

# Ascend C: 构建多级API，支撑多维场景算子开发

作者：黄金华

时间：2025/12/01

# 目录

**Part 1 Ascend C 介绍**

Part 2 Ascend C 编程模型&新增特性介绍

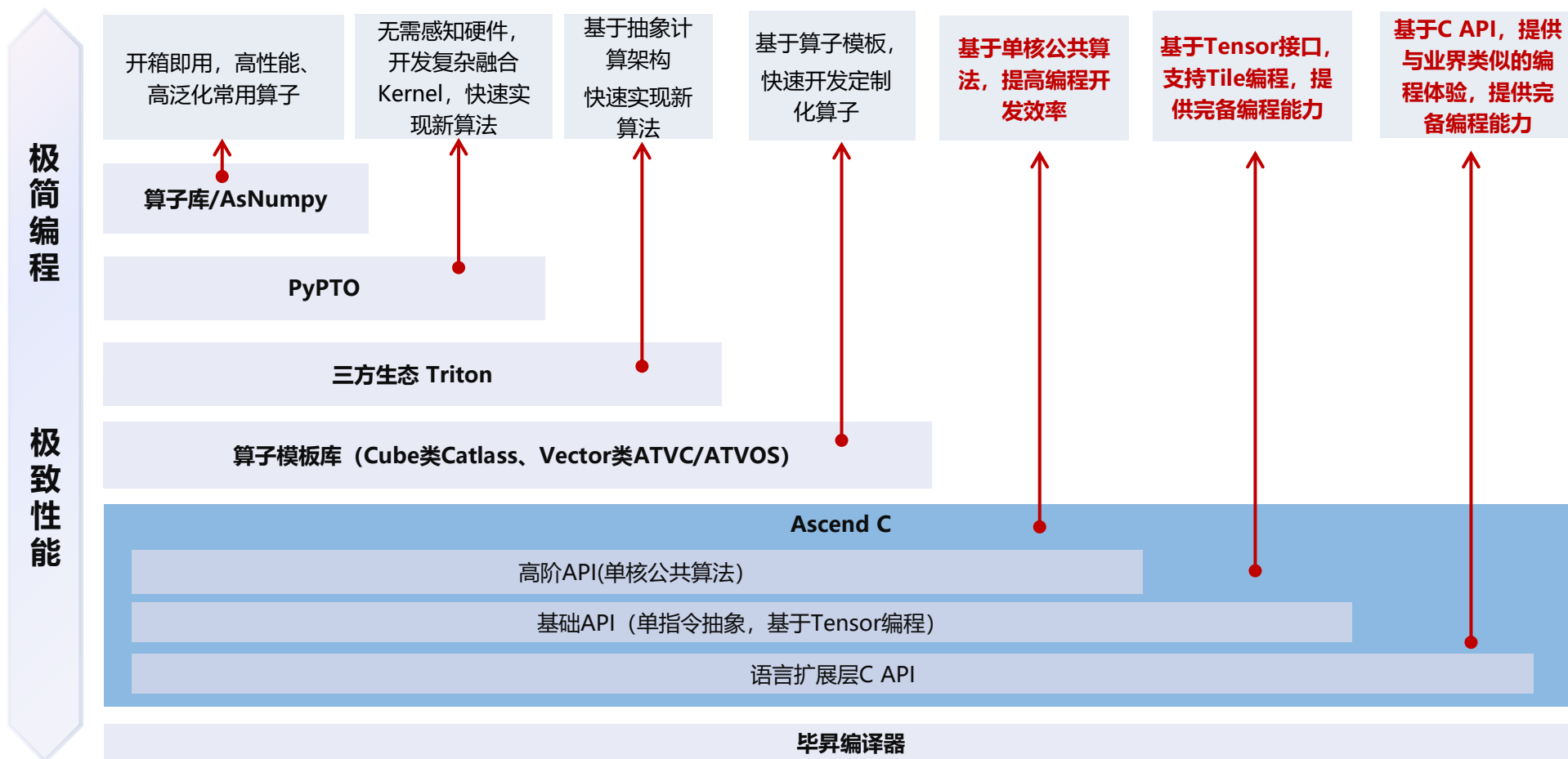
Part 3 Ascend C 下一代规划

# Ascend C: 能基于“手工”支撑实现极致性能

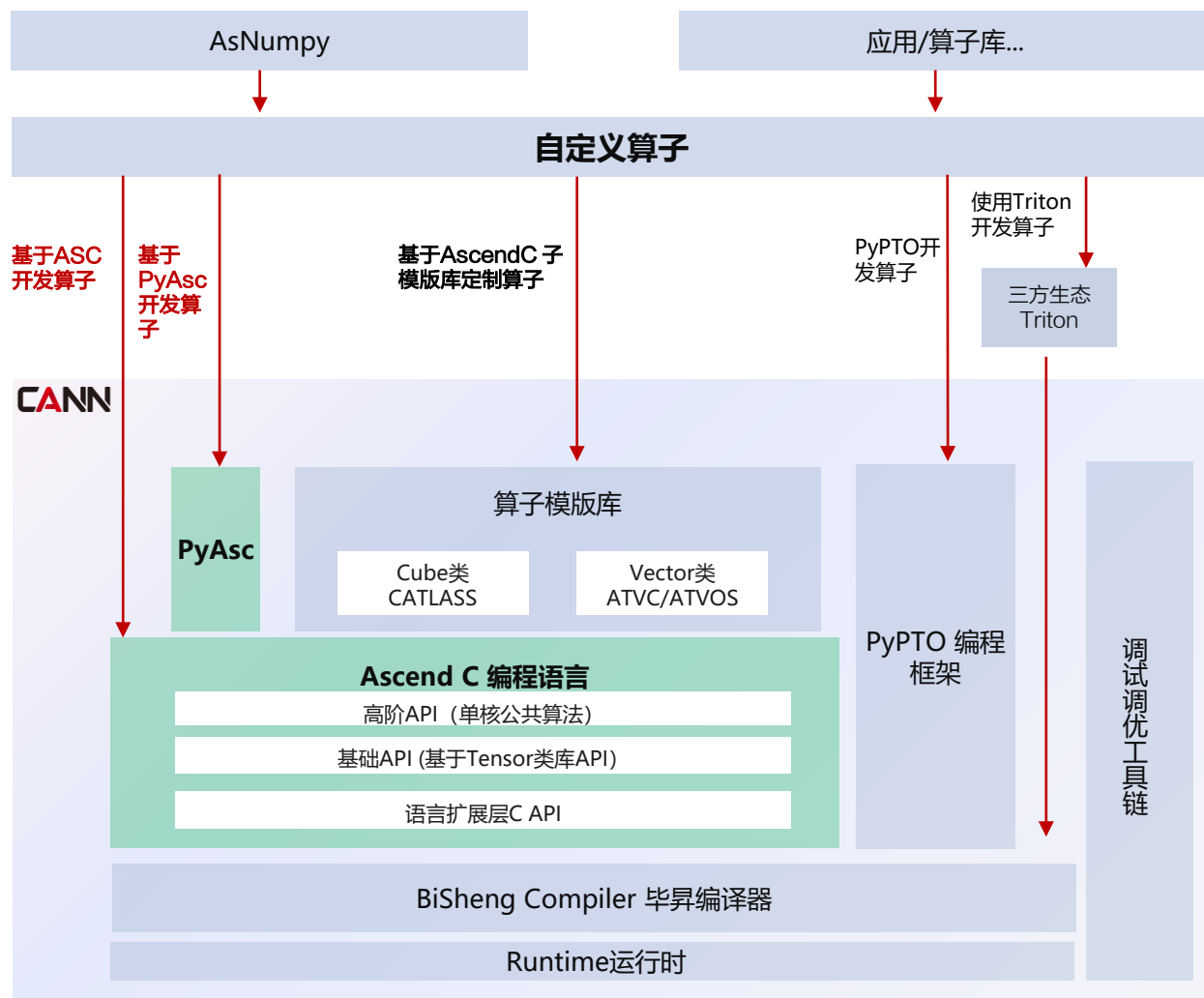
算法工  
程师



系统优化  
工程师



# Ascend C: 基于C/C++扩展的算子编程语言

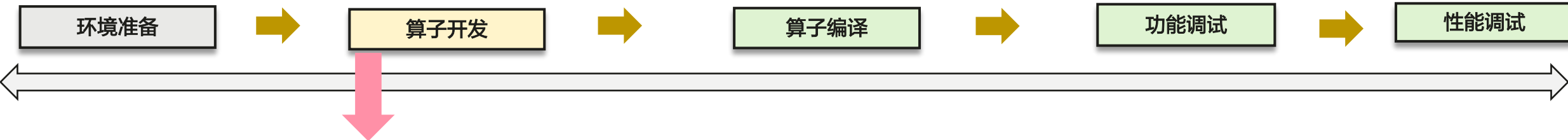


- 基于C/C++扩展，遵循C/C++标准规范；
- 提供底层芯片完备编程能力，支撑实现极致性能；
- 构建多级接口，满足多场景算子开发诉求，匹配业界开发习惯；
- 支持异构编程和<<<>>>直调
- 一套代码支持CPU/NPU孪生调试

