**Distinguished Paper Award**

# ReSym: Harnessing LLMs to Recover Variable and Data Structure Symbols from Stripped Binaries

**Danning Xie**, Zhuo Zhang, Nan Jiang, Xiangzhe Xu, Lin Tan, Xiangyu Zhang

**PURDUE**
UNIVERSITY

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11 }
```

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10  }
11 }
```

**Compiler**

**Binary File**
(with debugging info.)

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1   void ixp_pstrings
2       (IxpMsg *msg, ushort num){
3    ushort len;
4    uint size;
5    uchar *s;
6    if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11  }
```
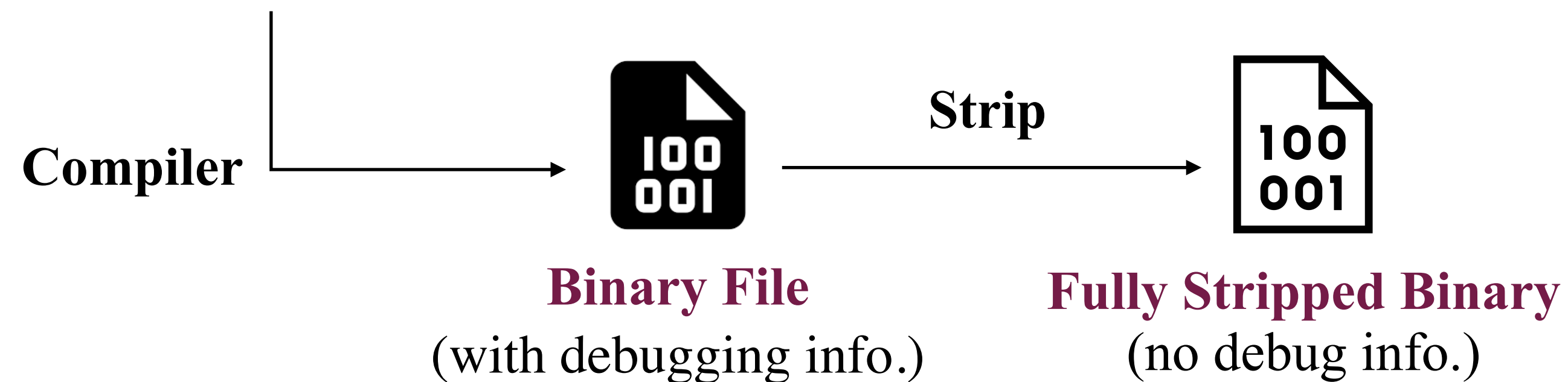
**Compiler**



**Binary File**
(with debugging info.)

Debugging Information: locations, sizes, and layout of functions and objects

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11  }
```

**Compiler**

**Strip**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

Debugging Information: locations, sizes, and layout of functions and objects

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11  }
```
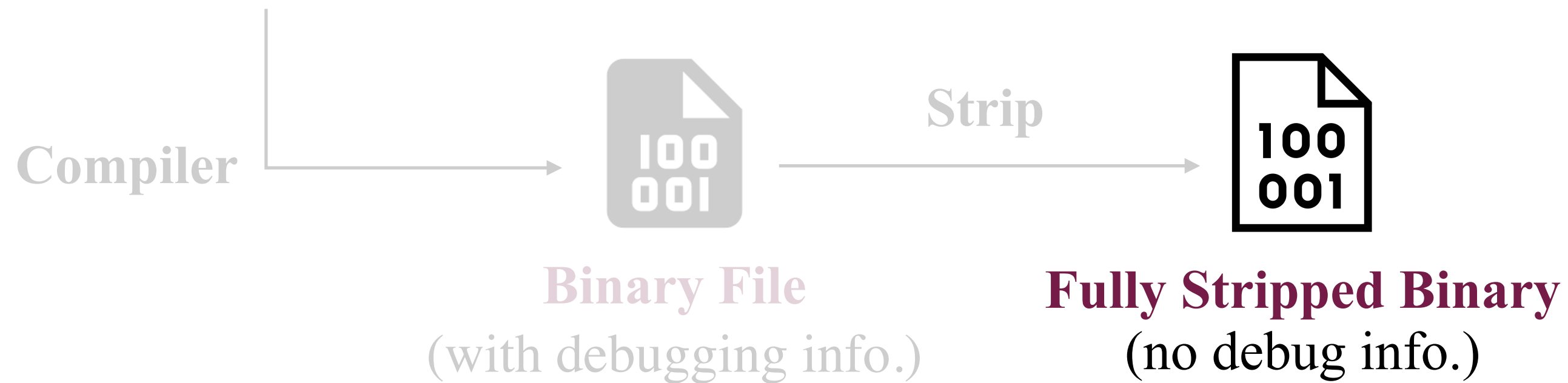
**Compiler**

**Strip**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

Debugging Information: locations, sizes, and layout of functions and objects

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2     (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11 }
```

