

# Artificial Intelligence applied to code repair after code static analysis verification

18th ESA Workshop on Avionics, Data, Control and Software  
Systems ~ ADCSS2024

Ainhoa López (TAS) / David de Fitero (UAH)

# Evolution of APR Methods and LLM Advances

- Traditional APR Methods
  - Heuristic-based: Uses predefined rules to generate patches (Saha et al., 2017)
  - Constraint-based: Reduces search space using logical restrictions (Xiong et al., 2017)
  - Template-based: Applies predefined correction patterns (J. Jiang et al., 2018)
- Modern ML Approaches
  - Frame APR as a translation task using RNNs, LSTMs, and Transformers
- Vulnerability Repair
  - General-purpose tools (e.g., Angelix (Mechtaev et al., 2016), concolic program repair (Shariffdeen et al., 2021)).
  - Security-specific methods (deductive reasoning, inductive reasoning, and static analysis)
- LLM developments in APR
  - LLMs with completion engines improve patch generation accuracy (Wei et al., 2023)
  - Fine-tuned LLMs surpass traditional APR tools (Huang et al., 2023; Silva et al., 2024; Xia et al., 2023)

# Followed Approach

- Datasets
  - CommitPackFt capturing real-world bug fixes
  - SonarQube rules dataset based on code quality rules
- Fine-tuning
  - QLoRA (Dettmers et al., 2023): Combines LoRA (Low-Rank Adaptation) with 4-bit quantization for resource-efficient fine-tuning
  - NEFTune (Jain et al., 2023): Adds noise to embeddings for improved performance
- Refinement after initial results
  - Expanded training data
    - Synthetic dataset to augment coverage of SonarQube rules, providing more diverse training examples
    - Added 650 manual corrections aligned with MISRA guidelines
  - Model upgrade - Llama 3 - 8B

# Description of the Datasets

- CommitPackFT
  - Filtered version of the CommitPack dataset
  - Contains **high-quality commit messages** resembling natural language instructions
  - Paired with **before and after code** corresponding to the commit
- Synthetic Dataset
  - Generated using **Llama 3 - 70B** model
  - Seed examples provided with **code before/after changes** based on **SonarQube** rules
  - Model generated **new samples** following similar patterns from seeds
- SonarQube Dataset
  - Extracted from SonarQube analysis
  - Aligned with **MISRA C** standards
  - Includes **human-generated fixes** for detected issues
  - **Smallest dataset** - upsampled for balance

