# From Engineer to Architect – Engineering Disciplines for High–Quality and Secure Code
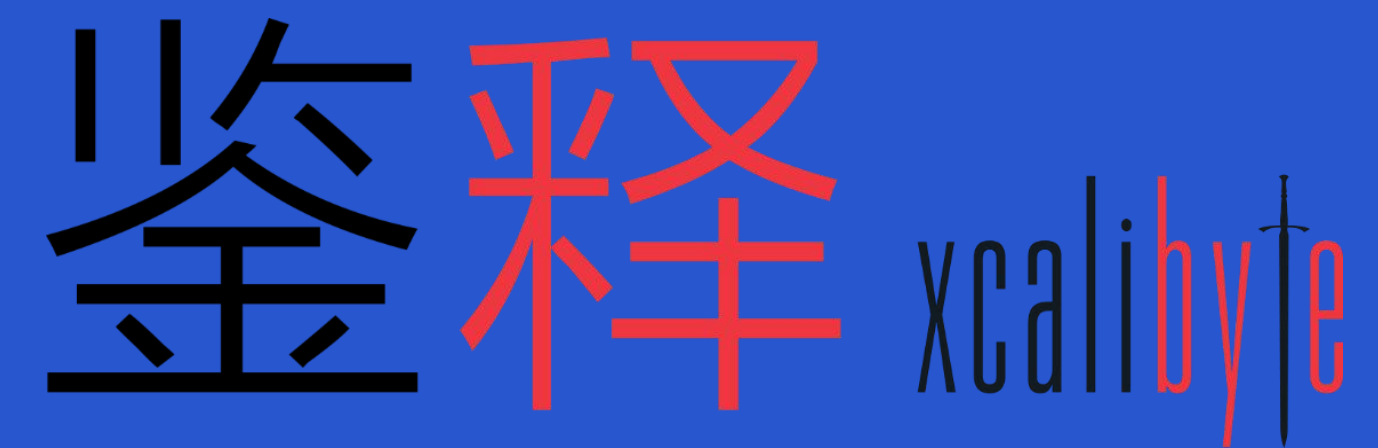
鉴释 Xcalibyte