CANN

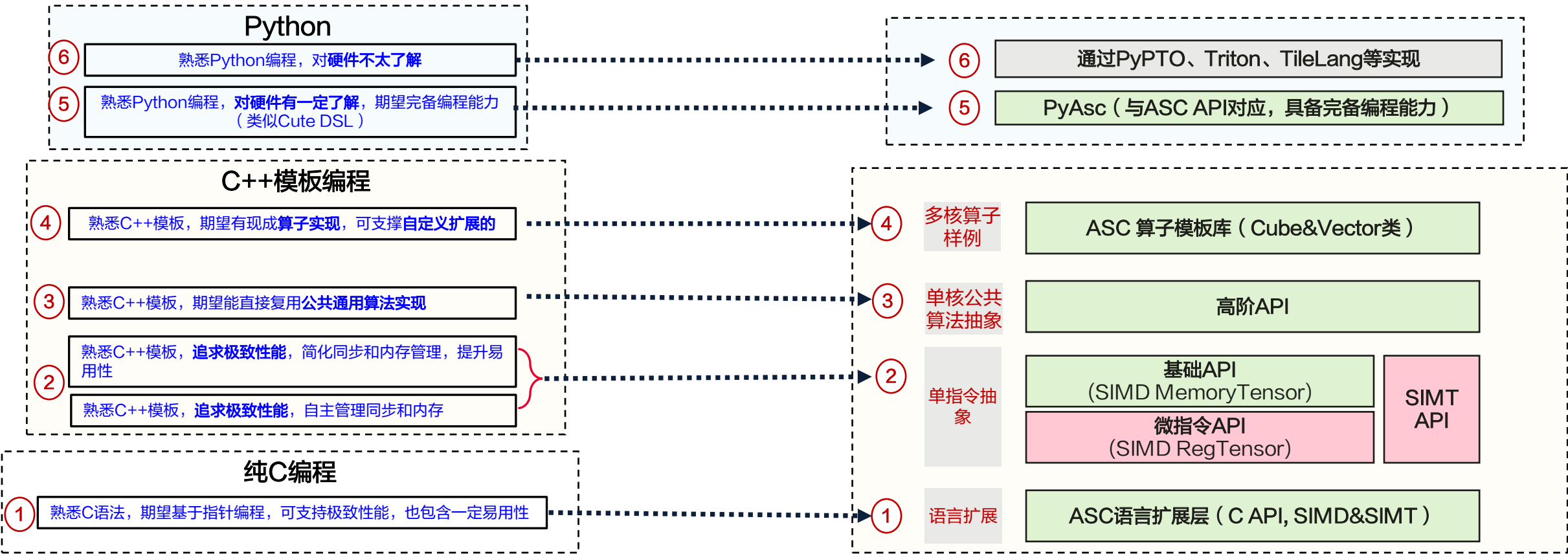
# Ascend C算子编程：构建多层级API，支撑多维场景算子开发诉求

典型算子开发流程

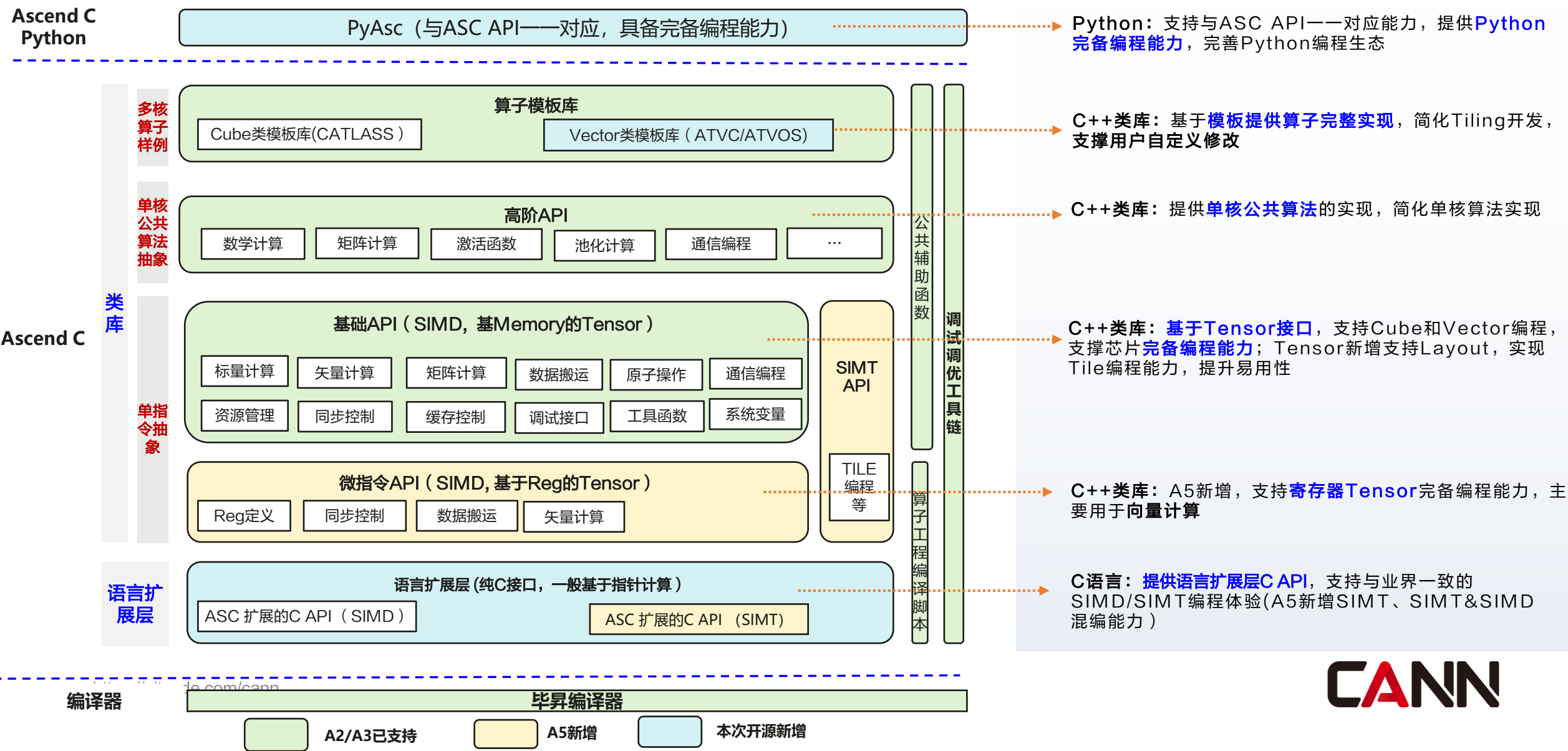


多维场景算子开发诉求

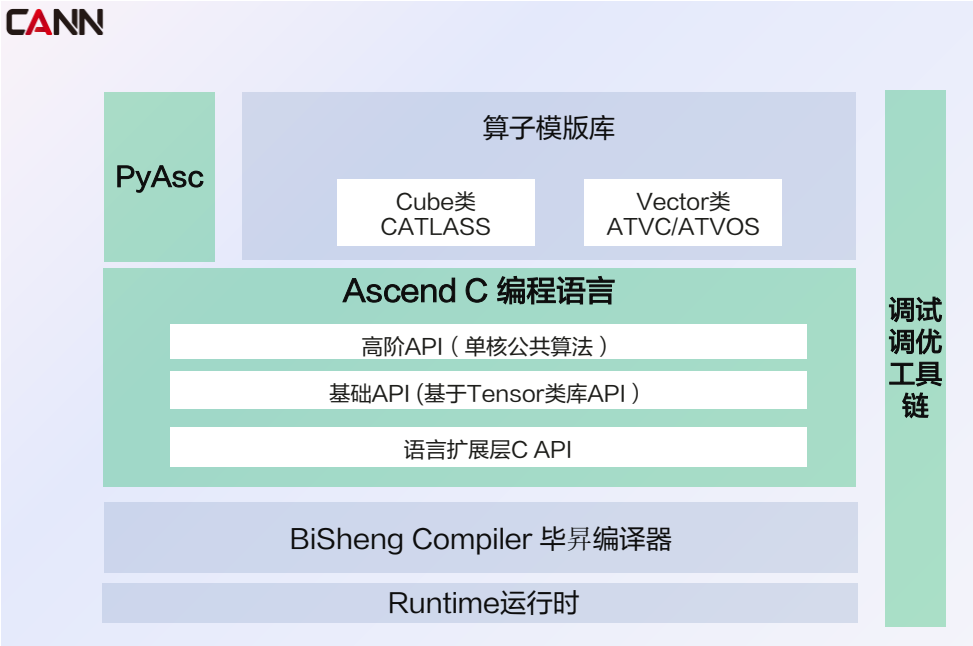
ASC 多层级API



## Ascend C算子编程架构：构建多层级API，支撑多维场景算子开发诉求



# Ascend C 分仓分包方案：支持按需安装，独立升级



No	仓	说明
1	CANN/asc-devkit	● ASC 算子开发的需要的 <b>最小集</b> （含API、编译脚本）；
2	CANN/asc-tools	● Ascend C 算子调试工具，包括CPU孪生调试等
3	CANN/pyasc	● Ascend C Python前端，与基础API/高阶API一一对应，完善Python编程生态；
4	CANN/atvos	● Vector类算子模板库
5	CANN/atvc	● Vector类Compute Only 模板库

