

格点 QCD 的计算方法与软件研发

宫明

中国科学院高能物理研究所
CLQCD 合作组

中国科学技术大学
2024.07.16

格点 QCD 是从第一性原理出发非微扰地系统研究强子和强相互作用的高能物理前沿领域

格点 QCD 是从第一性原理出发非微扰地系统研究强子和强相互作用的高能物理前沿领域

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

格点量子色动力学 (Lattice QCD)

第一性原理

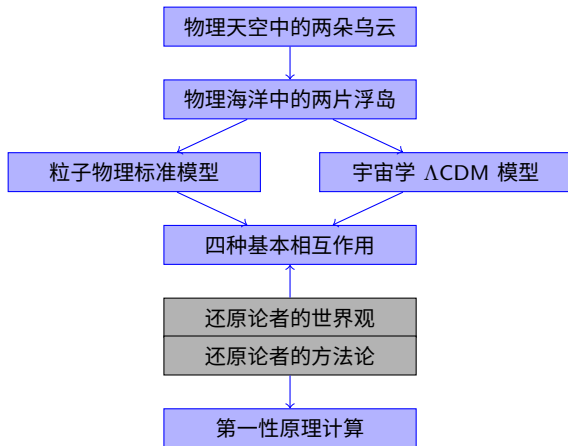
非微扰

系统研究

强相互作用

高能物理前沿

领域



第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

粒子物理标准模型



The Standard Model

□ Matter particles

■ 6 quarks

$$\begin{pmatrix} u \\ d \end{pmatrix} \quad \begin{pmatrix} s \\ c \end{pmatrix} \quad \begin{pmatrix} t \\ b \end{pmatrix}$$

■ 6 leptons

$$\begin{pmatrix} e \\ \nu_e \end{pmatrix} \quad \begin{pmatrix} \mu \\ \nu_\mu \end{pmatrix} \quad \begin{pmatrix} \tau \\ \nu_\tau \end{pmatrix}$$

Hadrons (proton, neutron, pion etc) are composites of 3 or 2 quarks



proton=uud

□ Particles mediating interactions

■ photon

$$\gamma \quad \text{EM}$$

■ Weak bosons

$$W, Z \quad \text{Weak interactions}$$

■ gluons

$$g \quad \text{Strong interaction}$$

Weinberg-Salam theory

Quantum Chromodynamics (QCD)

□ Gauge field theory based on

$$\underbrace{SU(3)}_{\text{QCD}} \otimes \underbrace{SU(2) \otimes U(1)}_{\text{EM+Weak}}$$

$$L_{QCD} = -\frac{1}{4} \text{Tr}(F^{\mu\nu} F_{\mu\nu}) + \sum_f \bar{q}_f \left[i\gamma^\mu \cdot (\partial_\mu - igA_\mu) - m_f \right] q_f$$

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

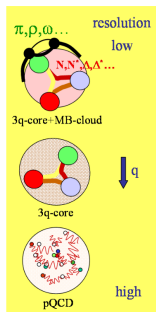
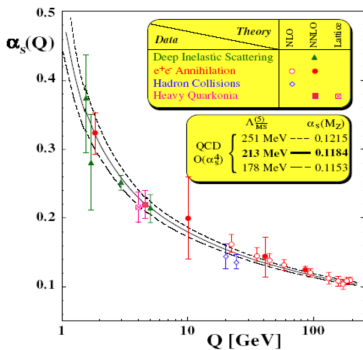
系统研究

强相互作用

高能物理前沿

领域

夸克禁闭与渐进自由



格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

Wick 转动

$$x_0 \equiv t \rightarrow -ix_4 \equiv -i\tau$$

$$p_0 \equiv E \rightarrow ip_4$$

$$e^{iS_M} \equiv e^{i \int d^4x_M L(x_M)} = e^{\int d^4x_E L(x_E)} \equiv e^{-S_E}$$

欧氏时空的场论语言和统计语言

The equivalences between a Euclidean field theory and Classical Statistical Mechanics.

Euclidean Field Theory	Classical Statistical Mechanics
Action	Hamiltonian
unit of action \hbar	units of energy $\beta = 1/kT$
Feynman weight for amplitudes	Boltzmann factor $e^{-\beta H}$
$e^{-S/\hbar} = e^{-\int \mathcal{L} dt/\hbar}$	
Vacuum to vacuum amplitude	Partition function $\sum_{conf.} e^{-\beta H}$
$\int \mathcal{D}\phi e^{-S/\hbar}$	
Vacuum energy	Free Energy
Vacuum expectation value $\langle 0 \mathcal{O} 0 \rangle$	Canonical ensemble average $\langle \mathcal{O} \rangle$
Time ordered products	Ordinary products
Green's functions $\langle 0 T[\mathcal{O}_1 \dots \mathcal{O}_n] 0 \rangle$	Correlation functions $\langle \mathcal{O}_1 \dots \mathcal{O}_n \rangle$
Mass M	correlation length $\xi = 1/M$
Mass-gap	exponential decrease of correlation functions
Mass-less excitations	spin waves
Regularization: cutoff Λ	lattice spacing a
Renormalization: $\Lambda \rightarrow \infty$	continuum limit $a \rightarrow 0$
Changes in the vacuum	phase transitions

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

配分函数与作用量

$$Z = \int DA_\mu D\psi D\bar{\psi} e^{-S}$$

$$S = \int d^4x \left(\frac{1}{4} F_{\mu\nu} F^{\mu\nu} - \bar{\psi} M \psi \right)$$

把费米场积分掉

$$Z = \int DA_\mu \det M e^{-\frac{1}{4} F_{\mu\nu} F^{\mu\nu}}$$

物理观测量

$$\langle O \rangle = \frac{1}{Z} \int DA_\mu O e^{-S}$$

问题的转化

在以规范场为自由度的统计系统里，求解观测量的系综平均值

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

科学的研究范式

- 1、实验观测
- 2、理论推演
- 3、**数值模拟**：格点 QCD 面向理论，但“以上帝的视角做实验”
- 4、数据推动？

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

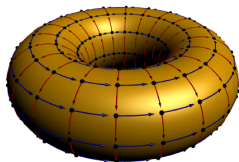
高能物理前沿

领域

科学的研究范式

- 1、实验观测
- 2、理论推演
- 3、数值模拟：格点 QCD 面向理论，但“以上帝的视角做实验”
- 4、数据推动？

离散化：四维环面



在计算机上的表示

- 规范场用群表示： $U_\mu(x) = e^{iagA_\mu(x)}$ ，是 3×3 的复矩阵
- 费米场只需夸克传播子： $G_q(U) = M^{-1}(U)$ ，由规范场决定