# Training Process Overview

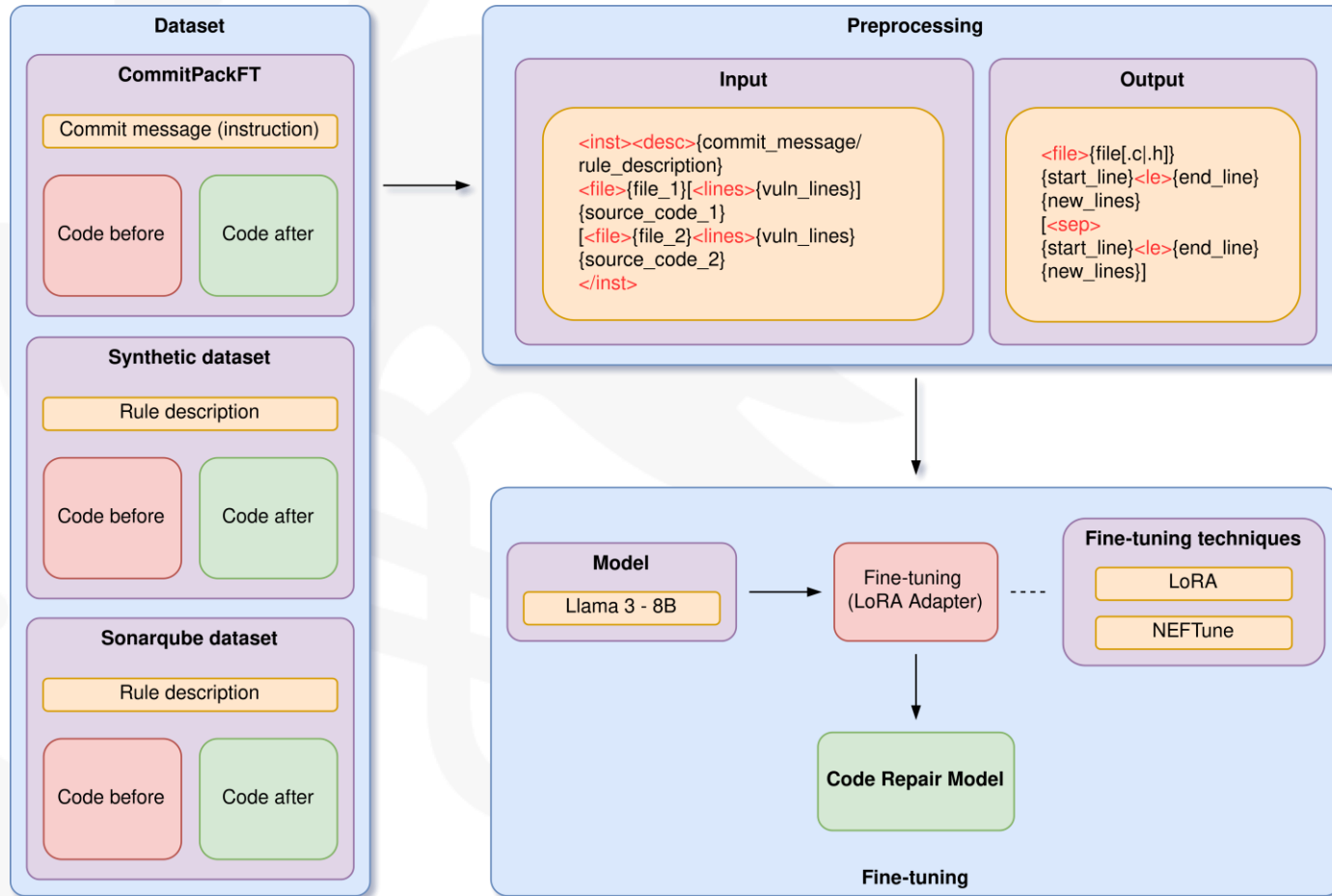
- Data Preprocessing
  - Input
    - Commit message or SonarQube rule
    - Source code with line numbers
  - Output
    - Line range of modification and the new lines to replace them
- Model & Techniques
  - CodeLlama / Llama 3 - 8B
  - Supervised Fine-tuning with QLoRA/LoRA and NEFTune
- Hardware
  - Trained on a NVIDIA H100

```

Input
-----
<inst><desc>Uppercase the literal suffix.
<file>file.c<line>424
420     if (dds_corr_magnitude > algo_wrapper_cmcu_2_man_cfg->fw_algo_fdir_cfg_4.dyn_fdir_dis_thr) {
421         algo_fdir_cycles = algo_wrapper_cmcu_2_man_cfg->fw_algo_fdir_cfg_4.dyn_fdir_dis_time;
422     }
423
424     if (algo_fdir_cycles > 0u) {
425         algo_fdir_cycles = algo_fdir_cycles - 1u;
426         (void) bsp_gpio_set_low(BSP_GPIO_FCA_EN);
427     }
428 </inst>

Output
-----
<file>file.c
423<le>426
    if (algo_fdir_cycles > 0U) {
        algo_fdir_cycles = algo_fdir_cycles - 1U;
  
```

# Approach diagram



# Refinement techniques

- Tailored contextual input
  - Use **tree-sitter** to reduce model input
  - Adjust context size based on the specific rule
- Additional Instructions
  - Adds specific instructions to certain rules to improve repair efficacy and ensure adherence to coding standards
- Simultaneous Fixing
  - Utilizes **batch inference** to accelerate processing when multiple repair rules are provided
- Quantization with EETQ
  - **Int8 quantization** reduces memory use with minimal accuracy loss
- Inference Engine
  - Uses **Text Generation Inference (TGI)** for fast and efficient model inference

