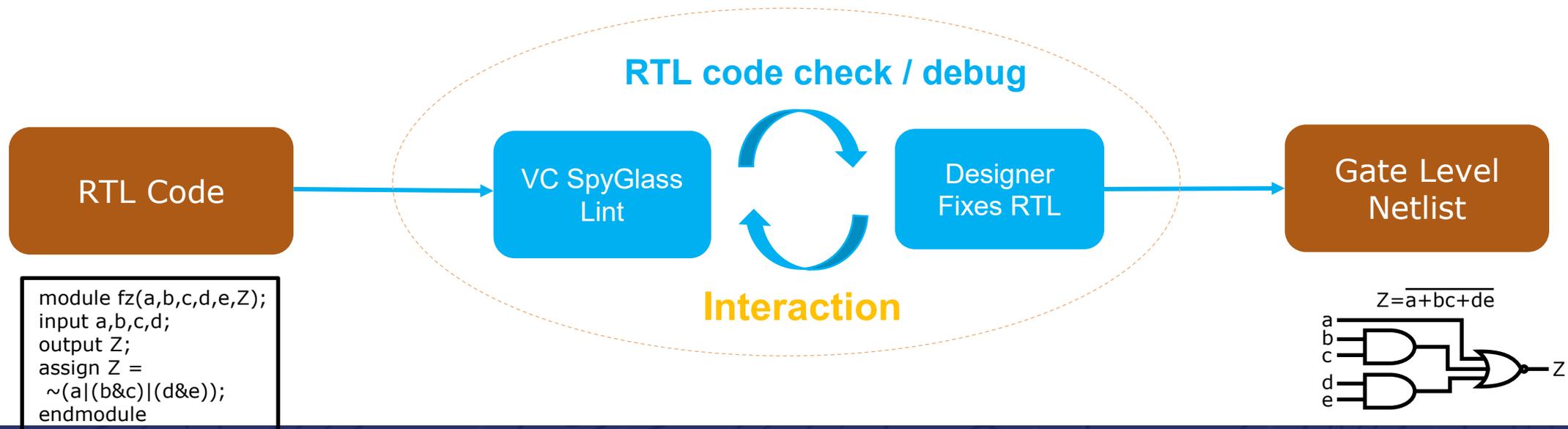# Agenda

- Problem Statement

- How to Automatically Fix RTL Code With AI?
  - Using Phison in-house solution
  - The results of using Phison in-house solution

