

# **Augmenting Decompiler Output with Learned Variable Names and Types**

USENIX Security Symposium

**Qibin Chen, Jeremy Lacomis, Edward J. Schwartz, Claire Le Goues,  
Graham Neubig, and Bogdan Vasilescu**

# Variable Naming is Essential

```
1| void read(char *a1) {  
2|     int v1 = 0;  
3|     char v2;  
4|     while (1) {  
5|         v2 = getchar();  
6|         if (v2 == EOF || v2 == '\n') {  
7|             a1[v1] = '\0';  
8|             break;  
9|         } else {  
10|             a1[v1] = v2;  
11|         }  
12|         v1++;  
13|     }  
14| }
```

# Variable Naming is Essential

```
1I void read(char *out_buf) {
2I   int loc = 0;
3I   char cur;
4I   while (1) {
5I     cur = getchar();
6I     if (cur == EOF || cur == '\n') {
7I       out_buf[loc] = '\0';
8I       break;
9I     } else {
10I       out_buf[loc] = cur;
11I     }
12I     loc++;
13I   }
14I }
```

# Variable Naming is Essential

```
1I void read(char *banana) {  
2I   int str = 0;  
3I   char moo;  
4I   while (1) {  
5I     moo = getchar();  
6I     if (moo == EOF || moo == '\n') {  
7I       banana[str] = '\0';  
8I       break;  
9I     } else {  
10I       banana[str] = moo;  
11I     }  
12I     str++;  
13I   }  
14I }
```

# DIRE

```
1I file *f_open(char **filename, char *mode, int create) {  
2I     int fd;  
3I     if (!create)  
4I         return fopen(*filename, mode);  
5I     if (*mode != 119)  
6I         assert_fail("fopen");  
7I     fd = open(*filename, 577, 384);  
8I     if (fd >= 0)  
9I         return reopen(fd, mode);  
10I    else  
11I        return 0;  
12I }
```

Hex-Rays	Developer
a1	filename
a2	mode
a3	is_private

# Types Are Also Useful

```
double func(int *a1, int *a2) {  
    double v1, v2;  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
    return sqrt(v1 + v2);  
}
```

# Types Are Also Useful

```
double func(int *a1, int *a2) {  
    double v1, v2;  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
double func(point *a1, point *a2) {  
    double v1, v2;  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
    return sqrt(v1 + v2);  
}
```

# Predicting Types should be easier than names!

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char); // 1  
sizeof(int); // 4  
sizeof(float); // 4  
sizeof(point); // 8  
sizeof(int[2]); // 8
```

# **Constraints Experiment**

# **Constraints Experiment**

- Mine GitHub, creating the same dataset as for DIRE

# Constraints Experiment

- Mine GitHub, creating the same dataset as for DIRE
- Also maintain a database of types including their substructure

# Constraints Experiment

- Mine GitHub, creating the same dataset as for DIRE
- Also maintain a database of types including their substructure
- Train a model to predict a type based on its context and constrain by size

# Constraints Experiment

- Mine GitHub, creating the same dataset as for DIRE
- Also maintain a database of types including their substructure
- Train a model to predict a type based on its context and constrain by size
- Results: Quite Bad!

# The Problem: Padding Bytes

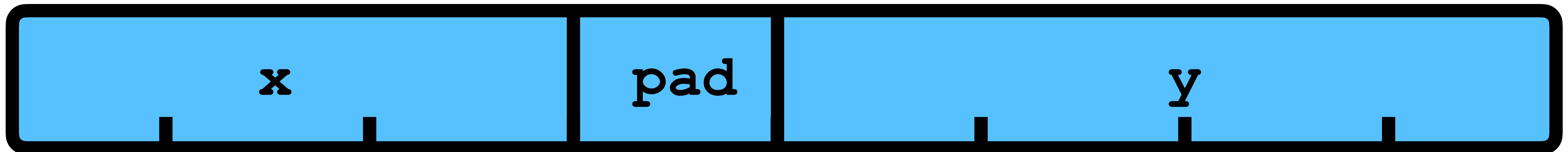
Original

```
void fun() {  
    char x[3];  
    int y;  
    // ...  
}
```

Decompiled

```
void fun() {  
    char x[4];  
    int y;  
    // ...  
}
```

Stack



# The Problem: Padding Bytes

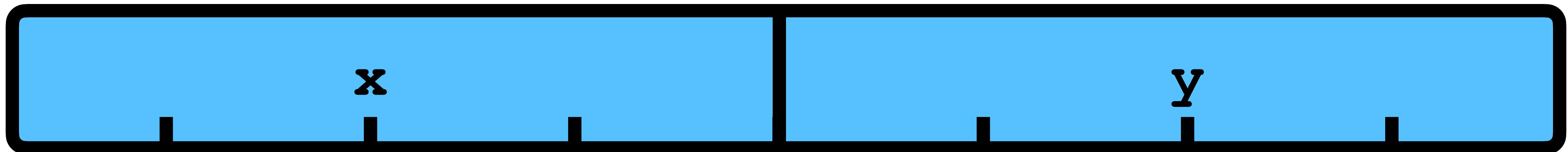
Original