# 目录

Part 1 Ascend C 介绍

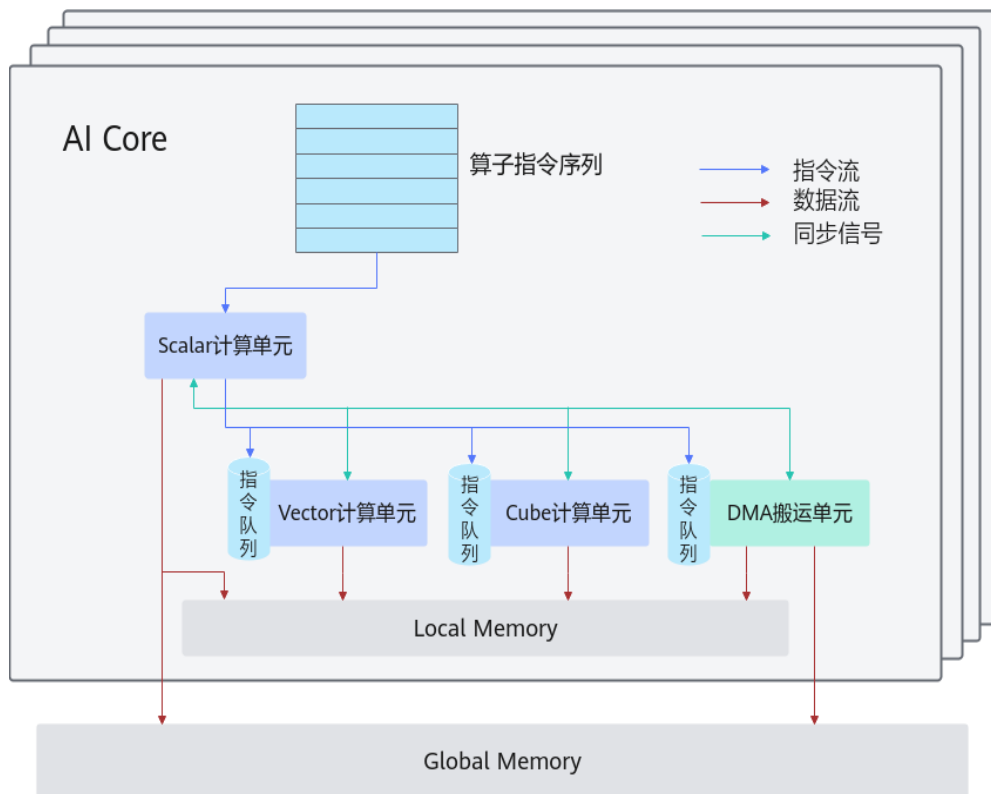
**Part 2 Ascend C 编程模型与新增特性介绍**

Part 3 Ascend C 下一代规划

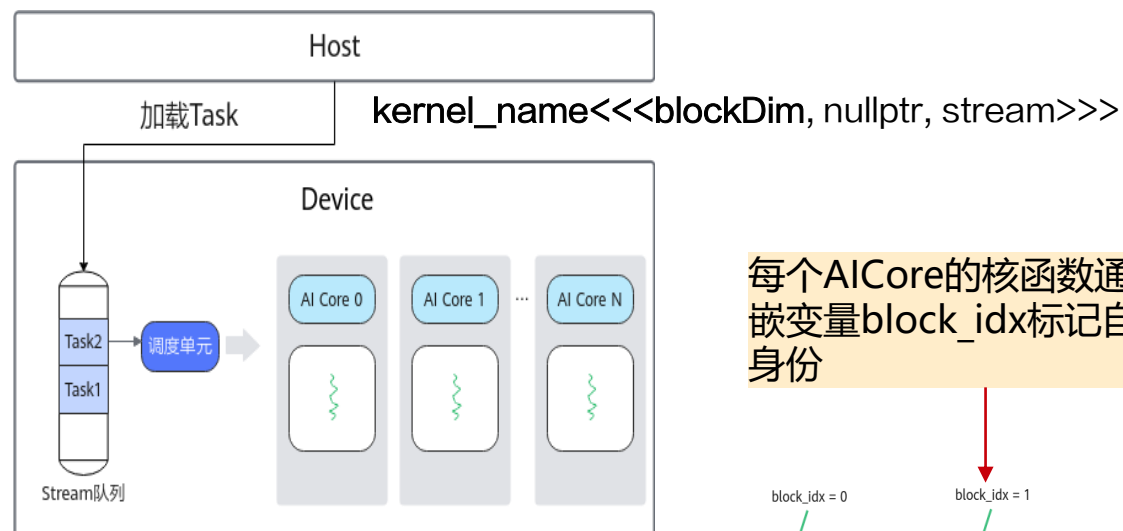


# Ascend C 编程模型简介

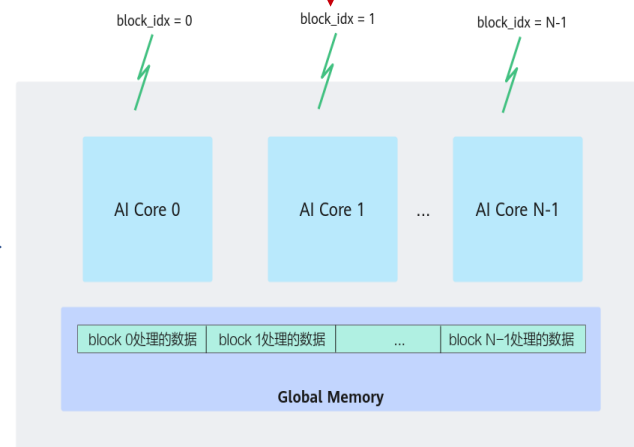
## NPU抽象硬件架构



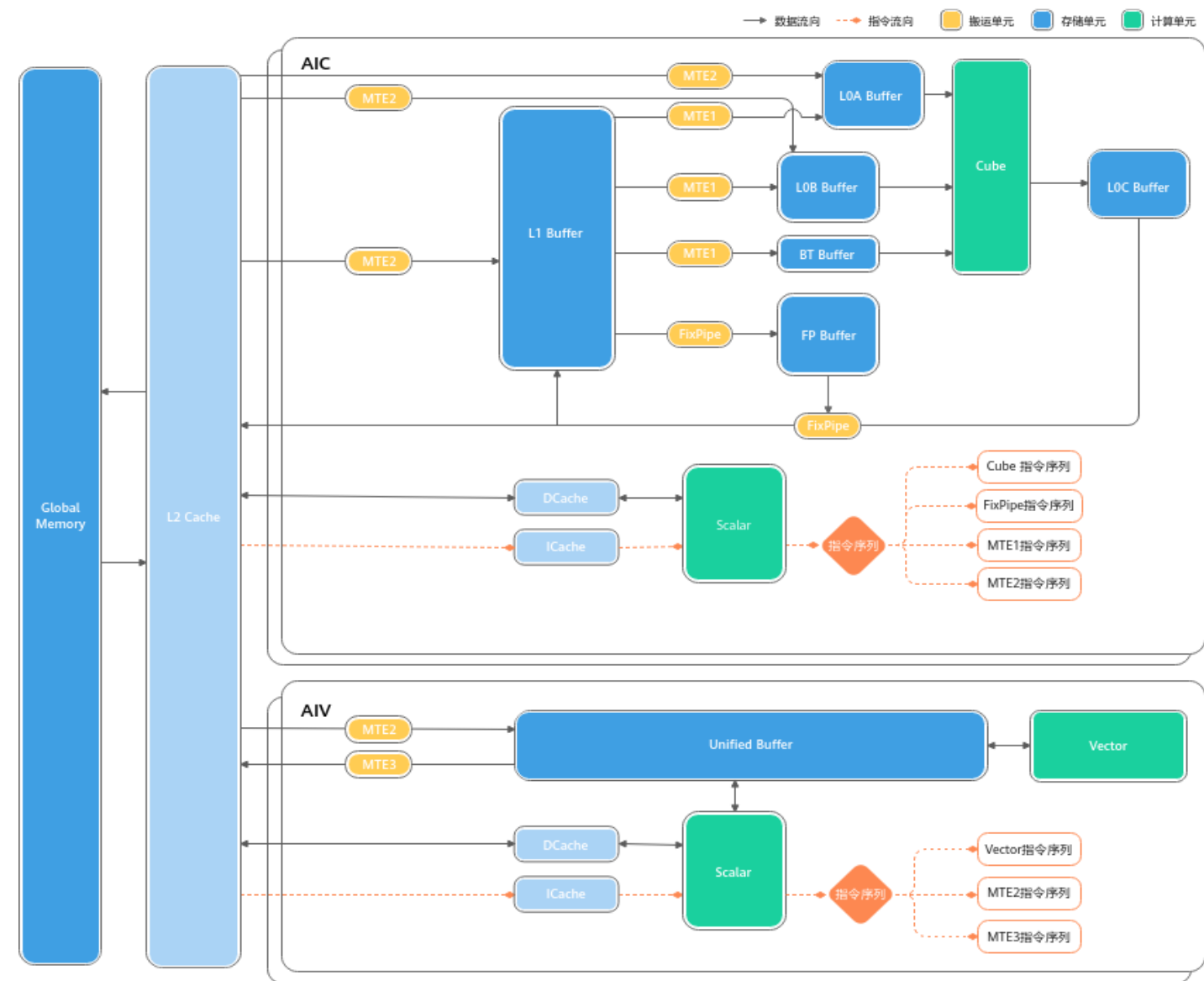
## Ascend C SPMD 编程模型 (SIMD)



每个AICore的核函数通过内嵌变量block\_idx标记自己的身份



# A2/A3 硬件架构介绍

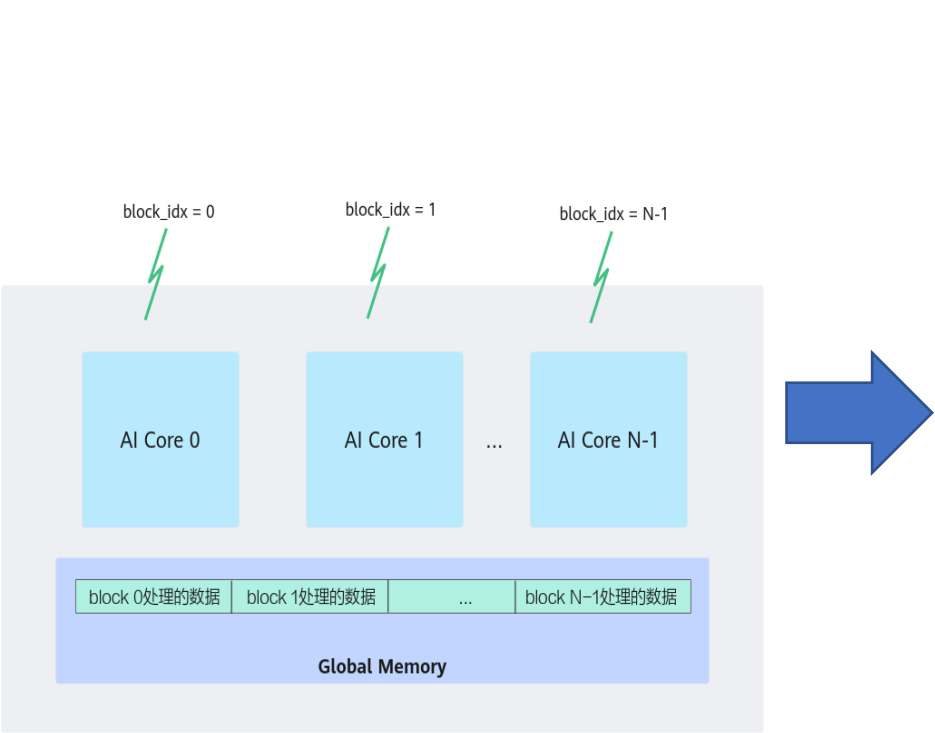


