

RioSet: A Design-Centric Open-Source Dataset for Enhanced EDA Applications in Machine Learning



Yifei Zhu, Zhenxuan Luan, Guohua Yin, Xinze and Zhangxi Tan*

RISC-V International open-source Laboratory, Tsinghua University

OVERVIEW

RioSet, a pioneering open-source Electronic Design Automation (EDA) dataset, is tailored specifically for various Machine Learning (ML) tasks. Compatible with OpenEDA tools, RioSet abstracts the EDA system into a universal set of features, organized into multiple ML-specific views. Each view encapsulates elements and features vital for training, enabling a structured and focused approach to optimize system performance.

PROBLEM FORMULATION

- The hardening process can be formalized as optimizing DB feature sets F through operations s_1 to s_N :

$$s_1(F_1) = F_{s1}, \dots, s_N(F_N) = F_{sN} = F_{\text{final}} \quad (1)$$

- Optimization methods o_i aim to minimize cost functions f_i for each task:

$$\arg \min_{o_i} f_i[o_i(F_i)] \quad (2)$$

- Finding suitable views for a task f_i is solving Eq.2's dual:

$$\arg \min_{F_i} f_i[o_i(F_i)] \quad (3)$$

FEATURE EXTRACTION FROM OPEN DB

Key contributions:

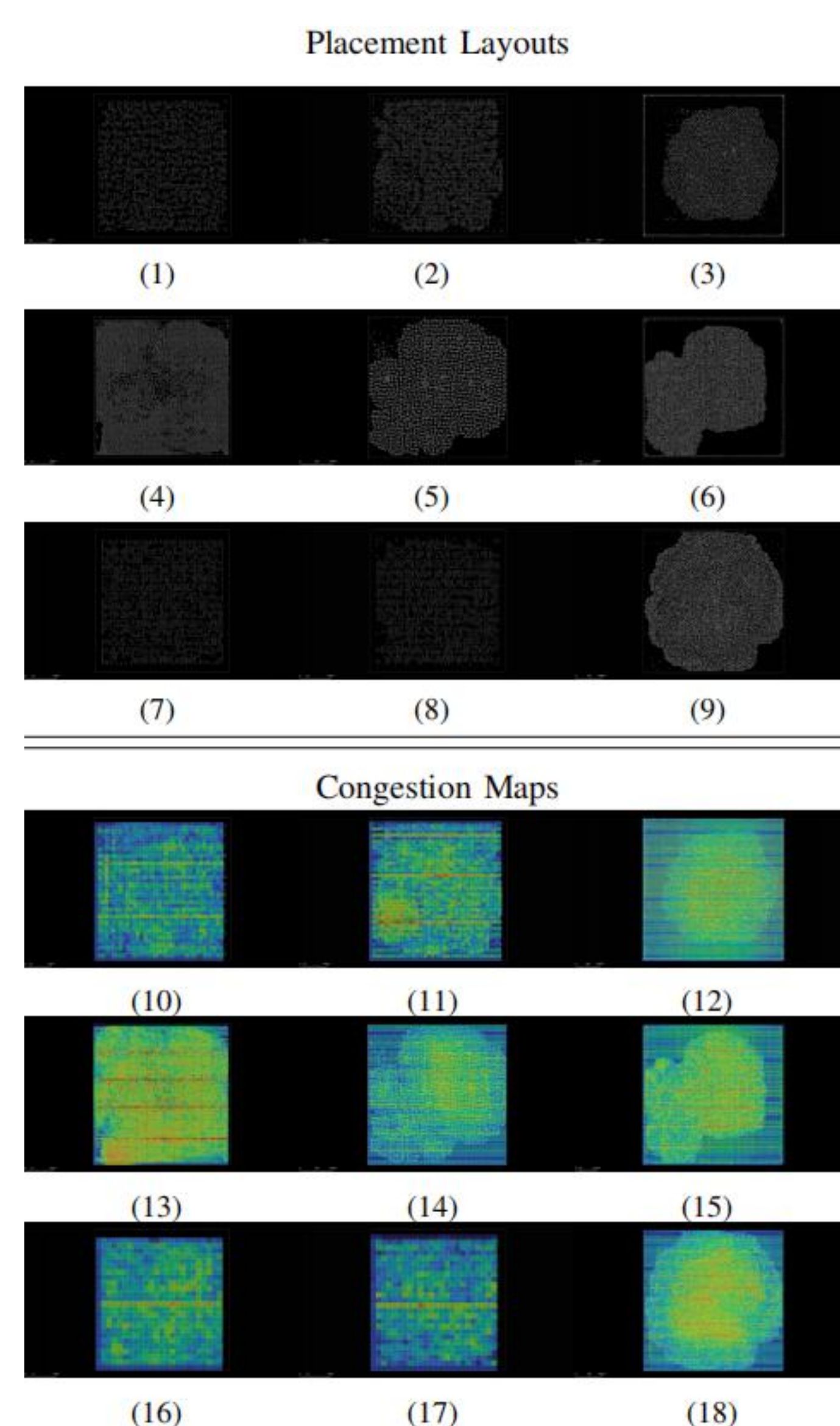
- ❖ 1. Predefined optimization views;
- ❖ 2. HDL classification using LLMs;
- ❖ 3. Full compatibility with open-source EDA tools.

```

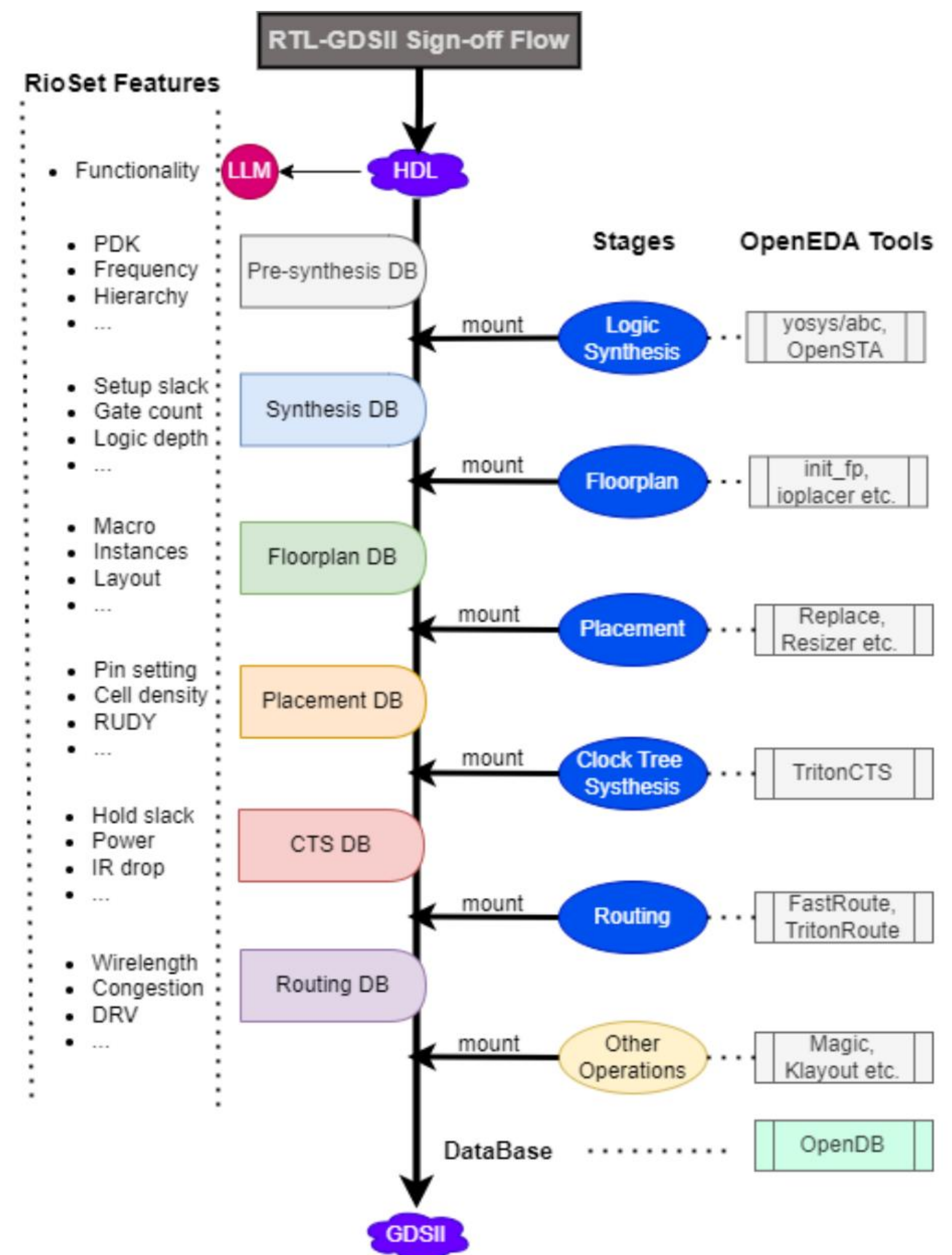
1 "CLASS": "
2   function_module",
3 "PDK": "sky130A",
4 "CELL": "fd_sc_hd",
5 "AREA": 128718,
6 "Hierarchy": 0,
7 "MACRO": 0,
8 "STRATEGY::AREA0": {
9   "GATE_COUNTS":
10    12229,
11   "SYN_AREA": 128718
12 },
13 "CLOCK_PERIOD::11.35": {
14   "SYN_SLACK_MAX":
15    "5.65",
16   "START": "Y87[1]",
17   "END": "_21258_",
18   "PATH_GROUP": "clk",
19   "INPUT_DELAY": 2.27,
20   "DEPTH": 9,
21   "STATUS": "MET",
22   "SYN_SLACK_MIN":
23    "0.17"
24 },
25 "CLOCK_PERIOD::3": {
26   "SYN_SLACK_MAX":
27    "-0.76",
28   "START": "_21366_",
29   "END": "_21380_",
30   "PATH_GROUP": "inout",
31   "INPUT_DELAY": 0,
32   "DEPTH": 12,
33   "STATUS": "VIOLATED",
34   "SYN_SLACK_MIN":
35    "0.17"
36 }
37 }
38 "NET": 14145,
39 "AREA": 143772,
40 "VIAS": 110173,
41 "DIE": [0.0,700,700],
42 "WIRELENGTH": {
43   "li1": 0,
44   "met1": 342190,
45   "met2": 354733,
46   "met3": 123665,
47   "met4": 108018,
48   "met5": 365
49 },
50 "PIN": {
51   "NUM": 749,
52   "MET2": 385,