```
void fun() {  
    char x[3];  
    int y;  
    // ...  
}
```

Decompiled

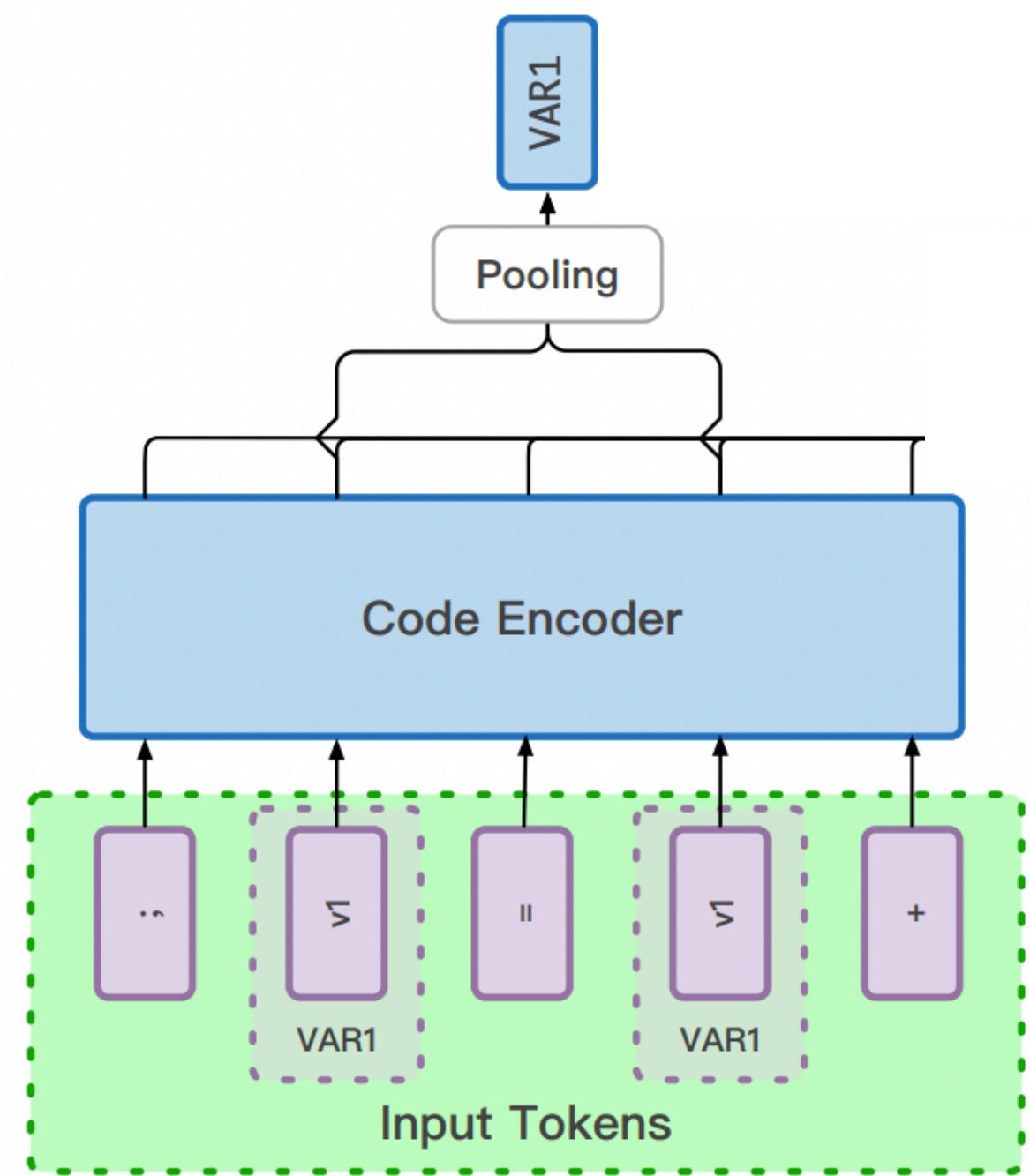
```
void fun() {  
    char x[4];  
    int y;  
    // ...  
}
```

Stack

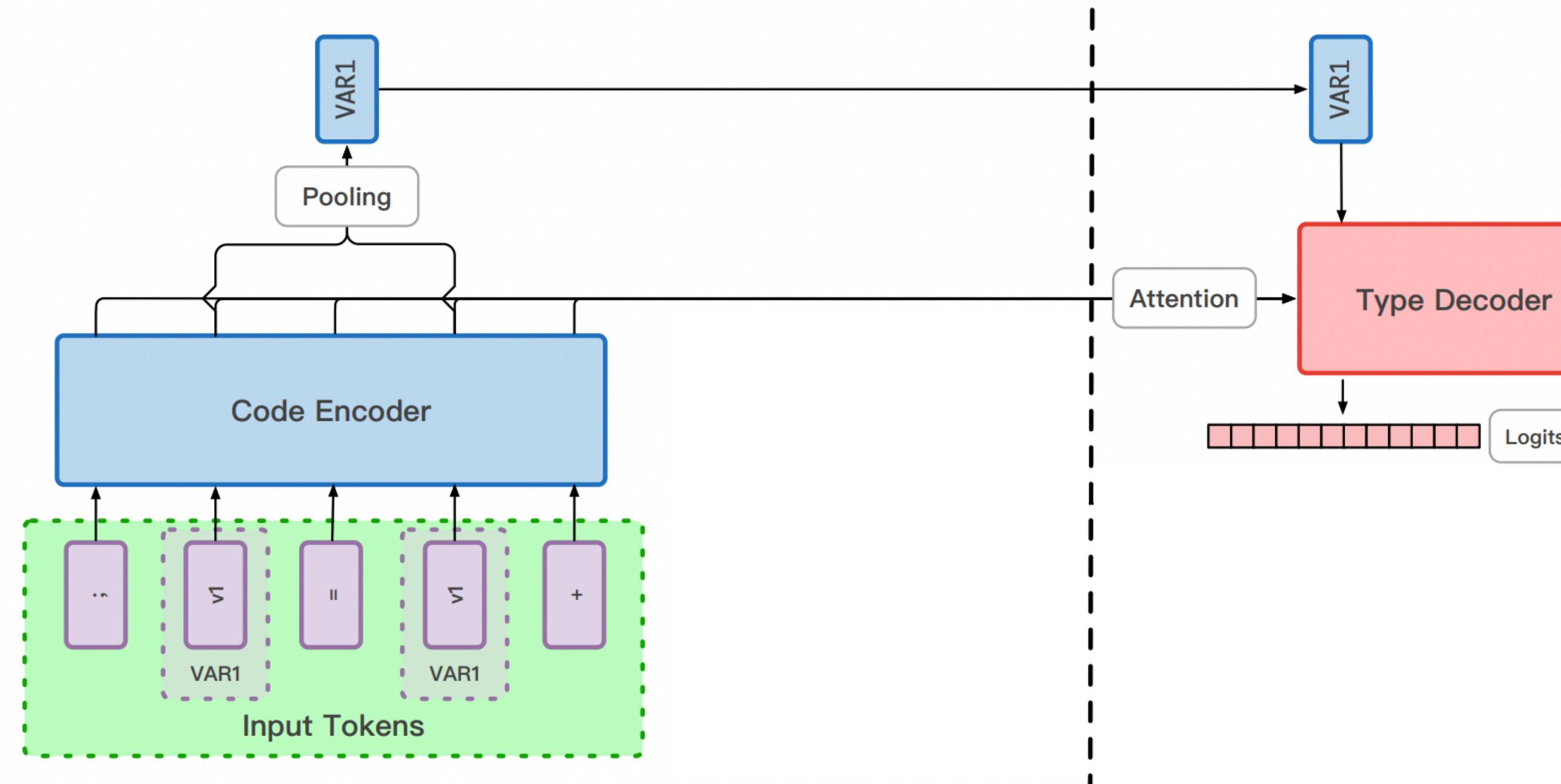


# Decompiled variable ReTYper (DIRTY) Architecture

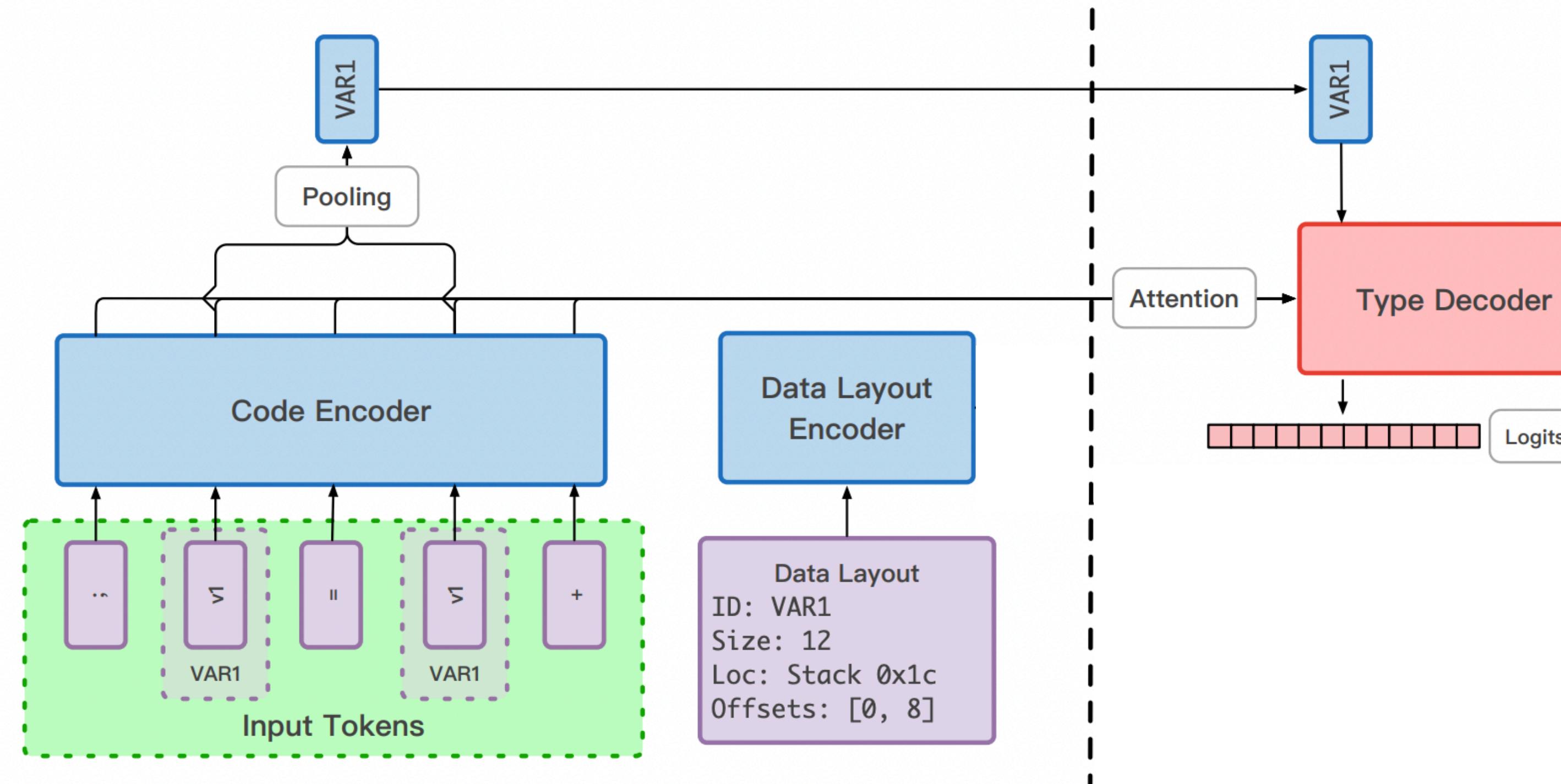
# Decompiled variable ReTYper (DIRTY) Architecture



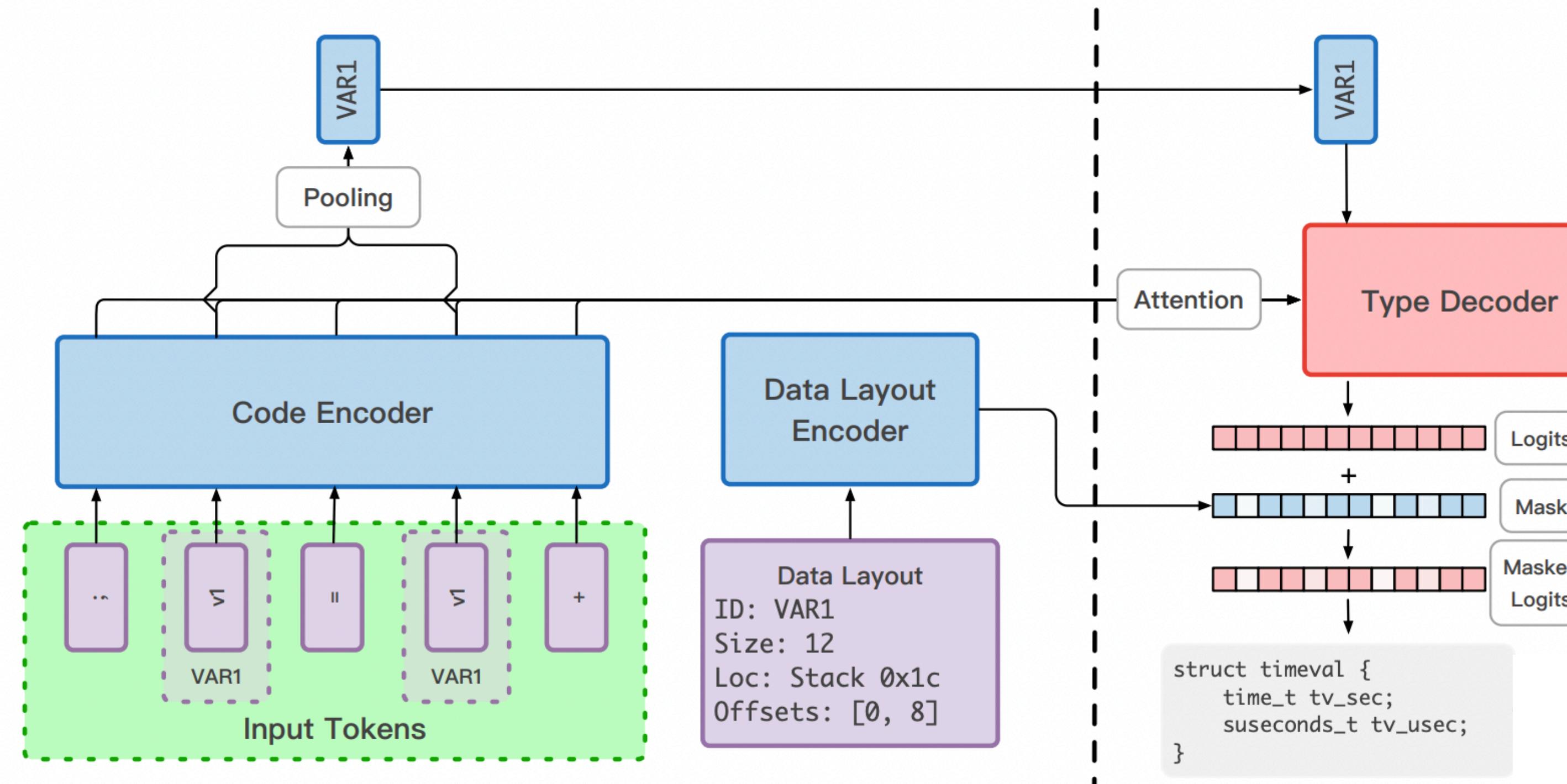