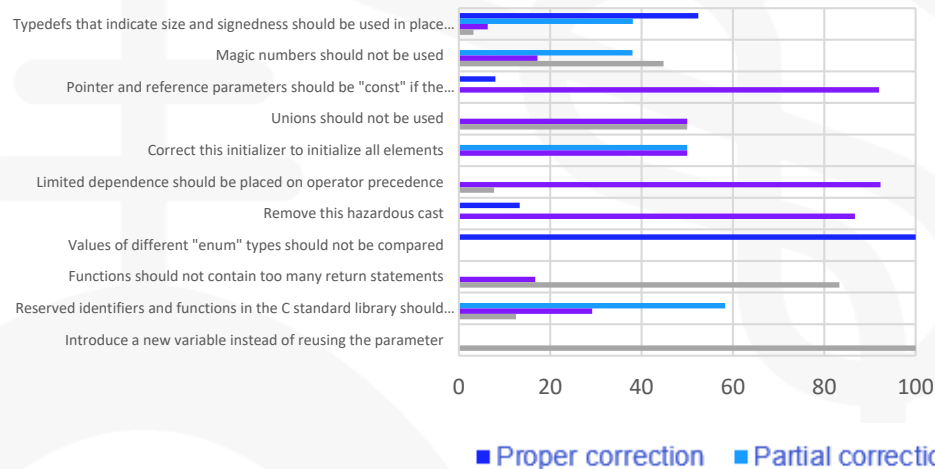




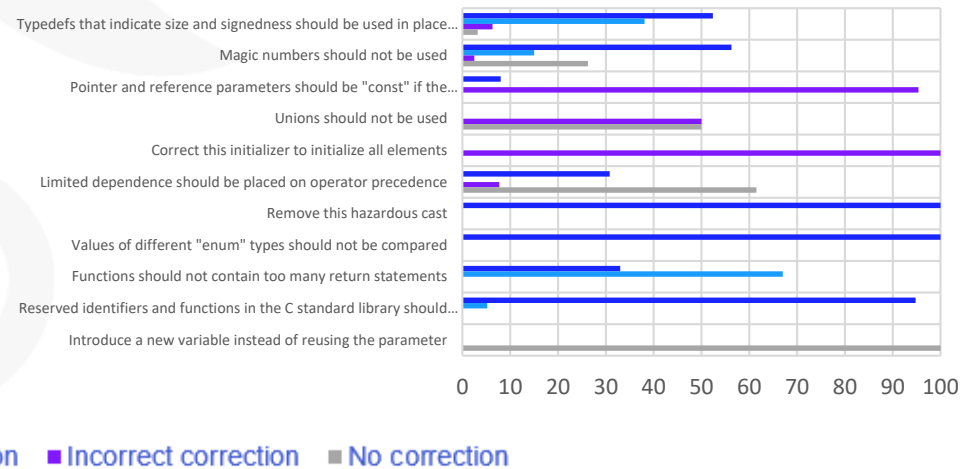
# Experiments

- Precision metrics
  - Results broken down for each rule: comparison CodeLlama vs Llama3
    - Variation in correction effectiveness
      - CodeLlama showed that the model's ability to handle different rules varied considerably.
      - Llama3 increased the proper and partial corrections but still shows variability in its correction ability

CodeLlama correction percentage by MISRAC 2012 Rule



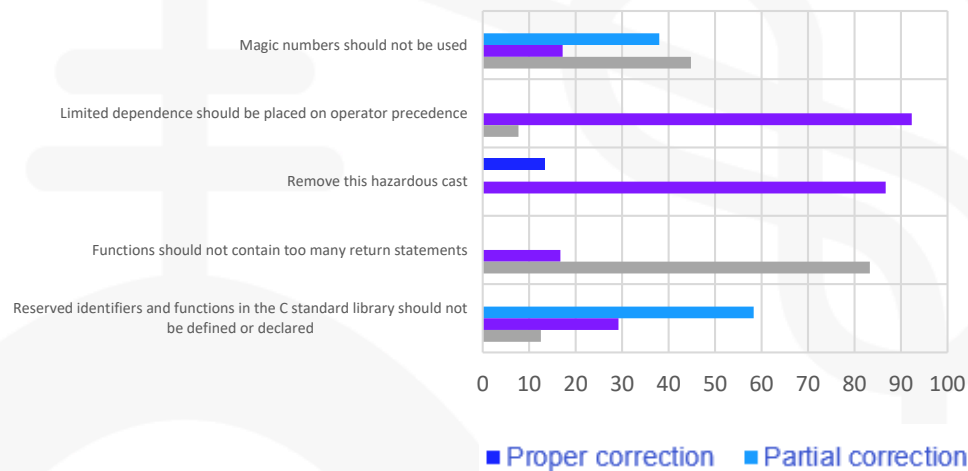
Llama3 correction percentage by MISRAC 2012 Rule