Shinming Liu

Chief Architect, Xcalibyte

ArchSummit 全球架构师峰会

鉴释 InfoQ

为一线互联网公司核心
技术人员提供优质内容

☑ 每日要闻  ☑ 技术干货

☑ 大咖访谈  ☑ 行业趋势

关注 InfoQ 官方微博

**InfoQ** 促进软件开发及相关领域知识与创新的传播

# Who Am I

刘新铭 **Shin-Ming Liu**

Chief Architect & Co–founder, Xcalibyte

- Compiler Scientist
- Former Director, Intel IOT Research Lab, China
- Former Director, HP Compiler Technology Lab
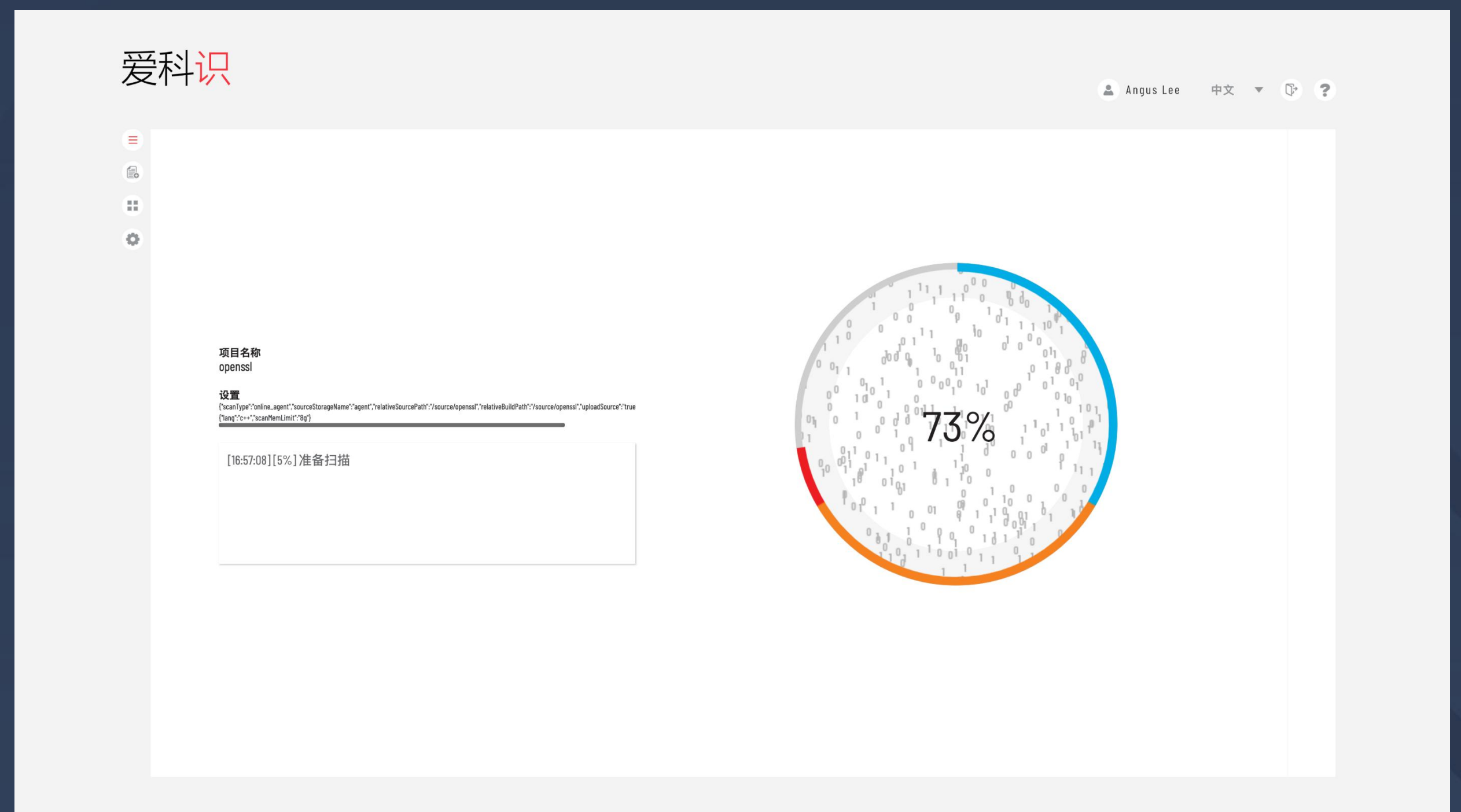- 10+ patents granted in program analysis and compiler optimization

