

LLM与安全风险治理 端到端解决校验缺失风险代码

刘彦南

字节跳动无恒实验室安全工程师



目 录

1. 安全风险治理的挑战与LLM机遇
2. 基于LLM的风险代码识别与修复
3. 实验结果与分析
4. 讨论 Q&A

The background is a dark, stylized illustration of a futuristic digital environment. In the center, a large, dark, curved screen or monitor dominates the upper half. Below it, two characters are depicted in a dynamic, action-oriented pose. On the left, a character in a red and white suit is crouching, holding a tablet. On the right, a character in a blue and red suit is also crouching, with one arm extended forward. The ground is covered in glowing, geometric patterns. In the background, there are silhouettes of futuristic buildings and a large, glowing blue rectangular shape in the upper left corner.

安全风险治理的挑战与LLM机遇

事前

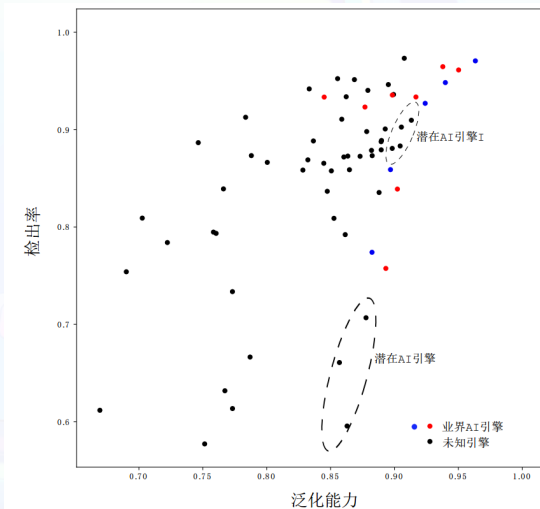
低成本

理解业务



但木桶理论让“攻守效率”不平等，使用相同的方法无法实现目标

1. 风险治理时机左移: 尽早解决风险
2. AI for 安全: “泛化能力”、“检测未知风险”



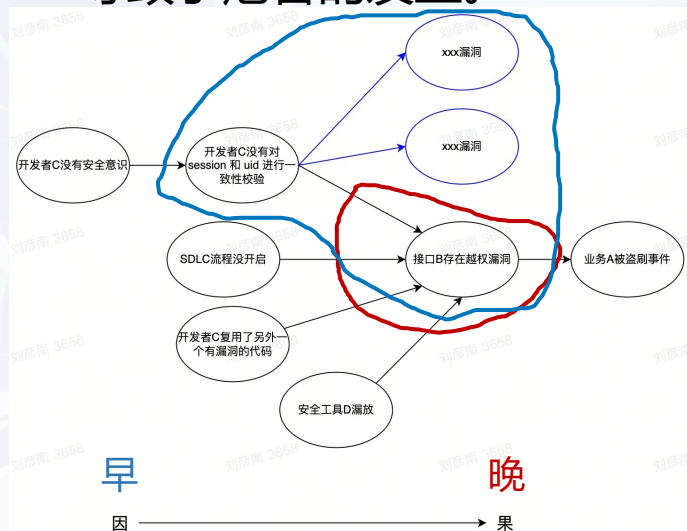
AI Detection

?

为什么能缓解木桶理论?

提升发现效率

- 安全风险其实是一个过程而非瞬时状态，暨多层脆弱性构成的因果链路最终导致了危害的发生。



```
Go 1 func GetOrder(ctx *gin.Context, req *OrderReq) (*OrderResp, error) {
2     if req.OrderId == nil {
3         logs.CtxError(ctx, "OrderId is nil")
4         return nil, common.NewError(common.ReqError, "OrderId is nil")
5     }
6
7     orderResp, err := order.GetOrder(req.OrderId)
8     if err != nil {
9         logs.CtxError(ctx, "Get Order error")
10        return nil, err
11    }
12    if orderResp == nil {
13        logs.CtxError(ctx, "Order is nil")
14        return nil, common.NewError(common.RespError, "Order is nil")
15    }
16
17    if orderResp.GetAccount().Id != GetUserId(ctx) {
18        logs.CtxError(ctx, "The current user does not have permission to obtain")
19        return nil, common.NewError(common.AuthError, "The current user does not")
20    }
21    return orderResp, nil
22 }
```

