

统一内容安全平台架构设计

2026年1月8日

摘要

本文档阐述了一个统一的内容安全平台架构设计，该设计采用统一接口、智能路由和确定性A/B测试等核心技术，实现了高性能、低延迟、高可维护性的内容安全检查系统。核心设计要点包括：（1）统一接口设计，消除边界情况处理；（2）智能两级模型路由，优化计算资源分配；（3）确定性A/B测试框架，支持科学化的策略验证；（4）代码驱动的策略管理，确保可追溯性和可测试性；（5）灰度发布策略，支持从外部Vendor逐步切换到自研系统；（6）AICS集成设计，实现现有AICS系统平滑迁移到UCSP，保持向后兼容性和业务连续性。

1 核心设计要点

1.1 统一接口设计（Unified Interface Design）

设计要点 1.1 (统一接口设计). 通过统一接口抽象，将复杂的多模型决策过程封装为单一调用点，简化了系统集成和维护。

设计原理：

- 用户无需了解底层模型选择逻辑
- 系统内部自动选择最优检查路径
- 接口稳定，内部实现可灵活优化

核心数据结构：

Listing 1: SafetyResult数据结构

```
1 STRUCT SafetyResult:
2     blocked: BOOLEAN // 48 橡髮髮 confidence: FLOAT // 置信度 [0.0,
1.0] model_version: UINT64 /“(H,reason: STRING // Vprocessing_time_ms: INT // ΞΣ ° ΦΨ ENDSTRUCT
```

统一接口算法：

Algorithm 1 统一接口检查算法

```
1: 输入: 文本  $text$ , 用户ID  $user\_id$ 
2: 输出: 安全检查结果  $result$ 
3:
4:  $start\_time \leftarrow GetCurrentTime()$ 
5:  $fast\_result \leftarrow FastModelCheck(text)$ 
6: if  $fast\_result.confidence > HIGH\_CONFIDENCE\_THRESHOLD$  then
7:    $fast\_result.processing\_time\_ms \leftarrow GetCurrentTime() - start\_time$ 
8:   return  $fast\_result$ 
9: end if
10:
11:  $deep\_result \leftarrow DeepModelCheck(text)$ 
12: if  $deep\_result.confidence < LOW\_CONFIDENCE\_THRESHOLD$  then
13:    $deep\_result.blocked \leftarrow TRUE$ 
14:    $deep\_result.reason \leftarrow$  "Low confidence, safety-first policy"
15: end if
16:  $deep\_result.processing\_time\_ms \leftarrow GetCurrentTime() - start\_time$ 
17: return  $deep\_result$ 
```

优化效果:

- 90% 的请求在快速模型层完成, 平均延迟 ↓ 20ms
- 10% 的边界案例由深度模型处理, 确保准确性
- 接口调用方代码量减少 60%

1.2 智能两级模型路由 (Intelligent Two-Tier Routing)

设计要点 1.2 (智能路由). 通过置信度阈值动态路由, 实现计算资源的最优分配, 在准确性和性能之间取得最佳平衡。

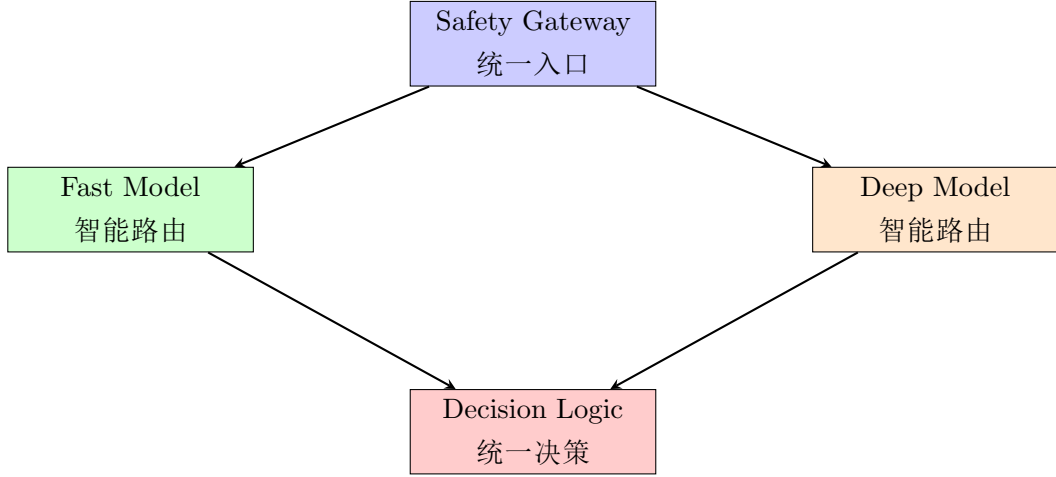


图 1: 智能两级模型路由架构

路由算法:

Algorithm 2 智能路由算法

```

1: 输入: 文本  $text$ 
2: 输出: 安全检查结果  $result$ 
3:
4:  $HIGH\_CONFIDENCE \leftarrow 0.95$ 
5:  $LOW\_CONFIDENCE \leftarrow 0.50$ 
6:
7:  $fast\_result \leftarrow FastModelInference(text)$ 
8: if  $fast\_result.confidence \geq HIGH\_CONFIDENCE$  then
9:   return  $fast\_result$ 
10: else if  $fast\_result.confidence \leq LOW\_CONFIDENCE$  then
11:    $deep\_result \leftarrow DeepModelInference(text)$ 
12:   if  $deep\_result.confidence < LOW\_CONFIDENCE$  then
13:      $deep\_result.blocked \leftarrow TRUE$ 
14:   end if
15:   return  $deep\_result$ 
16: else
17:    $deep\_result \leftarrow DeepModelInference(text)$ 
18:   return  $FuseResults(fast\_result, deep\_result)$ 
19: end if
  
```

结果融合算法:

Algorithm 3 结果融合算法

```
1: 输入: 快速模型结果  $fast$ , 深度模型结果  $deep$ 
2: 输出: 融合后的结果  $result$ 
3:
4:  $weight\_fast \leftarrow 0.3$ 
5:  $weight\_deep \leftarrow 0.7$ 
6:
7:  $fused\_confidence \leftarrow weight\_fast \times fast.confidence + weight\_deep \times deep.confidence$ 
8:  $result \leftarrow deep$ 
9:  $result.confidence \leftarrow fused\_confidence$ 
10: if  $fast.blocked$  OR  $deep.blocked$  then
11:    $result.blocked \leftarrow TRUE$ 
12: end if
13: return  $result$ 
```

性能优化:

