# Machine Learning Approach for Loop Unrolling Factor Prediction in High Level Synthesis
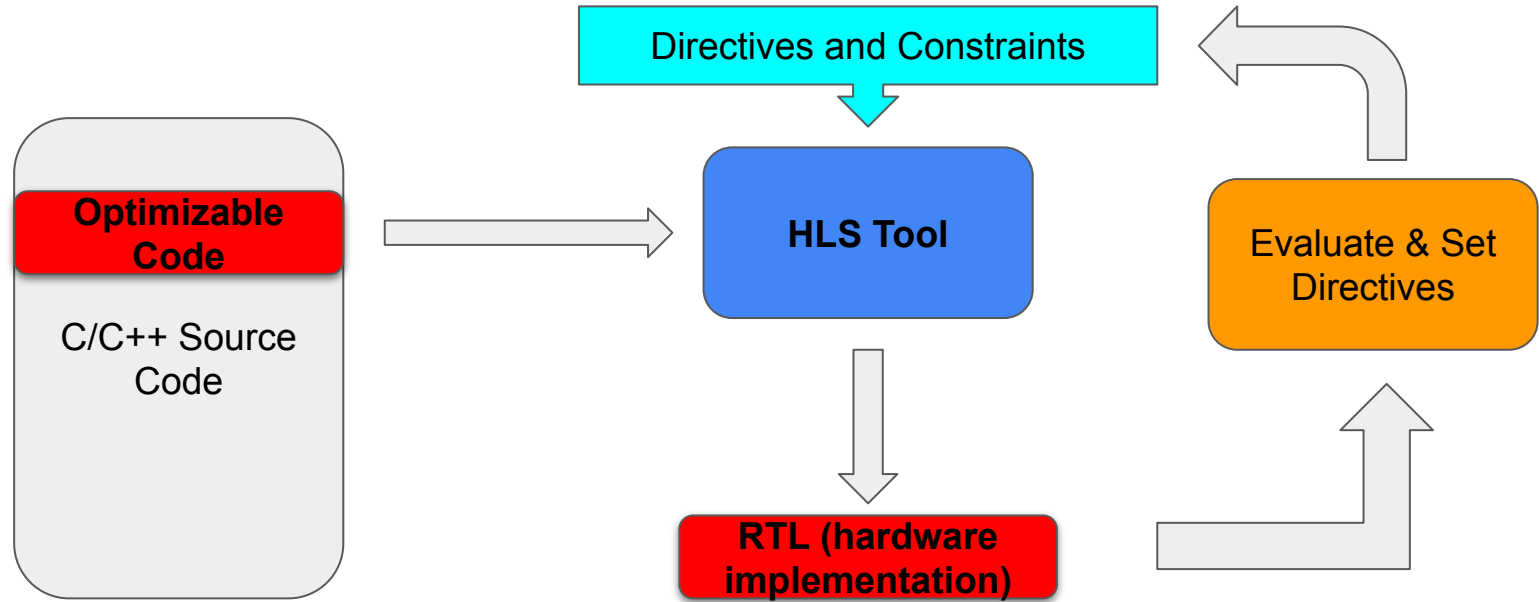
Authors: Georgios Zacharopoulos, Andrea Barbon, Giovanni Ansaloni and Laura Pozzi

Mehar Singh, Anisha Aggarwal, Parth Raut (Group #10)

# High Level Synthesis (HLS)

- HLS frameworks allow **hardware circuits** to be described at **higher abstraction** for customizable hardware
  - Languages like C or C++
- Differs from **Register Transfer Level** (RTL) methods which describe circuits in terms of registers, logical operations, and data movements
  - Languages like Verilog
- HLS allows engineers to deal with hardware design without having to know low-level details
  - Focus on **functionality**
  - Optimizing from code to hardware descriptions is done **automatically**
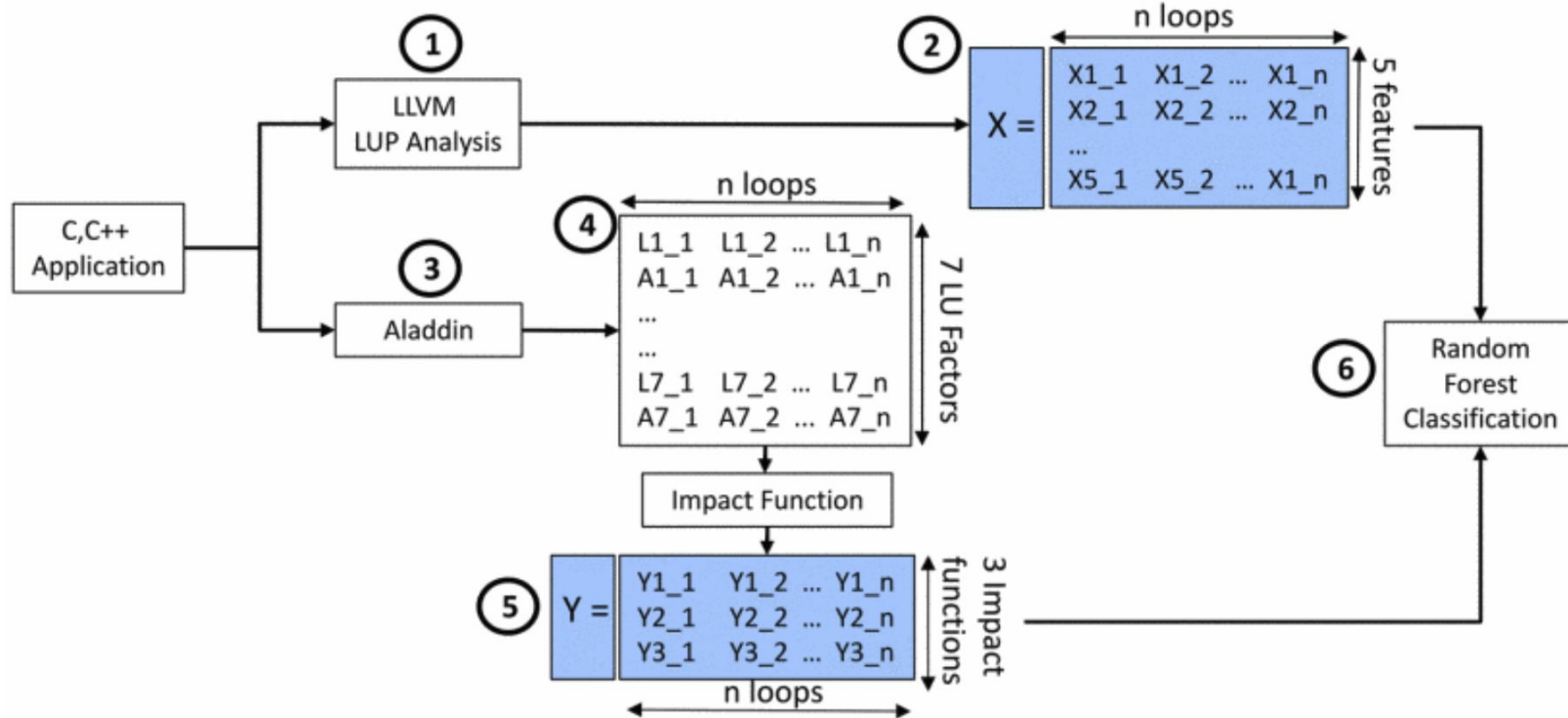
# High Level Synthesis (HLS)