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

作用量的设计

- 规范场能构成的规范不变的量，只有圈，最简单的示例：

$$S_g = \frac{6}{g^2} \sum_x \sum_{\mu < \nu} \text{ReTr} \frac{1}{3} \left(1 - W_{\mu\nu}^{1 \times 1} \right)$$

- 为了减小离散误差，其他定义按特定比例加入了较大的圈
- 费米场作用量示例：Wilson 费米子矩阵：

$$M_{x,y}^W[U] a = \delta_{x,y} - \kappa \sum_{\mu} \left[\left(r - \gamma_{\mu} \right) U_{x,\mu} \delta_{x,y-\mu} + \left(r + \gamma_{\mu} \right) U_{x-\mu,\mu}^{\dagger} \delta_{x,y+\mu} \right]$$

- 其他费米子作用量定义有：Staggered 费米子、Domain-wall 费米子、Overlap 费米子等

格点 QCD 计算的基本框架：蒙特卡洛法

- 1、随机生成大量规范场组态 $\{U_{\mu}(x)\}$ ，使得出现概率：

$$P(U) \propto e^{-S_g(U)} \det M(U)$$

- 2、计算可观测量：

$$\langle O \rangle = \frac{1}{N} \sum_{\{U\}} O(U, G_q(U))$$

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

误差预算

- 统计误差：随组态个数的平方根压低
 - 可通过增加统计量而稳步降低
- 系统误差：离散误差、有限体积效应
 - 可对格距和体积进行外推，得到连续的无限体积时空的结果
- 系统误差：参数偏差
 - 示例：夸克质量参数设置不精确，或故意设置较重的质量
 - 可通过多参数内插或外推消除偏差
- 系统误差：一些计算细节导致的其他误差或偏差
 - 示例：对关联函数进行拟合时选取拟合范围的模糊因素
 - 可通过改进算法消除或压低

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

误差预算

- 统计误差：随组态个数的平方根压低
 - 可通过增加统计量而稳步降低
- 系统误差：离散误差、有限体积效应
 - 可对格距和体积进行外推，得到连续的无限体积时空的结果
- 系统误差：参数偏差
 - 示例：夸克质量参数设置不精确，或故意设置较重的质量
 - 可通过多参数内插或外推消除偏差
- 系统误差：一些计算细节导致的其他误差或偏差
 - 示例：对关联函数进行拟合时选取拟合范围的模糊因素
 - 可通过改进算法消除或压低

所有误差，归根结底，都交由摩尔定律解决！

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

格点 QCD 的研究方向

- 强子谱学
 - 传统强子的高精度计算、奇特强子态、用散射方法研究共振态、……
- 强子结构
 - 各种形状因子、部分子分布函数、……
- 极端条件下的 QCD
 - 高温/高密/强磁场、QCD 临界点、手征相变、强磁场下的相结构、……
- 味物理和超出标准模型的新物理
 - CKM 矩阵元、中微子相关的强子矩阵元、……
- ……

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

高能物理的三个传统前沿

- 高能前沿
 - LHC(ATLAS、CMS、LHCb)、CEPC 等
- 高亮度前沿
 - BEPCII(BESIII)、Belle II 等
 - 中微子实验 (大亚湾、江门、TUNE 等)
- 宇宙学前沿
 - 宇宙线观测 (LHAASO 等)

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

高能物理的三个传统前沿

- 高能前沿
 - LHC(ATLAS、CMS、LHCb)、CEPC 等
- 高亮度前沿
 - BEPCII(BESIII)、Belle II 等
 - 中微子实验 (大亚湾、江门、TUNE 等)
- 宇宙学前沿
 - 宇宙线观测 (LHAASO 等)

第四个前沿

- 高精度前沿
 - 经过四十年的发展，从解释现象到预言现象
 - 很多探索新物理的瓶颈在强相互作用的高精度计算上

格点量子色动力学 (Lattice QCD)

第一性原理

非微扰

系统研究

强相互作用

高能物理前沿

领域

世界观

对量子场论的非微扰定义

方法论

可操作的计算模型

领域交叉

- 不同的研究方法：
 - 有效场论、粒子物理唯象学、少体核物理、……
- 拓展研究方法：
 - 量子计算、……
- 研究方法本身：计算机科学
 - 高性能计算、数值计算方法
 - 软件设计、计算机语言、计算机体系结构

国内外格点 QCD 计算软件现状

USQCD SciDAC 软件集 (Chroma、QUADA、QDP++、……)

- + 国际普遍使用
- + 模块化, 支持 intel、nvidia 等架构
- 原版不支持国产超算
- 单任务, 中小规模运行

Bridge++、Grid、tmLQCD、openQCD、……

- + 代码量较小, 易于扩展和修改
- + 适应各自特定需求
- 功能单薄
- 跨平台性较差
- 不支持国产超算

各研究组的其他软件等

- + 实现通用软件不支持的功能
- 规模小、功能单一、大多不跨平台
- 大多不开源

国产新一代软硬件环境的挑战

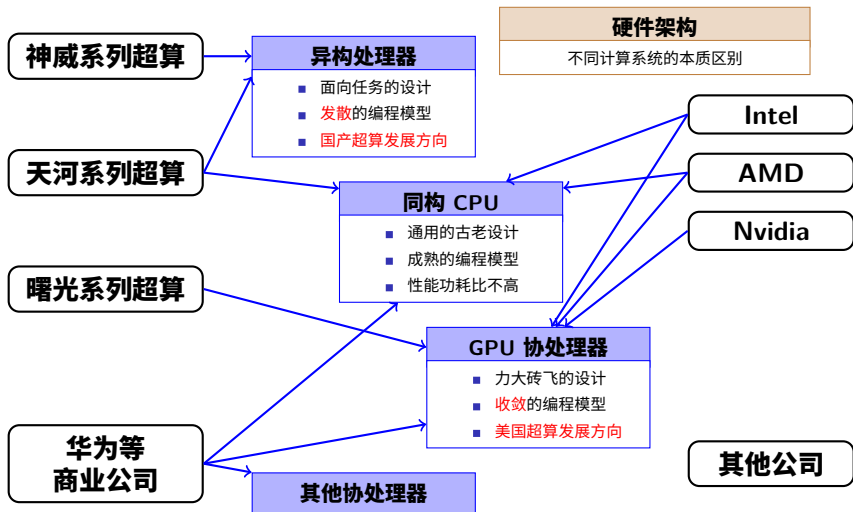
硬件的多样性

- 神威、天河、曙光三个架构都非常不同
 - 同系列下一代超算架构变化也非常大
 - 未来超算仍处于快速演化过程中
- + 所以软件设计必须挖掘共性特点、面向未来趋势。