# Decompiled variable ReTYper (DIRTY) Architecture



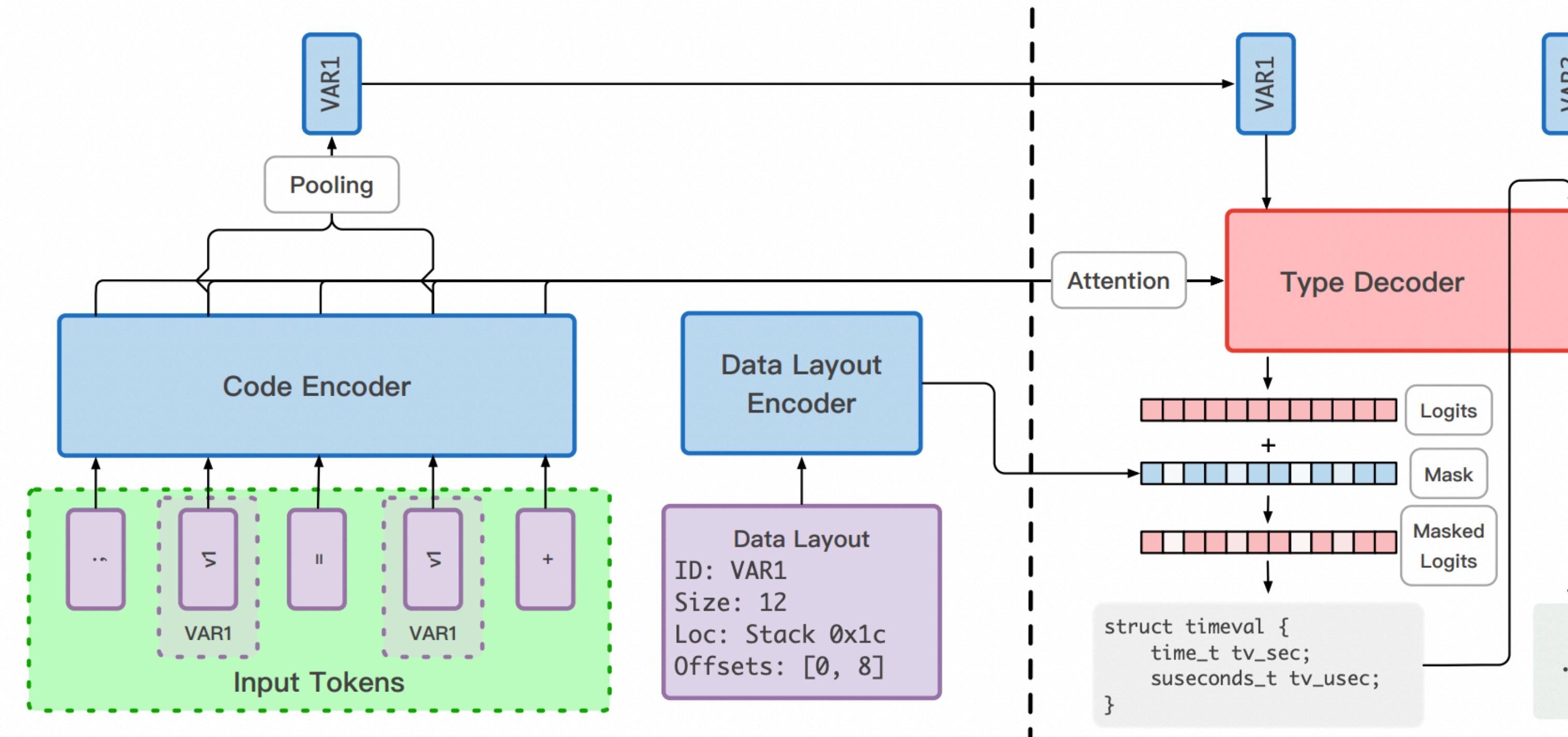
# Decompiled variable ReTYper (DIRTY) Architecture



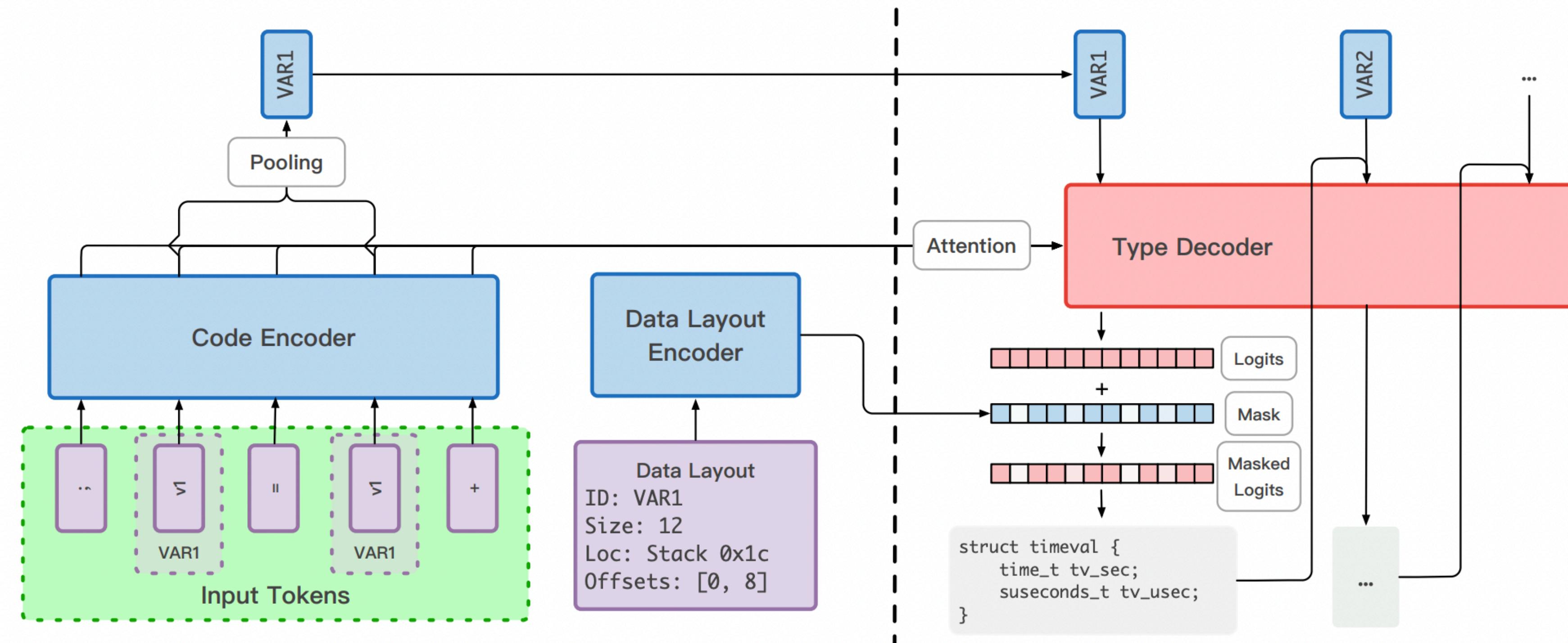
# Decompiled variable ReTYper (DIRTY) Architecture



# Decompiled variable ReTYper (DIRTY) Architecture



# Decompiled variable ReTYper (DIRTY) Architecture



# Evaluation

- DIRTY was evaluated on the Dataset for Idiomatic ReTyping (DIRT)