**Compiler**

**Strip**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

Debugging Information: locations, sizes, and layout of functions and objects

**Understanding stripped binaries is essential to ensure software security.**

# Background: Stripped Binary and Decompiled Code

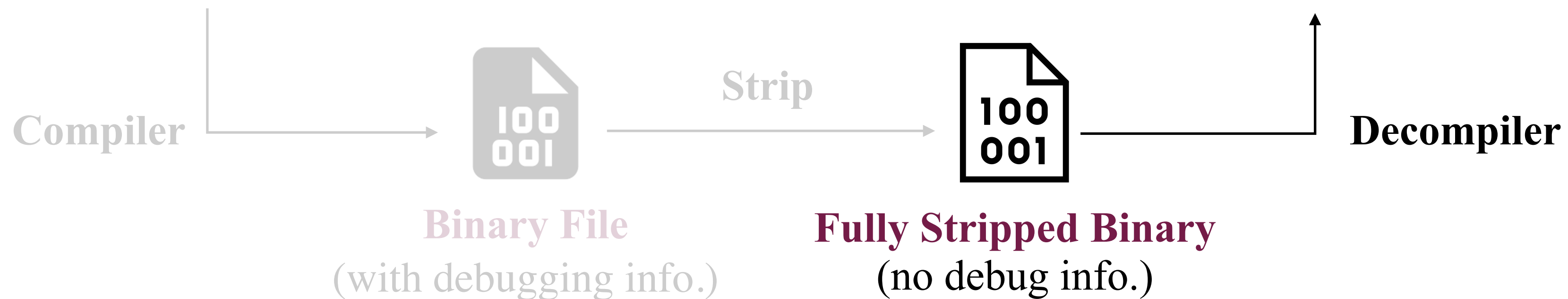**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11 }
```

**Compiler**

**Strip**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- Decompilation aims to recover the source code form of a binary executable.

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
 1  void ixp_pstrings
 2      (IxpMsg *msg, ushort num){
 3   ushort len;
 4   uint size;
 5   uchar *s;
 6   if(msg->mode == 1){
 7    s = msg->pos;
 8    size = len * num;
 9    emalloc(size);
10   }
11  }
```
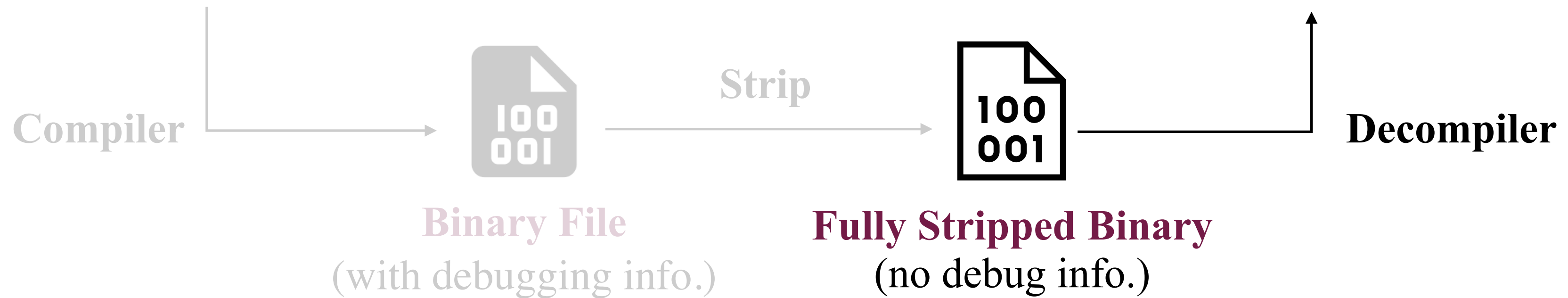
**Decompiled Code**

```
 1  unsigned int64 sub_404056
 2      (int64 a1, int16 a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if (*(int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 8);
 8    v4 = v3 * a2;
 9    sub_406BB9(v4);
10   }
11  }
```

**Compiler**

**Strip**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

**Decompiler**

- **Decompilation** aims to recover the source code form of a binary executable.

# Background: Stripped Binary and Decompiled Code
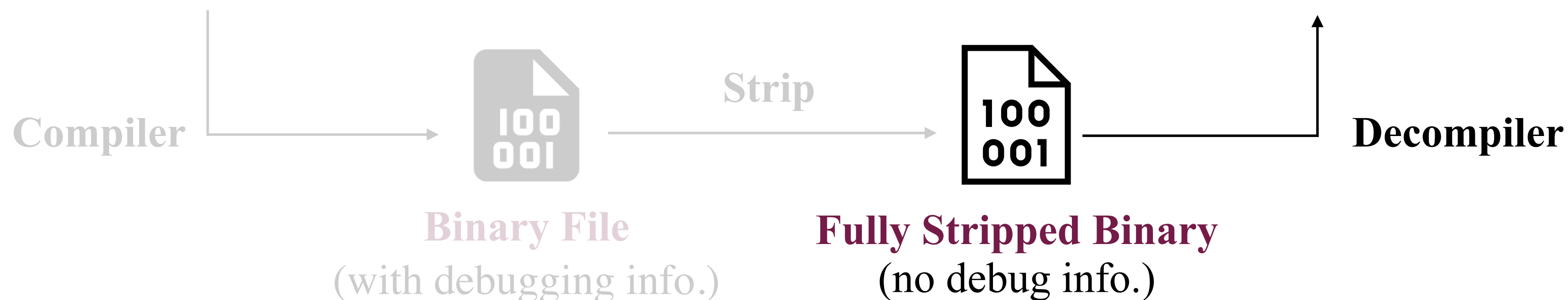
**Source Code**

```
 1  void ixp_pstrings
 2      (IxpMsg *msg, ushort num){
 3   ushort len;
 4   uint size;
 5   uchar *s;
 6   if(msg->mode == 1){
 7    s = msg->pos;
 8    size = len * num;
 9    emalloc(size);
10   }
11  }
```

**Decompiled Code**

```
 1  unsigned int64 sub_404056
 2      (int64 a1, int16 a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if (*(int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 8);
 8    v4 = v3 * a2;
 9    sub_406BB9(v4);
10   }
11  }
```

**Compiler**

**Strip**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- **Decompilation** aims to recover the source code form of a binary executable.

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10  }
11 }
```
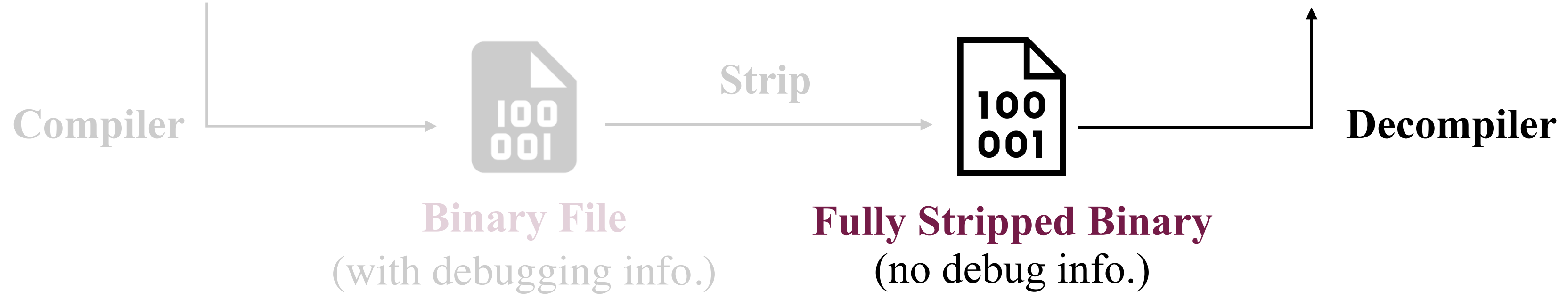
**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10  }
11 }
```

**Compiler**

**Strip**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- **Decompilation** aims to recover the source code form of a binary executable.