•**计算单元:** 包括Cube（矩阵）计算单元、Vector（矢量）计算单元和Scalar（标量）计算单元

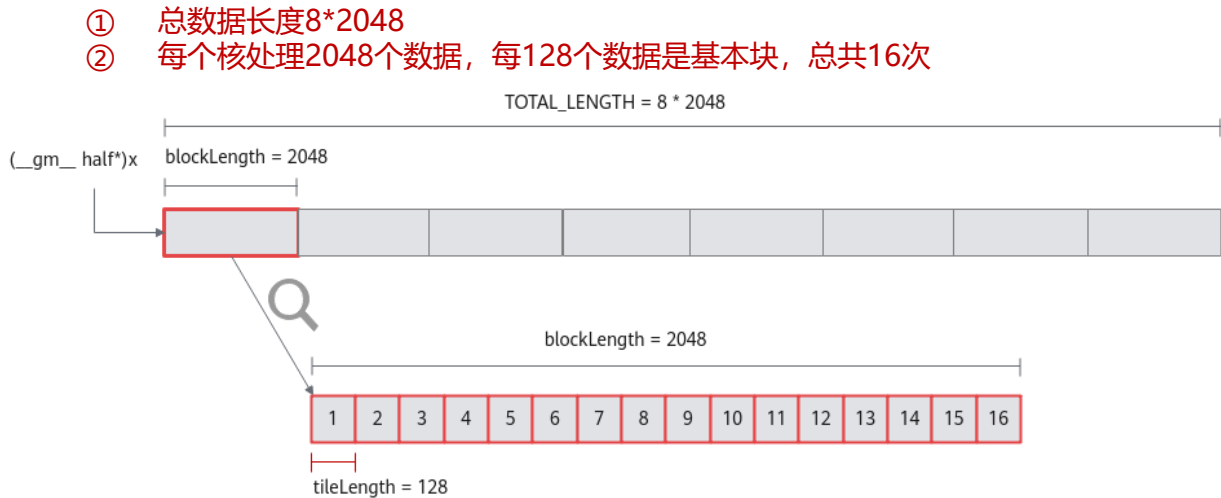
•**存储单元:** 包括L1 Buffer、L0A Buffer、L0B Buffer、L0C Buffer、Unified Buffer、BiasTable Buffer、Fixpipe Buffer等专为高效计算设计的存储单元。

•**搬运单元:** 包括MTE1、MTE2、MTE3和FixPipe，用于数据在不同存储单元之间的高效传输。

# Ascend C编程范式(SIMD): 以Add为例介绍数据切分策略

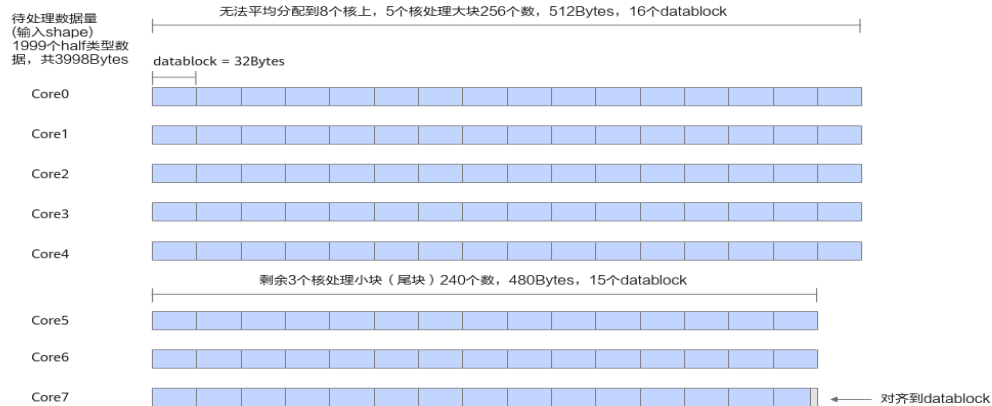


## ● 全对齐场景



- ① 尾块: 数据能平分到每个核, 但核内数据无法均分 (如上图: blockLength=2047)
- ② 尾核: 数据不能平分到每个核, 部分核处理数据较少;