# Hardware Loop Unrolling

- HLS is often used to customize hardware to **optimize loops** for **hardware accelerators**.
- Loop unrolling impacts the **performance** of hardware accelerators due to:
  - **High area cost** for duplicated logic
    - Area cost is the size of the chip occupied by the accelerator
  - **Loop-carried dependencies** & **frequent memory access**
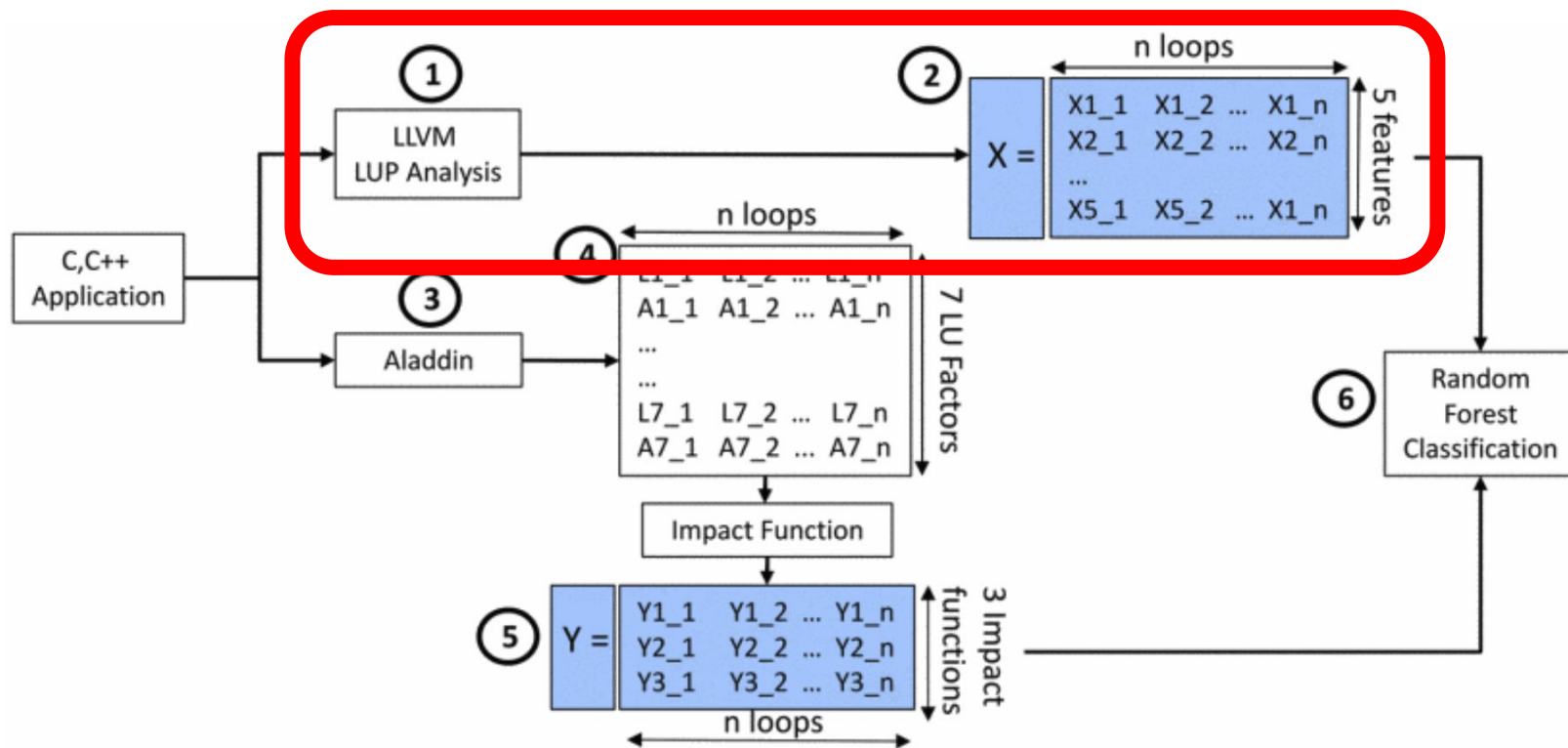    - Causes accelerators to act more sequentially