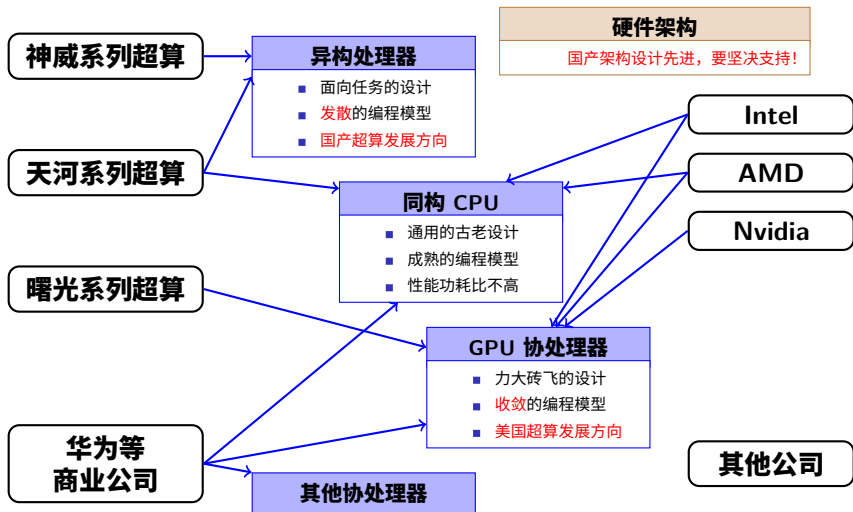
基础软件刚刚起步

- 新天河系统、神威系统仍不完善，急需大量基础软件支持。
 - 曙光系统有开源社区支持，但也缺少很多基础软件。
- + 所以软件研发需要从**基础软件**做起，但同时也抛弃了历史包袱：一张白纸绘新图。

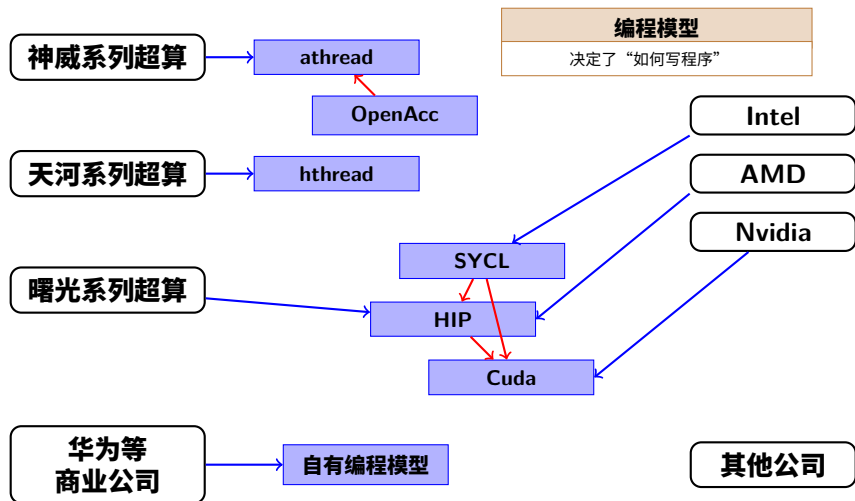
超级计算机的主要处理器平台



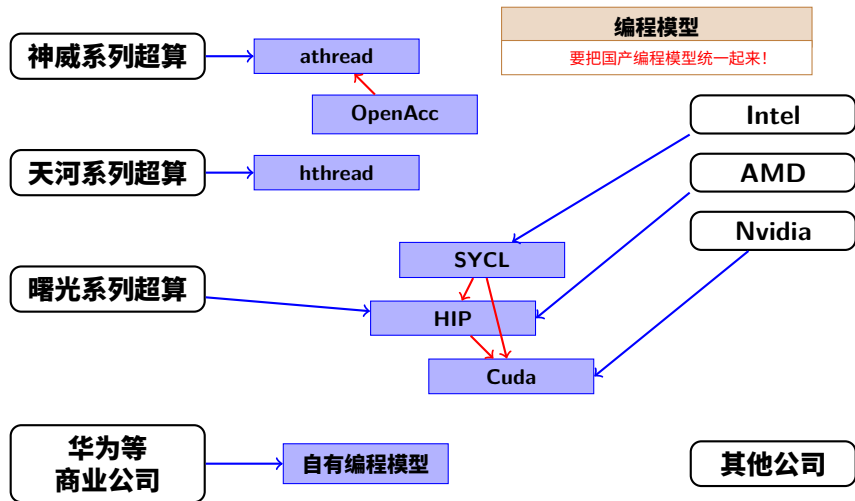
超级计算机的主要处理器平台



超级计算机的主要处理器平台



超级计算机的主要处理器平台



并行编程模型回顾

常用并行编程模型

- 格点 QCD 定制并行模型：QDP、Grid
- 进程并行：MPI
- 核心函数线程并行：CUDA、HIP、athread
- 匿名函数线程并行：SYCL、KOKKOS、RAJA
- 代码段线程并行：OpenMP、OpenAcc
- 指令并行：SIMD

- 以上编程模型在结构上分为两大类：对称、主从。
- 格点 QCD 软件通常混合采用对称和主从模型实现二级并行，代码结构很复杂。
- 以上模型中，除了 KOKKOS 外，大多数模型没有系统地处理多层次异构内存。
- 需要根据彻底的异构前提，设计**新的编程模型**。

未来科学计算技术的四个发展方向

自动代码生成

- 这是当前热点：元编程技术
- 路径：公式 → 中间表示 → 代码
- 是其他三个发展方向的基础设施
 - 参数化生成、面向平台生成、低代码界面

自动代码优化

- 人工智能相关技术正在跃进
 - 科学计算软件要准备好借东风
- 优化理论和智能算法有发展潜力
- 需要完整的支撑平台来测试和部署

跨平台部署和调度

- 大数据框架是关键基础设施
 - Spark、Map-Reduce、Pregel、……
- 传统科学计算软件在技术上已落后
- 急需“高性能 + 可移植”的两全设计