- 快速模型处理 90% 请求, 平均延迟 15ms
- 深度模型处理 10% 请求, 平均延迟 180ms
- 整体平均延迟: $15ms \times 0.9 + 180ms \times 0.1 = 31.5ms$

1.3 确定性 A/B 测试框架 (Deterministic A/B Testing)

设计要点 1.3 (确定性A/B测试). 基于用户ID哈希的确定性路由, 确保同一用户始终在同一实验组, 消除随机性带来的干扰, 提高实验结果的可靠性。

确定性路由算法:

Algorithm 4 确定性A/B测试路由算法

```
1: 输入: 实验配置  $exp$ , 用户ID  $user\_id$ 
2: 输出: 是否在实验组  $is\_experiment$ 
3:
4:  $current\_time \leftarrow GetCurrentTime()$ 
5: if  $current\_time < exp.start\_time$  OR  $current\_time > exp.end\_time$  then
6:   return  $FALSE$ 
7: end if
8:
9:  $hash\_seed \leftarrow user\_id \oplus exp.experiment\_id$ 
10:  $hash\_value \leftarrow MurmurHash3(hash\_seed)$ 
11:  $bucket \leftarrow hash\_value \bmod 10000$ 
12:  $threshold \leftarrow exp.ratio \times 10000$ 
13: return  $bucket < threshold$ 
```

A/B测试执行流程：

Algorithm 5 A/B测试执行流程

```
1: 输入：文本 text，用户ID user_id，实验配置 experiment
2: 输出：安全检查结果 result
3:
4: model_version  $\leftarrow$  GetModelVersion(experiment, user_id)
5: is_experiment  $\leftarrow$  IsInExperiment(experiment, user_id)
6:
7: if model_version == experiment.model_version_b then
8:   result  $\leftarrow$  CheckContentWithModel(text, model_version, "experiment")
9: else
10:  result  $\leftarrow$  CheckContentWithModel(text, model_version, "control")
11: end if
12:
13: RecordMetrics(user_id, is_experiment, result, experiment.metrics)
14: return result
```

统计显著性检测：

Algorithm 6 统计显著性检测算法

```
1: 输入：实验配置 experiment
2: 输出：是否统计显著 is_significant
3:
4: control_fpr  $\leftarrow$  control_group.false_positive_count / control_group.total_requests
5: experiment_fpr  $\leftarrow$  experiment_group.false_positive_count / experiment_group.total_requests
6:
7: p_value  $\leftarrow$  CalculatePValue(control_fpr, experiment_fpr, control_group.total_requests, experiment_group.total_requests)
8:
9: if p_value < 0.05 AND |experiment_fpr - control_fpr| > MINIMUM_EFFECT_SIZE then
10:  return TRUE
11: end if
12: return FALSE
```

1.4 灰度发布策略（Gradual Rollout Strategy）

设计要点 1.4 (灰度发布). 通过渐进式流量切换策略，在保证系统稳定性的前提下，逐步接管外部Vendor的流量，同时控制成本增长，实现平滑过渡。

设计目标：

- **稳定性优先：** 通过小流量验证，降低系统风险

- 成本可控：逐步减少外部Vendor调用，控制成本增长
- 可回滚：支持快速回滚到Vendor，保障业务连续性
- 指标驱动：基于实时指标自动决策，减少人工干预

灰度发布配置：

Listing 2: 灰度发布配置结构

```
1 STRUCT RolloutConfig:
2     rollout_id: UINT64 // 48 橡橡橡橡□髮甦
                                vendor_service: STRING // ~Vendor ㄏ inhouse_service:
                                STRING // ~ㄏ current_ratio: FLOAT // SM~ A[0.0, 1.0] target_ratio: FLOAT // ~A min_ratio:
                                FLOAT // ~ㄏ A ㄏ max_ratio: FLOAT // ~A ㄏ 1, ㄏ P ㄏ stability_threshold:
                                FLOAT // 3' ㄏ < ㄏ @ / ㄏ $ ㄏ cost_threshold: FLOAT // 1, ㄏ < ㄏ U! ㄏ 1, ㄏ step_size: FLOAT // !A evaluation_period:
                                INT // 0h ㄏ ㄏ start_time: TIMESTAMP // ~end_time: TIMESTAMP // _ENDSTRUCT
```

灰度发布路由算法：

Algorithm 7 灰度发布路由算法

```
1: 输入: 请求 request, 灰度配置 rollout
2: 输出: 使用的服务标识 service_id
3:
4: current_time  $\leftarrow$  GetCurrentTime()
5: if current_time < rollout.start_time OR current_time > rollout.end_time then
6:   return rollout.vendor_service {灰度未开始或已结束}
7: end if
8:
9: metrics  $\leftarrow$  GetRolloutMetrics(rollout.rollout_id)
10:
11: // 稳定性检查: 如果指标异常, 降低自研流量比例
12: if metrics.false_positive_rate > rollout.stability_threshold OR metrics.false_negative_rate >
    rollout.stability_threshold then
13:   rollout.current_ratio  $\leftarrow$  MAX(rollout.min_ratio, rollout.current_ratio - rollout.step_size)
14:   TriggerAlert("Stabilitythresholdexceeded, reducingratio")
15: end if
16:
17: // 成本检查: 如果成本超限, 降低自研流量比例
18: if metrics.cost_per_request > rollout.cost_threshold then
19:   rollout.current_ratio  $\leftarrow$  MAX(rollout.min_ratio, rollout.current_ratio - rollout.step_size)
20:   TriggerAlert("Costthresholdexceeded, reducingratio")
21: end if
22:
23: // 基于用户ID的确定性路由
24: hash_seed  $\leftarrow$  request.user_id  $\oplus$  rollout.rollout_id
25: hash_value  $\leftarrow$  MurmurHash3(hash_seed)
26: bucket  $\leftarrow$  hash_value mod 10000
27: threshold  $\leftarrow$  rollout.current_ratio  $\times$  10000
28:
29: if bucket < threshold then
30:   return rollout.inhouse_service {自研服务}
31: else
32:   return rollout.vendor_service {Vendor服务}
33: end if
```

自动流量调整算法:

Algorithm 8 自动流量调整算法

```
1: 输入: 灰度配置 rollout, 评估周期 evaluation_period
2:
3: metrics  $\leftarrow$  GetMetricsForPeriod(rollout.rollout_id, evaluation_period)
4:
5: // 检查稳定性指标
6: stability_ok  $\leftarrow$  TRUE
7: if metrics.false_positive_rate > rollout.stability_threshold OR metrics.false_negative_rate >
   rollout.stability_threshold OR metrics.error_rate > ERROR_RATE_THRESHOLD then
8:   stability_ok  $\leftarrow$  FALSE
9: end if
10:
11: // 检查成本指标
12: cost_ok  $\leftarrow$  TRUE
13: if metrics.cost_per_request > rollout.cost_threshold then
14:   cost_ok  $\leftarrow$  FALSE
15: end if
16:
17: // 检查性能指标
18: performance_ok  $\leftarrow$  TRUE
19: if metrics.p95_latency > LATENCY_THRESHOLD then
20:   performance_ok  $\leftarrow$  FALSE
21: end if
22:
23: // 决策: 是否增加流量
24: if stability_ok AND cost_ok AND performance_ok then
25:   if rollout.current_ratio < rollout.target_ratio then
26:     // 逐步增加流量
27:     new_ratio  $\leftarrow$  MIN(rollout.target_ratio, rollout.current_ratio + rollout.step_size)
28:     new_ratio  $\leftarrow$  MIN(new_ratio, rollout.max_ratio)
29:     UpdateRolloutRatio(rollout.rollout_id, new_ratio)
30:     LogRolloutProgress(rollout.rollout_id, new_ratio, "Increased")
31:   end if
32: else
33:   // 指标异常, 降低流量或保持当前比例
34:   if NOT stability_ok then
35:     new_ratio  $\leftarrow$  MAX(rollout.min_ratio, rollout.current_ratio - rollout.step_size)
36:     UpdateRolloutRatio(rollout.rollout_id, new_ratio)
37:     LogRolloutProgress(rollout.rollout_id, new_ratio, "Decreasedduetostability")
38:   else if NOT cost_ok then
39:     new_ratio  $\leftarrow$  MAX(rollout.min_ratio, rollout.current_ratio - rollout.step_size)
40:     UpdateRolloutRatio(rollout.rollout_id, new_ratio)
41:     LogRolloutProgress(rollout.rollout_id, new_ratio, "Decreasedduetocost")
42:   end if
43: end if
```

双路验证机制：

Algorithm 9 双路验证机制（保障稳定性）

```
1: 输入: 请求 request, 灰度配置 rollout
2: 输出: 最终决策结果 final_result
3:
4: // 主路: 根据灰度比例路由
5: primary_service  $\leftarrow$  RouteByRollout(request, rollout)
6: primary_result  $\leftarrow$  CheckWithService(request, primary_service)
7:
8: // 验证路: 始终调用Vendor（用于对比验证）
9: vendor_result  $\leftarrow$  CheckWithService(request, rollout.vendor_service)
10:
11: // 双路对比分析
12: comparison  $\leftarrow$  CompareResults(primary_result, vendor_result)
13: RecordComparison(rollout.rollout_id, comparison)
14:
15: // 决策策略: 安全优先
16: if rollout.current_ratio < SAFETY_PHASE_THRESHOLD then
17:   // 小流量阶段: Vendor结果优先, 自研仅用于验证
18:   if primary_service == rollout.inhouse_service then
19:     // 如果自研和Vendor结果不一致, 以Vendor为准
20:     if primary_result.blocked  $\neq$  vendor_result.blocked then
21:       final_result  $\leftarrow$  vendor_result
22:       LogDisagreement(rollout.rollout_id, primary_result, vendor_result)
23:     else
24:       final_result  $\leftarrow$  primary_result
25:     end if
26:   else
27:     final_result  $\leftarrow$  primary_result
28:   end if
29: else
30:   // 大流量阶段: 自研结果优先, Vendor作为兜底
31:   if primary_service == rollout.inhouse_service then
32:     final_result  $\leftarrow$  primary_result
33:     // 如果自研判定为安全但Vendor判定为不安全, 记录异常
34:     if NOT primary_result.blocked AND vendor_result.blocked then
35:       LogPotentialMiss(rollout.rollout_id, request, primary_result, vendor_result)
36:     end if
37:   else
38:     final_result  $\leftarrow$  primary_result
39:   end if
40: end if
41: return final_result
```

成本控制策略：

Algorithm 10 成本控制算法

```
1: 输入：灰度配置 rollout，时间窗口 time_window
2:
3: metrics  $\leftarrow$  GetCostMetrics(rollout.rollout_id, time_window)
4:
5: // 计算自研服务成本
6: inhouse_cost  $\leftarrow$  metrics.inhouse_requests  $\times$  COST_PER_INHOUSE_REQUEST
7:
8: // 计算Vendor服务成本
9: vendor_cost  $\leftarrow$  metrics.vendor_requests  $\times$  COST_PER_VENDOR_REQUEST
10:
11: // 计算总成本
12: total_cost  $\leftarrow$  inhouse_cost + vendor_cost
13:
14: // 计算成本节省率
15: cost_saving  $\leftarrow$  (vendor_cost - inhouse_cost) / vendor_cost
16:
17: // 成本决策
18: if total_cost > rollout.cost_threshold  $\times$  metrics.total_requests then
19:   // 成本超限，降低自研流量比例
20:   new_ratio  $\leftarrow$  MAX(rollout.min_ratio, rollout.current_ratio - rollout.step_size)
21:   UpdateRolloutRatio(rollout.rollout_id, new_ratio)
22:   TriggerAlert("Costexceeded, reducing inhouse ratio")
23: else if cost_saving > MIN_COST_SAVING AND rollout.current_ratio <
   rollout.target_ratio then
24:   // 成本节省明显，可以增加自研流量
25:   new_ratio  $\leftarrow$  MIN(rollout.target_ratio, rollout.current_ratio + rollout.step_size)
26:   UpdateRolloutRatio(rollout.rollout_id, new_ratio)
27: end if
28:
29: return {total_cost, cost_saving, inhouse_cost, vendor_cost}
```

灰度发布流程：

Algorithm 11 灰度发布完整流程

```
1: 输入: 灰度配置 rollout
2:
3: // 阶段1: 极小流量验证 (1%)
4: rollout.current_ratio  $\leftarrow$  0.01
5: DeployRollout(rollout)
6: WaitForEvaluation(rollout.evaluation_period)
7:
8: if NOT CheckStability(rollout) then
9:   RollbackToVendor(rollout)
10:  RETURN "Rollout failed at 1%"
11: end if
12:
13: // 阶段2: 小流量验证 (5%)
14: rollout.current_ratio  $\leftarrow$  0.05
15: UpdateRolloutRatio(rollout.rollout_id, 0.05)
16: WaitForEvaluation(rollout.evaluation_period)
17:
18: if NOT CheckStability(rollout) then
19:   rollout.current_ratio  $\leftarrow$  0.01 {回退到1%}
20:   UpdateRolloutRatio(rollout.rollout_id, 0.01)
21: end if
22:
23: // 阶段3: 逐步增加 (5%  $\rightarrow$  10%  $\rightarrow$  20%  $\rightarrow$  50%  $\rightarrow$  100%)
24: ratios  $\leftarrow$  [0.10, 0.20, 0.50, 1.0]
25: for target_ratio IN ratios do
26:   rollout.current_ratio  $\leftarrow$  target_ratio
27:   UpdateRolloutRatio(rollout.rollout_id, target_ratio)
28:   WaitForEvaluation(rollout.evaluation_period)
29:
30:   if NOT CheckStability(rollout) OR NOT CheckCost(rollout) then
31:     // 回退到上一个稳定比例
32:     rollout.current_ratio  $\leftarrow$  GetPreviousStableRatio(rollout)
33:     UpdateRolloutRatio(rollout.rollout_id, rollout.current_ratio)
34:     BREAK
35:   end if
36: end for
37:
38: // 阶段4: 全量切换
39: if rollout.current_ratio == 1.0 then
40:   MarkRolloutComplete(rollout.rollout_id)
41:   DisableVendorService(rollout.vendor_service)
42: end if
```