# Paper's Contributions

1. Trained a **Random Forest (RF) classifier** to predict unrolling factors for loops in HLS designs
2. Developed an **automated framework** in LLVM to **extract features** (for both training and inference) as input to the classifier
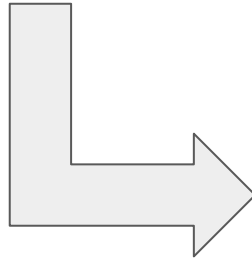
# Method

# Method

# Steps 1 & 2: Feature Selection for RF Classifier

1. Length of critical path
2. Loop trip count
3. The presence of loop carried dependencies
4. The # of load instructions
5. The # of store instructions

different features compared to original MIT paper

FEATURE VECTORS SELECTED BY STEPHENSON ET AL. [9].

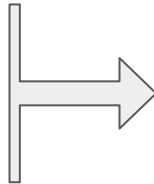| Features - X Vector 1 | Features - X Vector 2 |
|---|---|
| # Operands | # Floating Point Operations |
| Range Size | Loop Nest Level |
| Critical Path | # Operands |
| # Operations | # Branches |
| Loop Trip Count | # Memory Operations |

# Steps 1 & 2: LLVM Feature Extraction Pass

**Algorithm 1** LLVM Analysis Pass - Loop Unrolling Prediction Analysis
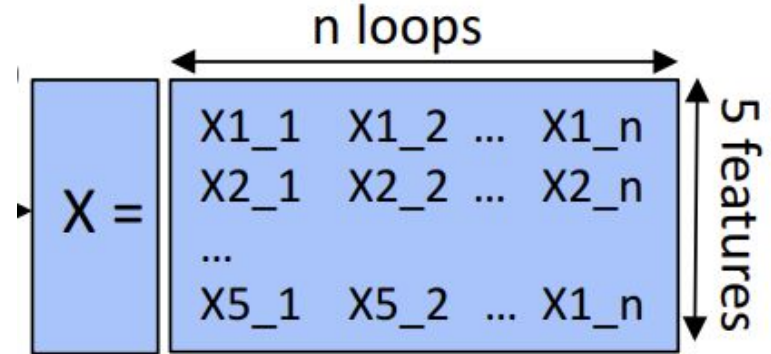
**Input:** Application written in C, C++
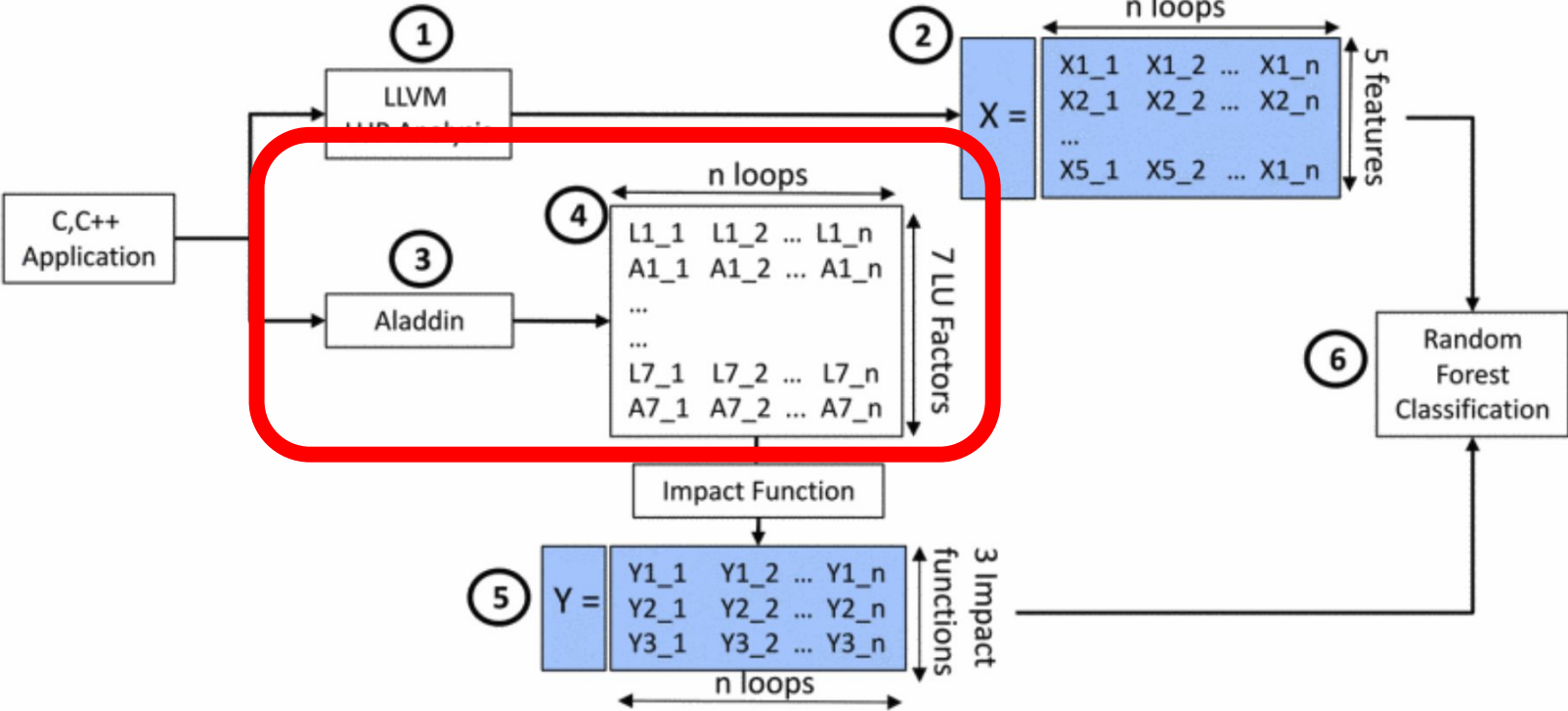**Output:** X (Feature Vector)

```
 1: function RunOnFunction( )
 2:     for BB in Function  do
 3:         if L=getLoopForBB()   then
 4:             LoopUnrollingPredictionAnalysis(BB,L)
 5:
 6: function  LoopUnrollingPredictionAnalysis(Basic  Block
    BB, Loop L)
 7:     LI=getLoopInfoAnalysis()
 8:     SE=getScalarEvolutionAnalysis()
 9:     DA=getDependenceAnalysis()
10:     /* Gather Features for X Vector */
11:     x1=getCriticalPath(BB)
12:     x2=getTripCountForLoop(L)
13:     x3=getLoopCarriedDependencies(BB)
14:     x4=getNumberOfLoadInstructions(BB)
15:     x5=getNumberOfStoreInstructions(BB)
```
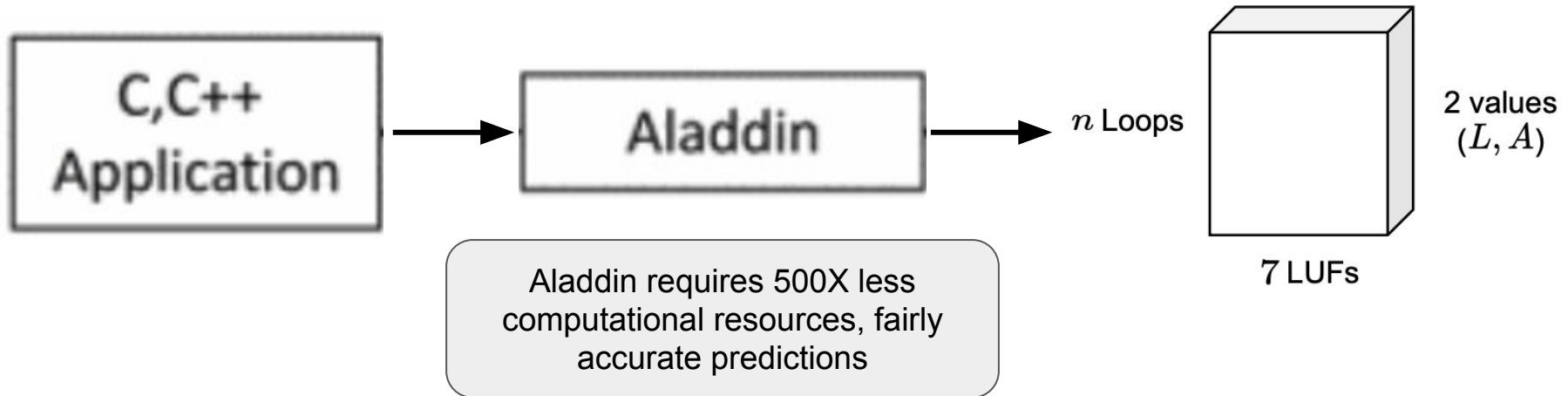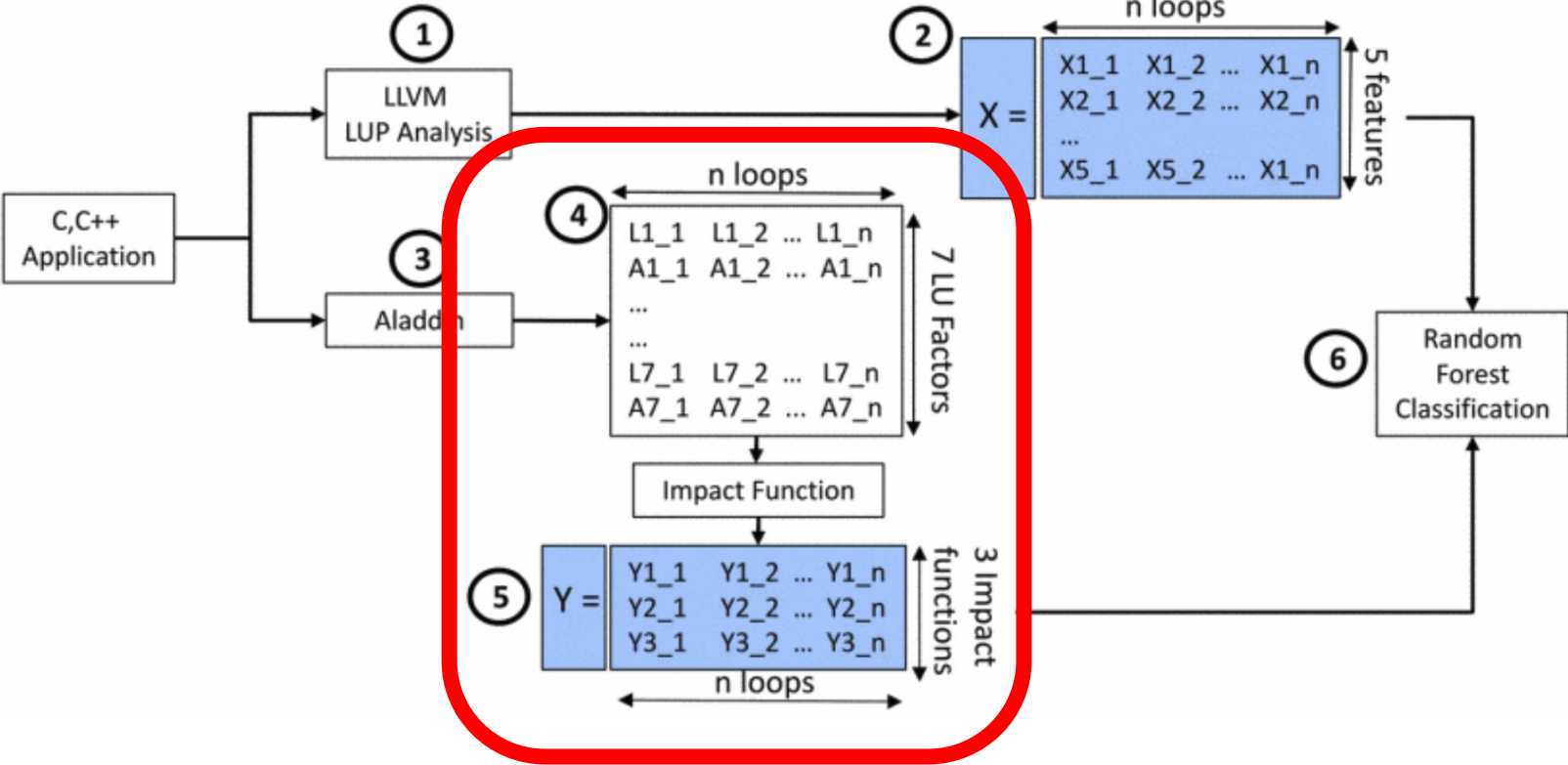