# Agenda

- 案例：数据库重构项目
- 用 First Principle 来分析，思考
- 用性能数据来佐证、决策，并管控
- 用程序内置工具来监控架构劣坏
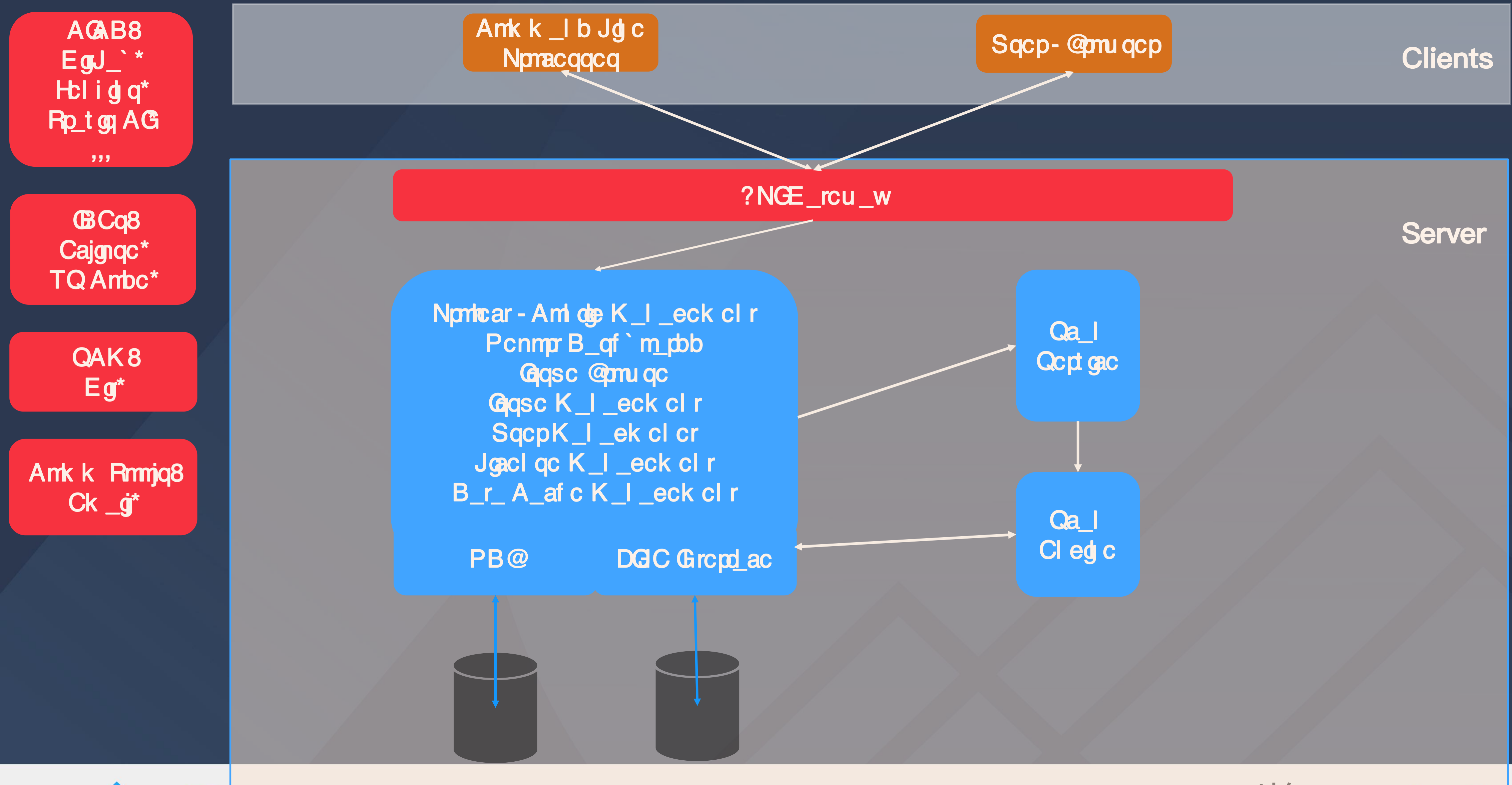- 为算法适配程序语言
- 重构要及时并持续进行
- 自我提升

# CASE — DATABASE REFACTOR PROJECT

## Pain Points

➤ Service Does Not Complete

➤ Service Timeout

➤ Random Service Restart

➤ UI Response Time Takes More than 25min for Extreme Cases

# Xcalscan was a Monolithic Java Application

# PERFORMANCE DATA COLLECTED

Measurement Made:

➢ Input Data Set Size

➢ DB Tables Size Collected For The Input

➢ Memory Footprint and Elapse Time

❑ Data Ingestion into DB for the Input

❑ IssueChangeAnalysis function for One Data Set

➢ Response Time for Dashboard Display

# PERFORMANCE DATA COLLECTED

## Measurement Made for MySQL5.7

- ➢ Input Data Set Size ─ 2.2 million LOC, 565MB Issue File

- ➢ DB Tables Size Collected For The Input ─ 2GB

- ➢ Memory Footprint and Elapse Time

  - ❏ Data Ingestion into DB for The Input ─ 22GB

  - ❏ IssueChangeAnalysis function for One Data Set ─ 12G

- ➢ Response Time for Dashboard Display ─ 180 Sec

# ISSUE IDENTIFIED

- Direct Correlation
  - ➢ Input Data Set Size ➔ Memory Footprint ➔ Elapse Time
- DB Tables Size Increases Accumulatively
  - ➢ This result in Slow Response Time DB Query
- Inefficient Query Function Used
  - ➢ API Involves Full Table Scan, Join, Select
- Interpretive Language Used in IssueChangeAnalysis Algorithm