- Fine-Tuning with Phison aiDAPTIV+ to Improve Correction Rate
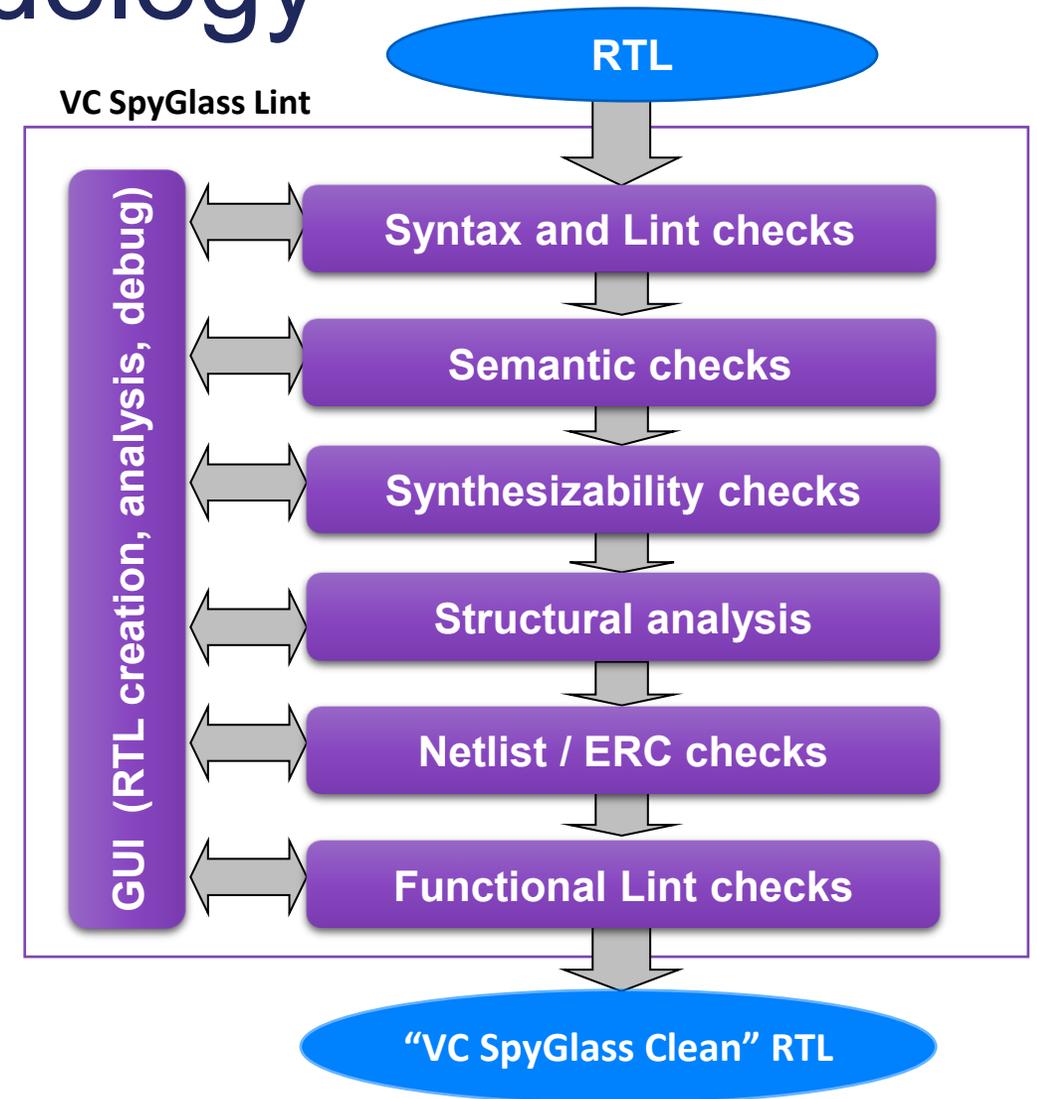  - Training LLM for better results

- Conclusions

- Future Works

# Problem Statement

- During initial stage of using VC SpyGlass Lint for RTL code development, debugging often results in numerous rule violations, which can consume a significant amount of time for designers.