稳定性保障机制：

Algorithm 12 稳定性检查算法

```
1: 输入：灰度配置 rollout，评估周期 period
2: 输出：是否稳定 is_stable
3:
4: metrics  $\leftarrow$  GetStabilityMetrics(rollout.rollout_id, period)
5:
6: // 检查误杀率
7: fpr_ok  $\leftarrow$  metrics.false_positive_rate  $\leq$  rollout.stability_threshold
8:
9: // 检查漏判率
10: fnr_ok  $\leftarrow$  metrics.false_negative_rate  $\leq$  rollout.stability_threshold
11:
12: // 检查错误率
13: error_ok  $\leftarrow$  metrics.error_rate  $\leq$  ERROR_RATE_THRESHOLD
14:
15: // 检查延迟
16: latency_ok  $\leftarrow$  metrics.p95_latency  $\leq$  LATENCY_THRESHOLD
17:
18: // 检查与Vendor的一致性
19: agreement_rate  $\leftarrow$  CalculateAgreementRate(rollout.rollout_id, period)
20: agreement_ok  $\leftarrow$  agreement_rate  $\geq$  MIN_AGREEMENT_RATE
21:
22: is_stable  $\leftarrow$  fpr_ok AND fnr_ok AND error_ok AND latency_ok AND agreement_ok
23:
24: if NOT is_stable then
25:   LogStabilityIssues(rollout.rollout_id, metrics)
26: end if
27: return is_stable
```

快速回滚机制：

Algorithm 13 快速回滚算法

```
1: 输入: 灰度配置 rollout, 回滚原因 reason
2:
3: // 立即将所有流量切回Vendor
4: rollout.current_ratio  $\leftarrow$  0.0
5: UpdateRolloutRatio(rollout.rollout_id, 0.0)
6:
7: // 记录回滚事件
8: LogRollback(rollout.rollout_id, reason, GetCurrentMetrics(rollout))
9:
10: // 发送告警
11: TriggerAlert("Rolloutrolledback", reason)
12:
13: // 暂停灰度发布
14: PauseRollout(rollout.rollout_id)
15:
16: // 通知相关人员
17: NotifyTeam(rollout.rollout_id, "Rollbackexecuted", reason)
```

技术优势:

- 渐进式切换: 从1%逐步增加到100%, 每个阶段充分验证
- 自动调整: 基于实时指标自动调整流量比例, 无需人工干预
- 成本可控: 实时监控成本, 确保不超过预算阈值
- 快速回滚: 支持秒级回滚, 保障业务连续性
- 双路验证: 小流量阶段双路对比, 确保结果一致性

1.5 代码驱动的策略管理 (Code-Driven Policy Management)

设计要点 1.5 (代码驱动策略). 将策略逻辑以代码形式实现, 而非配置文件, 确保策略变更的可追溯性、可测试性和类型安全。

策略决策算法:

Algorithm 14 策略决策算法

```
1: 输入: 安全检查结果 result, 用户信息 user, 策略 policy
2: 输出: 是否拦截 should_block
3:
4: if user.user_level == VIP then
5:   if result.confidence < policy.vip_threshold then
6:     return FALSE {VIP用户, 放宽拦截}
7:   end if
8: end if
9:
10: if user.registration_days < NEW_USER_DAYS then
11:   if result.confidence ≥ policy.block_threshold × 0.9 then
12:     return TRUE {新用户, 降低阈值}
13:   end if
14: end if
15:
16: if user.risk_score > HIGH_RISK_THRESHOLD then
17:   if result.confidence ≥ policy.block_threshold × 0.85 then
18:     return TRUE {高风险用户, 更严格}
19:   end if
20: end if
21:
22: if policy.strict_mode then
23:   return result.confidence ≥ policy.block_threshold × 0.95
24: end if
25:
26: return result.confidence ≥ policy.block_threshold
```