extracts all 5 features for each loop

$$X = \begin{matrix} X1\_1 & X1\_2 & ... & X1\_n \\ X2\_1 & X2\_2 & ... & X2\_n \\ ... & & & \\ X5\_1 & X5\_2 & ... & X1\_n \end{matrix}$$

n loops

5 features

# Method

# Steps 3 & 4: Approximating L & A with Aladdin

- Input: C/C++ programs to simulate
- Output: **latency** (L) and **area** (A) **values** for 7 loop unroll factors (LUFs)
  - LUFs: 1, 2, 4, 8, 16, 32, 64



C,C++ Application → Aladdin → $n$ Loops

7 LUFs

2 values $(L, A)$

Aladdin requires 500X less computational resources, fairly accurate predictions
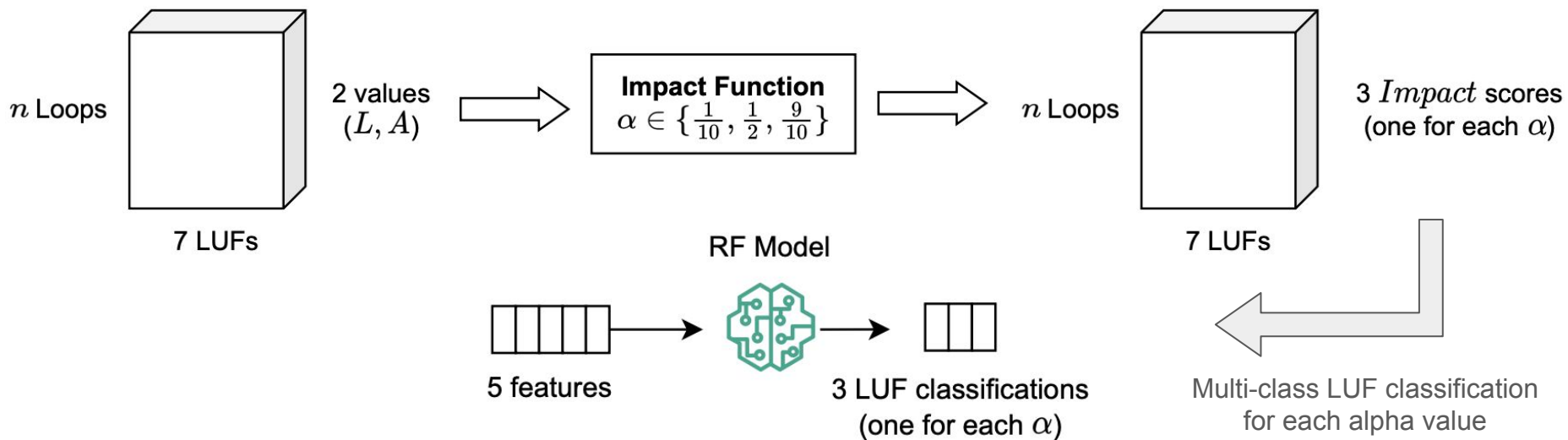
# Method