## ● 未对齐场景



# Ascend C编程范式(1): 基于ASC 语言扩展层C API的Add算子示例

基于语言扩展层接口C api, 自主管理内存和同步, 提供业界通用编程体验

示例一: 基于同步搬运/计算接口

```
#define TILE_LENGTH_PER_CORE 2048

__global__ __aicore__ void add (__gm__ float* x,
__gm__ float* y, __gm__ float* z)
{
    KERNEL_TASK_TYPE(KERNEL_TYPE_AIV_ONLY);
    asc_init ();

    __ubuf__ float xIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float yIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float zOut[TILE_LENGTH_PER_CORE ];

    asc_copy_gm2ub_sync(xIn, x + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);
    asc_copy_gm2ub_sync(yIn, y + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);

    asc_add_sync(zOut, xIn, yIn,
TILE_LENGTH_PER_CORE );

    asc_copy_ub2gm_sync(z + block_idx *
TILE_LENGTH_PER_CORE , zOut,
TILE_LENGTH_PER_CORE );
}
```

示例二: 基于异步搬运/计算+统一同步接口

```
__global__ __aicore__ void add(__gm__ float* x,
__gm__ float* y, __gm__ float* z)
{
    KERNEL_TASK_TYPE(KERNEL_TYPE_AIV_ONLY);
    asc_init ();

    __ubuf__ float xIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float yIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float zOut[TILE_LENGTH_PER_CORE];

    asc_copy_gm2ub(xIn, x + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);
    asc_copy_gm2ub(yIn, y + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);
    asc_sync();

    asc_add(zOut, xIn, yIn, TILE_LENGTH_PER_CORE );
    asc_sync();

    asc_copy_ub2gm(z + block_idx *
TILE_LENGTH_PER_CORE , zOut,
TILE_LENGTH_PER_CORE );
    asc_sync();
}
```

示例三: 基于异步搬运/计算+ 精细化同步接口

```
__global__ __aicore__ void add(__gm__ float* x, __gm__
float* y, __gm__ float* z)
{
    KERNEL_TASK_TYPE(KERNEL_TYPE_AIV_ONLY);
    asc_init ();

    __ubuf__ float xIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float yIn[TILE_LENGTH_PER_CORE ];
    __ubuf__ float zOut[TILE_LENGTH_PER_CORE ];

    asc_sync_notify(PIPE_V, PIPE_MTE2, EVENT_ID0);
    for (int i = 0; i < 2; i++) {
        asc_sync_wait(PIPE_V, PIPE_MTE2, EVENT_ID0);
        asc_copy_gm2ub(xIn, x + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);
        asc_copy_gm2ub(yIn, y + block_idx *
TILE_LENGTH_PER_CORE, TILE_LENGTH_PER_CORE);

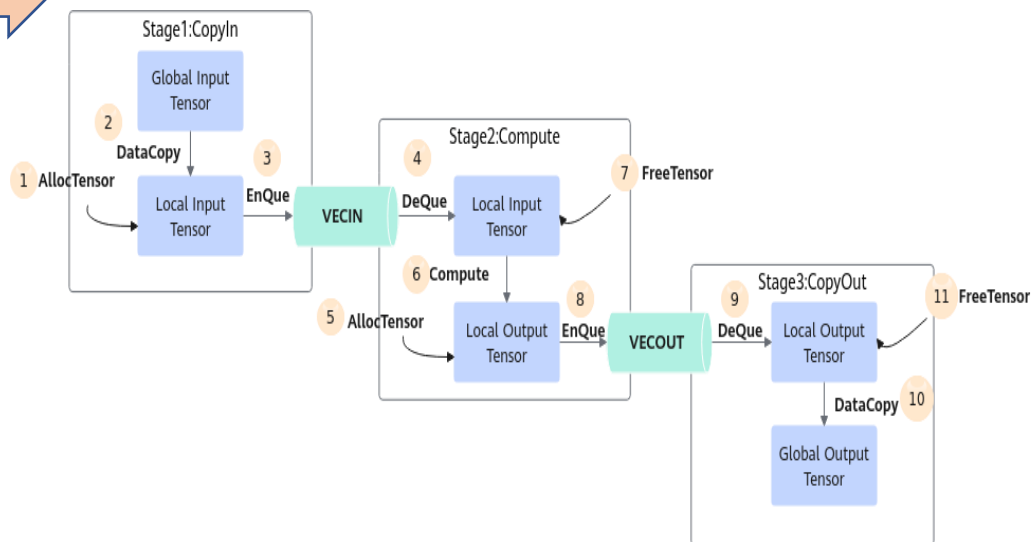
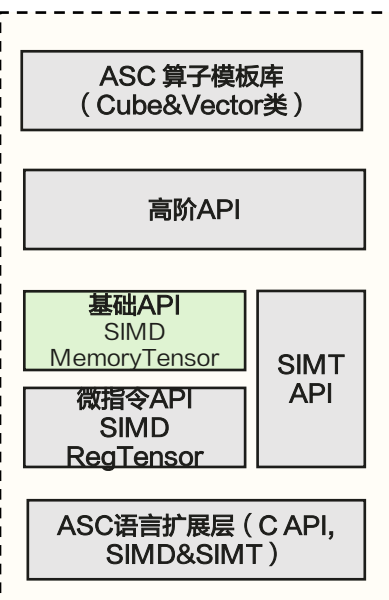
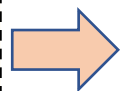
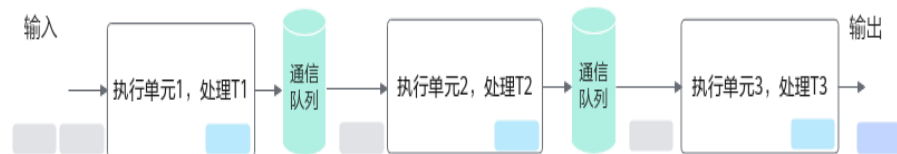
        asc_sync_notify(PIPE_MTE2, PIPE_V, EVENT_ID0);
        asc_sync_wait(PIPE_MTE2, PIPE_V, EVENT_ID0);

        asc_add(zOut, xIn, yIn, TILE_LENGTH_PER_CORE );
        asc_sync_notify(PIPE_V, PIPE_MTE3, EVENT_ID0);
        asc_sync_wait(PIPE_V, PIPE_MTE3, EVENT_ID0);

        asc_sync_notify(PIPE_3, PIPE_MTE2, EVENT_ID0);
        asc_copy_ub2gm(z + block_idx *
TILE_LENGTH_PER_CORE , zOut, TILE_LENGTH_PER_CORE );
    }
    asc_sync_wait(PIPE_V, PIPE_MTE2, EVENT_ID0);
}
```

# Ascend C编程范式(2): 基于基础API的Tque/Tpipe范式Add算子示例

基于Tque/Tpipe抽象硬件同步和内存管理, 简化编程



```
43  __aicore__ inline void Process()
44  {
45      int32_t loopCount = this->tileNum * BUFFER_NUM;
46      for (int32_t i = 0; i < loopCount; i++) {
47          CopyIn(i);
48          Compute(i);
49          CopyOut(i);
50      }
51  }
52
53  private:
54  __aicore__ inline void CopyIn(int32_t progress)
55  {
56      AscendC::LocalTensor<float> xLocal = inQueueX.AllocTensor<float>();
57      AscendC::LocalTensor<float> yLocal = inQueueY.AllocTensor<float>();
58      AscendC::DataCopy(xLocal, xGm[progress * this->tileLength], this->tileLength);
59      AscendC::DataCopy(yLocal, yGm[progress * this->tileLength], this->tileLength);
60      inQueueX.EnQue(xLocal);
61      inQueueY.EnQue(yLocal);
62  }
63  __aicore__ inline void Compute(int32_t progress)
64  {
65      AscendC::LocalTensor<float> xLocal = inQueueX.DeQue<float>();
66      AscendC::LocalTensor<float> yLocal = inQueueY.DeQue<float>();
67      AscendC::LocalTensor<float> zLocal = outQueueZ.AllocTensor<float>();
68      AscendC::Add(zLocal, xLocal, yLocal, this->tileLength);
69      outQueueZ.EnQue<float>(zLocal);
70      inQueueX.FreeTensor(xLocal);
71      inQueueY.FreeTensor(yLocal);
72  }
73  __aicore__ inline void CopyOut(int32_t progress)
74  {
75      AscendC::LocalTensor<float> zLocal = outQueueZ.DeQue<float>();
76      AscendC::DataCopy(zGm[progress * this->tileLength], zLocal, this->tileLength);
77      outQueueZ.FreeTensor(zLocal);
78  }
79  }
```