# DESIGN CHANGES MADE

- Redesign Data Format for Data Produced by Scan Engines
  - ➢ Use String Table ー Eliminate Duplicate String
  - ➢ Use Proprietary Data Format ー Internal Use Only
- Database Table Changes:
  - ➢ Source Code Issue Introduction Time (GIT commit ID)
  - ➢ Ingest New Issue Introduced Only, Minimize DB Size Increase
  - ➢ Remove Read–only Data from RDB
- Precompute Time Consuming Frequent Query
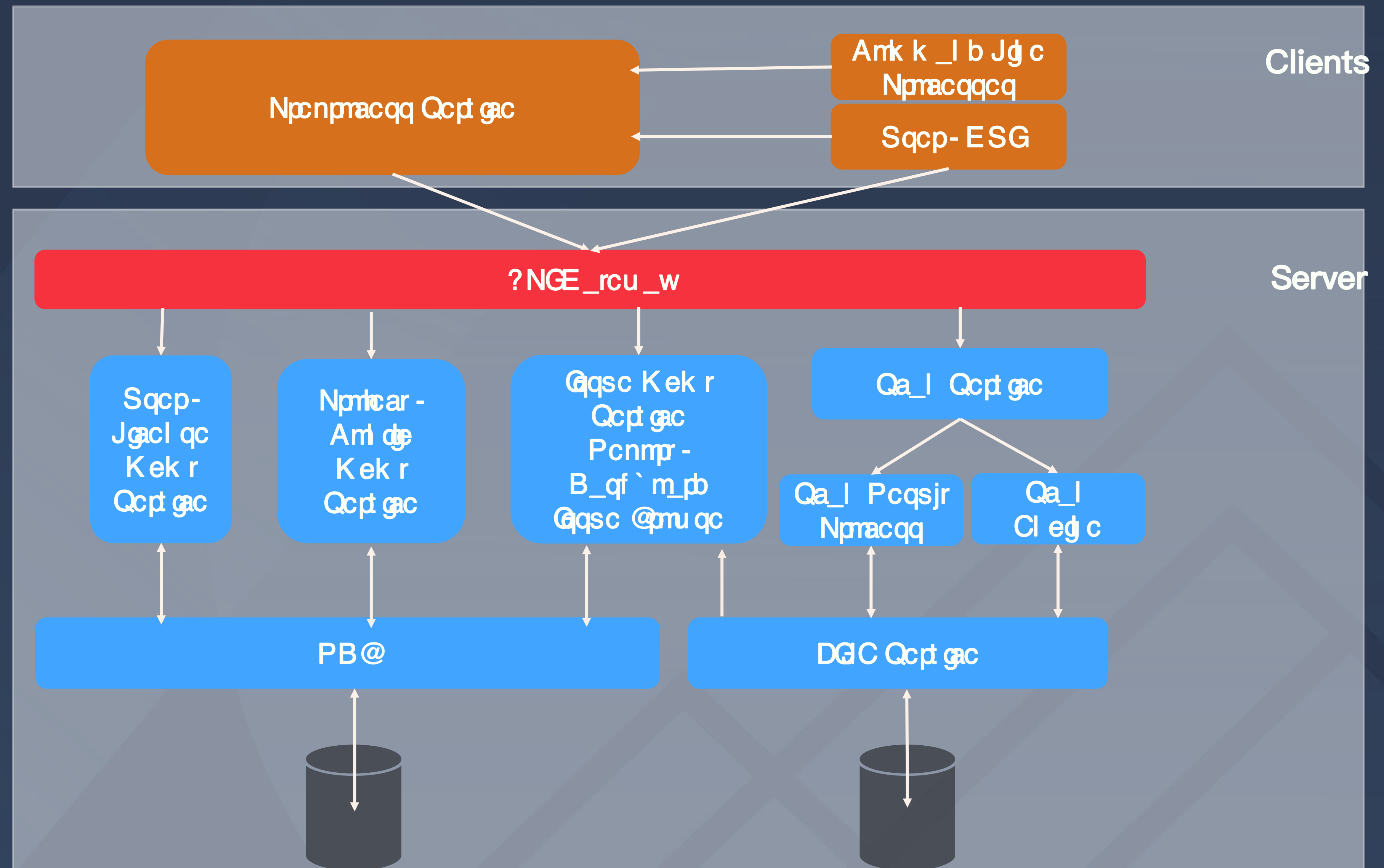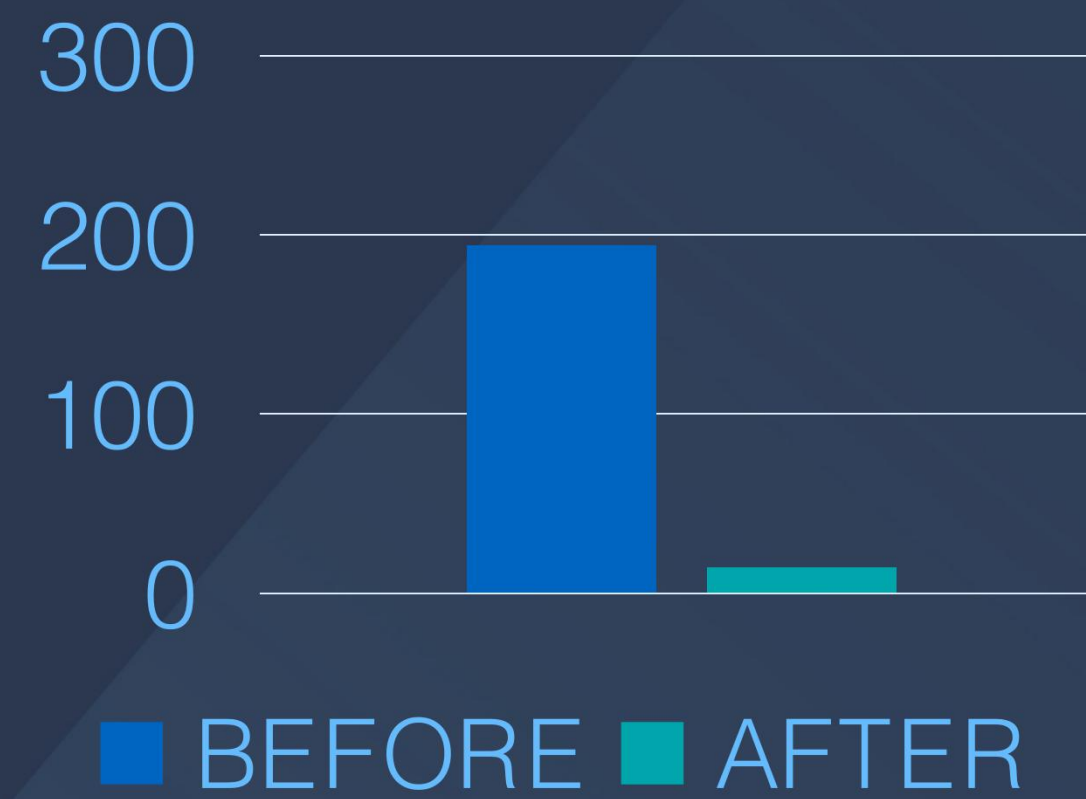- Use C++ for Complex Algorithms

