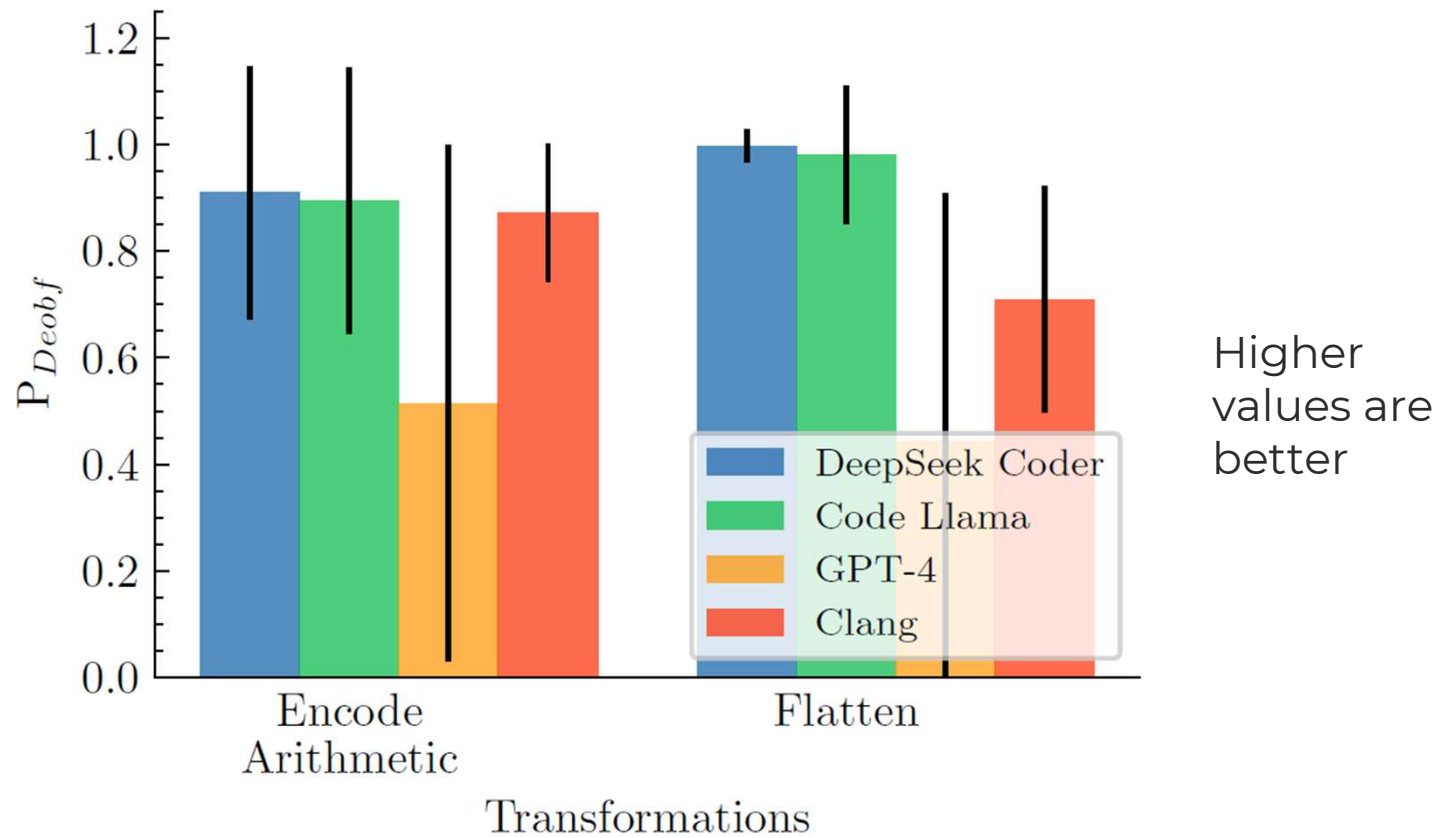
- **Only semantically correct samples are evaluated**