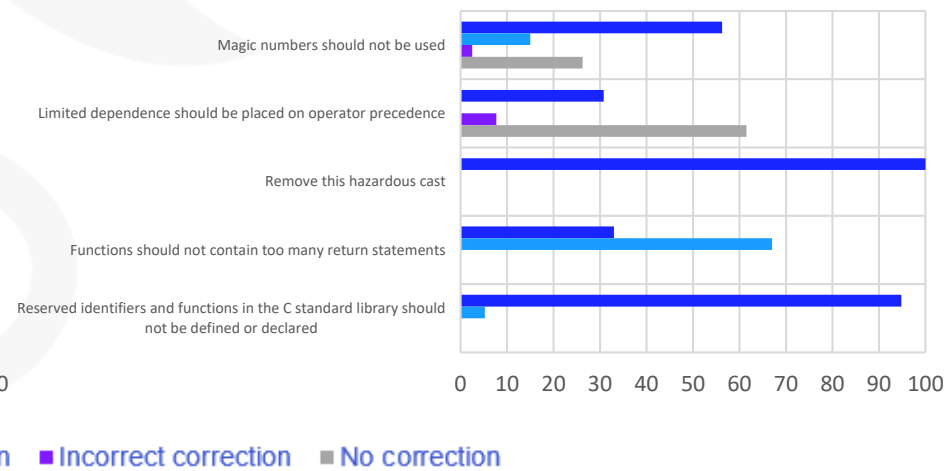
# Experiments

- Precision metrics
  - Results broken down for each rule: comparison CodeLlama vs Llama3
    - Variation in correction effectiveness
      - CodeLlama showed that the model's ability to handle different rules varied considerably.
      - Llama3 increased the proper and partial corrections but still shows variability in its correction ability

CodeLlama correction percentage by MISRAC 2012 Rule



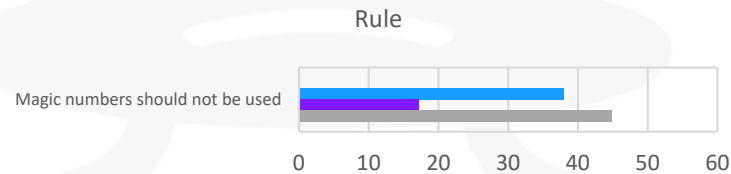
Llama3 correction percentage by MISRAC 2012 Rule



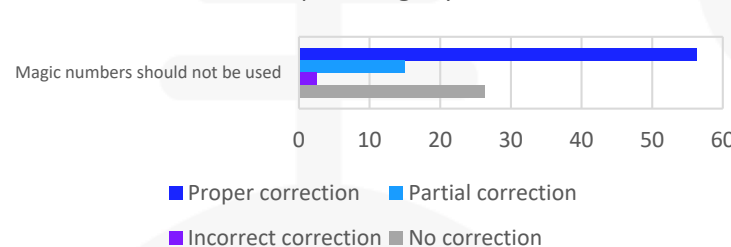
# Experiments

- Precision metrics: Variation in correction effectiveness
    - Improvement of Llama3 wrt CodeLlama
- Example: “Magic numbers should not be used”

CodeLlama correction percentage by MISRAC 2012



Llama3 correction percentage by MISRAC 2012 Rule



CodeLlama partial solution

```
#include <stddef.h>
#include <memctrl.h>
#include <01_054_ftmctrl.h>
#include <01_054_ftmctrl_hw_desc.h>
#include <01_054_ftmctrl_prom_io.h>
#include <01_054_ftmctrl_ram.h>
#include <01_054_ftmctrl_edac.h>
#include <01_054_ftmctrl_rs.h>

static int32_t ftmctrl_init_prom(ftmctrl_t *u, const mctrl_bus_t *bus)
{
    int32_t ret = 0;

    /* Configure bus width */
    switch (bus->width) {
    case PROM_WIDTH_8:
        /* Sets the data width of the PROM area ("00"=8) */
        ftmctrl_prom_set_width_8(u);
        break;
    case PROM_WIDTH_16:
        /* Sets the data width of the PROM area ("01"=16) */
        ftmctrl_prom_set_width_16(u);
        break;
    case PROM_WIDTH_32:
        /* Sets the data width of the PROM area ("10"=32) */
        ftmctrl_prom_set_width_32(u);
        break;
    default:
        ret = -1;
        break;
    }
}
```