- Therefore, we propose leveraging automatic generation capabilities of AI to automatically fix RTL, reducing time developers spend on manual corrections and thus improving efficiency.
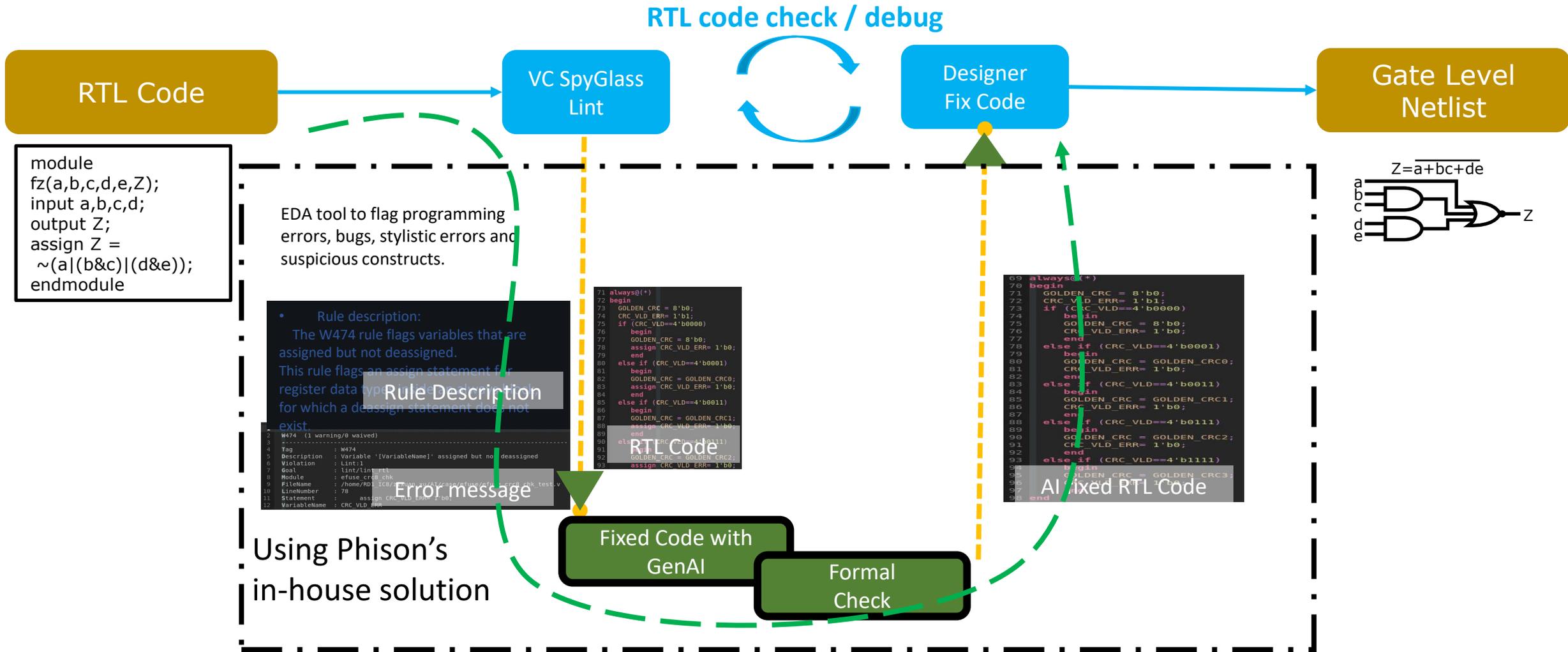
# VC SpyGlass Lint Methodology

- Find & fix bugs at source (correct-by-construction RTL)
  - Compliance to coding guidelines, STARC, OpenMore, Morelint, etc.
  - Synthesizability & simulation issues
  - Structural, logical and connectivity issues
  - Electrical rule checks
- Verilog, VHDL, SystemVerilog and mixed RTL support

- Structured methodology and templates help tackle design issues systematically

- Comprehensive waiver support

- Easy debug with cross-probe to RTL in Verdi GUI

**VC SpyGlass Lint**

**RTL**

**GUI (RTL creation, analysis, debug)**

- Syntax and Lint checks
- Semantic checks
- Synthesizability checks
- Structural analysis
- Netlist / ERC checks
- Functional Lint checks

**"VC SpyGlass Clean" RTL**

# How to Automatically Fix RTL Code With AI?

Using Phison in-house solution

# Automatically Fix RTL Lint Violations Using AI

# Prerequisites for Automatic Gen. of Lint Viol. Fixes

Selecting Correct LLM, Server Platform, Prompt, Token, and Rule

**Select a Suitable AI Model for Generating RTL Code**

LLM Selection

Server Platform

**Select a platform with reasonable pricing and high quality as the basis for AI-generated content**

**The size of the token and the accuracy of the prompt will significantly affect the accuracy and completeness of the AI-generated content**

Prompt & Token

Classification & Statistics

**Statistical Analysis of Commonly Waived Lint Rules and Categorization of Lint Rules Suitable for AI Correction**

# Violation Rule Classification & Statistics

- Choosing Lint rules to be targeted for automated generation of RTL Fixes
- Designer first categorizes violations to determine which types of violations have potential to be successfully corrected by AI and will also compile statistics on the most common violations in each project.