# Background: Stripped Binary and Decompiled Code
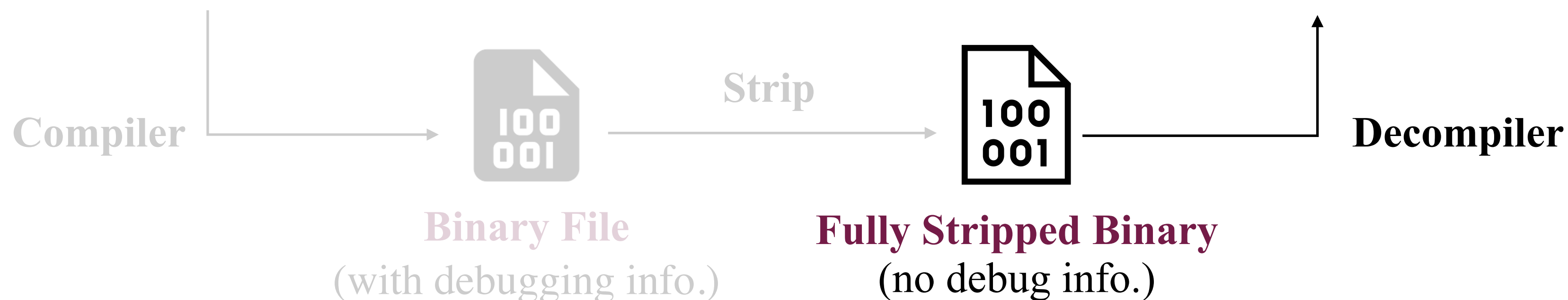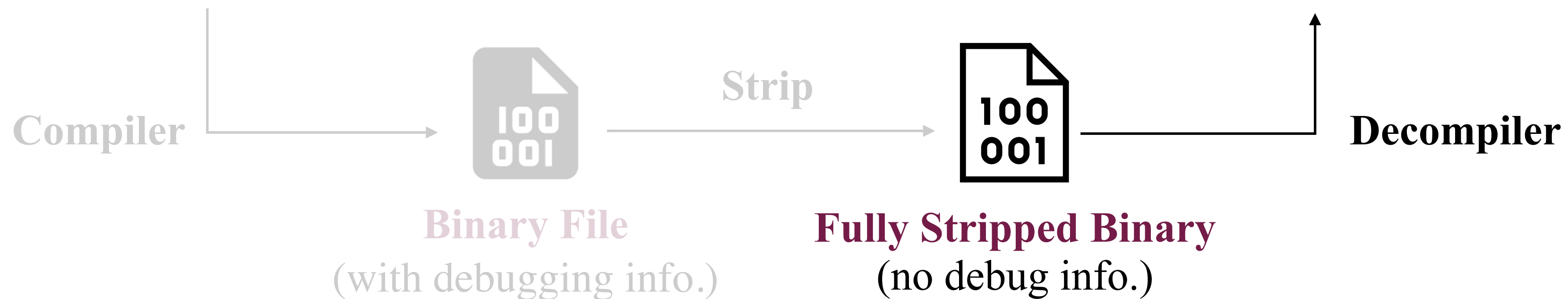
**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11 }
```

**Compiler**

**Strip**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- **Decompilation** aims to recover the source code form of a binary executable.

# Background: Stripped Binary and Decompiled Code
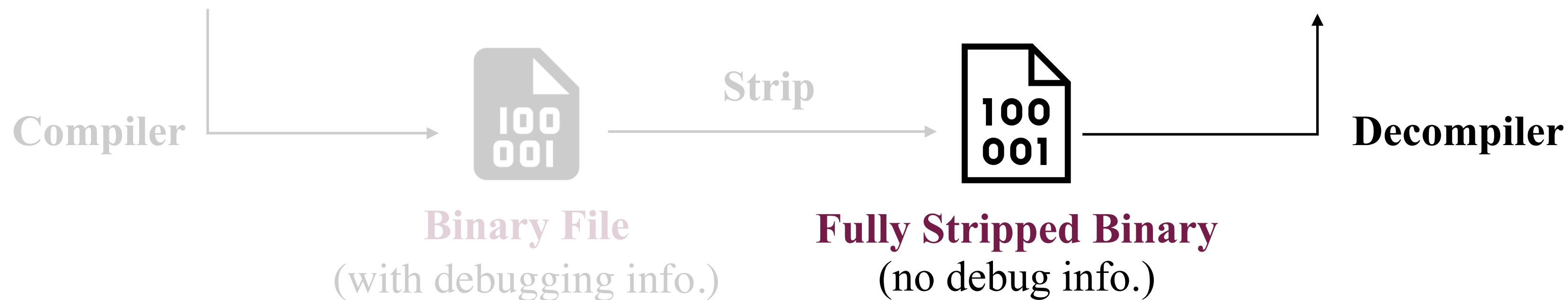
**Source Code**

```
1   void ixp_pstrings
2       (IxpMsg *msg, ushort num){
3    ushort len;
4    uint size;
5    uchar *s;
6    if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10    }
11  }
```

**Decompiled Code**

```
1   unsigned int64 sub_404056
2       (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10    }
11  }
```

**Strip**

**Compiler**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- Decompilation aims to recover the source code form of a binary executable

- Decompiled code from stripped binaries loses symbol information.

# Background: Stripped Binary and Decompiled Code

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10  }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10  }
11 }
```

**?**

**Compiler**

**Strip**

**Decompiler**

**Binary File**
(with debugging info.)

**Fully Stripped Binary**
(no debug info.)

- Decompilation aims to recover the source code form of a binary executable

- Decompiled code from stripped binaries loses symbol information.

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
 1  void ixp_pstrings
 2      (IxpMsg *msg, ushort num){
 3   ushort len;
 4   uint size;
 5   uchar *s;
 6   if(msg->mode == 1){
 7    s = msg->pos;
 8    size = len * num;
 9    emalloc(size);
10   }
11  }
```

**Decompiled Code**

```
 1  unsigned int64 sub_404056
 2      (int64 a1, int16 a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if (*(int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 8);
 8    v4 = v3 * a2;
 9    sub_406BB9(v4);
10   }
11  }
```

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10  }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10  }
11 }
```

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10  }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10  }
11 }
```

```
struct IxpMsg {
 char*    data;
 char*    pos;
 char*    end;
 _ixpuint size;
 _ixpuint mode;
};
```
**Ground Truth**

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
 1   void ixp_pstrings
 2      (IxpMsg *msg, ushort num){
 3    ushort len;
 4    uint size;
 5    uchar *s;
 6    if(msg->mode == 1){
 7     s = msg->pos;
 8     size = len * num;
 9     emalloc(size);
10    }
11   }
```

**Decompiled Code**

```
 1   unsigned int64 sub_404056
 2      (int64 a1, int16 a2){
 3    unsigned int16 v3;
 4    unsigned int v4;
 5    void *dest;
 6    if (*(int *)(a1 + 28) == 1){
 7     dest = *(void **)(a1 + 8);
 8     v4 = v3 * a2;
 9     sub_406BB9(v4);
10    }
11   }
```
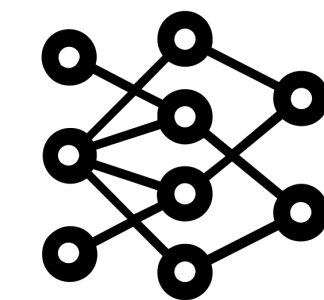
```
struct IxpMsg {
  char*    data;
  char*    pos;
  char*    end;
  _ixpuint size;
  _ixpuint mode;
};
```
**Ground Truth**

```
struct sha256_ctx {
  uint32_t H[8];
  uint32_t total[2];
  uint32_t buflen;
  char buffer[128];
};
```
**DIRTY**

```
struct struct0{
  int8*  s_0,
  int8*  s_1,
  int8*  s_2,
  int64  s_3,
  int64  s_4
};
```
**OSPREY**