# Evaluation

- DIRT<sup>Y</sup> was evaluated on the Dataset for Idiomatic ReTyping (DIRT)
- Baselines were Hex-Rays and Frequency by Size

# Evaluation

- DIRT<sup>Y</sup> was evaluated on the Dataset for Idiomatic ReTyping (DIRT)
- Baselines were Hex-Rays and Frequency by Size
- Accuracy: the percentage of types that exactly match the ground truth, including data layout, type name, and substructure, where applicable

# Evaluation

- DIRTY was evaluated on the Dataset for Idiomatic ReTyping (DIRT)
- Baselines were Hex-Rays and Frequency by Size
- Accuracy: the percentage of types that exactly match the ground truth, including data layout, type name, and substructure, where applicable

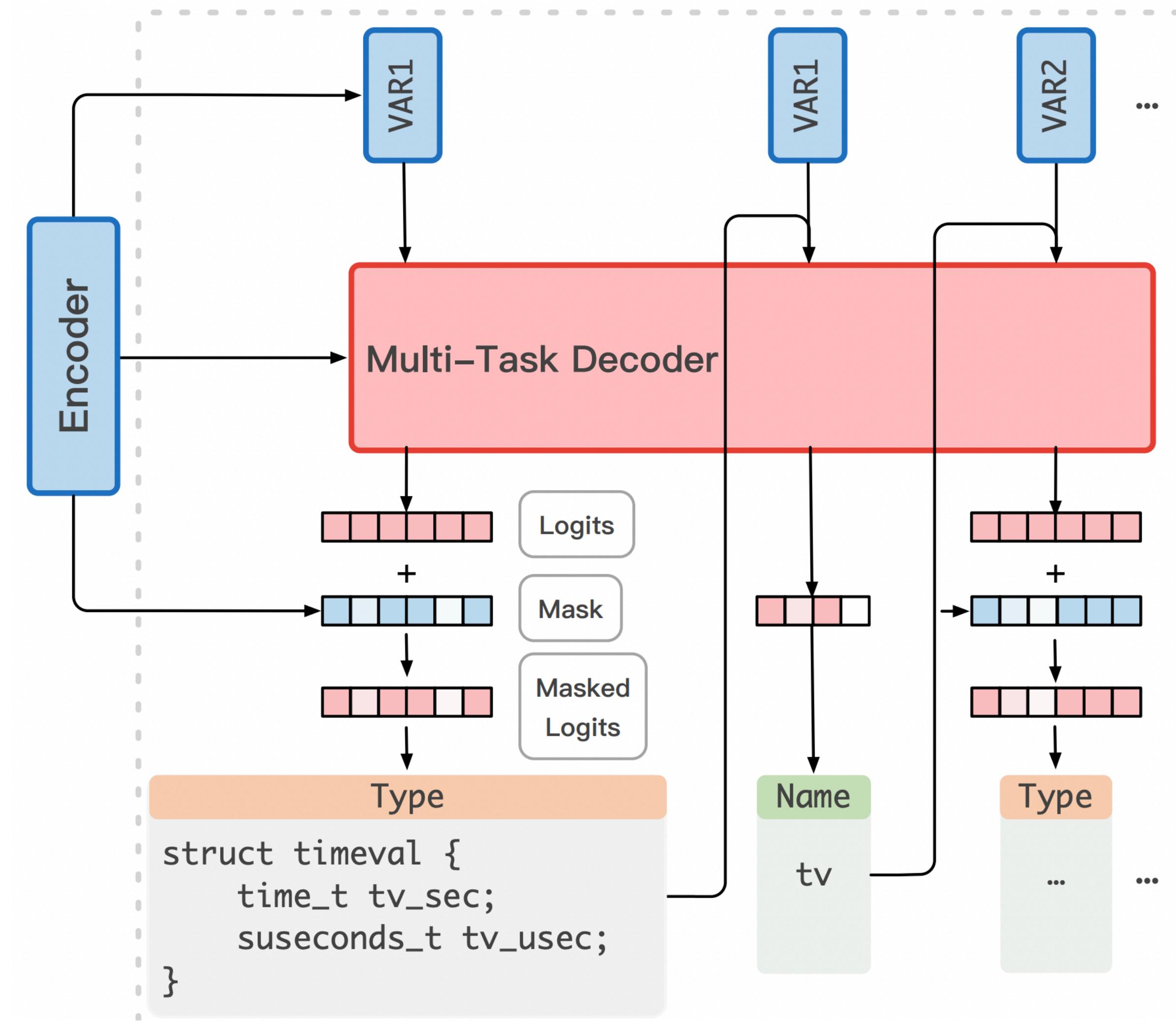
Method	All	Structures
Freq by Size	23.6	9.7
Hex-Rays	37.9	28.7
DIRTY	75.8	68.6

# Compiler Optimization Performance

Percent accuracy on 101 GNU Coreutils Programs

-00	-01	-02	-03
48.20	46.01	46.04	46.00

# Predicting types and names simultaneously



# Evaluation

- DIRTY compared with DIRE

# Evaluation

- DIRTY compared with DIRE
- Used the dataset from the original DIRE paper and also DIRT

# Evaluation

- DIRTY compared with DIRE
- Used the dataset from the original DIRE paper and also DIRT
- Accuracy is the percentage of names that exactly match the developer's choice.

# Evaluation

- DIRTY compared with DIRE
- Used the dataset from the original DIRE paper and also DIRT
- Accuracy is the percentage of names that exactly match the developer's choice.

Method	DIRE Data	DIRT
DIRE	72.8	57.5
DIRTY	81.4	66.4