| LINT RULE | Classification | Project 1 Violation Number |
|---|---|---|
| STARC05-3.3.1.4b | C | 9448 |
| W164b | B | 4397 |
| AsgnToOneBit-ML | A | 2553 |
| ImproperRangeIndex-ML | C | 2205 |
| W416 | B | 2050 |
| W164a | B | 1215 |
| W154 | A | 924 |
| FlopClockConstant | C | 753 |
| W362 | B | 700 |
| NoExprInPort-ML | A | 502 |
| STARC05-2.3.1.4 | C | 465 |
| ResetFlop-ML | C | 341 |
| W241 | C | 160 |

| LINT RULE | Classification | Project 2 Violation Number |
|---|---|---|
| STARC05-3.3.1.4b | C | 1171 |
| W164a | B | 906 |
| W287a | C | 718 |
| W362 | B | 440 |
| STARC05-2.2.3.1 | A | 385 |
| RptNegEdgeFF-ML | C | 354 |
| W164b | C | 329 |
| NoExprInPort-ML | B | 294 |
| ResetFlop-ML | C | 280 |
| W443 | C | 236 |
| W241 | C | 229 |

**A: Can be Corrected, Accuracy 90%~100%**
**B: Potentially Correctable, Accuracy 60%~80%**
**C: Not Suitable, Requires Human Intervention**

# Providing Prompt for GenAI to Fix RTL Code

RTL Code
```
71 always@(*)
72 begin
73   GOLDEN_CRC = 8'b0;
74   CRC_VLD_ERR= 1'b1;
75   if (CRC_VLD==4'b0000)
76     begin
77       GOLDEN_CRC = 8'b0;
78       assign CRC_VLD_ERR= 1'b0;
79     end
80   else if (CRC_VLD==4'b0001)
81     begin
```

Rule Description
- Rule description:
  The W474 rule flags variables that are assigned but not deassigned.
  This rule flags an assign statement for register data types inside an always block, for which a

Error Message
```
1  --------------------------------------------------
2  W474  (1 warning/0 waived)
3  --------------------------------------------------
4  Tag          : W474
5  Description  : Variable '[VariableName]' assigned but not deassigned
6  Violation    : Lint:1
7  Goal         : lint/lint_rtl
8  Module       : efuse_crc8_chk
9  FileName     : /home/RD1_IC8/zixuan_xu/AI/case/efuse/efuse_crc8_chk_test.v
10 LineNumber   : 78
```