# Ascend C编程范式(3): 基于基础API的自主管理内存/同步的Add算子示例

基于LocalMemoryAllocator、Barrier 自主管理内存和同步, 提供更底层控制能力

## As-Is: TQue 范式编程

```
TPipe tpipe;  
TQue<TPosition::VECIN, 1> vecInQue;  
TQue<TPosition::VECOUT, 1> vecOutQue;
```

```
constexpr uint32_t tileLength = 1024;  
// Initialize buffer size  
tpipe.InitBuffer(vecInQue, 1,  
tileLength*sizeof(half));  
tpipe.InitBuffer(vecOutQue, 1,  
tileLength*sizeof(half));
```

```
for(int i=0; i < Loop; i++) {  
    // Allocate Tensor  
    auto in = vecInQue.AllocTensor<half>();  
    DataCopy(in, gm+ i*1024, 1024);  
  
    // Insert sync  
    vecInQue.Enqueue(in);  
    vecInQue.DeQueue(in);
```

```
    auto out = vecOutQue.AllocTensor<half>();
```

```
    Abs(out, in, 1024);
```

```
    vecInQue.FreeTensor(in);
```

```
    // Insert sync  
    vecOutQue.Enqueue(out);  
    vecOutQue.DeQueue(out);  
    DataCopy(gm+i*1024, out, 1024);  
    vecOutQue.FreeTensor(out);  
}
```

## To Be: 用户底层控制, 同步与内存分离

```
LocalMemAllocator<MemType::UB> allocator;
```

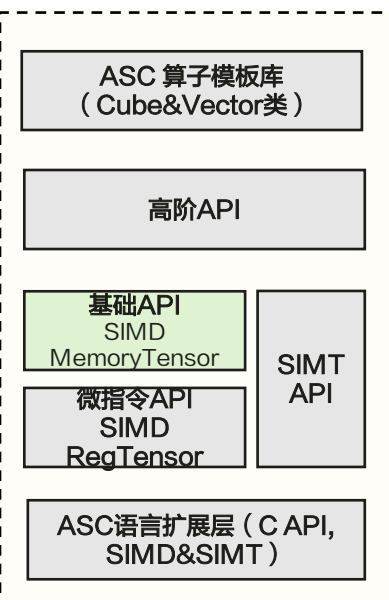
```
constexpr uint32_t tileLength = 1024;  
auto inLocalTensor = allocator.Alloc<half>(tileLength);  
auto outLocalTensor =  
allocator.Alloc<half>(tileLength);
```

```
Barrier<SyncScope::INTRA_CORE, mte3, mte2>  
firstBarrier;  
Barrier<SyncScope::INTRA_CORE, mte2, pipe_v>  
secBarrier;  
Barrier<SyncScope::INTRA_CORE, pipe_v, mte3>  
thirdBarrier;
```

```
firstBarrier.Arrive(0);
```

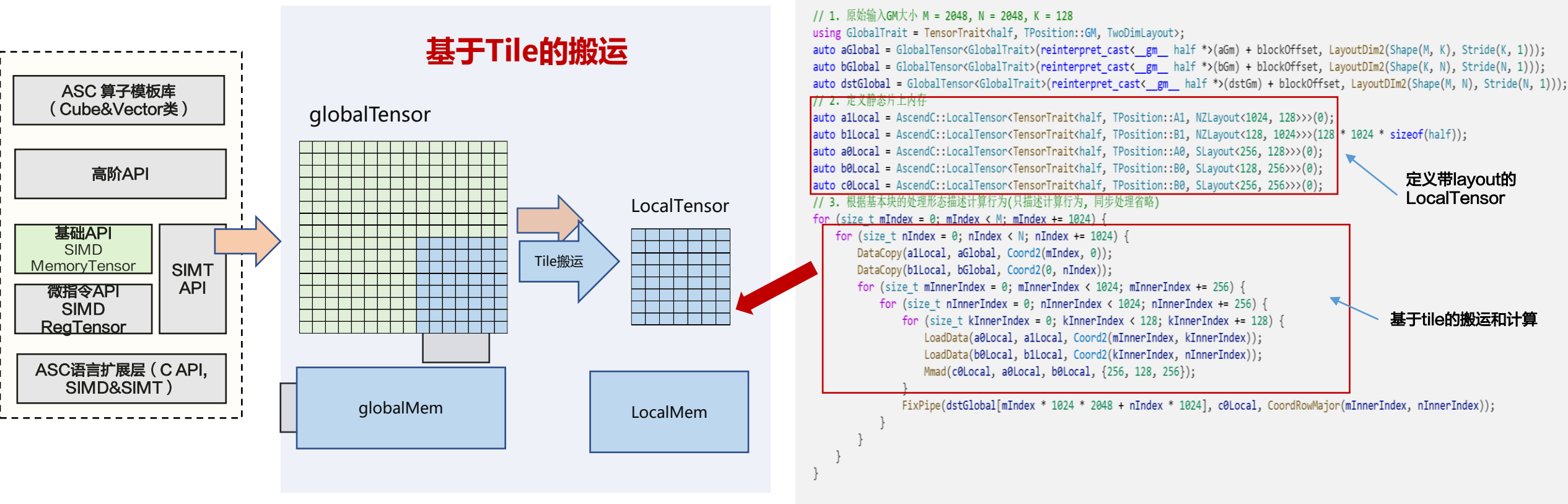
```
for(int i=0; i<loop; i++){  
    firstBarrier.wait(0);  
    DataCopy(in_ping, gm+ i *1024, 1024);  
    secBarrier.ArriveAndWait(0);  
    Abs(out_ping, in_ping, 1024);  
    thirdBarrier.ArriveAndWait(0);  
    DataCopy(gm+i*1024, out_ping, 1024);  
    firstBarrier.Arrive(0);
```

```
}  
firstBarrier.Wait(0);
```



# Ascend C编程范式(4): 基于基础API Tensor Tile API编程

Tensor增加对Layout的支持, 可提升编程的灵活性和可维护性, 同时保持高性能; 增强CUBE及VECTOR模板元编程能力;



# ASC支持异构编译+<<<>>>直调，提供与业界一致编程体验

- 对标业界，支持异构融合编译和<<<>>>直调；
- 通过后缀名“.asc”或编译选项“-x asc”使能异构编译；

```
#include "kernel_operator.h"
#include "acl/acl.h"
__global__ __aicore__ void
hello_world()
{
    AscendC::printf("Hello World!!!\n");
}
int32_t main(int argc, char const
*argv[])
{
    aclInit(nullptr);
    int32_t deviceId = 0;
    aclrtSetDevice(deviceId);
    aclrtStream stream = nullptr;
    aclrtCreateStream(&stream);
    constexpr uint32_t blockDim = 8;
    hello_world<<<blockDim, nullptr,
stream>>>());
    aclrtSynchronizeStream(stream);
    aclrtDestroyStream(stream);
    aclrtResetDevice(deviceId);
    aclFinalize();
    return 0;
}
```

Device侧代码



Host/Device编译归一



Host侧代码

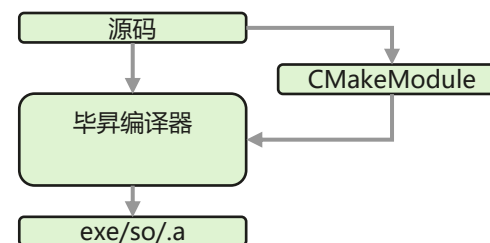


使用命令行一键编译:

```
bisheng main.asc --npu-soc
Ascend910B1
```

或基于标准cmake脚本

```
21 project(hello LANGUAGES ASC)
22 find_package(ASC)
23
24 add_executable(hello
25     hello_world.asc
26 )
```



<https://gitcode.com/cann/asc-devkit>

<https://gitcode.com/cann/asc-tools>

**CANN**



# ASC CPU孪生调试能力增强

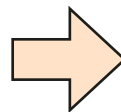
- 一套代码同时支持CPU和NPU调试，通过编译选项使能；
- CPU 孪生调试可解决大多数代码逻辑错误，包括同步指令不匹配校验等；