# Step 4 & 5: Impact Function

- Tradeoff between performance (latency) & required resources (area)

$$I(L,\ A) = \alpha \cdot \frac{(L_1 - L)}{L_1} + (1 - \alpha) \cdot \frac{(A_1 - A)}{A_1}, 0 \leq \alpha \leq 1$$
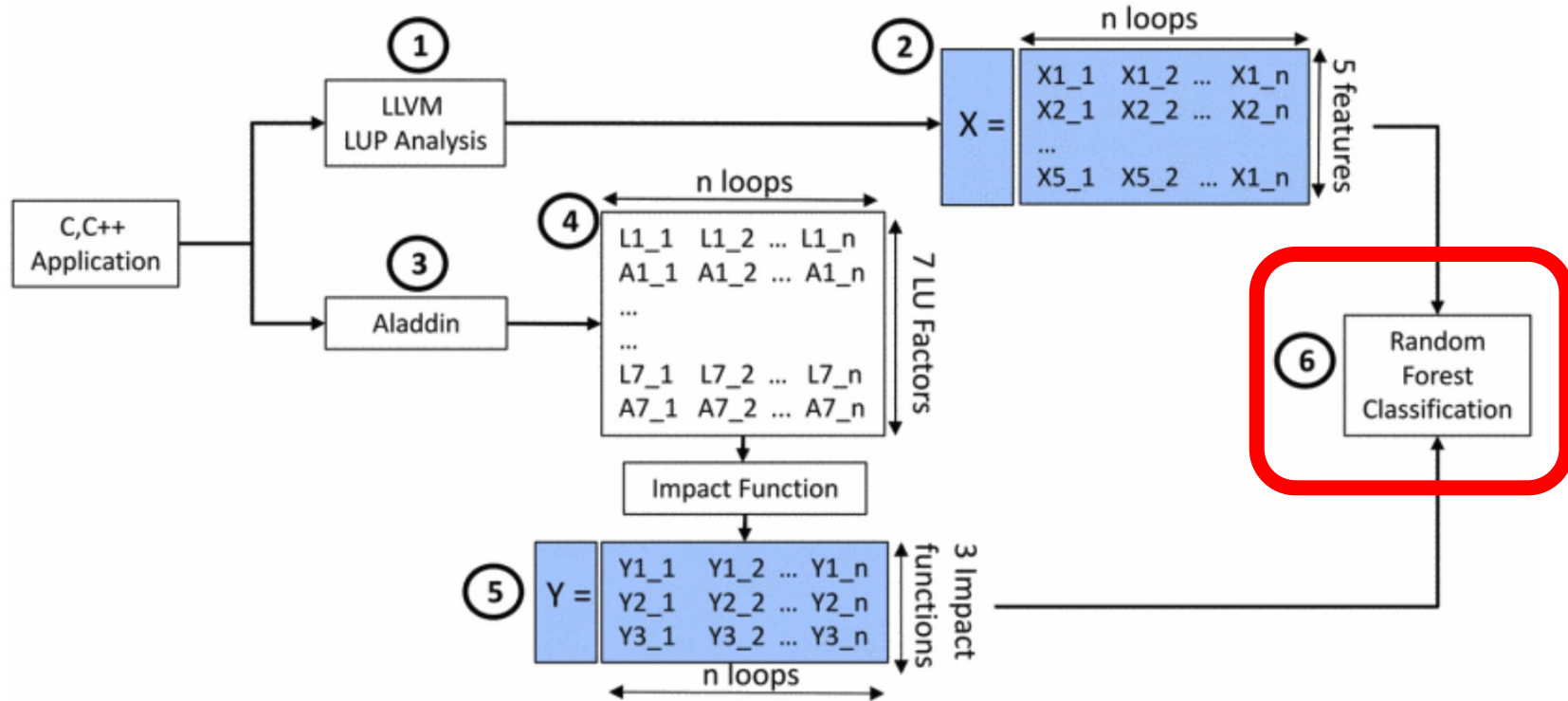
- L & A = latency & area of function synthesized as accelerator for chosen LUF
- $L_1$ & $A_1$ = latency & area of function synthesized as accelerator when LUF is 1
- α = relevance of latency & area
- The impact score is used to generate the ground truth for the RF model during training
  - The LUF that produces the highest impact score is considered the ground truth
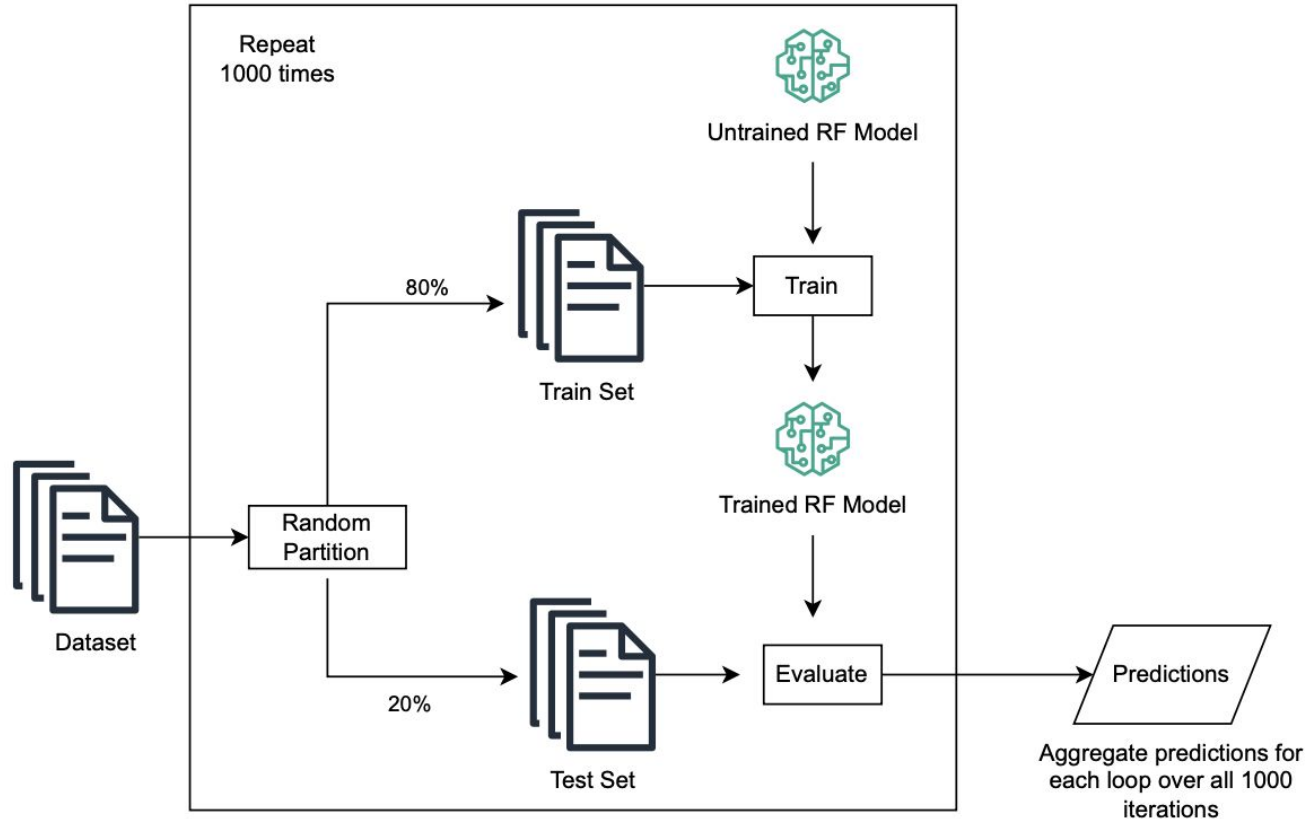
# Step 4 & 5: Impact Function

