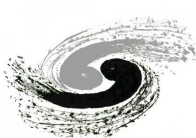


# 面向大数据的科学软件技术与实践

田浩来,邹佳恒,林韬,李卫东,张俊荣

高能物理研究所

2013 年 7 月 9 日



第十六届全国科学计算与信息化会议暨科研大数据论坛  
辽宁 大连

# 内容

- 1 背景和挑战
- 2 科学软件框架开发技术
- 3 科学软件框架的实践

- 1 背景和挑战
- 2 科学软件框架开发技术
- 3 科学软件框架的实践



# 思路

## 复杂性、抽象与耦合

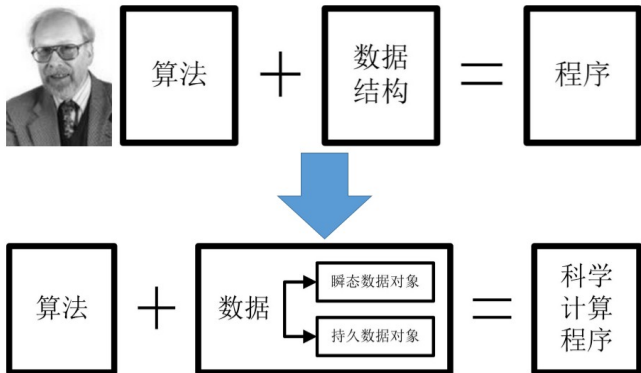
- 对功能的抽象带来耦合
- 解耦合带来复杂度
- 理解复杂需要抽象

## 解决方法

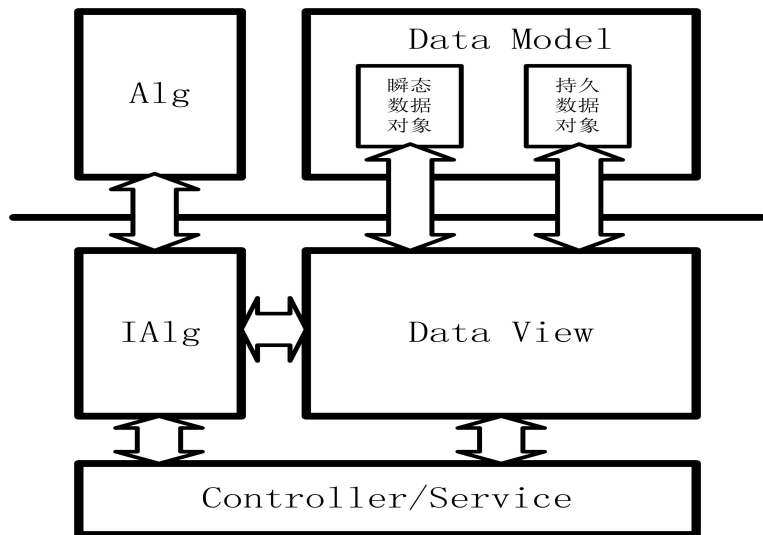
- 简洁轻量
- 面向接口编程
- 层次化模块化
- 插件式开发



# 以数据为中心算法与数据分离的软件架构



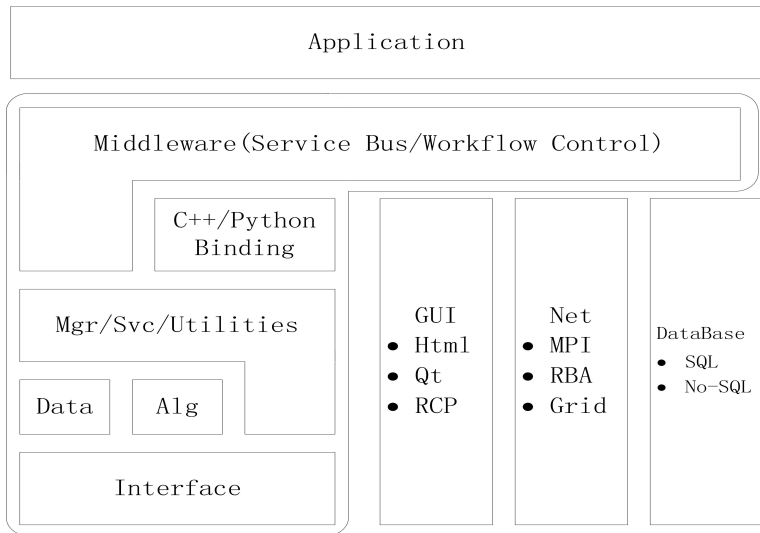
## 面向接口编程解除组件之间的耦合



- 1 背景和挑战
- 2 科学软件框架开发技术
- 3 科学软件框架的实践



# 整体框架



# 接口与接口管理

## 接口

- 一切皆为接口:所有对象的基类
- 编译时:类型与类型名注册
- 运行时:接口与接口名绑定，通过接口调用实例

## 接口管理

- 按Key-Value模式保存实例的名字和对应的指针
- 通过列表等数据结构管理实例
- 通过工厂模式创建实例

## 实例类型判断

- RTTI，运行时类型判断
- BOOST MPL 模版元编程，编译时类型计算判断

# 数据与算法

## 数据

- 树形结构
- Table/Matrix/List/Element 等简单数据类型
- 持久数据与瞬态数据的转换
  - 事例循环: 持续加载, 持续运算, 持续保存
  - 矩阵运算: 开始时一次性加载, 结束后一次性保存