# Methodology: Complexity Metrics

- Metrics used: **Halstead Program Length**
  - Halstead metrics has been shown to reflect human perceived program difficulty
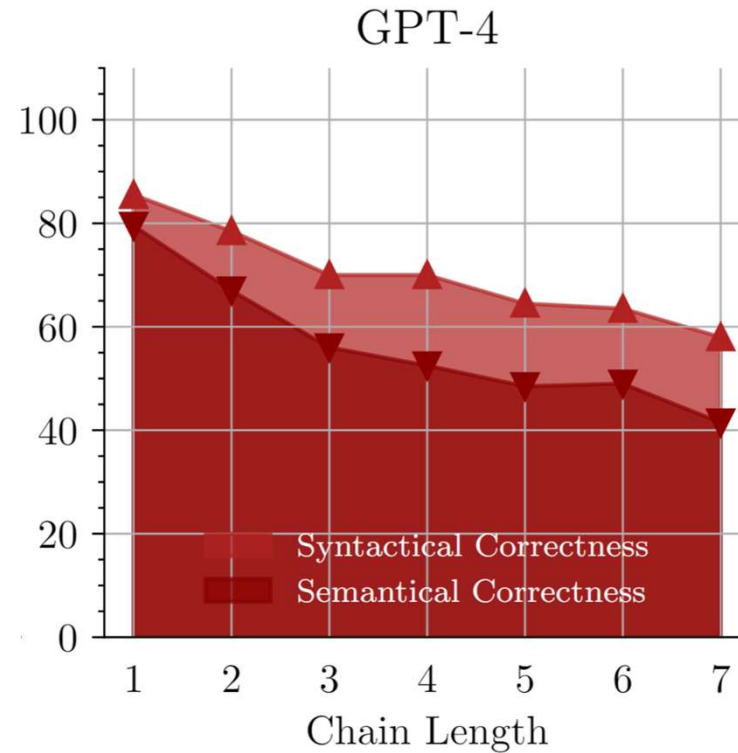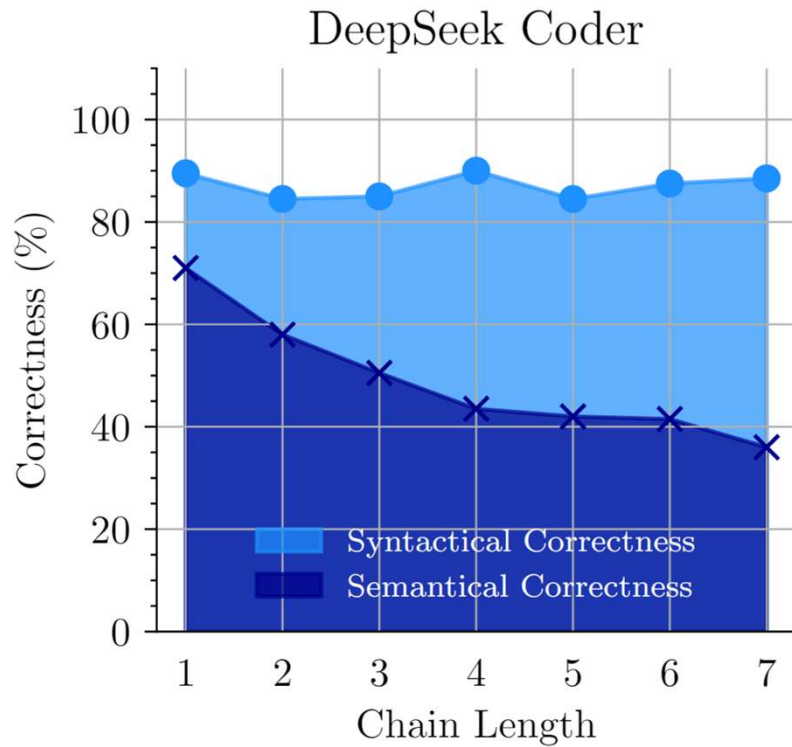- **Semantical correctness** check (I / O samples)