Uses a multi-classification model

Only recovers layout

8

# Opportunities: LLMs' Power of Code Comprehension

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10  }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10  }
11 }
```

# Opportunities: LLMs' Power of Code Comprehension

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11 }
```

# Opportunities: LLMs' Power of Code Comprehension

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10  }
11 }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10  }
11 }
```

**ChatGPT-4 Output**

```
1  unsigned long long sub_404056
2      (long long a1, unsigned short a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if(*(unsigned int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 28);
8
9    v4 = v3 * a2;
10   sub_406BB9(v4);
11  }
12 }
```

10

# Direct Prompting Pre-trained/Fine-tuned LLMs is not Ideal

### Source Code

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11  }
```

### Decompiled Code

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11  }
```

### ChatGPT-4 Output

```
1  unsigned long long sub_404056
2      (long long a1, unsigned short a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if(*(unsigned int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 28);
8
9    v4 = v3 * a2;
10   sub_406BB9(v4);
11   }
12  }
```

### Fine-tuned StarCoder-3B Model Output

```
1  void process
2      (state *s, unsigned short n){
3   unsigned short len;
4   int size;
5   char *tmp;
6   if (s->last == 1) {
7    tmp = s->tmp;
8    size = len * n;
9    min(len);
10   }
11  }
12
```

# Direct Prompting Pre-trained/Fine-tuned LLMs is not Ideal

## Source Code

```
 1  void ixp_pstrings
 2      (IxpMsg *msg, ushort num){
 3   ushort len;
 4   uint size;
 5   uchar *s;
 6   if(msg->mode == 1){
 7    s = msg->pos;
 8    size = len * num;
 9    emalloc(size);
10   }
11  }
```

## Decompiled Code

```
 1  unsigned int64 sub_404056
 2      (int64 a1, int16 a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if (*(int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 8);
 8    v4 = v3 * a2;
 9    sub_406BB9(v4);
10   }
11  }
```
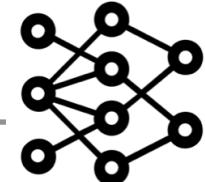
## ChatGPT-4 Output

```
 1  unsigned long long sub_404056
 2      (long long a1, unsigned short a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if(*(unsigned int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 28);
 8
 9    v4 = v3 * a2;
10    sub_406BB9(v4);
11   }
12  }
```

## Fine-tuned StarCoder-3B Model Output

```
 1  void process
 2      (state *s, unsigned short n){
 3   unsigned short len;
 4   int size;
 5   char *tmp;
 6   if (s->last == 1) {
 7    tmp = s->tmp;
 8    size = len * n;
 9    min(len);
10   }
11  }
12
```

11

# Direct Prompting Pre-trained/Fine-tuned LLMs is not Ideal

**Source Code**

```
1  void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
5   uchar *s;
6   if(msg->mode == 1){
7    s = msg->pos;
8    size = len * num;
9    emalloc(size);
10   }
11  }
```

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11  }
```

**ChatGPT-4 Output**

```
1  unsigned long long sub_404056
2    (long long a1, unsigned short a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if(*(unsigned int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 28);
8                    // should be `8`
9    v4 = v3 * a2;
10   sub_406BB9(v4);
11   }
12  }
```

**Fine-tuned StarCoder-3B Model Output**

```
1  void process
2    (state *s, unsigned short n){
3   unsigned short len;
4   int size;
5   char *tmp;
6   if (s->last == 1) {
7    tmp = s->tmp;
8    size = len * n;
9    min(len); // should be `size`
10   }
11  }
12  }
```

# Direct Prompting Pre-trained/Fine-tuned LLMs is not Ideal

Source Code

```
1  void ixp_pstrings
2    (IxpMsg *msg, ushort num){
3   ushort len;
4   uint size;
```

Decompiled Code

```
1  unsigned int64 sub_404056
2    (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
```

**Challenges:**

- A **general-purpose LLM** (e.g., ChatGPT) **struggles** to produce readable decompiled code.

- LLMs' direct outputs may have **incorrect semantics**.

- Accurate recovery requires a **global view**, while LLMs have token limits.

```
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if(*(unsigned int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 28);
8                   // should be `8`
9    v4 = v3 * a2;
10   sub_406BB9(v4);
11   }
12 }
```

```
3   unsigned short len;
4   int size;
5   char *tmp;
6   if (s->last == 1) {
7    tmp = s->tmp;
8    size = len * n;
9    min(len); // should be `size`
10   }
11 }
12 }
```
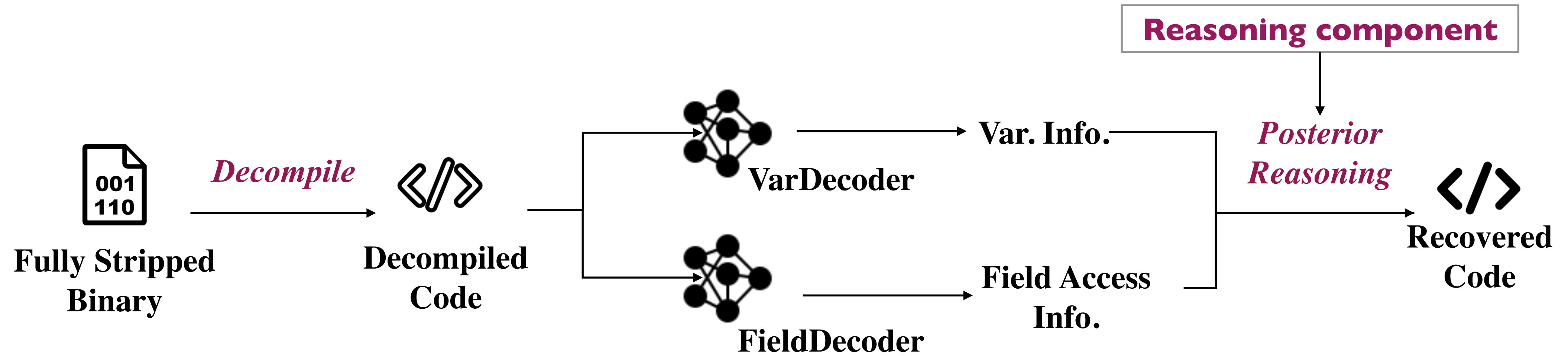
# Our Technique: **ReSym**

- A *hybrid approach* synergizes insights from LLMs and program analysis to recover variable and data structure symbols from stripped binaries