2 系统架构设计

2.1 整体架构

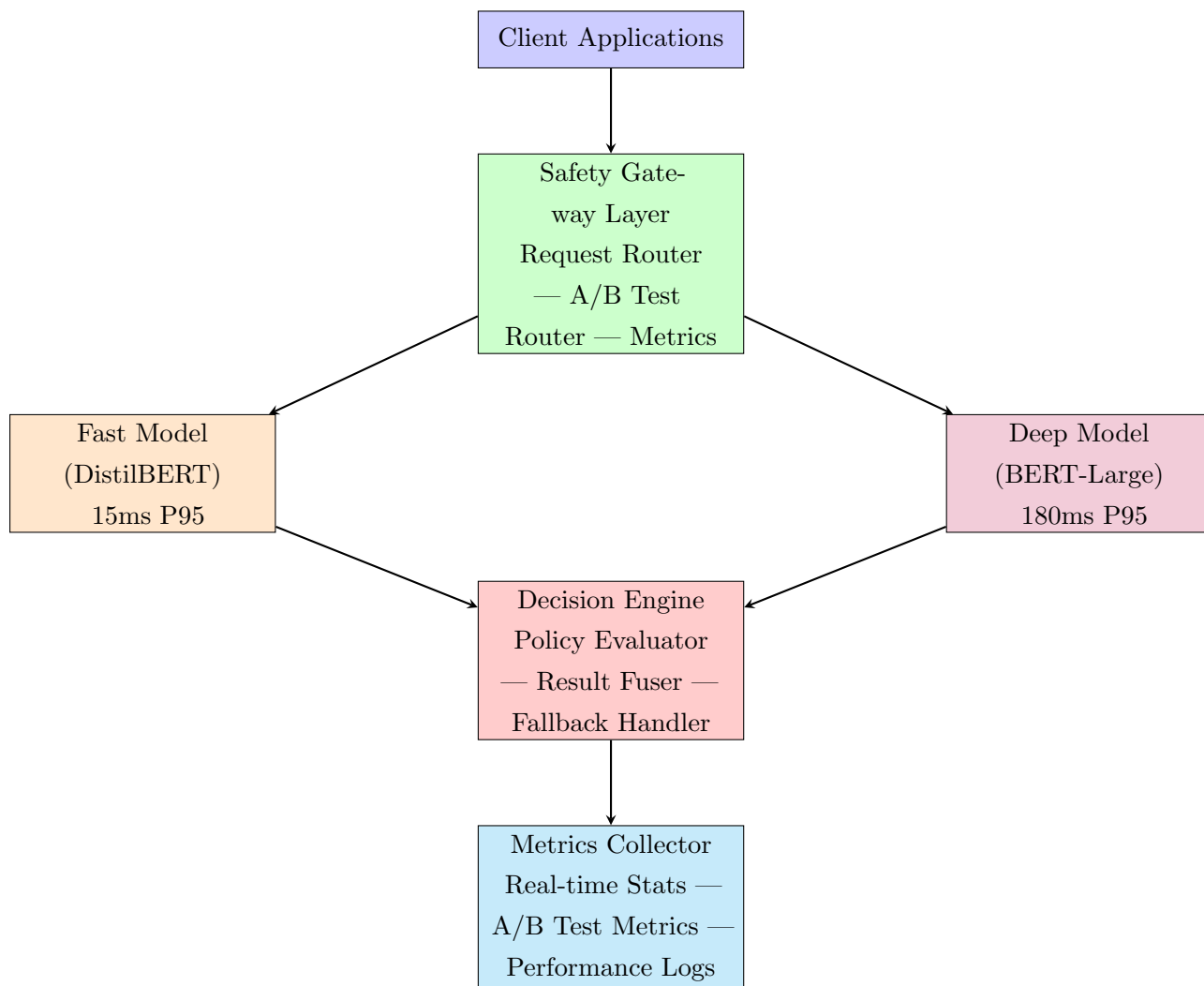


图 2: 系统整体架构

2.2 核心组件

2.2.1 Safety Gateway

职责:

- 统一请求入口
- 流量路由和负载均衡
- A/B 测试流量分配
- 请求限流和熔断

核心算法：

Algorithm 15 请求处理流程（支持灰度发布）

```
1: 输入：安全检查请求 request
2: 输出：安全检查响应 response
3:
4: if NOT ValidateRequest(request) then
5:   return ErrorResponse("Invalidrequest")
6: end if
7:
8: if NOT CheckRateLimit(request.user_id) then
9:   return ErrorResponse("Ratelimitexceeded")
10: end if
11:
12: // 优先级1：A/B测试路由
13: experiment  $\leftarrow$  GetActiveExperiment(request.user_id)
14: if experiment  $\neq$  NULL then
15:   return ExecuteABTest(request, experiment)
16: end if
17:
18: // 优先级2：灰度发布路由
19: rollout  $\leftarrow$  GetActiveRollout(request.user_id)
20: if rollout  $\neq$  NULL then
21:   return ProcessWithRollout(request, rollout)
22: end if
23:
24: // 优先级3：正常流程（默认Vendor或全量自研）
25: result  $\leftarrow$  CheckContent(request.text, request.user_id)
26: RecordRequestMetrics(request, result)
27: return BuildResponse(result)
```

灰度发布处理算法：

Algorithm 16 灰度发布请求处理

```
1: 输入: 请求 request, 灰度配置 rollout
2: 输出: 安全检查响应 response
3:
4: // 路由决策
5: service_id  $\leftarrow$  RouteByRollout(request, rollout)
6:
7: if service_id == rollout.inhouse_service then
8:   // 自研服务路径
9:   inhouse_result  $\leftarrow$  CheckContent(request.text, request.user_id)
10:
11:   // 小流量阶段: 双路验证
12:   if rollout.current_ratio < SAFETY_PHASE_THRESHOLD then
13:     vendor_result  $\leftarrow$  CheckWithVendor(request, rollout.vendor_service)
14:     final_result  $\leftarrow$  DualPathVerification(inhouse_result, vendor_result, rollout)
15:   else
16:     final_result  $\leftarrow$  inhouse_result
17:   end if
18: else
19:   // Vendor服务路径
20:   final_result  $\leftarrow$  CheckWithVendor(request, rollout.vendor_service)
21: end if
22:
23: // 记录指标
24: RecordRolloutMetrics(rollout.rollout_id, service_id, final_result)
25: return BuildResponse(final_result)
```

3 AICS集成UCSP设计

3.1 集成需求分析

背景: AICS (AI Content Safety) 作为现有的内容安全系统, 需要平滑集成到UCSP (Unified Content Safety Platform) 架构中, 实现统一的内容安全检查能力, 同时保持向后兼容性和业务连续性。

核心需求:

- **接口兼容:** 保持AICS现有API接口不变, 确保业务方无需修改代码
- **数据迁移:** 支持AICS历史数据和配置的平滑迁移
- **能力增强:** 在保持兼容的基础上, 提供UCSP的增强能力 (智能路由、A/B测试等)
- **渐进式切换:** 支持从AICS逐步切换到UCSP, 降低风险