代码1: 缺少if-condition校验代码(17-20行)会导致水平越权漏洞

(为方便本文阐述构造的测试代码, 不存在于模型数据集中)

左移的核心是解决更深层成因让治理风险更高效

缺失校验 (风险代码) \neq 越权漏洞 (漏洞)

修复成本成为瓶颈

1. 风险代码实例数 >> 漏洞数

及时且自动的修复风险来降低成本

- 风险代码是漏洞的必要条件

- 难以压缩：解决时机前移、确认‘可利用性’和‘危害’难度大

2. 现有面向漏洞的修复流程成本高昂

- 漏洞让“修复”与“发现”两部分割裂

漏洞工单

→ 安全人员确认

→ 业务确认

→ 复测

→ 解决

· 业务场景定义了安全风险成因与危害，以‘校验缺失’风险代码为例

1. 从业务代码中识别校验需求
 - 业务场景类型、函数功能语义
 - 变量结构定义、上下文函数定义
 - 支撑校验的Common Sense/客观关系知识
 - 值域/类型/模式校验需求
2. 理解历史上不同校验需求对应的实现方式
 - 交易、评论、订单、UG
 - 组件、中间件、自定义代码、三方函数

理解业务逻辑和历史实现方式

理解业务逻辑和历史实现方式

从风险代码更加高效治理风险

及时且自动的修复风险来降低成本

- **LLM的代码理解和分析能力**：LLM模型本身掌握通用的代码语法、文本语义分析、甚至是简单数据流分析，同时具备多种代码理解能力。
- **LLM对校验需求推断能力**：LLM涌现的逻辑演绎推理能力，可结合历史代码从代码语义、注释语义、数据流等多角度信息合理推断校验需求。
- **LLM的代码生成能力**：多种LLM代码补全应用证明了其理解开发者意图的能力，可利用其代码生成能力实现校验代码生成补全、实现自动修复。

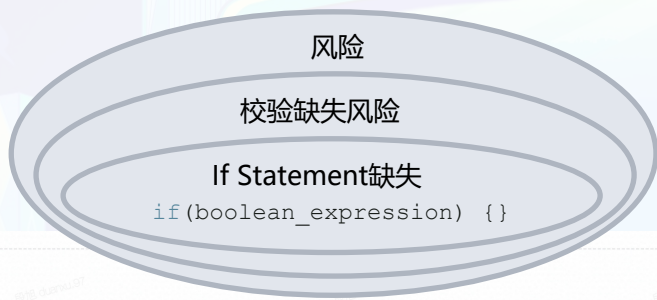
```
1 Q:<function>`func GetOrder(ctx *gin.Context, req *OrderReq) (*OrderResp, error) {
2   if req.OrderId == nil {
3     logs.CtxError(ctx, "OrderId is nil")
4     return nil, common.NewError(common.ReqError, "OrderId is nil")
5   }
6
7   orderResp, err := order.GetOrder(req.OrderId)
8   if err != nil {
9     logs.CtxError(ctx, "Get Order error")
10    return nil, err
11  }
12  if orderResp == nil {
13    logs.CtxError(ctx, "Order is nil")
14    return nil, common.NewError(common.RespError, "Order is nil")
15  }
16  return orderResp, nil
17 }`
18 <validate>
19
20 ---- 以下为Risky Code LLM补全结果----
21
22 A: if orderResp.UserId != common.GetUserID(ctx) {
23   logs.CtxError(ctx, "No permission to get order, OrderID:%v, UserID:%v", req.OrderId, c
24   return nil, common.NewError(common.RespError, "No permission to get order")
25 }
```

- 将原校验语句移除之后（不在训练集中），模型能够生成语义类似的校验语句，原代码为`orderResp.GetAccount().Id`，模型生成的为`orderResp.UserId`
- 说明模型在“语义层面”进行补全，成功识别到代码来自“获取订单信息”场景，关联‘需要对订单的属主与当前登录用户进行校验’需求，并生成校验代码