Prompt

**System Prompt:**
Clearly tell the AI what to do.

**User Prompt:**
Provide detailed items such as: analyze first, then fix the code, and provide examples for reference.
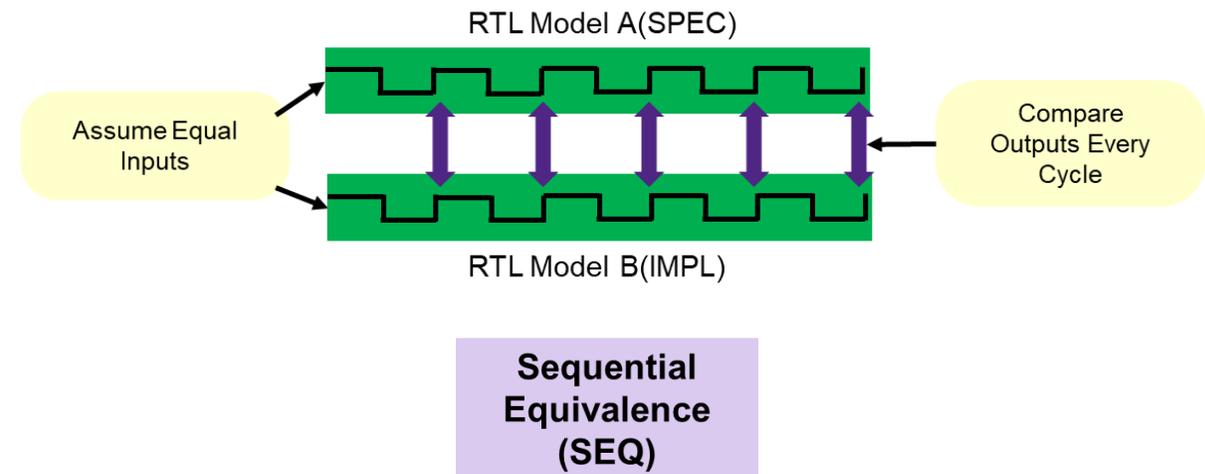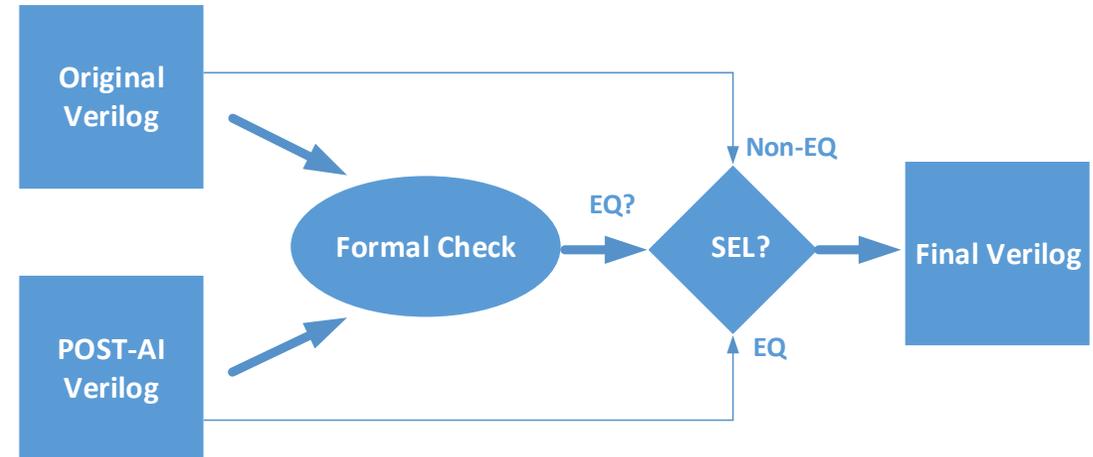
GenAI

Precise Inputs to GenAI is key.

The RTL code has been corrected by AI
```
69 always@(*)
70 begin
71   GOLDEN_CRC = 8'b0;
72   CRC_VLD_ERR= 1'b1;
73   if (CRC_VLD==4'b0000)
74     begin
75       GOLDEN_CRC = 8'b0;
76       CRC_VLD_ERR= 1'b0;
77     end
78   else if (CRC_VLD==4'b0001)
79     begin
80       GOLDEN_CRC = GOLDEN_CRC0;
81       CRC_VLD_ERR= 1'b0;
82     end
83   else if (CRC_VLD==4'b0011)
84     begin
85       GOLDEN_CRC = GOLDEN_CRC1;
86       CRC_VLD_ERR= 1'b0;
87     end
88   else if (CRC_VLD==4'b0111)
89     begin
90       GOLDEN_CRC = GOLDEN_CRC2;
91       CRC_VLD_ERR= 1'b0;
92     end
93   else if (CRC_VLD==4'b1111)
94     begin
95       GOLDEN_CRC = GOLDEN_CRC3;
96       CRC_VLD_ERR= 1'b0;
97     end
98 end
```

# Importance of Formal Verification in AI Generation Processes

- **Necessity of Formal Verification**
  - AI-generated results are probabilistic
  - There is a possibility of errors
- **Challenges in Automating RTL Code Correction**
  - Designers still need to confirm the correctness of AI corrections themselves
  - This affects the efficiency and purpose of automated RTL code correction
- **Solution**
  - Synopsys VC Formal Sequential Equivalence Checking (SEQ) can be used as an alternative tool
  - Adopt third-party Formal tools for post-generation verification