Higher values are better

# Methodology: Chained Transformations

- **Single** transformations and **chains**
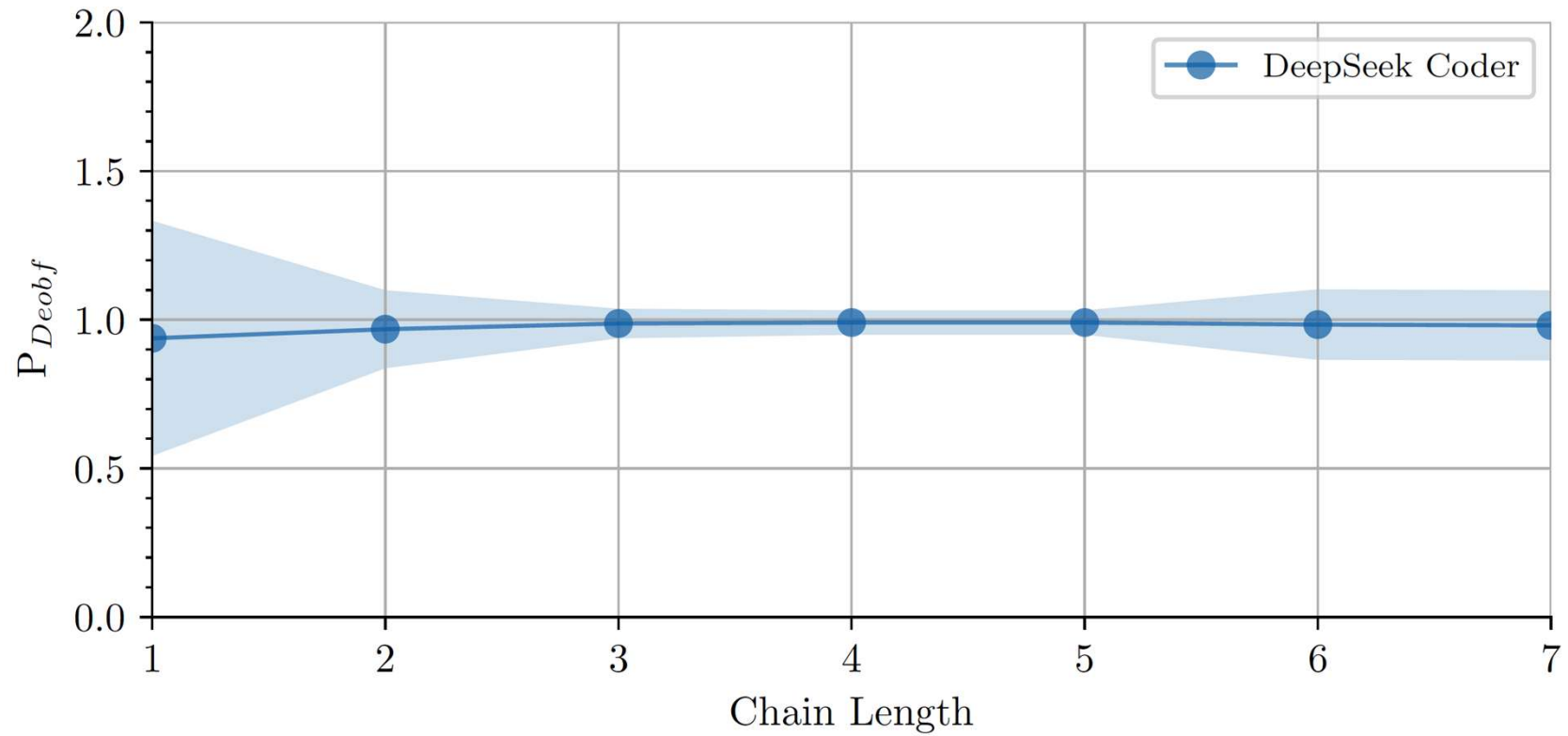  - Five for training
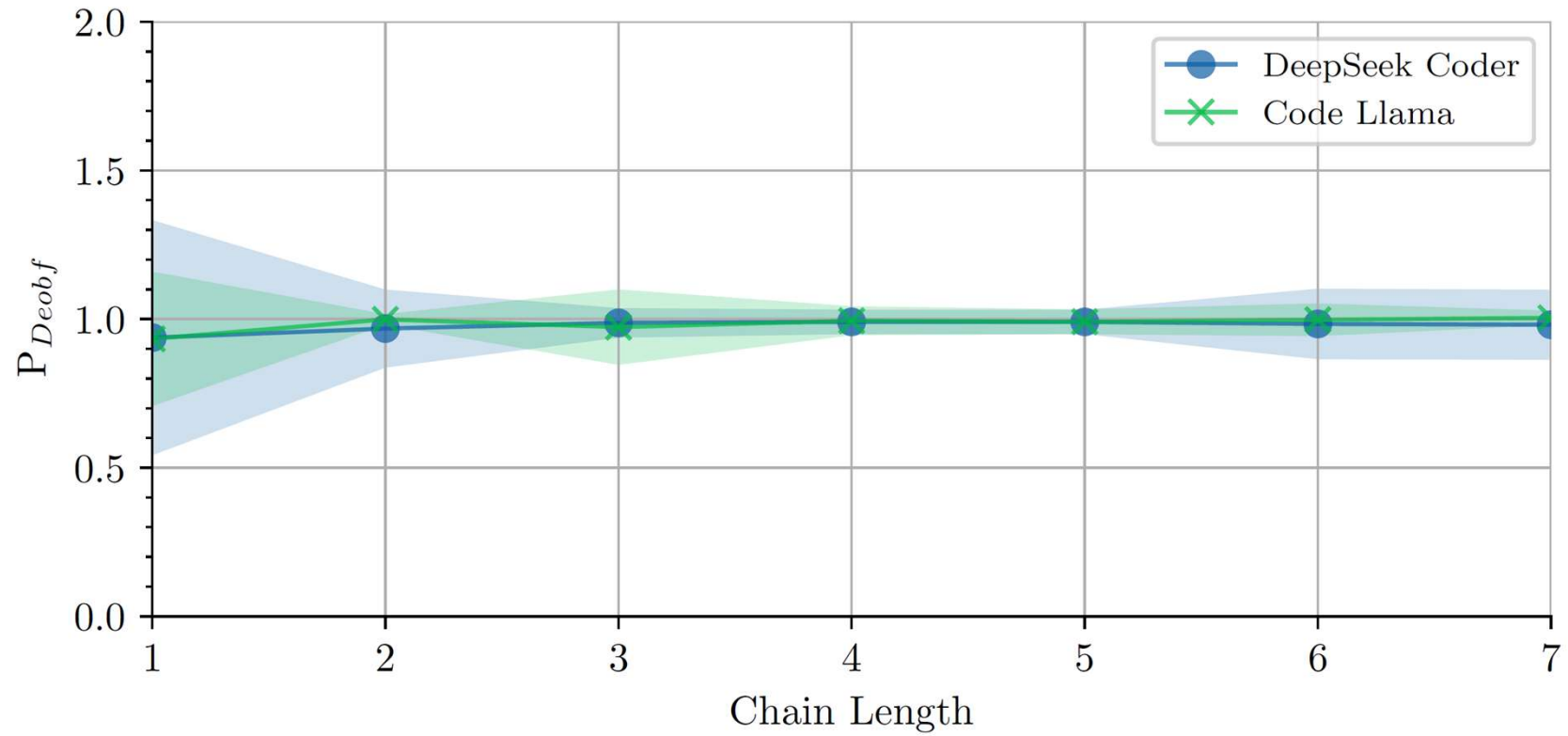  - Seven for evaluation