Risky Code LLM可在语义层面理解业务场景、分析校验需求并补全校验语句



基于LLM的风险代码识别与修复



- 统计, 在全部漏洞中, 约18.09%的漏洞均与If Statement校验语句缺失相关 (在漏洞修复MR中新增关键If Statement)

没有对订单属主鉴权，
导致可以越权获取任意订单

```

1 func GetOrder(ctx *gin.Context, req *OrderReq) (*OrderResp, error) {
2     if req.OrderId == nil {
3         logs.CtxError(ctx, "OrderId is nil")
4         return nil, common.NewError(common.ReqError, "OrderId is nil")
5     }
6
7     orderResp, err := order.GetOrder(req.OrderId)
8     if err != nil {
9         logs.CtxError(ctx, "Get Order failed")
10        return nil, err
11    }
12    if orderResp == nil {
13        logs.CtxError(ctx, "Order not found")
14        return nil, common.NewError(common.NotFound, "Order not found")
15    }
16
17    return orderResp, nil
18 }

```

```

1 func GetOrder(ctx *gin.Context, req *OrderReq) (*OrderResp, error) {
2     if req.OrderId == nil {
3         logs.CtxError(ctx, "OrderId is nil")
4         return nil, common.NewError(common.ReqError, "OrderId is nil")
5     }
6
7     orderResp, err := order.GetOrder(req.OrderId)
8     if err != nil {
9         logs.CtxError(ctx, "Get Order error")
10        return nil, err
11    }
12    if orderResp == nil {
13        logs.CtxError(ctx, "Order is nil")
14        return nil, common.NewError(common.RespError, "Order is nil")
15    }
16
17    if orderResp.GetAccount().Id != GetUserId(ctx) {
18        logs.CtxError(ctx, "The current user does not have permission to")
19        return nil, common.NewError(common.AuthError, "The current user")
20    }
21
22    return orderResp, nil
23 }

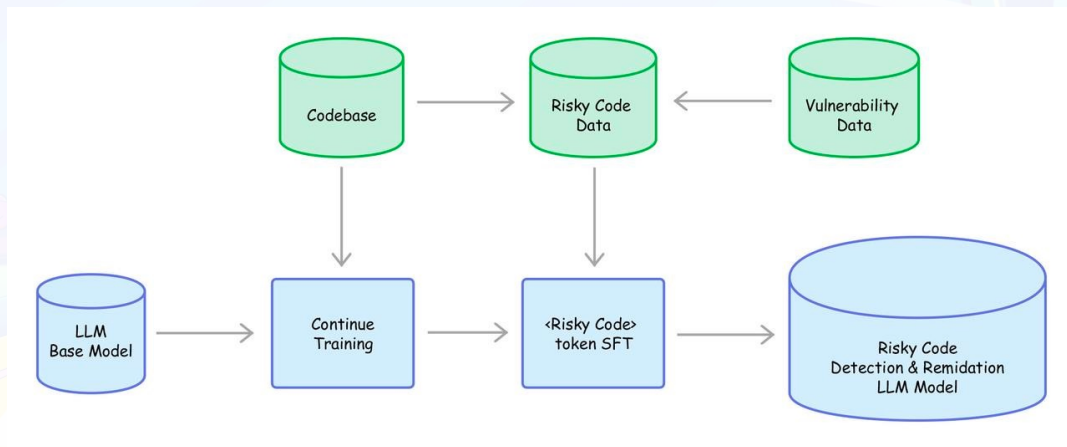
```

补全校验语句

补全鉴权校验，当请求订单不属于当前登录用户时，阻止获取订单

1. Pretraining: 使用具备基础代码能力的开源LLM模型作为Base Model。
2. Continue Training: 结合公司Codebase代码, 对Base Model进一步Continue Training, 填补开源LLM模型数据分布与公司代码分布存在的差异。
3. Supervised Fine-Tune (SFT) : 基于Codebase和历史漏洞数据, 对模型进行指令微调, 使模型了解“当前的任务/指令为补全校验语句”, 以及“在何种代码分布下需要补全校验语句”。