## As-Is

单独代码使能CPU孪生调试

```
108 int32_t main(int32_t argc, char *argv[])
109 {
110     uint32_t blockDim = 8;
111     size_t inputByteSize = 8 * 2048 * sizeof(uint16_t);
112     size_t outputByteSize = 8 * 2048 * sizeof(uint16_t);
113     #ifdef ASCEND_CPU_DEBUG
114         uint8_t *x = (uint8_t *)AscendC::GmAlloc(inputByteSize);
115         uint8_t *y = (uint8_t *)AscendC::GmAlloc(inputByteSize);
116         uint8_t *z = (uint8_t *)AscendC::GmAlloc(outputByteSize);
117
118         ReadFile("./input_x.bin", inputByteSize, x, inputByteSize);
119         ReadFile("./input_y.bin", inputByteSize, y, inputByteSize);
120
121         AscendC::SetKernelMode(KernelMode::AIV_MODE);
122         ICPU_RUN_KF(Add, blockDim, x, y, z); // use this macro for cpu debug
123
124         WriteFile("./output_z.bin", z, outputByteSize);
125
126         AscendC::GmFree((void *)x);
127         AscendC::GmFree((void *)y);
128         AscendC::GmFree((void *)z);
129     #else
130         CHECK_ACL(aclInit(nullptr));
131         int32_t deviceId = 0;
132         CHECK_ACL(aclrtSetDevice(deviceId));
133         aclrtStream stream = nullptr;
134         CHECK_ACL(aclrtCreateStream(&stream));
135
136         uint8_t *xHost, *yHost, *zHost;
```

孪生调试用户编写的特有代码



## To Be

一套代码，通过编译选项区分CPU孪生调试和NPU上板调试

```
104 int32_t main(int32_t argc, char *argv[])
105 {
106     uint32_t blockDim = 8;
107     size_t inputByteSize = 8 * 2048 * sizeof(uint16_t);
108     size_t outputByteSize = 8 * 2048 * sizeof(uint16_t);
109     CHECK_ACL(aclInit(nullptr));
110     int32_t deviceId = 0;
111     CHECK_ACL(aclrtSetDevice(deviceId));
112     aclrtStream stream = nullptr;
113     CHECK_ACL(aclrtCreateStream(&stream));
114
115     uint8_t *xHost, *yHost, *zHost;
116     uint8_t *xDevice, *yDevice, *zDevice;
117
118     CHECK_ACL(aclrtMallocHost((void **)&xHost, inputByteSize));
119     CHECK_ACL(aclrtMallocHost((void **)&yHost, inputByteSize));
120     CHECK_ACL(aclrtMallocHost((void **)&zHost, outputByteSize));
121     CHECK_ACL(aclrtMalloc((void **)&xDevice, inputByteSize, ACL_MEM_MALLOC_HUGE_FIRST));
122     CHECK_ACL(aclrtMalloc((void **)&yDevice, inputByteSize, ACL_MEM_MALLOC_HUGE_FIRST));
123     CHECK_ACL(aclrtMalloc((void **)&zDevice, outputByteSize, ACL_MEM_MALLOC_HUGE_FIRST));
124     ReadFile("./input_x.bin", inputByteSize, xHost, inputByteSize);
125     ReadFile("./input_y.bin", inputByteSize, yHost, inputByteSize);
126
127     CHECK_ACL(aclrtMemcpy(xDevice, inputByteSize, xHost, inputByteSize, ACL_MEMCPY_HOST_TO_DEVICE));
128     CHECK_ACL(aclrtMemcpy(yDevice, inputByteSize, yHost, inputByteSize, ACL_MEMCPY_HOST_TO_DEVICE));
129
130     Add<<blockDim, nullptr, stream>>(xDevice, yDevice, zDevice);
131     CHECK_ACL(aclrtSynchronizeStream(stream));
132
133     CHECK_ACL(aclrtMemcpy(zHost, outputByteSize, zDevice, outputByteSize, ACL_MEMCPY_DEVICE_TO_HOST));
134     WriteFile("./output_z.bin", zHost, outputByteSize);
135 }
```