- Leverages two fine-tuned LLMs

- Replicates the reverse engineering process used by human experts

# Our Technique: **ReSym**

- A *hybrid approach* synergizes insights from LLMs and program analysis to recover variable and data structure symbols from stripped binaries

- Leverages two fine-tuned LLMs

- Replicates the reverse engineering process used by human experts

  ① **Break** the task into smaller manageable subtasks

  ② Focus on one piece at a time and then **aggregate** insights from multiple code snippets
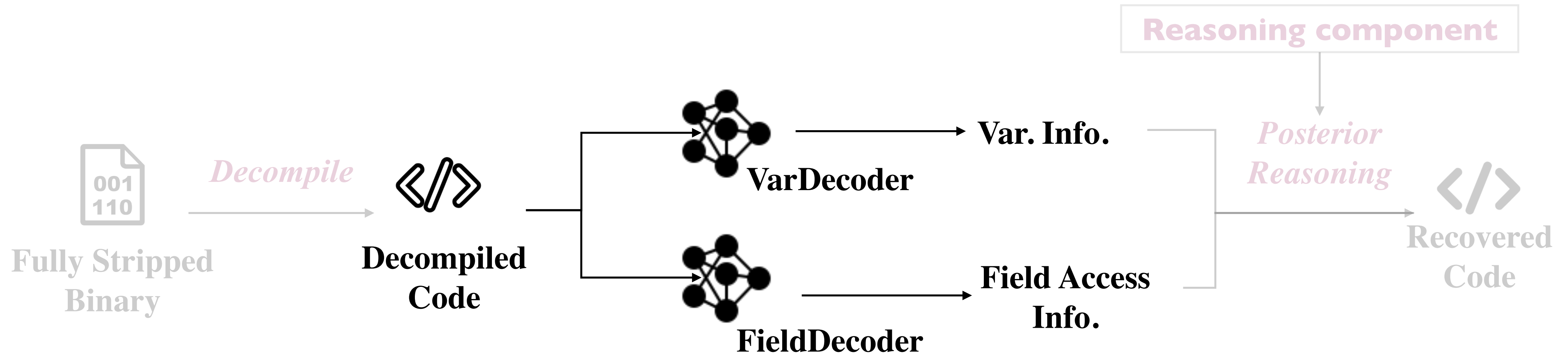
# ReSym Pipeline

Key idea: Break and Aggregate
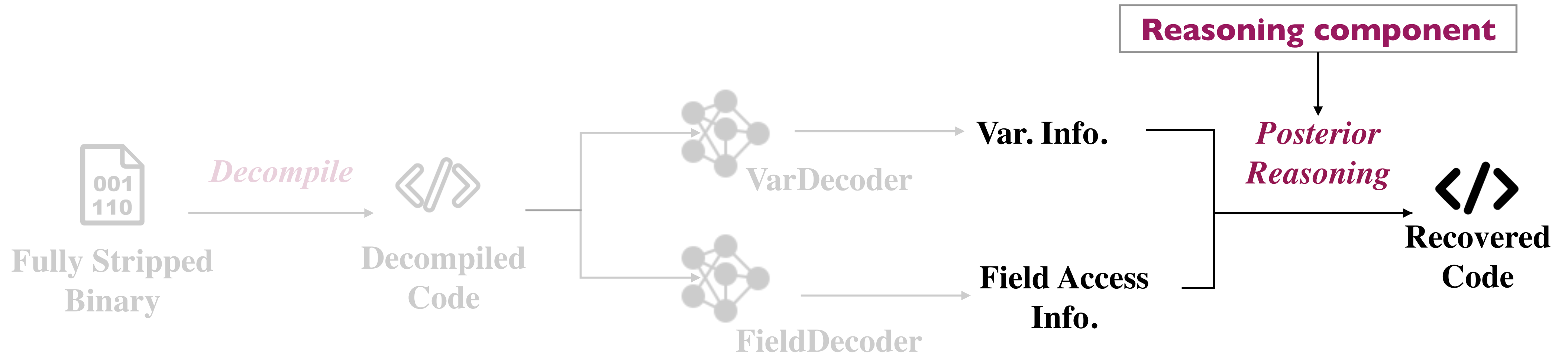
# ReSym Pipeline

Key idea: Break and Aggregate

# ReSym Pipeline

Key idea: Break and Aggregate

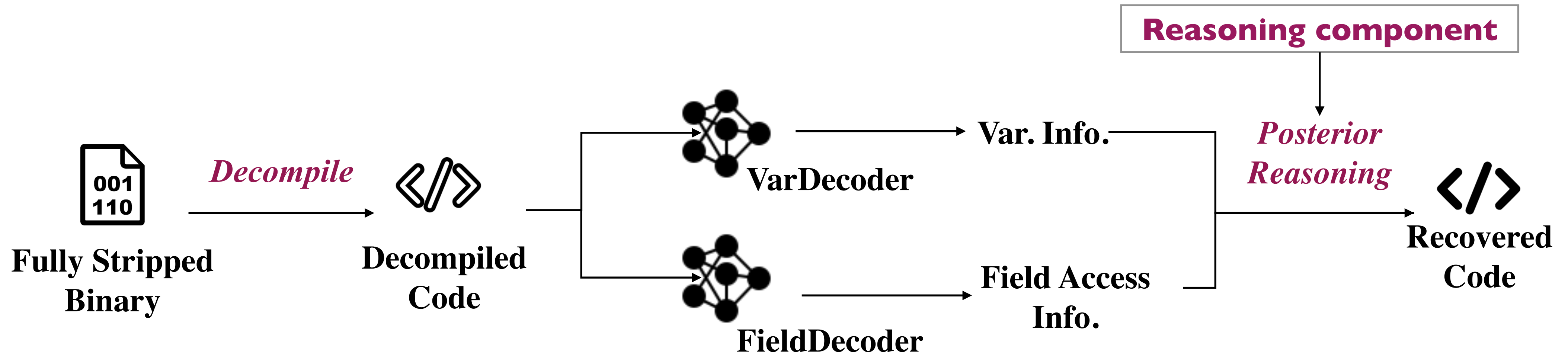# ReSym Pipeline

**Key idea: Break and Aggregate**



Fully Stripped Binary → *Decompile* → Decompiled Code → VarDecoder → **Var. Info.** / FieldDecoder → **Field Access Info.** → *Posterior Reasoning* → **Recovered Code**

**Reasoning component**

# ReSym Pipeline

Key idea: Break and Aggregate

# **VarDecoder**: Recover Variable Information

**Decompiled Code**

```
1   unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7      dest = *(void **)(a1 + 8);
8      v4 = v3 * a2;
9      sub_406BB9(v4);
10   }
11  }
```

# **VarDecoder**: Recover Variable Information

**Decompiled Code**

```
1   unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10    }
11  }
```

What are the original name and data type of
variables: `a1, a2, v3, v4, dest`? ⟨/⟩ ⟶ **VarDecoder**