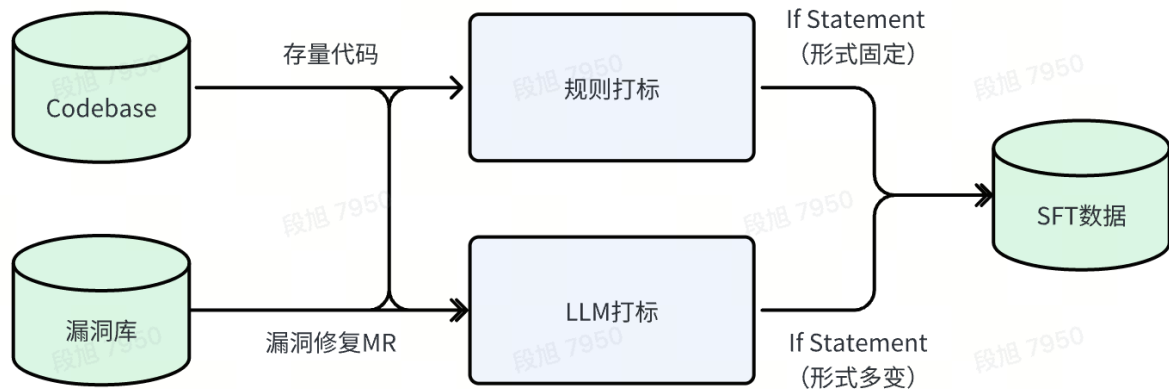
已知前提:
风险代码形成逻辑与
业务场景关系紧密



数据格式：采用了FIM (Fill in Middle) 的训练模型，通过引入Special Token，将函数方法上下文输入到模型中，模型需要输出该方法缺失的关键校验，具体格式如下：

Model Input: <bos><risk_check_context>[context]<risk_check_validation>

Model Output: <risk_check_result>[completed check]<eos> \longrightarrow 计算Loss



规则打标：AST抽取If Statement，通过简单规则过滤除关键的校验。打标出来的数据形式较为固定

LLM打标：部分关键校验形式多变，单纯规则无法识别（例如调鉴权语句再通过If判断结果，关键信息应是鉴权函数，而不是If语句），此类数据通过LLM识别。

- 在实践过程中，我们发现LLM有时会无中生有、强行编造一些结构体的成员变量，导致生成错误的校验语句。例如，如果orderResp中的属主是某个Group，同Group中的用户可以相互查看订单，那么对orderResp.UserId校验就是错误的。
- 此时可以在Prompt中补充orderResp的成员变量，从而生成正确的校验语句：

```
if orderResp.GroupId != common.GetGroupID(ctx)
```

```
1 Q:<context>`type OrderResp struct {  
2     GroupId      int  
3     OrderId      int  
4     OrderStatus  int  
5     OrderContent string  
6 }`  
7 <function>`func GetOrder(ctx *gin.Context, req *OrderReq) (*OrderResp, error) {  
8     if req.OrderId == nil {  
9         logs.CtxError(ctx, "OrderId is nil")  
10        return nil, common.NewError(common.ReqError, "OrderId is nil")  
11    }  
12  
13    orderResp, err := order.GetOrder(req.OrderId)  
14    if err != nil {  
15        logs.CtxError(ctx, "Get Order error")  
16        return nil, err  
17    }  
18    if orderResp == nil {  
19        logs.CtxError(ctx, "Order is nil")  
20        return nil, common.NewError(common.RespError, "Order is nil")  
21    }  
22    return orderResp, nil  
23 }`  
24 <validate>  
25  
26 ---- 以下为Risky Code LLM补全结果----  
27  
28 A: if orderResp.GroupId != common.GetGroupID(ctx) {  
29     logs.CtxError(ctx, "GroupId not match, req:%v, orderResp:%v", req, orderResp)  
30     return nil, common.NewError(common.RespError, "GroupId not match")  
31 }
```