53   "MET3": 362,
54   "POWER_PIN": "2
55   met4 met5"
56 },
57 "POWER": {
58   "INTERNAL": "2.34e
59   -02",
60   "SWITCHING": "8.31e
61   -03",
62   "LEAKAGE": "6.52e
63   -08",
64   "TOTAL": "3.17e-02"
65 },
66 "CLOCK_PERIOD::8": {
67   "PL_WNS_MAX":
68    "3.63",
69   "PL_WNS_MIN":
70    "0.10",
71   "RT_WNS_MAX":
72    "3.60",
73   "RT_WNS_MIN":
74    "0.09"
75 }
76 }

```

(a) logic synthesis (b) physical implementation



Layouts and Congestion Maps in RioSet (partial)



RioSet: available features and extraction stages

Textual Representation of Features (excerpts)

LLM FOR HDL CLASSIFICATION:

(a)-(c) and (d)-(e) are obfuscate and enhance operations respectively. A key insight is the crucial role of top module names in classification outcomes.

Correct: 12	core:2, storage: 1 communication: 2 function module: 7	Operation (a)	100% Correct
		Operation (b)	83.3% Correct
		Operation (c)	41.7% Correct
Incorrect: 5	core: 1 storage: 1 communication: 1 function module: 2	Operation (d)	80% Incorrect
		Operation (e)	10% Incorrect

Relevant dataset and auto-generation scripts are open-source and available on:

GitHub: <https://github.com/b224hisl/RioSet>

Email: zhuyf20@mails.tsinghua.edu.cn