# **VarDecoder**: Recover Variable Information

**Decompiled Code**

```
1   unsigned int64 sub_404056
2       (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11  }
```
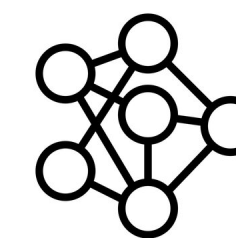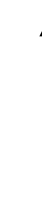
**Recovered Code**

```
1   unsigned int64 sub_404056
2       (Buffer *context, uint16 len){
3    uint16 chunk_len;
4    uint32 total_len;
5    char *temp_str;
6    if (*(int *)(a1 + 28) == 1) {
7     temp_str = *(void **)(a1 + 8);
8     total_len = chunk_len * len;
9     sub_406BB9(total_len);
10   }
11  }
```

What are the original name and data type of variables: `a1, a2, v3, v4, dest`? ⟨/⟩

**VarDecoder**

# **VarDecoder**: Recover Variable Information
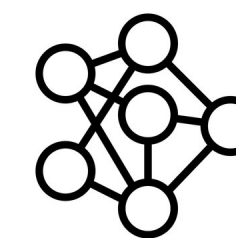
**Decompiled Code**

```
 1  unsigned int64 sub_404056
 2        (int64 a1, int16 a2){
 3   unsigned int16 v3;
 4   unsigned int v4;
 5   void *dest;
 6   if (*(int *)(a1 + 28) == 1){
 7    dest = *(void **)(a1 + 8);
 8    v4 = v3 * a2;
 9    sub_406BB9(v4);
10   }
11  }
```

**Recovered Code**

```
 1  unsigned int64 sub_404056
 2        (Buffer *context, uint16 len){
 3   uint16 chunk_len;
 4   uint32 total_len;
 5   char *temp_str;
 6   if (*(int *)(a1 + 28) == 1) {
 7    temp_str = *(void **)(a1 + 8);
 8    total_len = chunk_len * len;
 9    sub_406BB9(total_len);
10   }
11  }
```

What are the original name and data type of variables: `a1, a2, v3, v4, dest`? ⟨/⟩

**VarDecoder**

# **FieldDecoder**: Recover Field Access Information

**Decompiled Code**

```
1   unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7      dest = *(void **)(a1 + 8);
8      v4 = v3 * a2;
9      sub_406BB9(v4);
10   }
11  }
```

**Recovered Code**

```
1   unsigned int64 sub_404056
2        (Buffer *context, uint16 len){
3    uint16 chunk_len;
4    uint32 total_len;
5    char *temp_str;
6    if (*(int *)(a1 + 28) == 1) {
7      temp_str = *(void **)(a1 + 8);
8      total_len = chunk_len * len;
9      sub_406BB9(total_len);
10   }
11  }
```

What are the original name and data type of
variables: `a1, a2, v3, v4, dest`? ⟨/⟩

**VarDecoder**

21

# **FieldDecoder**: Recover Field Access Information
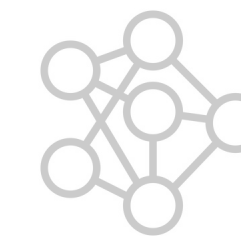
**Decompiled Code**

```
1  unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11 }
```

**Recovered Code**

```
1  unsigned int64 sub_404056
2        (Buffer *context, uint16 len){
3   uint16 chunk_len;
4   uint32 total_len;
5   char *temp_str;
6   if (*(int *)(a1 + 28) == 1) {
7     temp_str = *(void **)(a1 + 8);
8     total_len = chunk_len * len;
9     sub_406BB9(total_len);
10   }
11 }
```

What are the original name and data type of
variables: `a1, a2, v3, v4, dest`?  </>

**VarDecoder**

# **FieldDecoder**: Recover Field Access Information
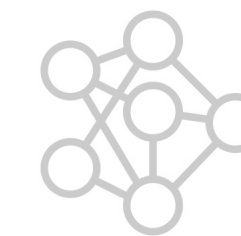
**Decompiled Code**

```
1  unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10  }
11 }
```

**Recovered Code**

```
1  unsigned int64 sub_404056
2        (Buffer *context, uint16 len){
3   uint16 chunk_len;
4   uint32 total_len;
5   char *temp_str;
6   if (context->type == 1) {
7    temp_str = context->pos;
8    total_len = chunk_len * len;
9    sub_406BB9(total_len);
10  }
11 }
```

What are the original name and data type of variables: `a1, a2, v3, v4, dest`?  〈/〉

**VarDecoder**

What are the variable name and type for the following field accesses:
`(int *)(a1 + 28), (void **)(a1 + 8)`?  〈/〉

**FieldDecoder**

# Posterior Reasoning:
## Aggregating Field Access from Multiple Functions

**Decompiled Code**

```
1  unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11  }
```

**Recovered Code**

```
1  unsigned int64 sub_404056
2      (Buffer *context, uint16 len){
3   uint16 chunk_len;
4   uint32 total_len;
5   char *temp_str;
6   if (context->type == 1) {
7    temp_str = context->pos;
8    total_len = chunk_len * len;
9    sub_406BB9(total_len);
10   }
11  }
```

**Recovered Data Structure**

```
struct Buffer {
    ?                    // 0-7
    uint8_t* pos;        // 8-15
    ?                    ?
    uint32_t type;       // 28-31
};
```

# Posterior Reasoning:
## Aggregating Field Access from Multiple Functions

**Decompiled Code**

```
1  unsigned int64 sub_404056
2        (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11 }
```
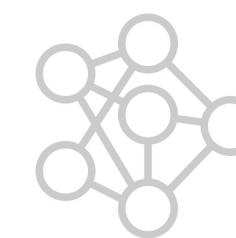
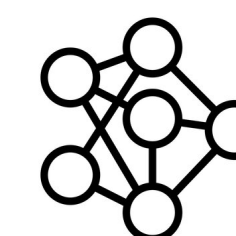**Recovered Code**

```
1  unsigned int64 sub_404056
2        (Buffer *context, uint16 len){
3   uint16 chunk_len;
4   uint32 total_len;
5   char *temp_str;
6   if (context->type == 1) {
7    temp_str = context->pos;
8    total_len = chunk_len * len;
9    sub_406BB9(total_len);
10   }
11 }
```

**Other Functions**

```
int64 sub_404362 (...){
   *(int64 *)(a1 + 8)
   (int64 *)(a1 + 16)
}
```

**Recovered Data Structure**

```
struct Buffer {
   ?                        // 0-7
   uint8_t* pos;            // 8-15
   uint8_t* streamPos;      // 16-23
   ?                        // 24-27
   uint32_t type;           // 28-31
};
```