# How to Automatically Fix RTL Code With AI?

The results of using Phison in-house solution

# Automatically Fix Code on Real Project

The following table presents real project cases where the flow was used to:

| Lint Rule | Initial Violations (A) | POST-AI Remaining[1] Violations at once (B) | Correction Rate (C) |
|---|---|---|---|
| AsgnToOneBit-ML | 2553 | 1167 | **54.29%** |
| W154 | 920 | 152 | **83.48%** |
| STARC05-2.2.3.1 | 9 | 2 | **77.78%** |
| NoExprInPort-ML | 492 | 340 | **30.89%** |
| W164b | 4395 | 4294 | **2.3%** |
| W164a | 1211 | 1132 | **6.52%** |
| W416 | 2050 | 1970 | **3.9%** |
| W362 | 700 | 611 | **11.9%** |

# Post-AI Correction Violations Handling: At once vs In segments

Correcting multiple rules simultaneously in a single block has been shown to affect the correction rate (RTL Fixes). Therefore, executing corrections in segments can effectively improve accuracy.

| Lint Rule | Post-AI Correction Rate at once | Post-AI Correction[1] Rate in segments | Once vs SEG. DIFF. Rate | Segments |
|---|---|---|---|---|
| AsgnToOneBit-ML | 54.29% | **63.22%** | +8.93% | Segment 1 |
| W154 | 83.48% | **88.26%** | +4.78% | |
| NoExprInPort-ML | 30.89% | **42.48%** | +11.59% | |
| STARC05-2.2.3.1 | 77.78% | **77.78%** | +0% | |
| W164b | 2.3% | **34.74%** | +32.44% | Segment 2 |
| W164a | 6.52% | **15.85%** | +9.33% | |
| W416 | 3.9% | **75.9%** | +72% | |
| W362 | 11.9% | **21.39%** | +9.49% | |

# RTL Correction Rate via Different AI Model Using Same Prompt

The same prompt might affect different model generation results as shown in following table:

| Lint Rule | Initial Violations | Mistral Large[1] POST-AI Remaining Violations | LLAMA[2] POST-AI Remaining Violations | QWEN[3] POST-AI Remaining Violations |
|---|---|---|---|---|
| AsgnToOneBit-ML | 2553 | 1167 | 2005 | 2026 |
| W154 | 920 | 152 | 174 | 173 |
| STARC05-2.2.3.1 | 9 | 2 | 2 | 2 |
| NoExprInPort-ML | 492 | 340 | 384 | 381 |
| W164b | 4395 | 4294 | 4239 | 4246 |
| W164a | 1211 | 1132 | 1275 | 1214 |
| W416 | 2050 | 1970 | 1975 | 1982 |
| W362 | 700 | 611 | 608 | 607 |
| Correction Rate (Post-AI Cor. / Initial Vio.) | - | **21.58%** | **13.52%** | **13.77%** |

# Analysis of the Relationship between AI Correction Time and Token Size

The size of tokens indeed affects the time it takes for AI to generate content. Filtering out large token files can help shorten the generation time.

| Lint Rule | Initial Violations | POST-AI Remaining[1] Violations | POST-AI Remaining[1] Violations |
|---|---|---|---|
| AsgnToOneBit-ML | 2553 | 1167 | 1454 |
| W154 | 920 | 152 | 268 |
| STARC05-2.2.3.1 | 9 | 2 | 2 |
| NoExprInPort-ML | 492 | 340 | 316 |
| W164b | 4395 | 4294 | 4376 |
| W164a | 1211 | 1132 | 1138 |
| W416 | 2050 | 1970 | 1997 |
| W362 | 700 | 611 | 640 |
| Token | - | <4000 | <2500 |
| Correction Rate (Post-AI Cor. / Initial Vio.) | - | **21.58%** | **17.34%** |
| Run time | - | **7H 29min** | **1H 55min** |