- Apply the impact function for 3 α values
    - α = 0.5: Optimize L & A equally
    - α = 0.9: Optimize L over A
    - α = 0.1: Optimize A over L

$n$ Loops

7 LUFs

2 values
$(L, A)$

**Impact Function**
$\alpha \in \{\frac{1}{10}, \frac{1}{2}, \frac{9}{10}\}$

$n$ Loops

7 LUFs

$3\ Impact$ scores
(one for each $\alpha$)

RF Model

5 features

3 LUF classifications
(one for each $\alpha$)

Multi-class LUF classification
for each alpha value

# Method

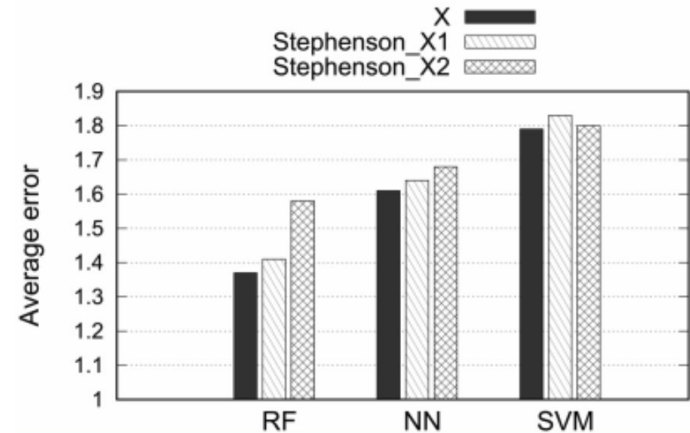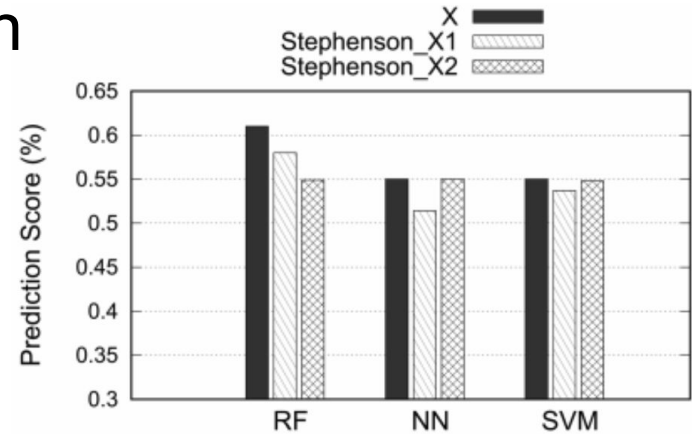# Step 6: Random Forest (RF) Classification



**Evaluation Metrics**

- *Prediction score:* percentage of optimal LUFs correctly identified on test set