# Xcalscan Architecture After Refactor

# PERFORMANCE IMPROVEMENTS



Ave. Scan Result
Import Time

Scan Result
Retrieval Time

BEFORE ■ AFTER

BEFORE ■ AFTER

- Precompute Summary Data, Reduce UI Dashboard Load Time

# PERFORMANCE IMPROVEMENTS

(cont)



- Identify and Skip Redundant Data, Enables Large Project Scan Such as MySQL5.7

# PERFORMANCE IMPROVEMENTS
(cont)

**ICA:IssueChangeAnalysis**



GB Mem. Required for ICA

■ BEFORE ■ AFTER

Time (sec.) Spent for ICA

■ BEFORE ■ AFTER

- Rewrite Complex Algorithm with C++, Save Both Time and Space

# PERFORMANCE IMPROVEMENTS

(cont)



- Remove Rule Information From RDB

  - ➤ Lazy Load of Rule Information

  - ➤ Reduced Memory Footprint During Query

- Time to retrieve rule info (Read Only Data) – **>10x less**

# DESIGN w/ FIRST PRINCIPLE

- Steps Involved
  - ➢ Collect Requirements and Facts
  - ➢ Partition Issues till Atomic
  - ➢ Reduce to Minimum: Discard Irrelevant & Unimportant
- What Mattered In The Refactor Effort:
  - ➢ There were Memory & I/O Bound Issues

# QUANTITATIVE OBJECTIVES & MEASURE

- Performance Data Played a Key Role for Architecture Decision
  - ➢ Architecture Choices are Subjective
  - ➢ Measured Data are Objective
- Measurements on Data, Resource and Logic
  - ➢ Data Size 一 Footprint, Persistent
  - ➢ Data Processing Logic 一 Access Pattern Matters