24

# Posterior Reasoning:
## Aggregating Field Access from Multiple Functions

**Decompiled Code**

```
1  unsigned int64 sub_404056
2       (int64 a1, int16 a2){
3   unsigned int16 v3;
4   unsigned int v4;
5   void *dest;
6   if (*(int *)(a1 + 28) == 1){
7    dest = *(void **)(a1 + 8);
8    v4 = v3 * a2;
9    sub_406BB9(v4);
10   }
11  }
```
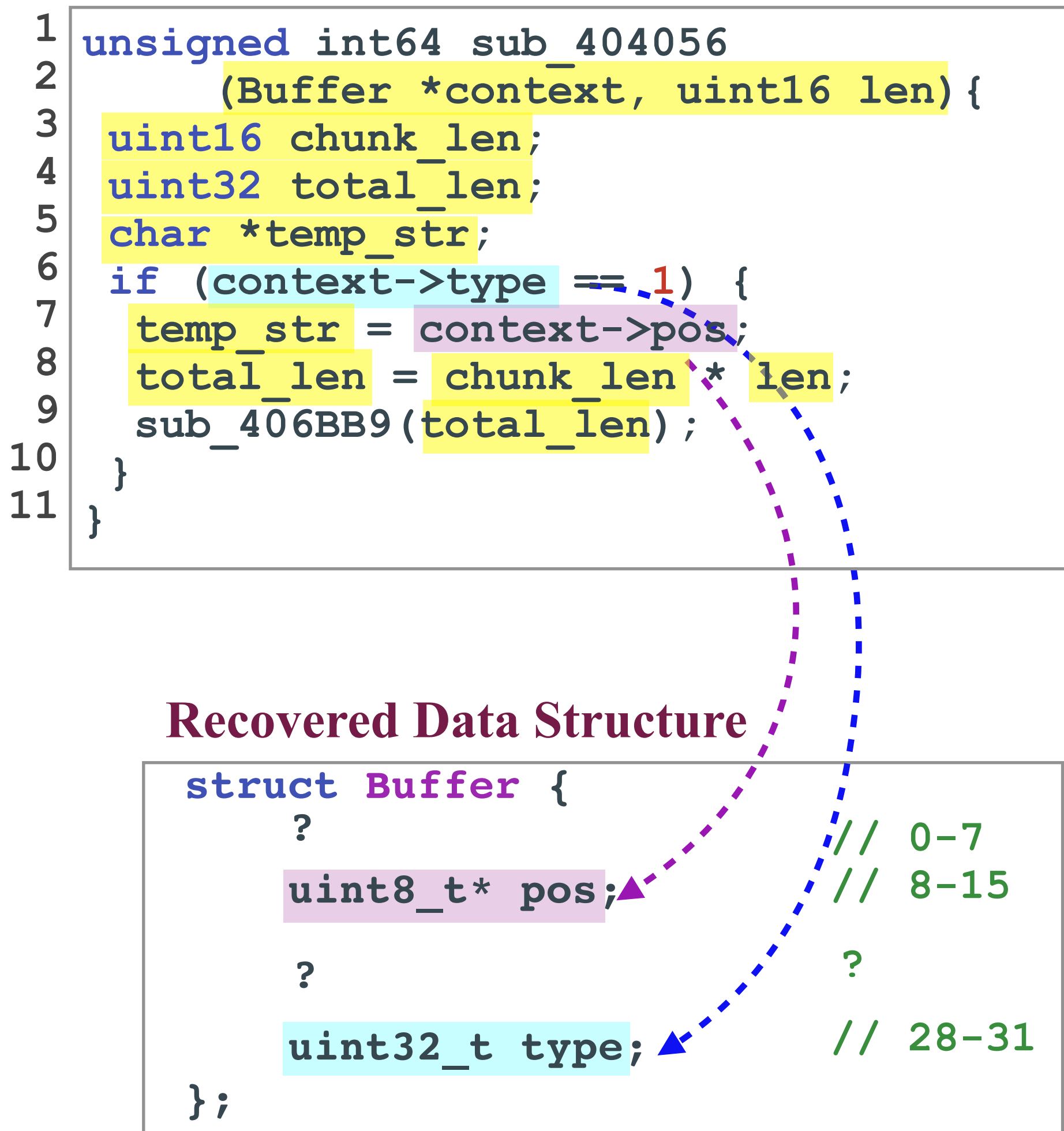
**Recovered Code**

```
1  unsigned int64 sub_404056
2       (Buffer *context, uint16 len){
3   uint16 chunk_len;
4   uint32 total_len;
5   char *temp_str;
6   if (context->type == 1) {
7    temp_str = context->pos;
8    total_len = chunk_len * len;
9    sub_406BB9(total_len);
10   }
11  }
```

**Other Functions**

```
int64 sub_404362 (...){
  *(int64 *)(a1 + 8)
  (int64 *)(a1 + 16)
}
```

**Recovered Data Structure**

```
struct Buffer {
    uint8_t* buffer;       // 0-7
    uint8_t* pos;          // 8-15
    uint8_t* streamPos;    // 16-23
    uint32_t bufferSize;   // 24-27
    uint32_t type;         // 28-31
};
```

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
1   void ixp_pstrings
2      (IxpMsg *msg, ushort num){
3    ushort len;
4    uint size;
5    uchar *s;
6    if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11  }
```

**Decompiled Code**

```
1   unsigned int64 sub_404056
2      (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11  }
```

```
struct IxpMsg {
 char*    data;
 char*    pos;
 char*    end;
 _ixpuint size;
 _ixpuint mode;
};
```
**Ground Truth**

```
struct sha256_ctx {
 uint32_t H[8];
 uint32_t total[2];
 uint32_t buflen;
 char buffer[128];
};
```
**DIRTY**

```
struct struct0{
 int8*  s_0,
 int8*  s_1,
 int8*  s_2,
 int64  s_3,
 int64  s_4
};
```
**OSPREY**

Uses a multi-classification model

Only recovers layout

26

# Existing Techniques are Limited on Recovering User-defined Data Structures

**Source Code**

```
1   void ixp_pstrings
2       (IxpMsg *msg, ushort num){
3    ushort len;
4    uint size;
5    uchar *s;
6    if(msg->mode == 1){
7     s = msg->pos;
8     size = len * num;
9     emalloc(size);
10   }
11  }
```

**Decompiled Code**

```
1   unsigned int64 sub_404056
2       (int64 a1, int16 a2){
3    unsigned int16 v3;
4    unsigned int v4;
5    void *dest;
6    if (*(int *)(a1 + 28) == 1){
7     dest = *(void **)(a1 + 8);
8     v4 = v3 * a2;
9     sub_406BB9(v4);
10   }
11  }
```

```
struct IxpMsg {
 char*    data;
 char*    pos;
 char*    end;
 _ixpuint size;
 _ixpuint mode;
};
```
**Ground Truth**