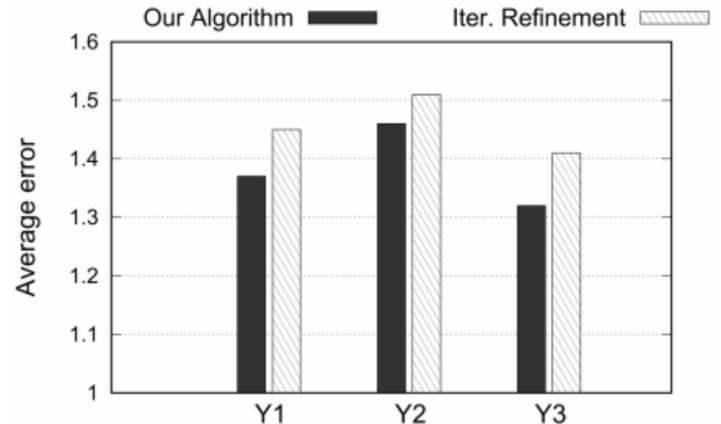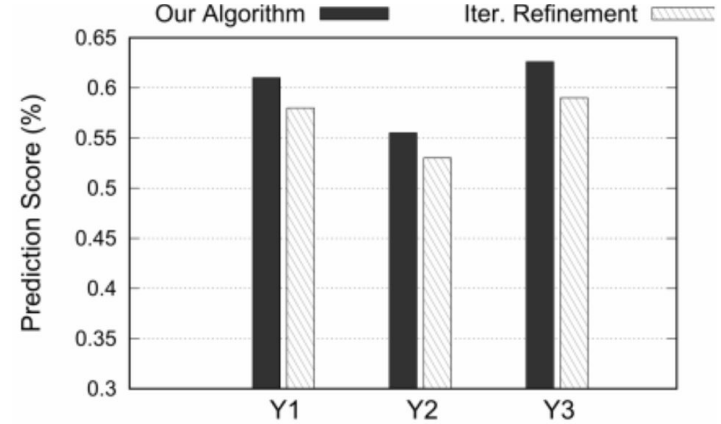- *Average error:* average distance between correct & predicted LUFs

# Results: Model & Feature Selection



- Compared feature selection against Stephenson et. al. (original MIT paper)
- α = 0.5
- Used 3 different models: RF, Nearest Neighbor (NN), and Support Vector Machine (SVM)
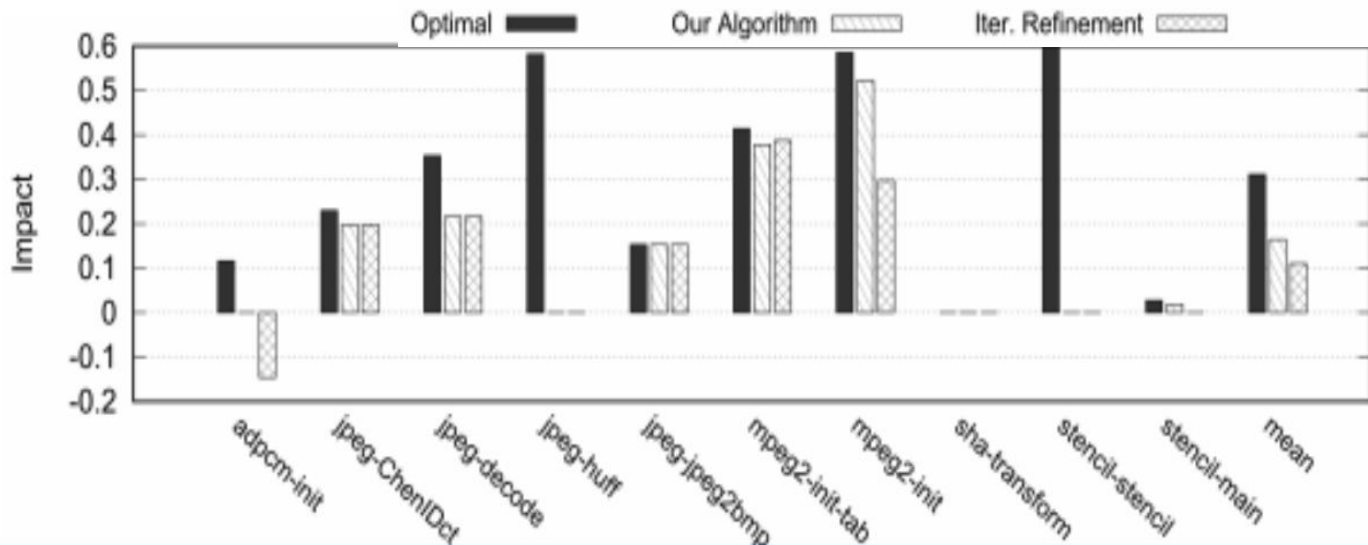- Paper presents **superior feature selection** and **choice of model**

# Results: Training Method

- Compared against an Iterative Refinement (IR) approach
  - IR first trains the model on a dataset, and then trains it again on a disjoint dataset
- Y1: α = 0.5
- Y2: α = 0.9
- Y3: α = 0.1
- Paper algorithm **outperforms** Iterative Refinement

# Results: Impact Score

- Compared against an Iterative Refinement (IR) approach
- Improvement or comparable impact against IR in all cases
- Poor performance on jpeg-huff and stencil-stencil, tests with complex loop structure

# Group Commentary

- Positives:
  - Tested 3 α values to consider the latency/area tradeoff
  - Considered some loop dependencies (MIT considered none)
  - Using Aladdin over full synthesis for computational efficiency
- Limitations:
  - Model may be too simplistic to understand complex looping structures
    - Evidenced by poor performance on jpeg-huff & stencil-stencil
  - Their 80-20% train/test split is random across all 1000 iterations
  - Feature extraction may be too simplistic
    - "Has Loop Carried Dependencies" feature - does not consider dependency distance
  - Regression to predict impact scores may be more suitable than classification
    - The LUF that maximizes predicted impact score should be selected
    - In classification, if ground truth LUF = 6, a predicted LUF of 1 and 5 would be considered equally wrong

# Thank you!

# Q&A