在Prompt内丰富上下文信息可缓解幻觉 (hallucination) 问题

[context] 构成： 给定待扫描函数，则该函数的 [context] 包含：该函数代码 + 该函数 Callee + 该函数中对象的成员变量 + 成员方法，四部分拼接在一起

```
229 func GetInvoiceOrderDetail(ctx context.Context, r *core.GetInvoiceOrderDetailRequest) (*  
230     response := core.NewGetInvoiceOrderDetailResponse()  
231     setDefaultRespParam(response)  
232     // 1. 参数校验  
233     if err := GetInvoiceOrderDetailCheckParams(ctx, r); err != nil {  
234         return response, err  
235     }  
236     // 2. 根据userId, 关联的订单号等信息在DB中查询订单详情  
237     invoiceApp := &invoiceApp{voiceDB().ReadDB()}.G  
238     if err := invoiceApp.  
239     lo.  
240     re.  
241 }  
242 // 只有科意。  
243 if consts.I  
244 ta  
245 us  
246 ex  
247 pa  
248 if  
249  
250     ckPermission block"  
251     chCustomMessage("che  
252 }  
253 // check 账单详情的用户Id和请求Userid是否一致  
254 if r.GetUserId() != 0 && invoiceApplyOrder.UserId != r.GetUserId() {  
255     logs.CtxError(ctx, "GetInvoiceOrderDetailCheckParams userid not match db  
256     return response, werror.RequestParamsError.WithCustomMessage("userid is  
257 }
```

当前函数代码

```
154 func CheckPermissionForSwitch(ctx context.Context, uType int64, uID string, targetID int64  
155     pass, err := wcc.GetPermissionSwitch()  
156     if err !=  
157         re  
158     }  
159     return CheckPermission(ctx, uType, uID, targetID, extraStr) || !pass  
160 }  
161 func SyncQueryInvoiceResultByInvoiceOrderNo(ctx context.Context, invoiceOrder *model.Inv  
162     defer func() {  
163         if err := recover(); err != nil {  
164             v", err)  
165         }  
166     }()  
167     QueryInvoiceResultByInvoiceOrder(ctx, invoiceOrder)  
168 }  
169 }  
170 func buildGetInvoiceOrderDetailResponse(ctx context.Context, invoiceOrder *model.Invoice  
171     processingProgressItemList, invoiceOrderDetailDetailItemList := buildProcessInfo  
172     invoiceDetailItemList := make([]*core.InvoiceDetailItem, 0)  
173     for _, item :=  
174         inv  
175         t,  
176         t,  
177     }  
178 }  
179     invoiceDetailItemList = append(invoiceDetailItemList, invoiceDetails  
180 }
```

Callee函数代码

```
49 func (p *GetInvoiceOrderDetailRequest) InitDefault()  
50 func (p *GetInvoiceOrderDetailRequest) GetInvoiceOrderNo() (v string)  
51 func (p *GetInvoiceOrderDetailRequest) GetUserId() (v int64)  
52 func (p *GetInvoiceOrderDetailRequest) GetBase() (v *base.Base)  
53 func (p *GetInvoiceOrderDetailRequest) SetInvoiceOrderNo(val string)  
54 func (p *GetInvoiceOrderDetailRequest) SetUserId(val *int64)  
55 func (p *Ge  
56 func (p *Ge  
57 func (p *Ge  
58 func (p *Ge  
59 func (p *Ge  
60 func (p *Ge  
61 func (p *Ge  
62 func (p *Ge  
63 func (p *Ge  
64 func (p *Ge  
65 func (p *Ge  
66 func (p *Ge  
67 func (p *GetInvoiceOrderDetailResponse) GetInvoiceType() (v InvoiceTypeEnum)  
68 func (p *GetInvoiceOrderDetailResponse) GetInvoiceDetails() (v []*InvoiceDetailItem)  
69 func (p *GetInvoiceOrderDetailResponse) GetBusinessOrderDesc() (v string)  
70 func (p *GetInvoiceOrderDetailResponse) GetBusinessLicense() (v string)  
71 func (p *GetInvoiceOrderDetailResponse) GetBankAccountPermission() (v string)
```