```
struct sha256_ctx {
 uint32_t H[8];
 uint32_t total[2];
 uint32_t buflen;
 char buffer[128];
};
```
**DIRTY**

```
struct struct0{
 int8*  s_0,
 int8*  s_1,
 int8*  s_2,
 int64  s_3,
 int64  s_4
};
```
**OSPREY**

```
struct Buffer {
 uint8_t* buffer;
 uint8_t* pos;
 uint8_t* streamPos;
 uint32_t bufferSize;
 uint32_t type;
};
```
**ReSym**

Uses a multi-classification model

Only recovers layout

# Experimental Setup

- **3,058** C/C++ real-world projects collected from GitHub

- **16,217** binary files

  - average size: **116 KB**; maximum size: **8.9 MB**

- Split train/test set **by project** with a ratio of 0.95

- Fine-tune two **StarCoder 3B** models for VarDecoder and FieldDecoder

# Research Questions

- How good is ReSym at recovering variable names and types?

- How good is ReSym at recovering user-defined data structures?

# Evaluation: Name and Type Recovery

## Perfect match accuracy (%)

| Method | Overall | |
|--------|---------|------|
| | name | type |
| ReSym | 56.7  8.0↑ | 60.7  4.9↑ |
| DIRTY | 48.7 | 55.8 |

# Evaluation: Name and Type Recovery

## Perfect match accuracy (%)

| Method | Overall | |
|--------|---------|---|
| | name | type |
| ReSym | 56.7   8.0↑ | 60.7   4.9↑ |
| DIRTY | 48.7 | 55.8 |

ReSym is effective in recovering variable names and types.
ReSym outperforms DIRTY by 4.9 — **8.0%**.

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|--------|---------------|---|---|------------------------------|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|--------|-----------|--------|------|-------------|------------|------------|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

**Recovered Data Structure**

```
struct Buffer {
 uint8_t* buffer;     // 0-7
 uint8_t* pos;        // 8-15
 size_t*  streamPos;  // 16-23
 uint32_t bufferSize; // 24-27
 uint32_t type;       // 28-31
};
```

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

**Recovered Data Structure**

```
struct Buffer {
 uint8_t* buffer;      // 0-7
 uint8_t* pos;         // 8-15
 size_t*  streamPos;   // 16-23
 uint32_t bufferSize;  // 24-27
 uint32_t type;        // 28-31
};
```

predicted offsets and sizes

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|--------|---------------|--------|------|-----------------|------------|------------|
|  | Precision | Recall | F1 | **Struct Type** | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

**Recovered Data Structure**

```
struct Buffer {
 uint8_t* buffer;      // 0-7
 uint8_t* pos;         // 8-15
 size_t*  streamPos;   // 16-23
 uint32_t bufferSize;  // 24-27
 uint32_t type;        // 28-31
};
```

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|--------|-----------|--------|------|-------------|------------|------------|
| | Precision | Recall | F1 | Struct Type | **Field Name** | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

**Recovered Data Structure**

```
struct Buffer {
 uint8_t* buffer;      // 0-7
 uint8_t* pos;         // 8-15
 size_t*  streamPos;   // 16-23
 uint32_t bufferSize;  // 24-27
 uint32_t type;        // 28-31
};
```

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|--------|-----------|--------|------|-------------|------------|------------|
| | Precision | Recall | F1 | Struct Type | Field Name | **Field Type** |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |

**Recovered Data Structure**

```
struct Buffer {
 uint8_t* buffer;     // 0-7
 uint8_t* pos;        // 8-15
 size_t*  streamPos;  // 16-23
 uint32_t bufferSize; // 24-27
 uint32_t type;       // 28-31
};
```

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |
| ReSym_Light | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |

ReSym without posterior reasoning

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | **48.6** | 44.4 | 14.4 | 15.5 |
| ReSym<sub>Light</sub> | 73.3 | 29.8 | **42.4** | 40.2 | 7.9 | 8.3 |

ReSym without posterior reasoning

Posterior Reasoning is effective and improves the F1 score by 6.2%.

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | 81.9 | 34.6 | 48.6 | 44.4 | 14.4 | 15.5 |
| ReSym Light | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |
| OSPREY | 38.1 | 60.2 | 46.7 | - | - | - |
| DIRTY | 54.6 | 3.3 | 6.2 | 0.5 | 0.8 | 0.9 |

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | **81.9** | 34.6 | **48.6** | **44.4** | **14.4** | **15.5** |
| ReSym<sub>Light</sub> | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |
| OSPREY | 38.1 | **60.2** | 46.7 | - | - | - |
| DIRTY | 54.6 | 3.3 | 6.2 | 0.5 | 0.8 | 0.9 |

**ReSym's lower recall**:

1. Discarding functions due to token limit

2. Data-flow analysis is inherently undecidable

39

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | **81.9** | 34.6 | **48.6** | **44.4** | **14.4** | **15.5** |
| ReSym<sub>Light</sub> | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |
| OSPREY | 38.1 | **60.2** | 46.7 | - | - | - |
| DIRTY | 54.6 | 3.3 | 6.2 | 0.5 | 0.8 | 0.9 |

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | **81.9** | 34.6 | **48.6** | **44.4** | **14.4** | **15.5** |
| ReSym<sub>Light</sub> | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |
| OSPREY | 38.1 | **60.2** | 46.7 | - | - | - |
| DIRTY | 54.6 | 3.3 | 6.2 | 0.5 | 0.8 | 0.9 |

ReSym achieves the highest F1 score and accurately recovers structure annotation.

# Evaluation: User-Defined Data Structure Recovery

| Method | Struct Layout | | | Struct Annotation (Accuracy) | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Struct Type | Field Name | Field Type |
| ReSym | **81.9** | 34.6 | **48.6** | **44.4** | **14.4** | **15.5** |
| ReSym_Light | 73.3 | 29.8 | 42.4 | 40.2 | 7.9 | 8.3 |
| OSPREY | 38.1 | **60.2** | 46.7 | - | - | - |
| DIRTY | 54.6 | 3.3 | 6.2 | 0.5 | 0.8 | 0.9 |

ReSym achieves the highest F1 score and accurately recovers structure annotation.

ReSym analyzes each binary file in 3.4s, while OSPREY takes 528.24s.

# Conclusions

- Propose a prototype, **ReSym**, that harnesses LLMs to effectively recover variable and data structure symbols from stripped binaries.

- **Divide** the difficult symbol recovery problem into two manageable sub-problems.

- Develop a rigorous reasoning component to **aggregate** and cross-check local results, enabling the recovery of **comprehensive data structures**.

- ReSym **outperforms** state-of-the-art techniques, DIRTY and OSPREY.

- Build a **large-scale public dataset** of C/C++ projects containing binaries annotated with corresponding symbols, together with the automatic annotation pipeline, to facilitate future research.