# QUANTITATIVE OBJECTIVES & MEASURE

- **Strategies Applied** in DB Refactor Project

  ➢ Program Stability Issues Happened on Cases with Large Data Set

  ➢ Systematically Measure Data Size, Memory Footprint

  ➢ Track Response Time for APIs that is > 2 Seconds

- **Policy Changed** during and after DB Refactor Project

  ➢ Performance Measured and Summarized Nightly

# PROGRAM TO DETECT ISSUES INTELLIGENTLY

- **Strategies Applied** in DB Refactor Project

  ➢ Log Data Characteristics at Building Block Level

  ➢ Instrument Performance Monitor ─ White Box Approach

  ➢ Nightly Regression Test Accelerates Project Integration

  ➢ Use Option to Enable / Disable Functionality for Triage Purpose

  ➢ Unified Log Format Facilitate Issue Reproduction

# CODE WITH RIGHT PROGRAMMING LANGUAGE IN RIGHT PLACE

- Strategy Applied in DB Refactor Project
  - ➢ Easy DB manipulation, Java Stayed w/ Reduced Scope
  - ➢ Performance intensive, C++ : Pre–calculate Data, ICA function
  - ➢ Wrapper Functions, Added JavaScript for React Programming

# REFACTOR OFTEN AND EARLY

- Recommend Book:

  ➢ http://stepanovpapers.com/notes.pdf by Alex Stepanov

  ➢ He invented Generic Programming and C++ STL

  ➢ Decomposing An Application Into A Collection Of General–purpose Algorithms And Data Structures Makes It Robust

- Strategy Adopted: Incremental Phase–in Changes

# RENEW KNOWLEDGE PERIODICALLY

- Bottom Line: You Own Your Own Career !

  ➢ Get Ready for Your Next Job 一 In The Same Company Or Not

- Drive the Learning To Enhance Your Knowledge Framework

  ➢ When And Where Will You Need New Techniques To Improve Your Project Architecture

  ➢ Learning the Design Rationale Behind Open–Source Package

- Dedicate 8 Hours Per Week

# 严谨 — 架构师的必要特质

- 用科学方法来工作：实验尝试，小心求证
- 用 First Principle 来分析、思考
- 用性能数据来佐证、决策，并管控
- 用程序内置工具来监控架构劣坏
- 为算法适配程序语言
- 重构要及时并持续进行
- Stay Hungry, Stay Foolish

# THANKS

—

Global

Architect Summit

# AGENDA

- Case 一 Database Refactor (DR) Project
- Design with First Principle
- Define Quantitative Objectives and Measure Nightly
- Program to Detect Issues Intelligently
- Code with Right Programming Language in Right Place
- Refactor Often and Early
- Renew Knowledge Periodically
- Summary

# IMPROVEMENTS AFTER DB REFACTOR

- DB Ingestion Time — **Now 4x faster**

- DB Query Time — **Now 20x faster**

- Response Time for UI dashboard (50 proj) — **2–3 sec, was 3 — 25 min**

- Enable Lage Project Scans
  - ➤ **File size 10x smaller; DB Ingestion Time 10x less; Smaller Memory Footprint, e.g. MySQL5.7 — 565MB vs. 63MB; 1807s vs. 43s; 22G vs. < 1G**

- Time and memory required to generate IssueChange
  - ➤ **Time 10x less; Memory 25x less, e.g. SQLite 20+min vs. 2 min; 10G vs. 0.4G.**

- Time to retrieve rule info (Read Only Data) – **>10x less**

# QCon<sup>+</sup> 案例研习社

## 学习前沿案例，向行业领先迈进

**40** 个
热门专题
—
行业专家把关内容筹备，
助你快速掌握最新技术发展趋势

**200** 个
实战案例
—
了解大厂前沿实战案例，
为 200 个真问题找到最优解

**40** 场
直播答疑
—
40 位技术大咖，每周分享最新
技术认知，互动答疑

**365** 天
持续学习
—
视频结合配套 PPT
畅学 365 天

# SUMMARY

- 用科学方法来工作：实验尝试，小心求证
- 用First Principle 来思考：Connecting the Dot
- 用 OKR 来自我管理
- 用流水线来安排工作
- 在制高点审视全局
- 时时准备交班
- 时时刻刻提升自己的市场价值
- 每三年检讨自己在 Maslow's Hierarchy 的位置，步步为营
- 失败为成功之母，失之东隅，收之桑榆
- Stay Hungry, Stay Foolish