成员方法

```
1 type GetInvoiceOrderDetailRequest struct {  
2     InvoiceOrderNo string  
3  
4  
5 }  
6 type GetInvoiceOrderDetailResponse struct {  
7     InvoiceTypeDesc string  
8     InvoiceAmount int64  
9     InvoiceAmtTime int64  
10    In  
11    In  
12    Pr  
13    It  
14    In  
15    In  
16    Bu  
17    Bu  
18    Ba  
19    GeneralTaxpayerCertificate *string  
20    BaseResp *base.BaseResp  
21 }  
22 type InvoiceBusinessOrderDAO struct {  
23     db *gorm.DB  
24 }
```

成员变量

The background is a dark, stylized illustration of a futuristic city at night. A large, glowing screen dominates the center, displaying the title. Two characters are positioned in the foreground: on the left, a character in a white and red suit is looking at a tablet; on the right, a character in a blue and red suit is pointing towards the screen. The scene is filled with various architectural elements and a sense of high-tech activity.

实验结果与分析

细粒度评价指标：

- BLEU：机器翻译领域的常见指标，衡量参考文本和生成文本在1-gram至n-gram的相似性；
- 加权BLEU：在BLEU基础上提高编程语言关键字的权重；
- AST相似性：衡量模型生成代码与原代码的语法相似性，即参考代码和生成代码中相同AST子树的比例；
- DFG相似性：衡量模型生成代码与原代码的语义相似性，即参考代码和生成代码中相同DFG边的比例；
- 校验变量的Jaccard相似度：IF语句布尔表达式中被校验变量的Jaccard相似度；

Codebleu：加权平均

参数量、Continue Training以及在Callee、成员变量、成员方法均能够带来有效提升

模型	参数量	codebleu	bleu	加权bleu	AST相似性	校验变量Jaccard相似度
Customized	1.7b	0.49	0.45	0.45	0.58	0.40
Codegen	6b	0.51	0.48	0.50	0.56	0.53
Starcoder	15b	0.58	0.55	0.57	0.65	0.55

模型	训练方式	codebleu	bleu	加权bleu	AST相似性	校验变量Jaccard相似度
Codegen 6b	SFT	0.49	0.47	0.49	0.54	0.50
Codegen 6b	Continue Training - SFT	0.51	0.48	0.50	0.56	0.53

模型	Context构成	codebleu	bleu	加权bleu	语法匹配程度	校验变量Jaccard相似度
Codegen 6b	No Additional Context	0.58	0.55	0.56	0.66	0.63
Codegen 6b	with Additional Context	0.60	0.57	0.59	0.68	0.65

- 上传漏洞、SSRF等“校验形式单一”的传统web洞，及涉及到鉴权的逻辑漏洞具有较高的准确性

漏洞类型	bleu	codebleu	加权N-gram	语法匹配程度	校验变量Jaccard相似度
服务端请求伪造(SSRF)	0.78	0.85	0.79	0.92	0.77
上传漏洞	0.78	0.85	0.78	0.93	0.79
跨站请求伪造(CSRF)	0.75	0.83	0.75	0.93	0.73
潜在信息泄露	0.78	0.83	0.78	0.91	0.75
未授权访问	0.77	0.85	0.77	0.93	0.76
逻辑漏洞	0.69	0.79	0.70	0.84	0.67
权限绕过	0.75	0.82	0.76	0.90	0.76
敏感字段泄露	0.71	0.79	0.72	0.85	0.70
SQL注入	0.68	0.76	0.68	0.81	0.71
跨站脚本攻击(XSS)	0.34	0.49	0.37	0.49	0.31
非法图床漏洞	0.37	0.52	0.39	0.54	0.33

测试数据来源：

历史漏洞数据： 筛选历史漏洞修复MR中包含新增校验语句的漏洞。将Merge之前的代码作为[context]，将MR中新增的校验语句作为[completed check]；



Q&A

THANKS