# Fine-Tuning with Phison aiDAPTIV+ to Improve Correction Rate

Training LLM for better results

# Fine-Tuning a Model Using LoRA to Improve Correction Rate

We first fine-tuned W164b to confirm whether it could improve the correction rate.
The experimental results showed that the correction rate indeed increased.

| Item | Initial Violation | POST-AI Remaining Violations | POST-AI Remaining Violations |
|---|---|---|---|
| W164b | 4395 | 2868 | **2746** |
| Correction Rate | - | 34.74% | **37.51%** |
| Fine-Tuning | - | X | **V** |

Extend GPU RAM: aiDAPTIV+
AI Model : Mistral Large 123b
Dataset : 1000 patterns
EPOCH : 3
Run time : 21H
GPU:A6000*8 ,RAM: 384GB

# Phison's aiDAPTIVCache Removes HBM Limitation

## Current AI Computing Architecture

### Limited by GPU HBM

18 H100 GPU to run LLaM2 (70B)

(Requires 3 DGX Chassis)



## Phison's aiDAPTIV+ AI Computing Architecture

### Unlimited Model Size[1]

Scale # GPU to match budget

(Requires 4 GPU and 2 SSD)

Extend GPU RAM

AI Application

PyTorch
AI Framework

*aiDAPTIVLink*
Middleware Library

GPU
Cost Efficient | *aiDAPTIVCache*
100 DWPD

*ai100E* aiDAPTIVCache
*w/ 100 DWPD*

1: Based on aiDAPTIVCatch capacity

# Conclusions

- **Automated Correction Using AI Generation Technology**:
  - By leveraging AI generation technology in conjunction with Lint violation reports, automated corrections were achieved, significantly improving work efficiency.
  - Formal tools were used to automatically verify consistency before and after corrections, ensuring the accuracy and reliability of each correction.

- **Significant Correction of Violations**:
  - Experimental results showed that this process can significantly correct at least 10% of violations, and in some cases, the correction rate is even higher.

- **Fine-tuning with Phison aiDAPTIV+**:
  - Using Phison aiDAPTIV+ with limited hardware support for fine-tuning further enhanced the accuracy of corrections, ensuring optimal correction outcomes.

# Future Work

- We also plan to explore GenAI powered **VC SpyGlass Lint Advisor with Phison's aiDAPTIV+** AI Computing Architecture for
  - Automated RTL Lint violation fixes, verifying design specific custom bug finding and automated waiver generation