```

int find_unused_picture(int a1, int a2, int a3) {
    int i, j, v1;
    if (a3) {
        for (i = <NUM>; ++i) {
            if (i > <NUM>)
                goto LABEL_13;
            if (!*(*(<NUM> * i + a2) + <NUM>))
                break;
        }
        v1 = i;
    } else {
        for (j = <NUM>; ++j) {
            if (j > <NUM>) {
                LABEL_13:
                av_log(a1, <NUM>, <STR>);
            }
        }
    }
}

```

ID	Developer	DIRTY
a1	AVCodecContext *avctx	MpegEncContext *
a2	Picture *picture	Picture *pic
a3	int shared	int shared
v1	int result	int result

**BONUS  
SLIDES**

# NUMERICAL RECIPES in C

The Art of Scientific Computing

Second Edition

William H. Press      Saul A. Teukolsky

William T. Vetterling      Brian P. Flannery

```
void scrsho(float (*fx)(float))
```

For interactive CRT terminal use. Produce a crude graph of the function  $y=f(x)$  for interval  $x_1, x_2$ . Query for another plot until the user signs off.

```
{
```

```
    int jz,j,i;
```

```
    float ysml,ybig,x2,x1,x,dyj,dx,y[ISCR+1];
```

```
    char scr[ISCR+1][JSCR+1];
```

```
    for (;;) {
```

```
        printf("\nEnter x1 x2 (x1=x2 to stop):\n");
```

```
        scanf("%f %f",&x1,&x2);
```

```
        if (x1 == x2) break;
```

```
        for (j=1;j<=JSCR;j++)
```

```
            scr[1][j]=scr[ISCR][j]=YY;
```

```
        for (i=2;i<=(ISCR-1);i++) {
```

```
            scr[i][1]=scr[i][JSCR]=XX;
```

```
            for (j=2;j<=(JSCR-1);j++)
```

```
                scr[i][j]=BLANK;
```

```
}
```

```
dx=(x2-x1)/(ISCR-1);
```

Fill vertical

Fill top

Fill intermediate

books/scrsho.cpp at master · b · +

github.com/burakbayramli/books/blob/master/NumericalRecipes2/recipes/scrsho.cpp

Pulls Issues Marketplace Explore

burakbayramli / books

Code Issues Pull requests

master · books / NumericalRecipes2 / recipes / scrsho.cpp

BB Initial commit · Latest commit d32c78b on Jan 3, 2017 · History

0 contributors

55 lines (53 sloc) | 1.25 KB

```
1 #include <string>
2 #include <iostream>
3 #include <iomanip>
4 #include "nr.h"
5 using namespace std;
6
7 void NR::scrsho(DP fx(const DP))
8 {
9     const int ISCR=60, JSCR=21;
10    const char BLANK=' ', ZERO='-', YY='l', XX='-', FF='x';
11    int jz,j,i;
12    DP ysm1,ybig,x2,x1,x,dyj,dx;
13    Vec_DP y(ISCR);
14    string scr[JSCR];
15
16    for (;;) {
17        cout << endl << "Enter x1 x2 (x1=x2 to stop):" << endl;
18        cin >> x1 >> x2;
19        if (x1 == x2) break;
20        scr[0]=YY;
```