## 算法

- 算法模版: initialize(), execute(), finalize()
- 算法嵌套调用

# 公共模块

## 管理模块

负责算法、数据区、服务和应用等运行实例的创建和管理

## 服务模块

与数据处理算法相关，但不是算法和数据的其他运行组件

## 公用模块

Timing、Log和Message等独立公共应用集合

# C++ vs. Python

## C++

算法和数据区等内存和效率要求高的场合

## Python

workflow引擎，调用接口，组件和应用配置等灵活的模块，性能要求不高的组件的快速开发。

## BOOST\_Python

C++和Python的绑定: 算法和数据区等C++模块的Python调用接口实现

# 中间层

## 总线

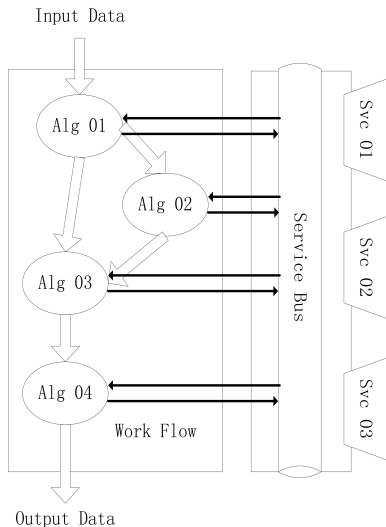
本进程，跨进程和跨主机的服务访问机制

- 按名访问
- 命名规则

## workflow

对数据进行处理的数据调用序列

- 分支以及条件判断
- 继承以及组合



# 高性能计算

- Many Cores(GPUs) → OpenCL/CUDA → Data/Algorithm
- Multi Cores(CPUs) → Multi-thread → Mgr/Svc
- Multi Cores(CPUs) → Multi-process → Middle ware
- Distributed Computation → MPI/Grid/CORBA → Middle ware

- 1 背景和挑战
- 2 科学软件框架开发技术
- 3 科学软件框架的实践

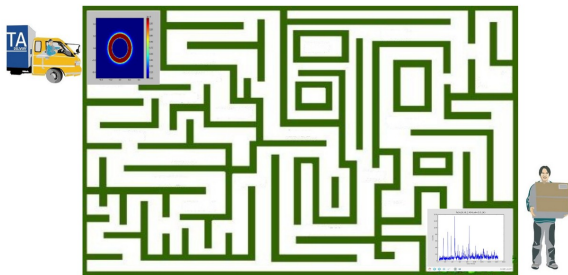


# Gaudi: 对撞机物理实验数据处理框架



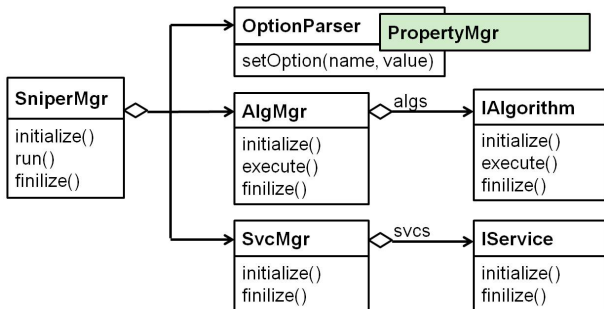
- 数据：树形结构
- 事例循环，事例相互独立
- 复杂的数据IO服务
- 简单的算法调用序列
- 业务焦点：数据格式，算法实现
- 业务难点：计算效率

# Mantid: 中子散射实验数据处理框架



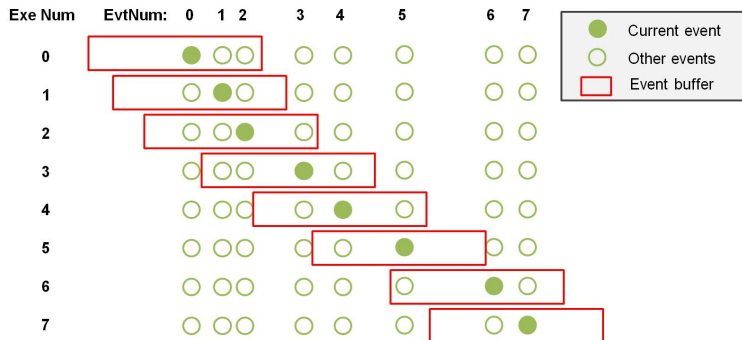
- 数据：列表和矩阵
- 列表处理和矩阵运算
- 通过Load/Save算法对数据进行IO
- 复杂的工作流引擎
- 业务焦点：数据加载，工作流实现
- 业务难点：内存占用

# SNiPER: 非对撞机物理实验数据处理框架



- 基于开放架构，通过插件实现Gaudi和Mantid的大部分功能
- 支持事例循环、事件驱动、列表处理和矩阵运算
- 业务焦点: 组件管理，数据管理，框架的可扩展性
- 业务难点: 数据管理的灵活性

# 关联事例分析



# 总结

**谢谢!**