适合科研用户的人机界面

- “物理的归物理，计算的归计算”
 - 这是目前科学计算的最大痛点
- 面向物理公式的编程模型创新
- 面向科研工作的敏捷开发平台

格点 QCD 在四个方向上的需求和解决方案

MetaTensor: 自动代码生成

- 物理计算基于复杂的公式
- 所有公式都可以转写为广义的“张量表达式”

张量表达式计算

DDQ: 高效跨平台调度器

- 格点 QCD 的算法结构主要是循环迭代
- 需要大规模并行、需要高效数据交换

循环迭代

拆解算法流程

格点 QCD 计算软件

拆解计算程序

纯计算部分

- 与物理结果无关，即为优化空间
- 参数化此空间，开展自动优化研究

纯物理部分

- 隔离计算细节，面向物理公式
- 适合做科研用户界面的编程模型

Trene: 通用自动优化示例

x 语言: 图形化编程语言

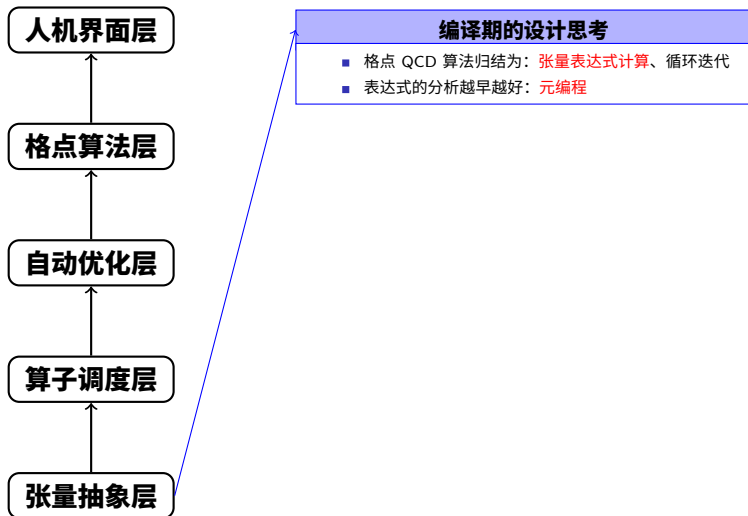
研发贡献者

- 中国科学院高能物理研究所
宫明、陈莹、刘朝峰、毕玉江、孙玮、施春江、蒋翔宇
- 北京航空航天大学
栾钟治、韩斌、王御臣、肖敏毅、房歆哲、龚煜涵、马世清、刘星宇、李根、王逸杰、李亦白
- 中国科学院计算机网络信息中心
徐顺、张克龙、韩秉豫、张术飞
- 中国科学院理论物理研究所
王建成