# Evaluation: Chained Correctness
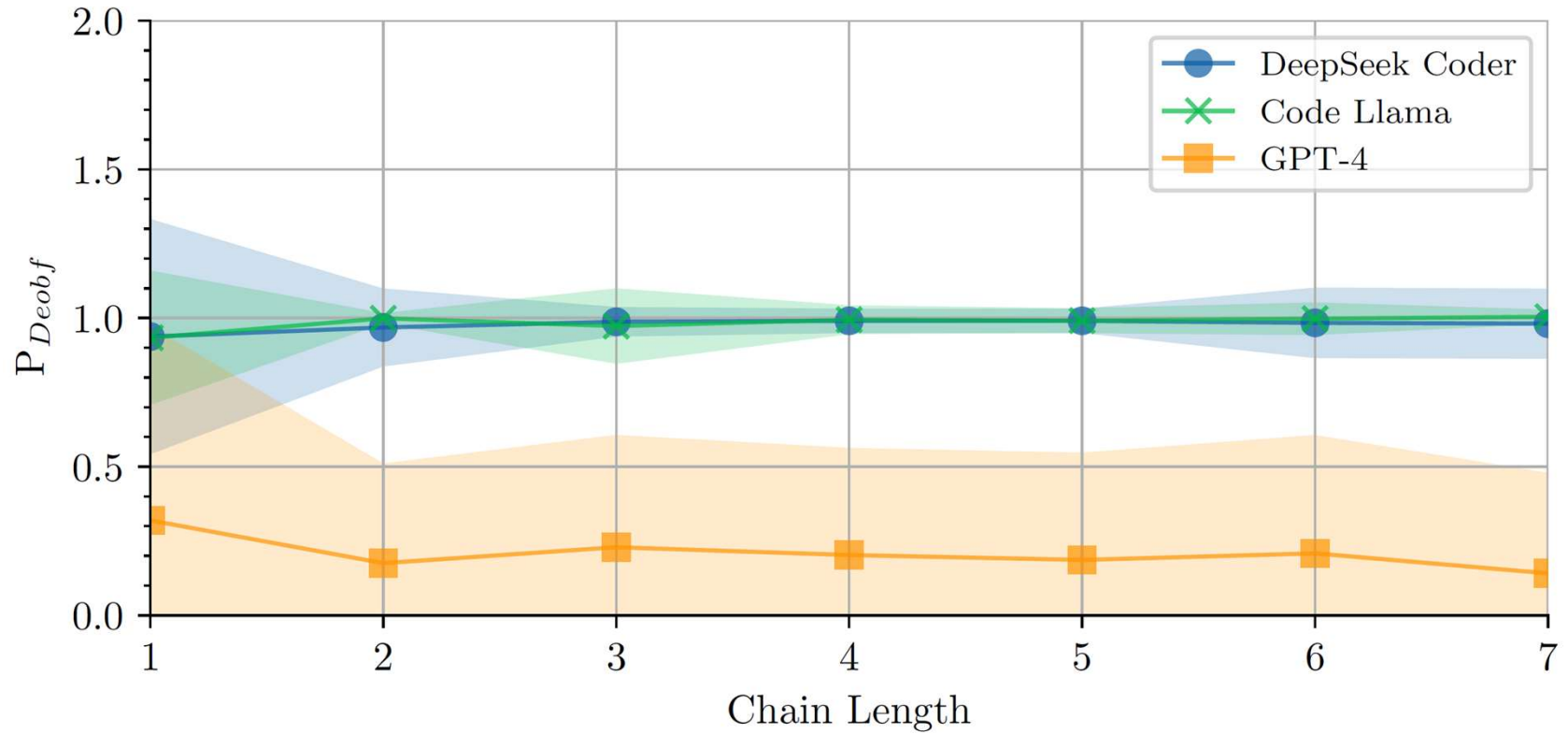
# Evaluation: Chained Deobfuscation Performance