Llama3 correct solution

```
#include <stddef.h>
#include <memctrl.h>
#include <01_054_ftmctrl.h>
#include <01_054_ftmctrl_hw_desc.h>
#include <01_054_ftmctrl_prom_io.h>
#include <01_054_ftmctrl_ram.h>
#include <01_054_ftmctrl_edac.h>
#include <01_054_ftmctrl_rs.h>

#define PROM_WIDTH_8    8U
#define PROM_WIDTH_16   16U
#define PROM_WIDTH_32   32U

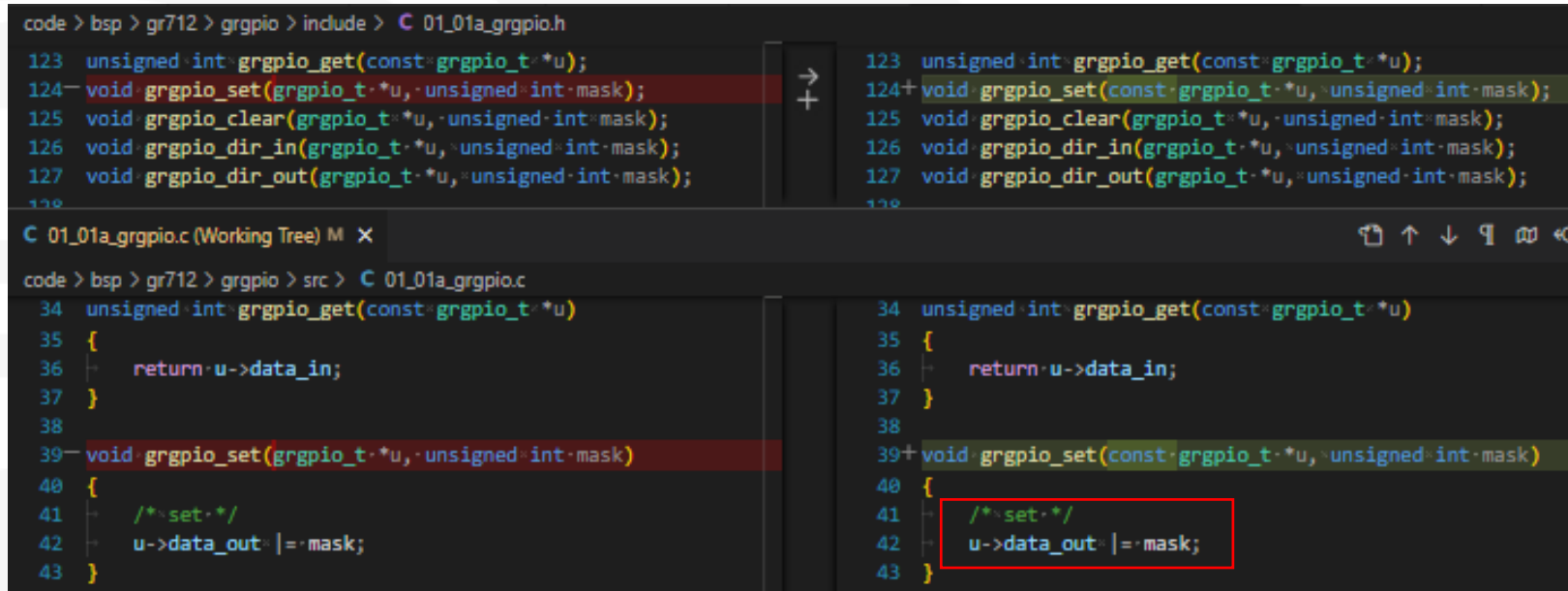
static int32_t ftmctrl_init_prom(ftmctrl_t *u, const mctrl_bus_t *bus)
{
    int32_t ret = 0;

    /* Configure bus width */
    switch (bus->width) {
    case PROM_WIDTH_8:
        /* Sets the data width of the PROM area ("00"=8) */
        ftmctrl_prom_set_width_8(u);
        break;
    case PROM_WIDTH_16:
        /* Sets the data width of the PROM area ("01"=16) */
        ftmctrl_prom_set_width_16(u);
        break;
    case PROM_WIDTH_32:
        /* Sets the data width of the PROM area ("10"=32) */
        ftmctrl_prom_set_width_32(u);
        break;
    default:
        ret = -1;
        break;
    }
}
```