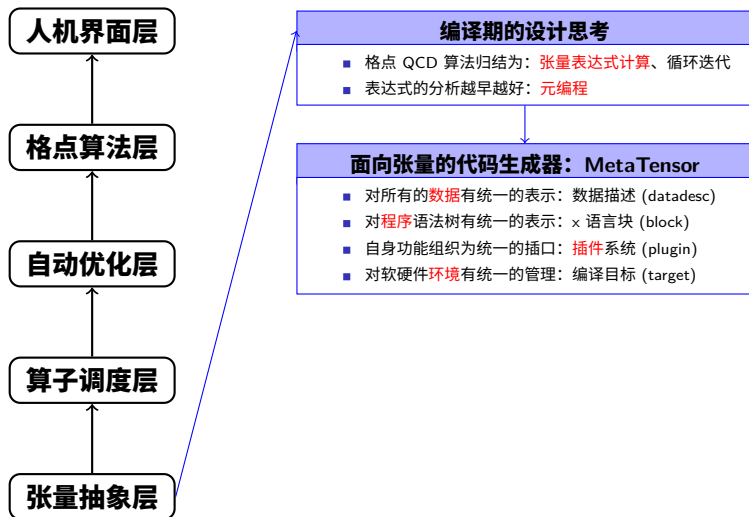
格点 QCD 新软件框架



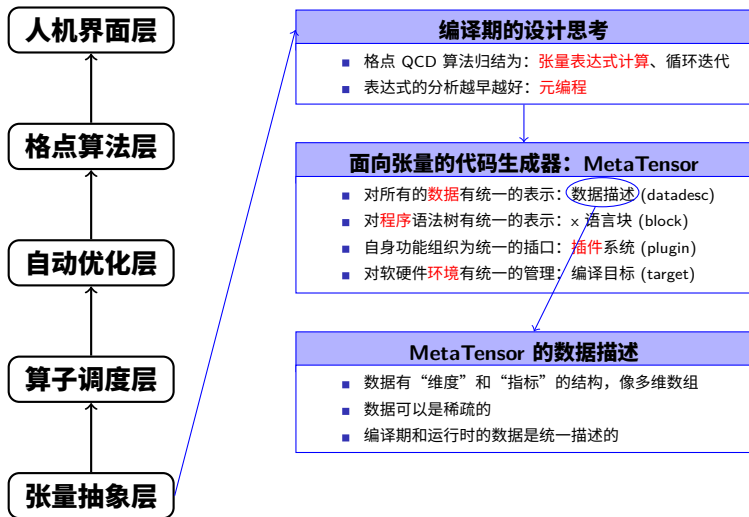
格点 QCD 新软件框架



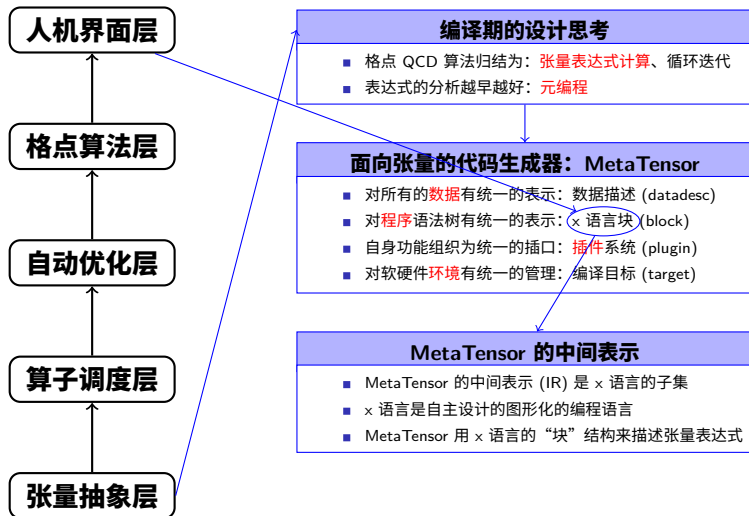
格点 QCD 新软件框架



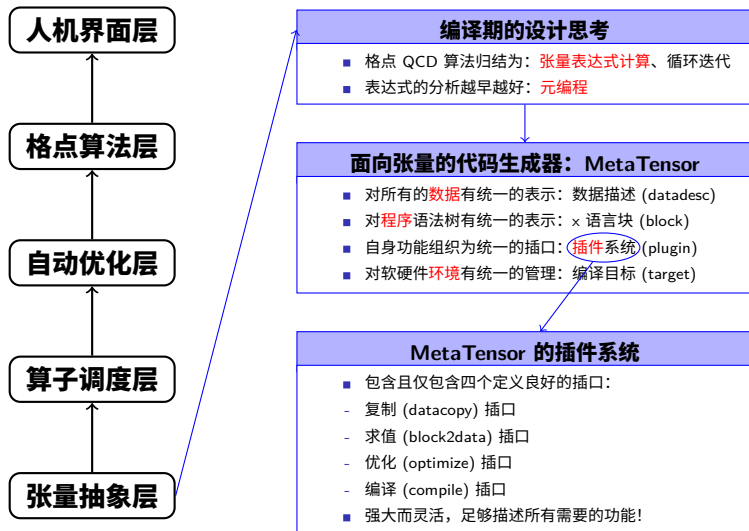
格点 QCD 新软件框架

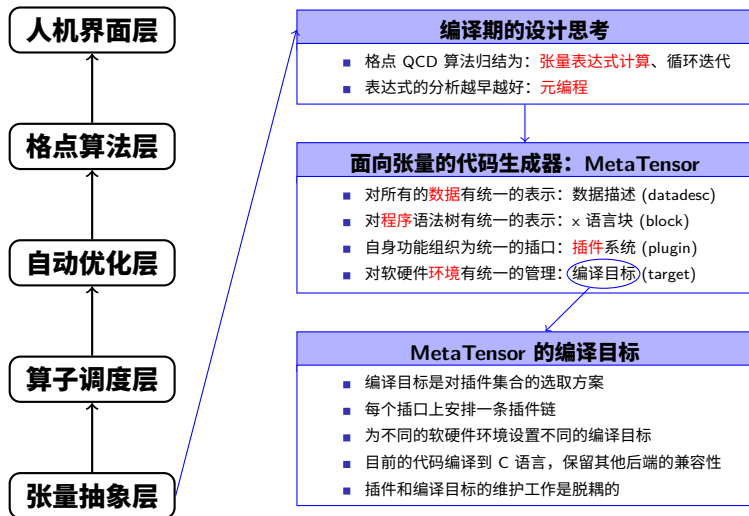


格点 QCD 新软件框架

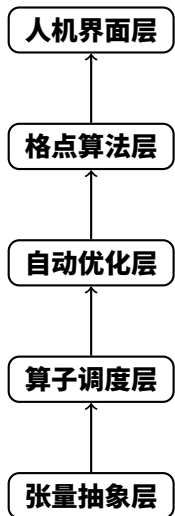


格点 QCD 新软件框架





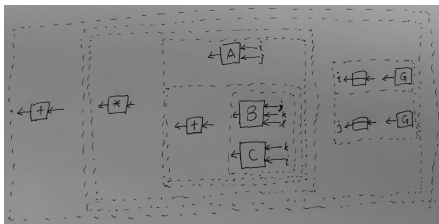
格点 QCD 新软件框架



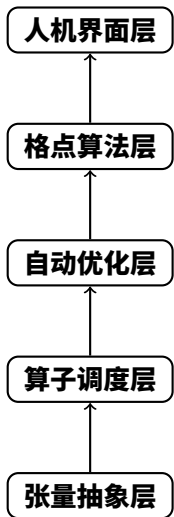
需要计算的张量表达式:

$$result_{k,l} = \sum_{i,j} A_{i,j} * (B_{j,k,l} + C_{k,i})$$

用 x 语言转写为:



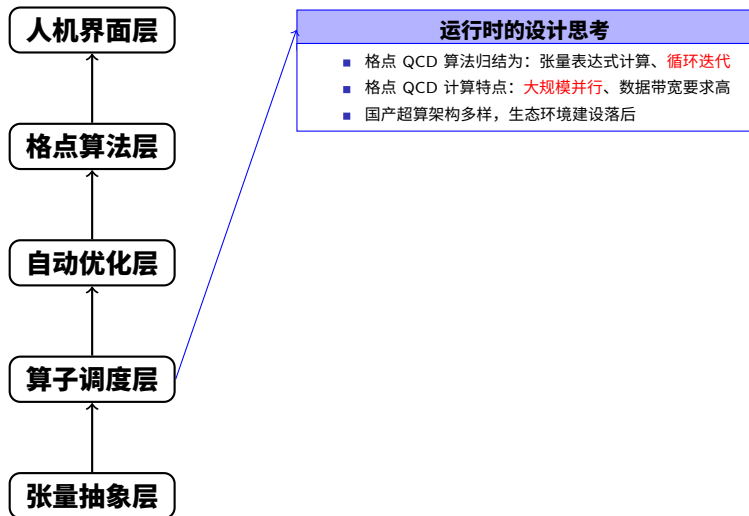
格点 QCD 新软件框架



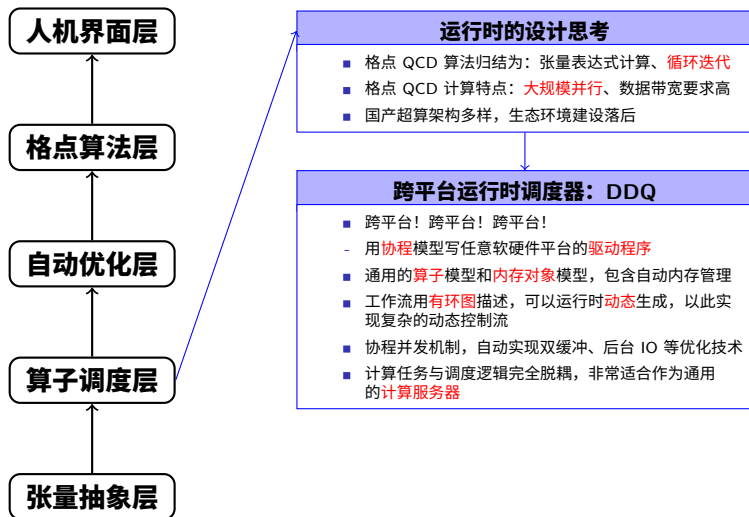
MetaTensor 把它编译到 C 语言:

```
#include <stdint.h>
#include <stdlib.h>
void calc(void *res, void *a, void *b, void *c){
    double*var_5 =malloc(2048);
    {
        double*var_2 =(double*)((char *) (b)+0);
        double*var_3 =(double*)((char *) (c)+0);
        double*var_4 =(double*)((char *) (var_5 )+0);
        int32_t var_7 =0;
        int32_t var_8 =0;
        int32_t var_9 =0;
        const int32_t var_10 [4]={ 0,16,32,40};
        const int32_t var_11 [4]={ 0,64,128,192};
        for(int32_t var_12 =0; var_12 <4; var_12 ++){
            int32_t var_13 =var_7 +var_10 [var_12 ];
            int32_t var_14 =var_9 +var_11 [var_12 ];
            const int32_t var_15 [4]={ 0,4,8,12};
            const int32_t var_16 [4]={ 0,4,8,12};
            const int32_t var_17 [4]={ 0,16,32,40};
            for(int32_t var_18 =0; var_18 <4; var_18 ++){
                int32_t var_19 =var_13 +var_15 [var_18 ];
                int32_t var_20 =var_8 +var_16 [var_18 ];
                int32_t var_21 =var_14 +var_17 [var_18 ];
                const int32_t var_22 [4]={ 0,1,2,3};
                const int32_t var_23 [4]={ 0,4,8,12};
                for(int32_t var_24 =0; var_24 <4; var_24 ++){
                    int32_t var_25 =var_19 +var_22 [var_24 ];
                    int32_t var_26 =var_21 +var_23 [var_24 ];
                    const int32_t var_27 [4]={ 0,1,2,3};
                    const int32_t var_28 [4]={ 0,1,2,3};
                    for(int32_t var_29 =0; var_29 <4; var_29 ++){
                        int32_t var_30 =var_20 +var_27 [var_29 ];
                        int32_t var_31 =var_26 +var_28 [var_29 ];
                        var_4 [var_31 ]+=var_2 [var_25 ]+var_3 [var_30 ];
                    }
                }
            }
        }
        double*var_37 =malloc(2048);
        {
            double*var_34 =(double*)((char *) (a)+0);
            double*var_35 =(double*)((char *) (var_5 )+0);
            double*var_36 =(double*)((char *) (var_37 )+0);
            int32_t var_39 =0;
            int32_t var_40 =0;
            int32_t var_41 =0;
            const int32_t var_42 [4]={ 0,4,8,12};
            const int32_t var_43 [4]={ 0,1,2,3};
            const int32_t var_44 [4]={ 0,64,128,192};
            const int32_t var_45 =0; var_45 <4; var_45 ++){
                for(int32_t var_46 =0; var_46 <4; var_46 ++){
                    int32_t var_47 =var_39 +var_42 [var_45 ];
                    int32_t var_48 =var_40 +var_43 [var_45 ];
```