# Evaluation: Chained Deobfuscation Performance

# Evaluation: Chained Deobfuscation Performance

# Methodology: Memorization

- Training on public code makes LLMs prone to memorization
- Are deobfuscated samples memorized due to bias?
- **Experiment:** Change constants, then obfuscate again
- Check if LLM deobfuscates with the changed constants
    - If not, sample likely memorized

# Evaluation: Memorization

- Semantical plausibility unimportant, only if the LLM correctly identifies the correct constants

- **Results:** Memorization was not a significant issue

```
void temp_init(double *temps )
{
  int t ;
  double dT ;

  {
  t = 0;
  while (t < 10) {
    dT = 5.0 / (double )10;
    *(temps + t) = 5.0 - (double )t * dT;
    t ++;
  }
  return;
  }
}
```

```
void temp_init(double *temps )
{
  int t ;
  double dT ;

  {
  t = -2;
  while (t < 28) {
    dT = 49.37 / (double )848.88;
    *(temps + t) = 22.88 - (double )t * dT;
    t ++;
  }
  return;
  }
}
```
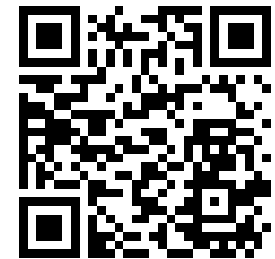
# Summary

- Trained and evaluated two **LLMs for deobfuscation** tasks
- **Fine-tuning** small coding models shows **promising results** for deobfuscation
- Challenges with **functional correctness -> larger models** very likely to reduce this problem
- **Memorization** was **non-significant** in our test -> Indication of genuine code understanding capabilities of LLMs

- **Thank you for your attention!**

https://github.com/DavidBeste/llm-code-deobfuscation

# Evaluation