# Final

# IC Design Flow Illustration: House-Building



Focus of this session:
RTL code check / debug

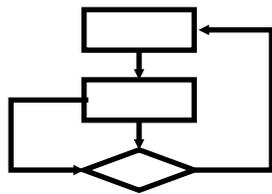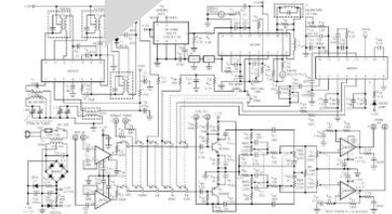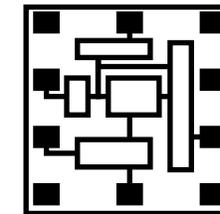| IC Specification & Functional Design | RTL Code | Pre-sim Synthesis | Gate Level Netlist | Placement Routing | Layout | Post-sim Verification | Tape-out |

```
module fz(a,b,c,d,e,Z);
input a,b,c,d;
output Z;
assign Z =
 ~(a|(b&c)|(d&e));
endmodule
```
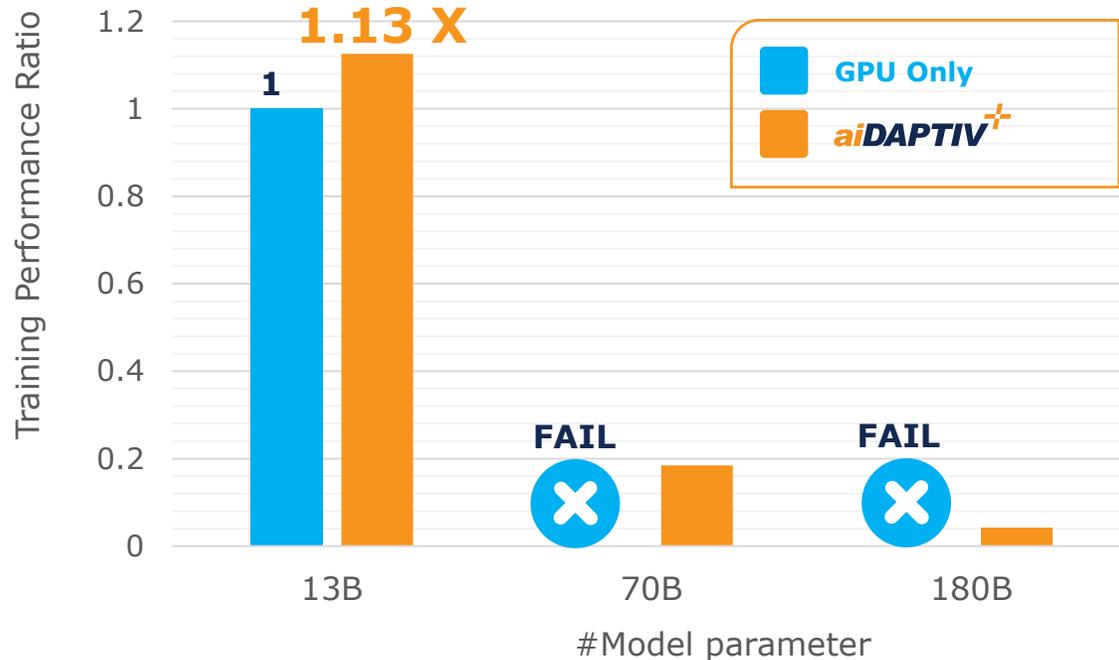
$Z = \overline{a + bc + de}$

Use VC SpyGlass Lint to check RTL code

# Performance Improvement with aiDAPTIV+

**w/ aiDAPTIV+ AI Training Server**

- The computing power and performance comparison for fine-tuning the LLAMA-2 13B/70B/ FALCON 180B model:
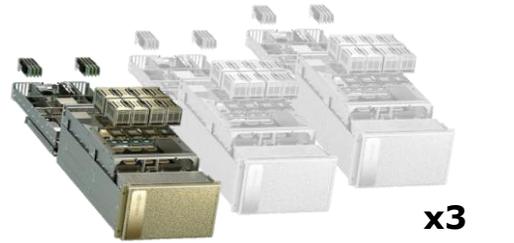


**aiDAPTIV+**
**Make it possible for larger model training**
**&**
**Enhances training performance in limit #GPU condition**

**System Configuration**
- CPU : Intel Xeon Gold6430
- GPU: **A6000*8**
- RAM: 64GB*16

# aiDAPTIV+ Benefit

Fine-tuning the **Llama 3.1 70B** model with **10M** tokens:



| | DGX H100 | Tower Workstation aiDAPTIV+ / 4000 Ada | Rack Workstation aiDAPTIV+ / A6000 | Server aiDAPTIV+ / A6000 |
|---|---|---|---|---|
| | • H100 GPU x24<br>• GPU Memory 1920GB<br>• Without SSD | • 4000 Ada GPU x4<br>• GPU Memory 80GB<br>• 2TB AI100E x2 | • A6000 GPU x4<br>• GPU Memory 192GB<br>• 2TB AI100E x2 | • A6000 GPU x8<br>• GPU Memory 384GB<br>• 2TB AI100E x2 |
| Price (USD) | $1,920K | $40K ⬇ 97.9% | $80K ⬇ 95.8% | $100K ⬇ 94.7% |
| Training Time (Hour) | 0.3 H | 7.6 H | 4.3 H | 2.2 H |