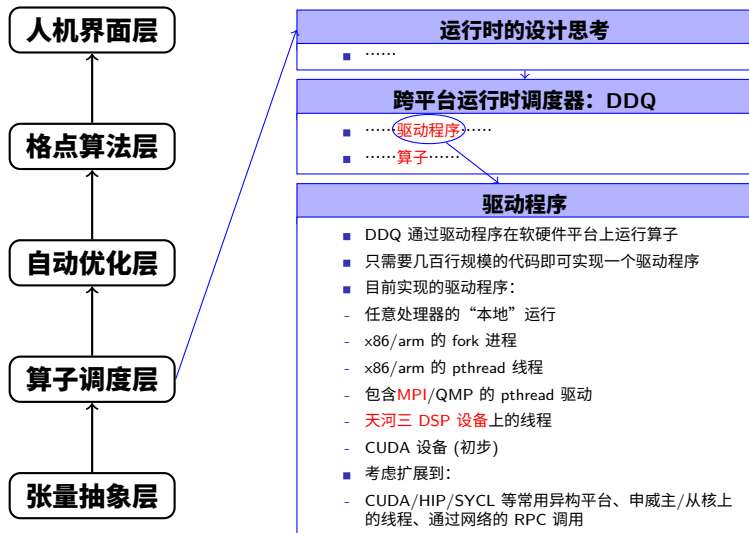
格点 QCD 新软件框架



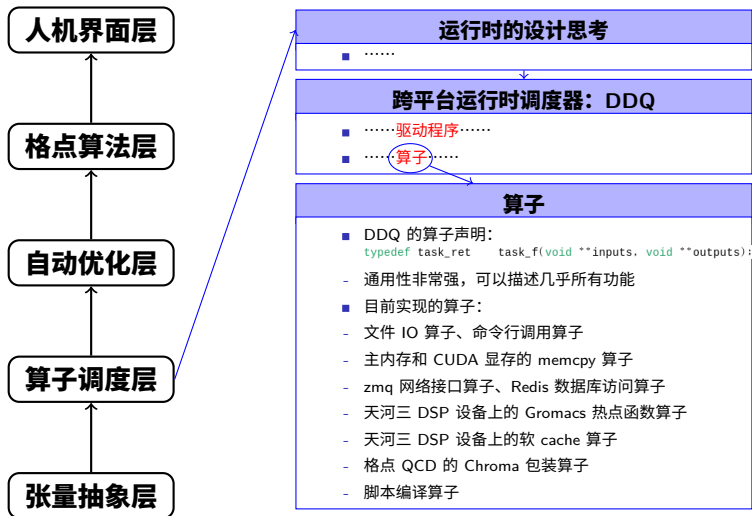
格点 QCD 新软件框架



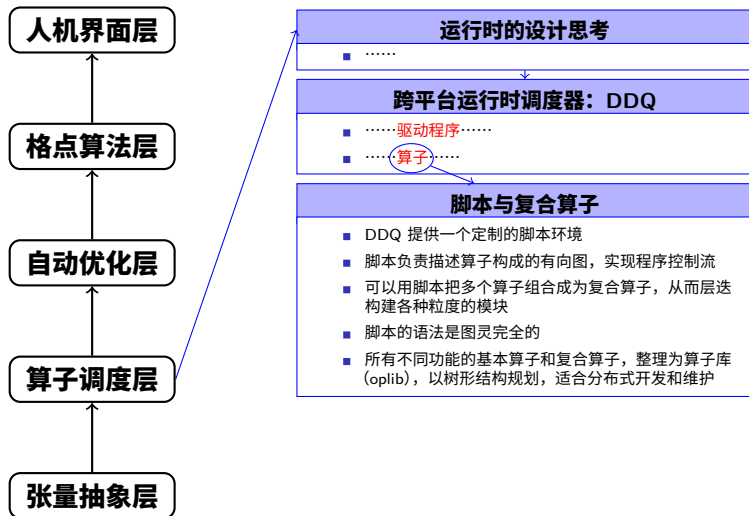
格点 QCD 新软件框架



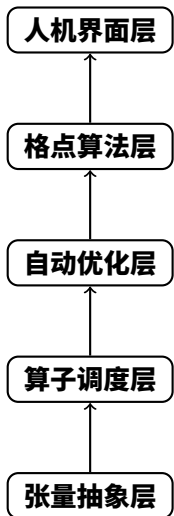
格点 QCD 新软件框架



格点 QCD 新软件框架



格点 QCD 新软件框架



举例：计算圆周率 π

$$\frac{1}{4}(\pi - 3) = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{2k(2k+1)(2k+2)} = \frac{1}{2 \cdot 3 \cdot 4} - \frac{1}{4 \cdot 5 \cdot 6} + \frac{1}{6 \cdot 7 \cdot 8} - \dots$$

用 DDQ 脚本转写为：

```
tests/pi.ddq : |
-----
print = load_op(std/real/print)
add := load_op(std/real/add)
mul4 := load_op(std/real/mul4)
mul := load_op(std/real/mul)
inv := load_op(std/real/inv)
init := load_op(std/real/init)
sleep := load_op(std/time/sleep)

print < [pi] < add < [pi term]
[term] < mul4 < [s ia ib ic]

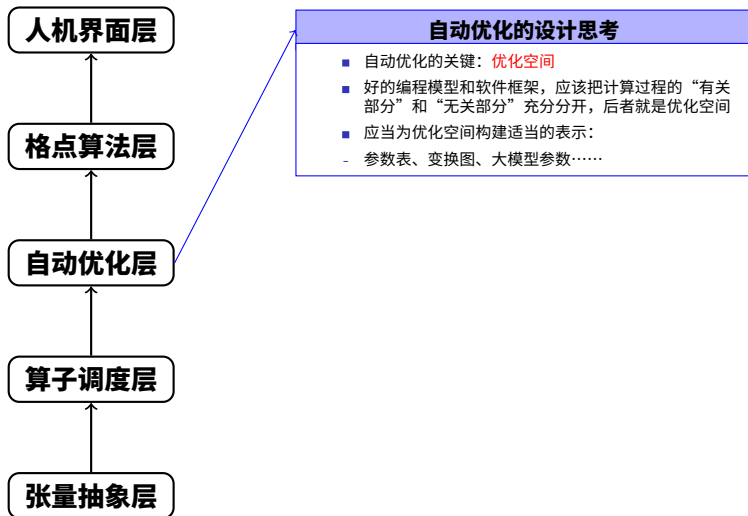
[s] < mul < [s -1.0]
[ia] < inv < [a] < add < [a 2.0]
[ib] < inv < [b] < add < [b 2.0]
[ic] < inv < [c] < add < [c 2.0]

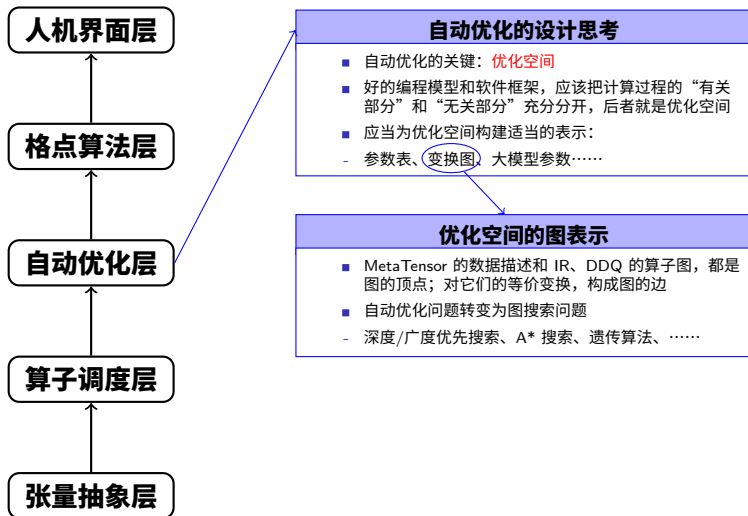
[pi] < init < [3.0]
[s] < init < [4.0]
[a] < init < [2.0]
[b] < init < [3.0]
[c] < init < [4.0]

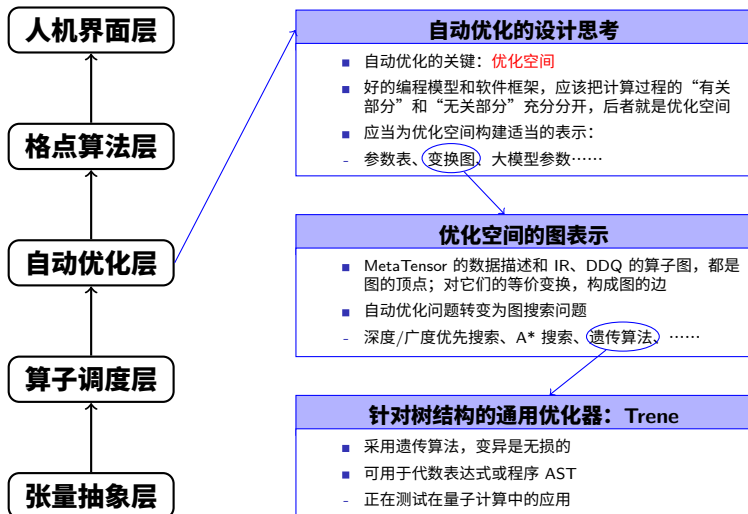
std/real/mul.ddq : |
-----
newop = load_so(std/real/mul)
newop.processor = "direct"
[load_type(std/real)] < newop < [load_type(std/real)]

std/real/mul4 : |
[s] < newop < [a b c d]
[s] < load_op(std/real/mul) < [abc d]
[abc] < load_op(std/real/mul) < [ab c]
[ab] < load_op(std/real/mul) < [a b]

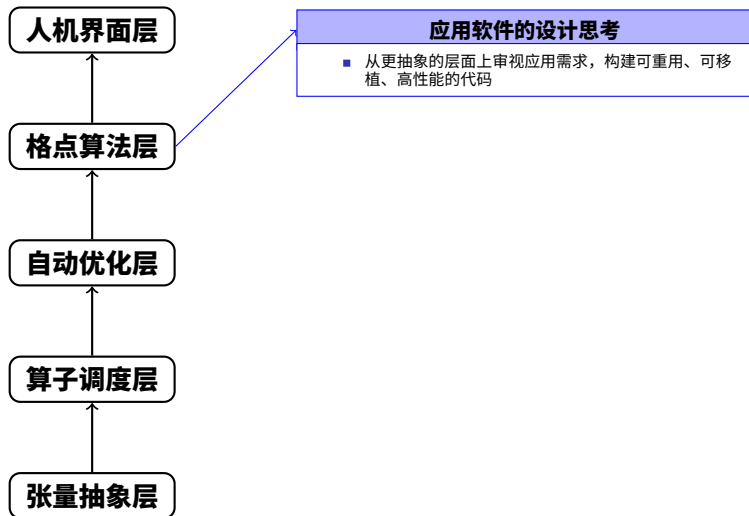
std/real/inv : |
[out] < newop < [in]
[out] < load_op(std/real/div) < [1.0 in]
```