# Experiments

- Precision metrics: False positive and code context
  - Still a weakness in the refined model

Example: "Pointer and reference parameters should be 'const' if the corresponding object is not modified":  
92% of corrections were incorrect.

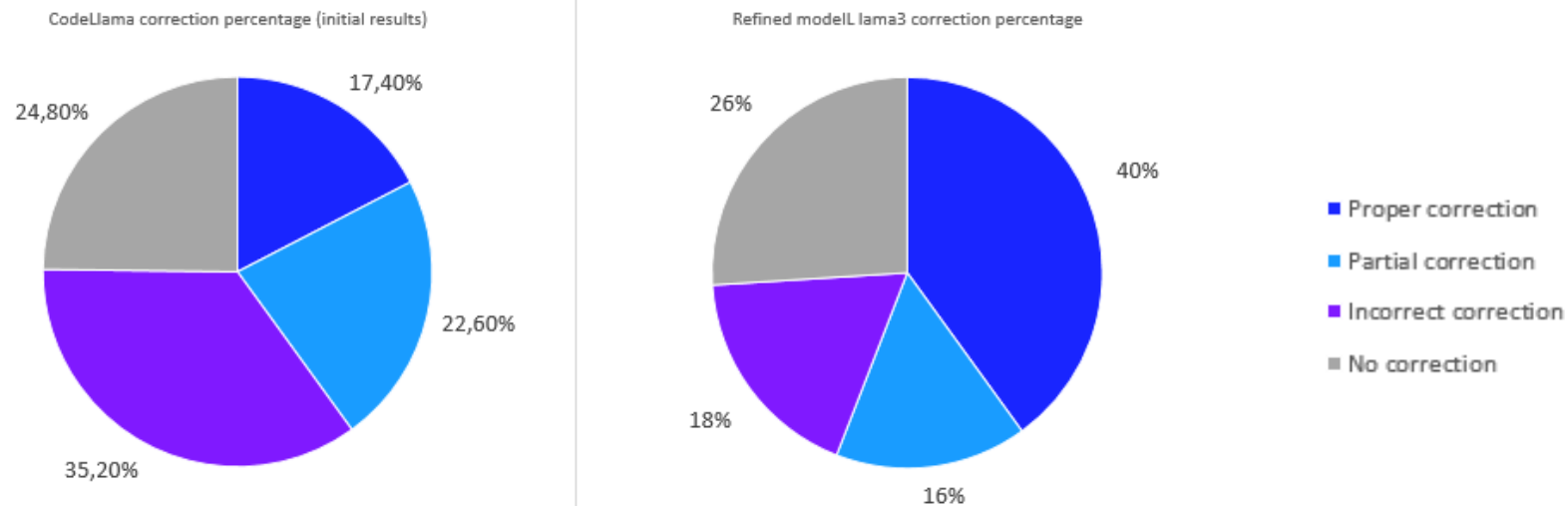


```
code > bsp > gr712 > grgpio > include > C 01_01a_grgpio.h
123 unsigned int grgpio_get(const grgpio_t *u);
124- void grgpio_set(grgpio_t *u, unsigned int mask);
125 void grgpio_clear(grgpio_t *u, unsigned int mask);
126 void grgpio_dir_in(grgpio_t *u, unsigned int mask);
127 void grgpio_dir_out(grgpio_t *u, unsigned int mask);
128

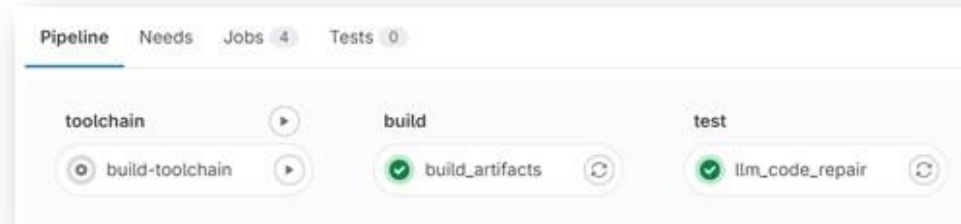
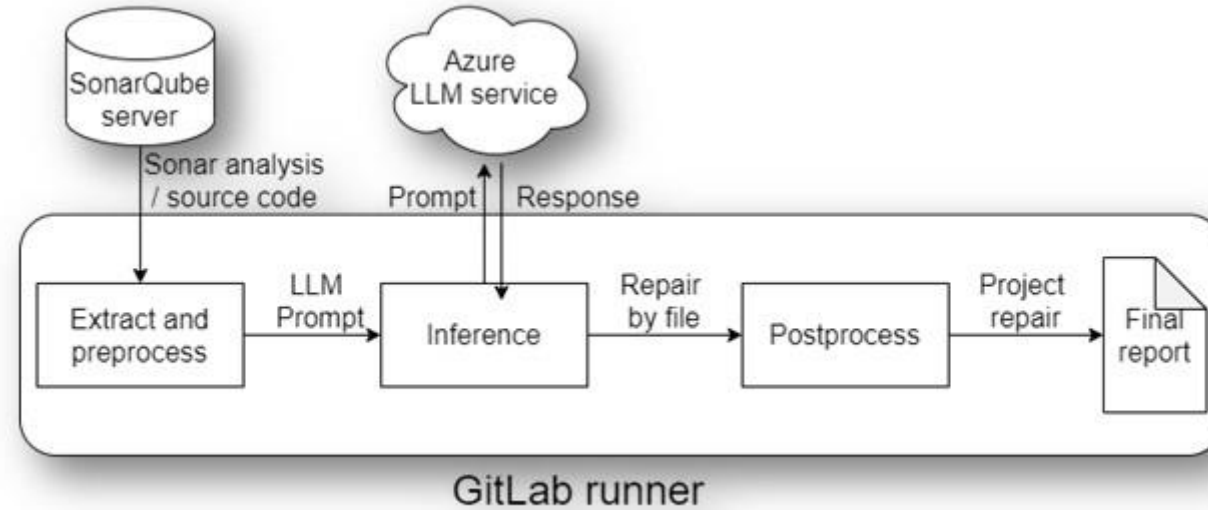
C 01_01a_grgpio.c (Working Tree) M X
code > bsp > gr712 > grgpio > src > C 01_01a_grgpio.c
34 unsigned int grgpio_get(const grgpio_t *u)
35 {
36     return u->data_in;
37 }
38
39- void grgpio_set(grgpio_t *u, unsigned int mask)
40 {
41     /* set */
42     u->data_out |= mask;
43 }
```

# Experiments

- Precision metrics: overall results
  - Fully accurate corrections: increased up to 40% → reduces manual workload and saves developer's time
  - Partial corrections: slightly decrease to 15,9% → provides useful insights to guide developers in resolving complex issues
  - Incorrect corrections: dropped to 18,1% → increased model reliability and less time spent reviewing errors



# Final Deployment



# Conclusions

- **Still not a full replacement for developers**
  - Despite its improvements, the model cannot yet replace developers entirely. Human oversight is still essential for more intricate, context-specific coding decisions and fine-tuning.
- **Better code quality and compliance**
  - The system's improved alignment with coding standards like MISRA ensures that the code is not only more reliable but also compliant with industry regulations. This contributes to long-term maintainability and reduces the risk of defects.
- **Increased confidence in automated suggestions**
  - The sharp decrease in incorrect corrections means developers can rely more on the model's suggestions, reducing time spent double-checking or correcting faulty outputs, thus streamlining the entire review process.
- **Support for complex issue resolution**
  - While full corrections have increased, the value of partial corrections lies in guiding developers through complex or edge-case issues, providing a foundation upon which they can make finer adjustments with confidence.