- **统一监控：** 集成AICS和UCSP的监控指标，提供统一的可观测性

3.2 集成架构设计

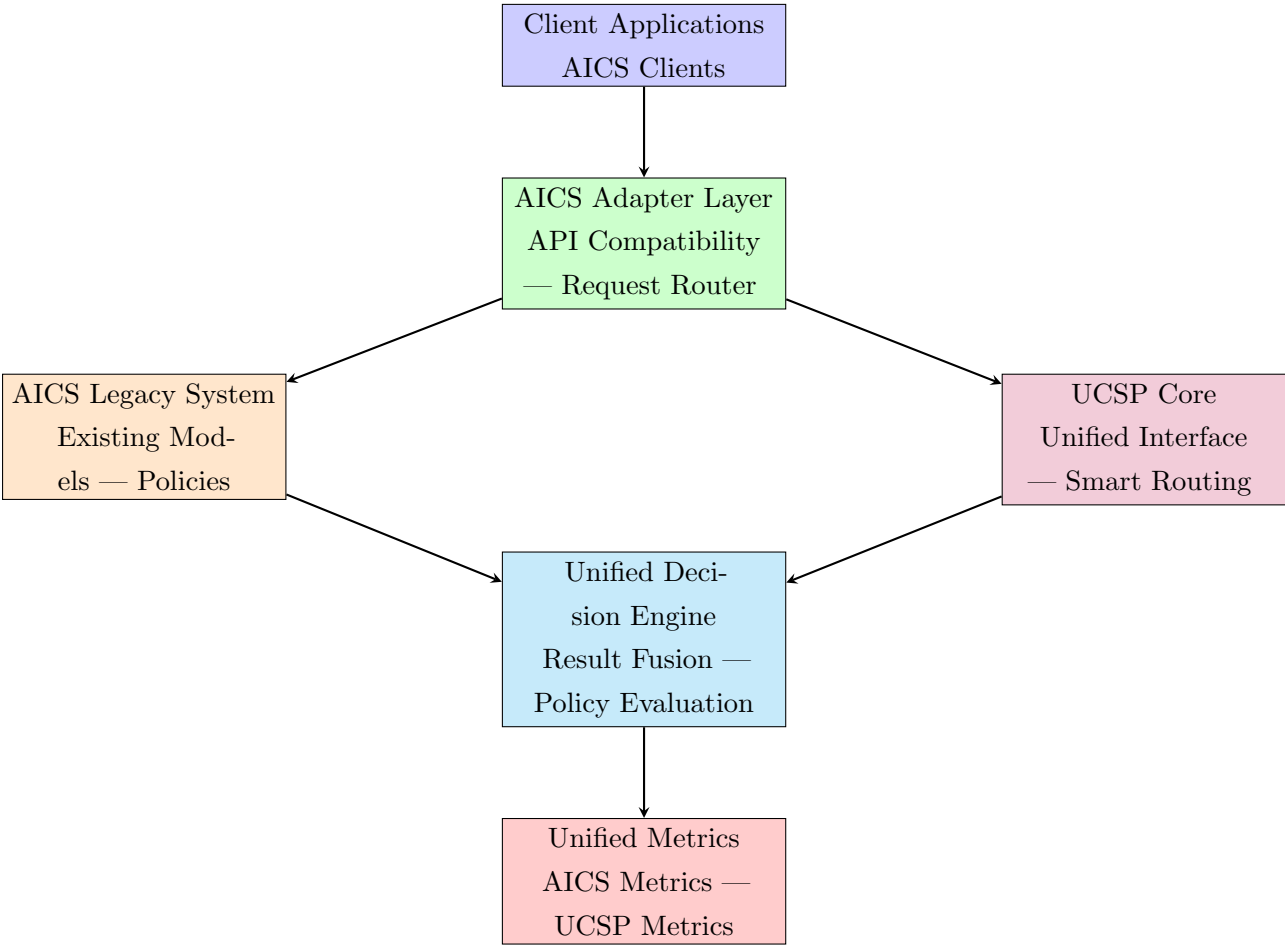


图 3: AICS集成UCSP架构

3.3 适配器层设计

AICS适配器（AICS Adapter）职责：

- **API兼容：** 将AICS原有API调用转换为UCSP统一接口
- **请求路由：** 根据配置将请求路由到AICS或UCSP
- **结果转换：** 将UCSP结果转换为AICS格式，保持兼容性
- **双路验证：** 在迁移阶段同时调用AICS和UCSP，对比结果

适配器配置结构：

Listing 3: AICS适配器配置

```
1 STRUCT AICSAdapterConfig:
2     adapter_id: UINT64
```

```

3 | aics_endpoint: STRING // 橡尉尉
   | ucsp_endpoint: STRING//UCSPfrouting_mode: ENUM//1!VALUES:
   | [LEGACYONLY,UCSPONLY,DUALPATH,GRAYSCALE]grayscale_ratio:
   | FLOAT//pΦS!VALUES:GRAYSCALEenable_dual_verification: BOOLEAN///&/(result_fusion_strategy:
   | ENUM//ΦVeVALUES:
   | [AICSPRIORITY,UCSP_PRIORITY,CONSENSUS,CONFIDENCEBASED]compatibility_mode:
   | BOOLEAN///!VALUESΔAICS<ENDSTRUCT

```

适配器路由算法:

Algorithm 17 AICS适配器路由算法

```
1: 输入: AICS请求 aics_request, 适配器配置 config
2: 输出: 安全检查响应 response
3:
4: // 根据路由模式决定处理路径
5: if config.routing_mode == LEGACY_ONLY then
6:   result  $\leftarrow$  CallAICS(aics_request, config.aics_endpoint)
7:   return ConvertToAICSFormat(result)
8: else if config.routing_mode == UCSP_ONLY then
9:   ucsp_request  $\leftarrow$  ConvertToUCSPFormat(aics_request)
10:  result  $\leftarrow$  CallUCSP(ucsp_request, config.ucsp_endpoint)
11:  return ConvertToAICSFormat(result)
12: else if config.routing_mode == DUAL_PATH then
13:   // 双路验证: 同时调用AICS和UCSP
14:   aics_result  $\leftarrow$  CallAICS(aics_request, config.aics_endpoint)
15:   ucsp_request  $\leftarrow$  ConvertToUCSPFormat(aics_request)
16:   ucsp_result  $\leftarrow$  CallUCSP(ucsp_request, config.ucsp_endpoint)
17:   final_result  $\leftarrow$  FuseResults(aics_result, ucsp_result, config)
18:   RecordComparison(aics_result, ucsp_result)
19:   return ConvertToAICSFormat(final_result)
20: else if config.routing_mode == GRAYSCALE then
21:   // 灰度发布: 基于用户ID路由
22:   user_id  $\leftarrow$  ExtractUserID(aics_request)
23:   hash  $\leftarrow$  MurmurHash3(user_id  $\oplus$  config.adapter_id)
24:   bucket  $\leftarrow$  hash mod 10000
25:   threshold  $\leftarrow$  config.grayscale_ratio  $\times$  10000
26:   if bucket < threshold then
27:     // 走UCSP路径
28:     ucsp_request  $\leftarrow$  ConvertToUCSPFormat(aics_request)
29:     result  $\leftarrow$  CallUCSP(ucsp_request, config.ucsp_endpoint)
30:     if config.enable_dual_verification then
31:       aics_result  $\leftarrow$  CallAICS(aics_request, config.aics_endpoint)
32:       RecordComparison(aics_result, result)
33:     end if
34:   else
35:     // 走AICS路径
36:     result  $\leftarrow$  CallAICS(aics_request, config.aics_endpoint)
37:   end if
38:   return ConvertToAICSFormat(result)
39: end if
```