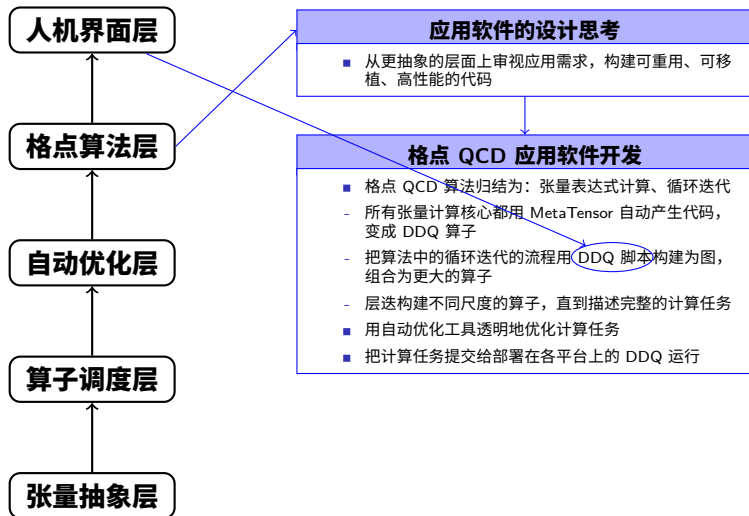




格点 QCD 新软件框架

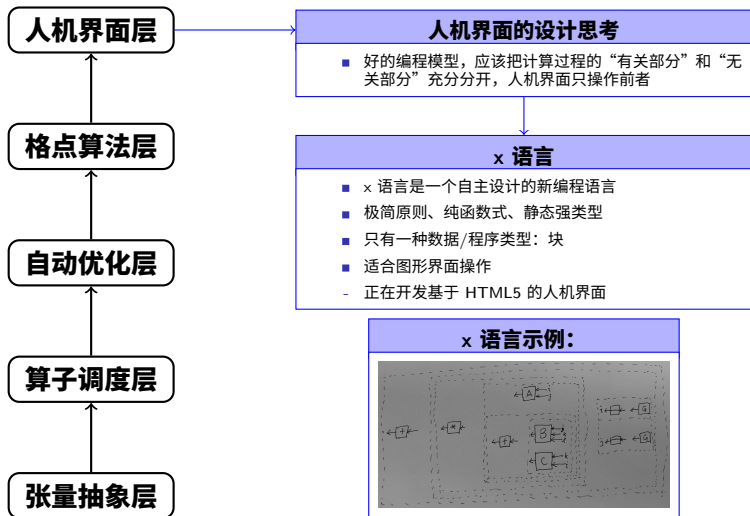


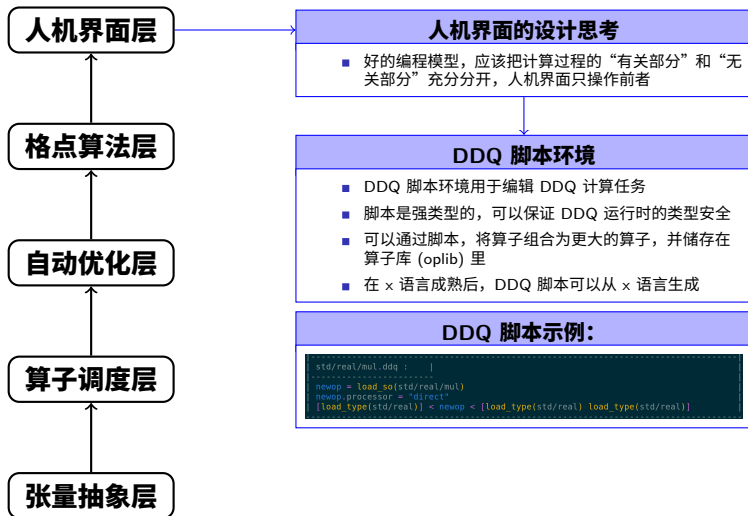
格点 QCD 新软件框架





格点 QCD 新软件框架





- 格点 QCD 是从第一性原理出发非微扰地系统研究强子和强相互作用的高能物理前沿领域
- 我们正在为格点 QCD 编写新的软件框架：
 - 它的研发目标是通用于各类科学计算的
 - 它的技术设计是面向未来发展的

谢谢!