执行代码一致，用户不用写而外代码

简化CPU孪生调试使能

# 目录

Part 1 Ascend C 介绍

Part 2 Ascend C 编程模型与新增特性介绍

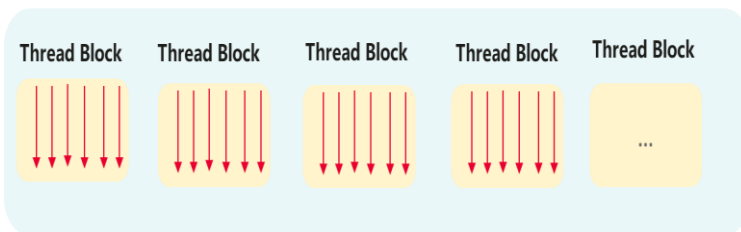
**Part 3 Ascend C 下一代规划**

# ASC 基于新一代硬件支持SIMT编程

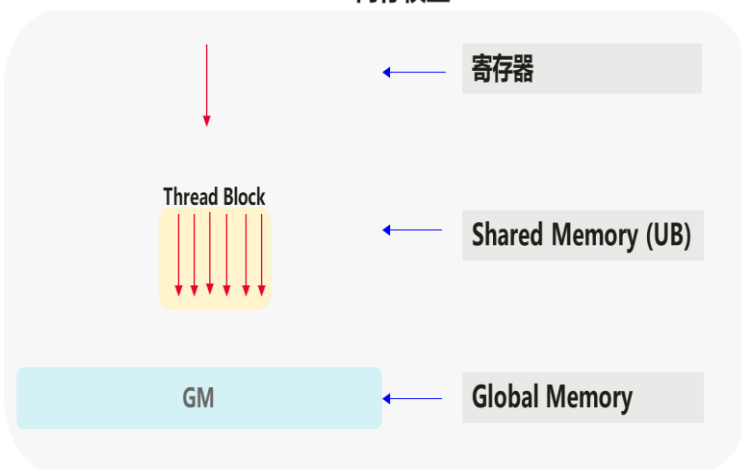
支持与业界类似的SIMT编程体验，优化离散和复杂逻辑处理等场景算子开发

## SIMT 编程模型

A5 SIMT 线程架构



A5 内存模型



## SIMT 算子示例

```
template <typename T>
__global__ LAUNCH_BOUND(512) inline void VectorAdd(T* dst, T*
src0, T* src1, uint32_t len)
{
    // Just show how to allocate share memory
    __ubuf__ T sharedBuf[512];

    // blockDim3 will be supported soon
    int idx = blockDim.x * blockIdx.x + threadIdx.x;

    // Just show how to use share memory
    sharedBuf[threadIdx.x] = src0[idx] + src1[idx];

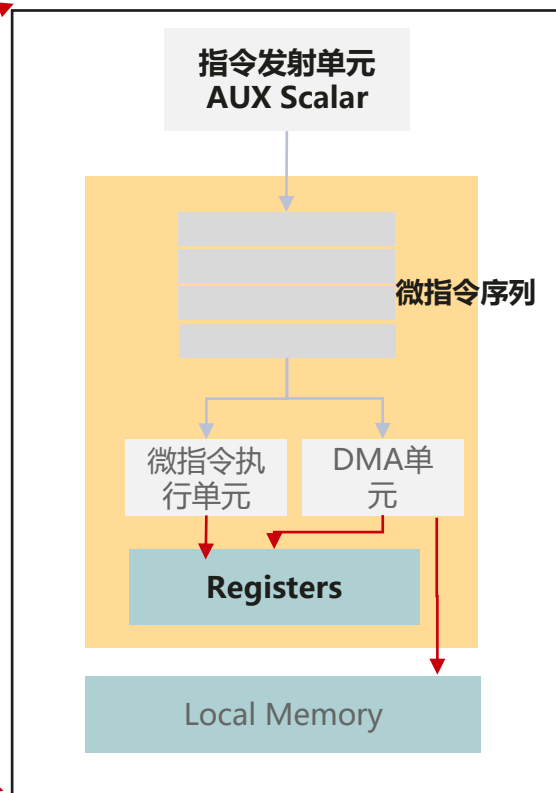
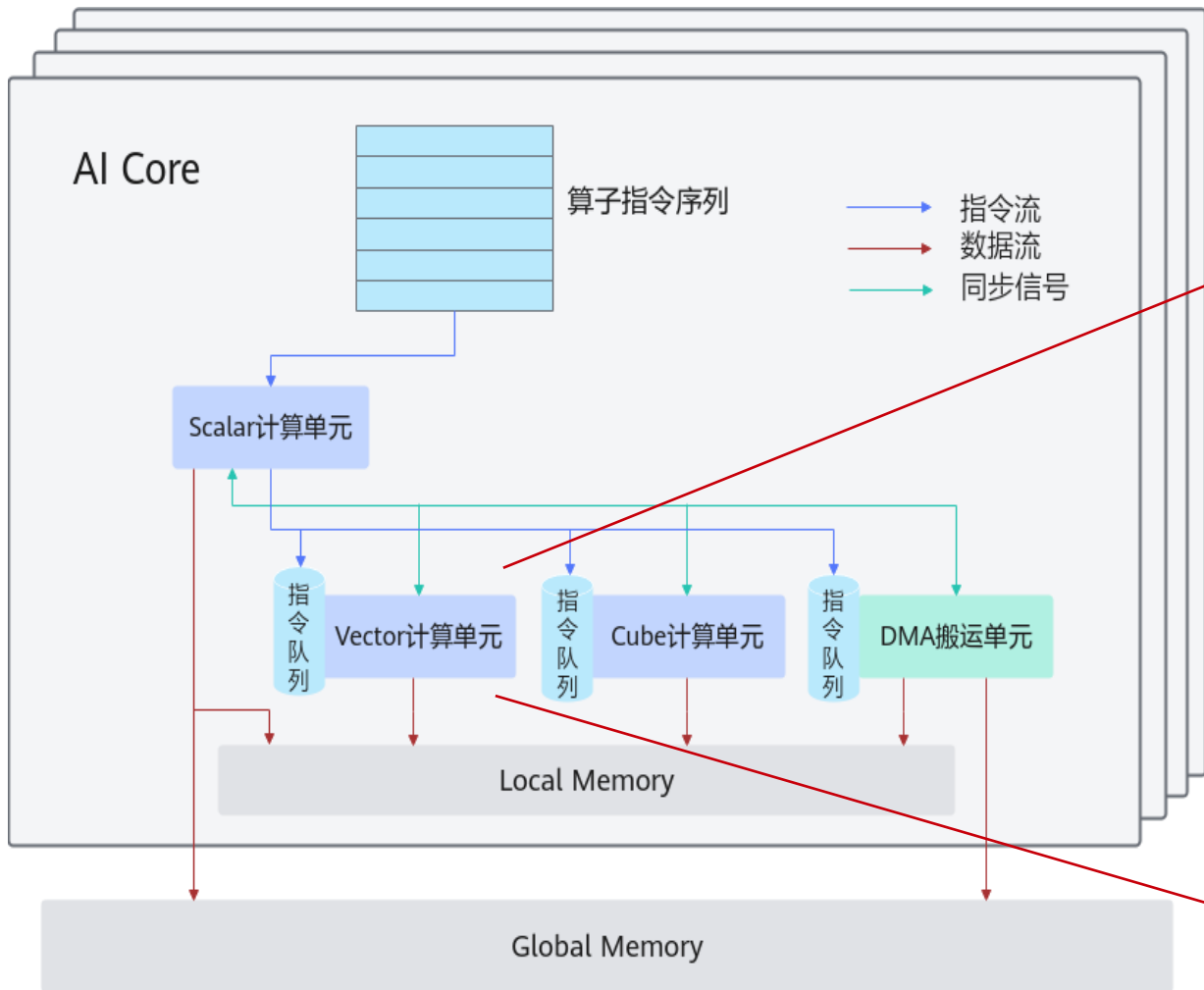
    if (idx < len) {
        dst[idx] = sharedBuf[threadIdx.x];
    }
    ...
    asc_syncthreads();
    ...
}
```

开放SIMT编程，适合处理离散访问、复杂控制逻辑处理等场景的优势

# 下一代SIMD 抽象架构的演进

持续优化SIMD架构，基于寄存器架构更好发挥SIMD优势

当前向量计算单元内部计算逻辑为固化的实现，无法实现可编程。下一代引进了向量计算单元可编程能力



Vector可编程计算单元

# ASC 基于新一代支持SIMT& SIMD新同构编程

同一个算子Kernel同时支持SIMD&SIMT，充分发挥各自优势

## 新同构SIMT示例

```
template <typename T>
__simt_vf__ __aicore__
LAUNCH_BOUND(512) inline void
SimtReduce(
__ubuf__ T* dst, __ubuf__ T* src0, uint32_t
len)
{
__ubuf__ T shareData;
int idx = blockDim.x * blockIdx.x +
threadIdx.x;

if (idx < len) {
shareData[threadIdx.x] = src0[idx];
}
asc_syncthreads();

...

if (threadIdx.x == 0) {
dst[blockIdx.x] = shareData[0];
}
}
```

开放SIMT编程，适合处理离散访问、复杂控制逻辑处理等场景的优势

## 新同构SIMD示例

```
__simd_vf__ __aicore__ inline void
SimdAdd(__ubuf__ half* dstPtr, __ubuf__
half * src0Ptr, __ubuf__ half * src1Ptr, uint32_t len)
{
constexpr uint16_t oneRepeatSize =
asc_get_vf_len() / sizeof(half);
uint16_t repeatTimes = ceil(len, oneRepeatSize);

vector_half dstReg0, srcReg0, srcReg1;
vector_bool maskReg;

for (uint16_t i = 0; i < repeatTimes; i++) {
maskReg = asc_update_mask(len);
asc_load(srcReg0, src0Ptr + i *
oneRepeatSize);
asc_load(srcReg1, src1Ptr + i *
oneRepeatSize);
asc_add(dstReg0, srcReg0, srcReg1, maskReg);
asc_store(dstPtr + i * oneRepeatSize, dstReg0,
maskReg);
}
}
```

开放SIMD微指令API，直接操作VEC向量寄存器编程，适合连续规整计算

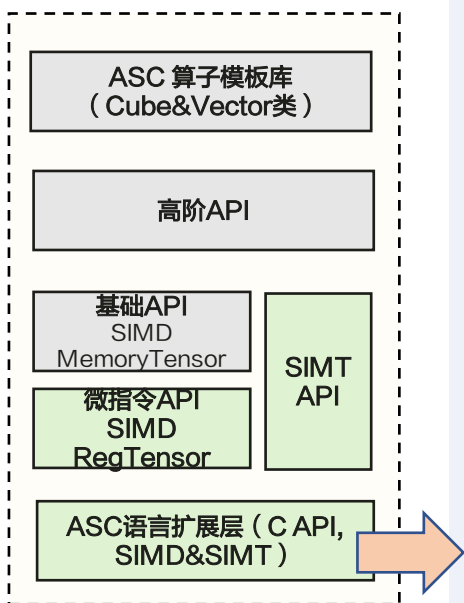
## SIMT&SIMD混合调用示例

通过VF\_CALL调用SIMT和SIMD函数

```
template<typename T>
extern "C" __global__ __aicore__ inline void
MixKernel(__gm__ T* x, __gm__ T* y, __gm__
T* z, uint32_t len)
{
...
// Execute SIMT Vector Function(VF) B
asc_call_vf<SimtReduce<T>>(blockDim,
dst0Local, src0Local, src1Local, len);

// Execute SIMD Vector Function(VF)
asc_call_vf<SimdAdd>(dst0Local,
src0Local, src1Local, len);
}
```

基于SIMT&SIMD新同构编程，充分发挥各自的优势



# Ascend C 规划路线图

## Ascend C 全面开源

- 支持A2/A3 语言扩展层C API初版
- 支持Tensor Tile API初稿
- 基础API全面开源

11/28

2025/12/30

- 支持A2/A3 语言扩展层C API，支撑向量开发
- 支持Tensor Tile API
- 完善样例等

2026下半年

- 完善设备侧类库的支持
- 完善Python前端编程能力

2026  
上半年

- 全量支持A2/A3 C API
- 基于下一代支持SIMT编程
- 基于下一代支持SIMD&SIMT新同构编程
- 基础API 完善Tensor Tile API
- 提升调试调优效率
- Python前端支持SIMT、Tensor 编程能力

- 适配新硬件，提供更优编程体验

# 欢迎交流 & 贡献



**asc-devkit 代码仓**

<https://gitcode.com/cann/asc-devkit>  
<https://gitcode.com/cann/asc-tools>



**asc-tools 代码仓**

# Thank you.

社区愿景：打造开放易用、技术领先的AI算力新生态

社区使命：使能开发者基于CANN社区自主研究创新，构筑根深叶茂、跨产业协同共享共赢的CANN生态

Vision: Building an Open, Easy-to-Use, and Technology-leading AI Computing Ecosystem

Mission: Enable developers to independently research and innovate based on the CANN community and build a win-win CANN ecosystem with deep roots and cross-industry collaboration and sharing.



上CANN社区获取干货



关注CANN公众号获取资讯