3.4 结果融合策略

结果融合算法：

Algorithm 18 AICS与UCSP结果融合算法

```
1: 输入: AICS结果 aics_result, UCSP结果 ucsp_result, 融合策略 strategy
2: 输出: 融合后的结果 fused_result
3:
4: if strategy == AICS_PRIORITY then
5:   // AICS优先: 仅在AICS不确定时使用UCSP结果
6:   if aics_result.confidence < UNCERTAIN_THRESHOLD then
7:     return ucsp_result
8:   else
9:     return aics_result
10:  end if
11: else if strategy == UCSP_PRIORITY then
12:   // UCSP优先: 仅在UCSP不确定时使用AICS结果
13:   if ucsp_result.confidence < UNCERTAIN_THRESHOLD then
14:     return aics_result
15:   else
16:     return ucsp_result
17:   end if
18: else if strategy == CONSENSUS then
19:   // 共识策略: 两者都拦截或都放行时采用, 否则采用更严格的结果
20:   if aics_result.blocked == ucsp_result.blocked then
21:     return aics_result {一致时采用AICS结果保持兼容}
22:   else
23:     // 不一致时, 采用更严格的结果 (拦截优先)
24:     if aics_result.blocked then
25:       return aics_result
26:     else
27:       return ucsp_result
28:     end if
29:   end if
30: else if strategy == CONFIDENCE_BASED then
31:   // 置信度策略: 选择置信度更高的结果
32:   if aics_result.confidence > ucsp_result.confidence then
33:     return aics_result
34:   else
35:     return ucsp_result
36:   end if
37: end if
```

3.5 数据迁移方案

迁移策略:

Algorithm 19 AICS数据迁移流程

```
1: 输入: AICS数据源 aics_data
2:
3: // 阶段1: 数据导出
4: exported_data  $\leftarrow$  ExportAICSData(aics_data)
5: ValidateDataFormat(exported_data)
6:
7: // 阶段2: 数据转换
8: converted_data  $\leftarrow$  ConvertToUCSPFormat(exported_data)
9: ValidateUCSPFormat(converted_data)
10:
11: // 阶段3: 数据导入 (分批进行)
12: batches  $\leftarrow$  SplitIntoBatches(converted_data, BATCH_SIZE)
13: for batch IN batches do
14:   ImportToUCSP(batch)
15:   VerifyImport(batch)
16:   if NOT VerifyImport(batch) then
17:     RollbackBatch(batch)
18:     LogError("Batchimport failed", batch)
19:   end if
20: end for
21:
22: // 阶段4: 数据一致性验证
23: VerifyDataConsistency(aics_data, ucsp_data)
24:
25: // 阶段5: 切换流量
26: SwitchTrafficToUCSP()
```

配置迁移:

Algorithm 20 AICS配置迁移算法

```
1: 输入: AICS配置 aics_config
2: 输出: UCSP配置 ucsp_config
3:
4: // 模型配置迁移
5: ucsp_config.fast_model  $\leftarrow$  MapAICSModel(aics_config.primary_model)
6: ucsp_config.deep_model  $\leftarrow$  MapAICSModel(aics_config.fallback_model)
7:
8: // 策略配置迁移
9: ucsp_config.policy  $\leftarrow$  ConvertAICSPolicy(aics_config.policy)
10:
11: // 阈值配置迁移
12: ucsp_config.block_threshold  $\leftarrow$  aics_config.block_threshold
13: ucsp_config.warn_threshold  $\leftarrow$  aics_config.warn_threshold
14:
15: // 验证配置有效性
16: ValidateUCSPConfig(ucsp_config)
17: return ucsp_config
```

3.6 渐进式迁移方案

迁移阶段规划:

1. 阶段1: 影子模式 (Shadow Mode)

- 所有请求走AICS, 同时异步调用UCSP
- 对比AICS和UCSP结果, 收集差异数据
- 评估UCSP性能和准确性
- 持续时间: 2-4周

2. 阶段2: 小流量验证 (1%-5%)

- 1%流量走UCSP, 99%走AICS
- 双路验证, 确保结果一致性
- 监控UCSP稳定性指标
- 持续时间: 2-3周

3. 阶段3: 逐步增加 (5%-50%)

- 每周增加5-10%流量
- 持续监控和调整
- 处理发现的问题

- 持续时间：6-8周

4. 阶段4：全量切换（50%-100%）

- 逐步增加到100%
- AICS作为兜底备用
- 最终完全切换到UCSP
- 持续时间：4-6周

迁移监控指标：

Algorithm 21 迁移监控算法

```

1: 输入：迁移配置 migration_config
2:
3: // 实时监控指标
4: metrics  $\leftarrow$  GetMigrationMetrics(migration_config)
5:
6: // 一致性指标
7: agreement_rate  $\leftarrow$  CalculateAgreementRate(metrics.aics_results, metrics.ucsp_results)
8:
9: // 性能指标
10: ucsp_latency  $\leftarrow$  CalculateAverageLatency(metrics.ucsp_requests)
11: aics_latency  $\leftarrow$  CalculateAverageLatency(metrics.aics_requests)
12:
13: // 准确性指标
14: ucsp_accuracy  $\leftarrow$  CalculateAccuracy(metrics.ucsp_results, metrics.ground_truth)
15: aics_accuracy  $\leftarrow$  CalculateAccuracy(metrics.aics_results, metrics.ground_truth)
16:
17: // 决策：是否继续迁移
18: if agreement_rate > MIN_AGREEMENT_RATE AND ucsp_latency <
    LATENCY_THRESHOLD AND ucsp_accuracy  $\geq$  aics_accuracy then
19:   IncreaseMigrationRatio(migration_config)
20: else
21:   PauseMigration(migration_config)
22:   TriggerAlert("Migration paused due to metrics degradation")
23: end if

```

3.7 统一监控与可观测性

统一指标聚合：

Algorithm 22 统一指标聚合算法

```
1: 输入: AICS指标 aics_metrics, UCSP指标 ucsp_metrics
2: 输出: 统一指标 unified_metrics
3:
4: // 聚合请求量
5: unified_metrics.total_requests  $\leftarrow$  aics_metrics.requests + ucsp_metrics.requests
6:
7: // 聚合拦截量
8: unified_metrics.total_blocked  $\leftarrow$  aics_metrics.blocked + ucsp_metrics.blocked
9:
10: // 计算平均延迟 (加权平均)
11: aics_weight  $\leftarrow$  aics_metrics.requests / unified_metrics.total_requests
12: ucsp_weight  $\leftarrow$  ucsp_metrics.requests / unified_metrics.total_requests
13: unified_metrics.avg_latency  $\leftarrow$  aics_weight  $\times$  aics_metrics.avg_latency + ucsp_weight  $\times$ 
    ucsp_metrics.avg_latency
14:
15: // 聚合错误率
16: unified_metrics.error_rate  $\leftarrow$  (aics_metrics.errors + ucsp_metrics.errors) / unified_metrics.total_requests
17:
18: // 计算迁移进度
19: unified_metrics.migration_progress  $\leftarrow$  ucsp_metrics.requests / unified_metrics.total_requests
20:
21: return unified_metrics
```

4 性能优化

4.1 缓存策略

Algorithm 23 带缓存的内容检查

```
1: 输入: 文本 text, 用户ID user_id
2: 输出: 安全检查结果 result
3:
4: cache_key  $\leftarrow$  HashText(text)
5: cached_result  $\leftarrow$  Cache.Get(cache_key)
6: if cached_result  $\neq$  NULL then
7:   return cached_result
8: end if
9:
10: result  $\leftarrow$  CheckContent(text, user_id)
11: if result.confidence  $>$  0.90 then
12:   Cache.Set(cache_key, result, TTL = 3600) {1小时TTL}
13: end if
14: return result
```

4.2 批量处理优化

Algorithm 24 批量检查优化

```
1: 输入: 请求列表 requests[]
2: 输出: 结果列表 results[]
3:
4: fast_batch  $\leftarrow$  []
5: deep_batch  $\leftarrow$  []
6:
7: for request IN requests do
8:   quick_check  $\leftarrow$  QuickPreCheck(request.text)
9:   if quick_check.confidence  $>$  0.95 then
10:    fast_batch.APPEND(request)
11:   else
12:    deep_batch.APPEND(request)
13:   end if
14: end for
15:
16: fast_results  $\leftarrow$  PARALLELFastModel.BatchInference(fast_batch)
17: deep_results  $\leftarrow$  PARALLELDeepModel.BatchInference(deep_batch)
18: return MERGE(fast_results, deep_results)
```

5 监控与可观测性

5.1 实时指标监控

Algorithm 25 指标更新算法

```
1: 输入: 安全检查结果 result, 处理时间 processing_time
2:
3: ATOMIC INCREMENT metrics.total_requests
4: if result.blocked then
5:   ATOMIC INCREMENT metrics.blocked_count
6: end if
7:
8: UPDATE latency_histogram(processing_time)
9: metrics.avg_latency_ms  $\leftarrow$  CalculateAverage(latency_histogram)
10: metrics.p95_latency_ms  $\leftarrow$  CalculatePercentile(latency_histogram, 0.95)
11: metrics.p99_latency_ms  $\leftarrow$  CalculatePercentile(latency_histogram, 0.99)
12:
13: if result.model_version == FAST_MODEL_VERSION then
14:   INCREMENT fast_model_hits
15: else
16:   INCREMENT deep_model_hits
17: end if
18:
19: metrics.fast_model_hit_rate  $\leftarrow$  fast_model_hits/metrics.total_requests
20: metrics.deep_model_hit_rate  $\leftarrow$  deep_model_hits/metrics.total_requests
```

5.2 告警机制

Algorithm 26 告警检查算法

```
1: 输入: 系统指标 metrics
2:
3: if metrics.p95_latency_ms > LATENCY_THRESHOLD then
4:   TriggerAlert("Highlatencydetected", metrics.p95_latency_ms)
5: end if
6:
7: false_positive_rate  $\leftarrow$  metrics.false_positive_count/metrics.total_requests
8: if false_positive_rate > FALSE_POSITIVE_THRESHOLD then
9:   TriggerAlert("Highfalsepositiverate", false_positive_rate)
10: end if
11:
12: false_negative_rate  $\leftarrow$  metrics.false_negative_count/metrics.total_requests
13: if false_negative_rate > FALSE_NEGATIVE_THRESHOLD then
14:   TriggerAlert("Highfalsenegativerate", false_negative_rate)
15: end if
16:
17: if metrics.fast_model_hit_rate < MIN_HIT_RATE then
18:   TriggerAlert("Fastmodelperformancedegraded")
19: end if
```

6 实施路线图

6.1 阶段一：核心功能实现（1-2个月）

目标：实现统一接口和智能路由

任务：

1. 实现 Safety Gateway
2. 部署 Fast Model 和 Deep Model
3. 实现统一接口 CheckContent
4. 实现智能路由算法

验收标准：

- 90% 请求在快速模型完成
- 平均延迟 \leq 50ms
- 接口可用性 \geq 99.9%

6.2 阶段二：灰度发布框架（2-3个月）

目标：实现灰度发布，逐步接管Vendor流量

任务：

1. 实现灰度发布路由算法
2. 实现双路验证机制
3. 实现自动流量调整
4. 实现成本控制算法
5. 实现稳定性检查与快速回滚

灰度发布计划：

- 第1周： 1% 流量验证，双路对比，确保稳定性
- 第2-3周： 5% 流量验证，持续监控指标
- 第4-5周： 10% 流量，验证成本控制
- 第6-7周： 20% 流量，验证系统负载
- 第8-9周： 50% 流量，验证大规模场景
- 第10-11周： 80% 流量，接近全量
- 第12周： 100% 流量，完全接管

验收标准：

- 支持1%-100%流量比例动态调整
- 稳定性指标（误杀率、漏判率） \leq 阈值
- 成本控制在预算范围内
- 支持秒级回滚
- 与Vendor结果一致性 \geq 95%

6.3 阶段三：A/B 测试框架（3-4个月）

目标：实现确定性 A/B 测试

任务：

1. 实现确定性路由算法
2. 实现指标收集系统
3. 实现统计显著性检测

4. 实现实验管理界面

验收标准：

- 支持多实验并行
- 指标实时更新
- 自动决策支持

6.4 阶段四：策略管理系统（4-5个月）

目标：实现代码驱动的策略管理

任务：

1. 实现策略定义框架
2. 实现策略版本管理
3. 实现策略测试框架
4. 实现策略部署流程

验收标准：

- 策略变更可追溯
- 策略可单元测试
- 支持策略回滚

7 总结

本设计文档阐述了一个统一的内容安全平台架构，通过统一接口设计、智能路由、确定性A/B测试和代码驱动的策略管理等核心技术，实现了高性能、高可维护性的内容安全检查系统。该架构具有以下优势：

1. **高性能：** 90% 请求在 20ms 内完成，整体平均延迟 $\leq 35\text{ms}$
2. **高可维护性：** 统一接口设计，代码量减少 60%
3. **科学验证：** 确定性 A/B 测试框架，支持策略科学验证
4. **可追溯性：** 策略版本管理，所有变更可追溯
5. **可扩展性：** 模块化设计，易于扩展新功能
6. **平滑迁移：** 灰度发布策略，支持从Vendor逐步切换到自研系统
7. **向后兼容：** AICS集成设计，实现现有系统平滑迁移，保持业务连续性

该架构为构建企业级内容安全平台提供了坚实的技术基础，同时支持从现有系统（如AICS）平滑迁移，降低了系统升级的风险和成本。