



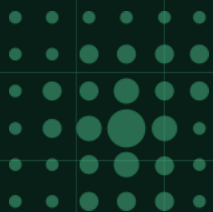
高性能向量计算编程实践

提高深度学习业务响应速度，降低部署成本

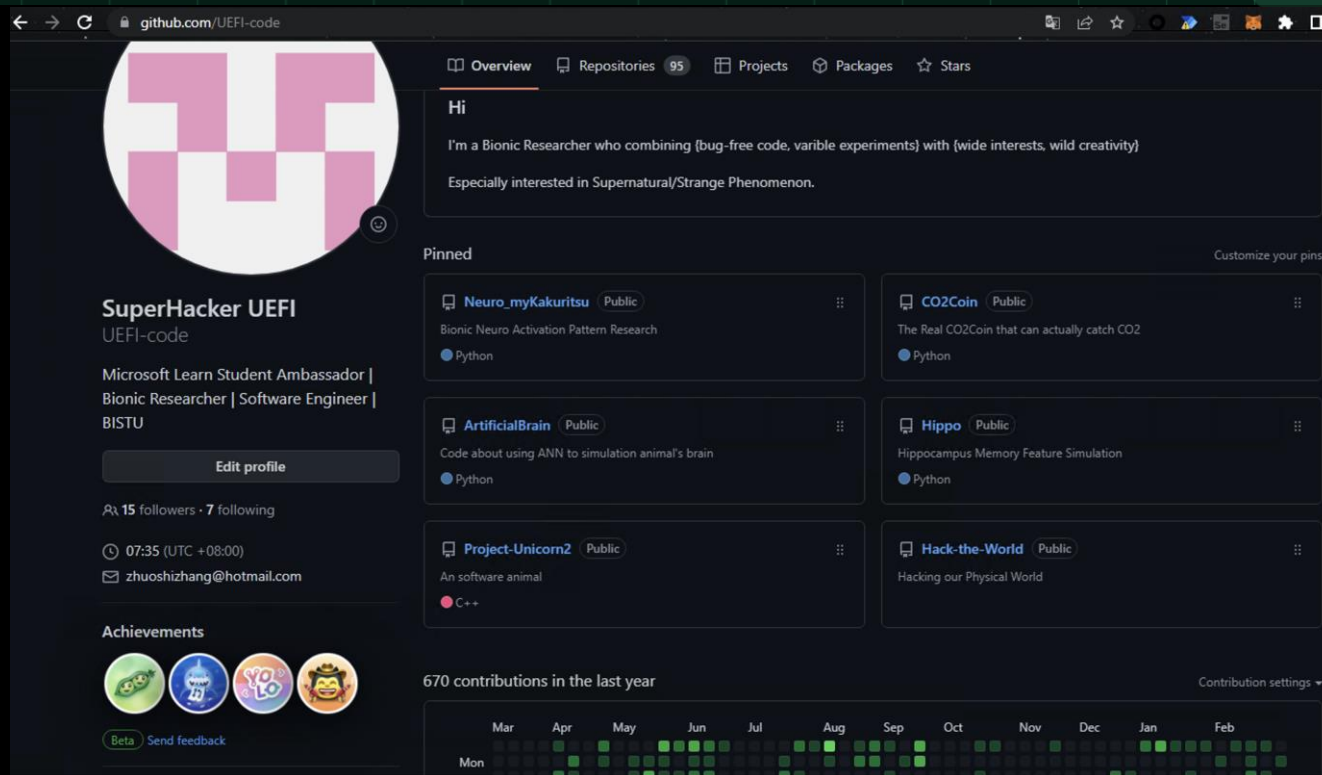
张世卓

微软学生大使 & 开源爱好者

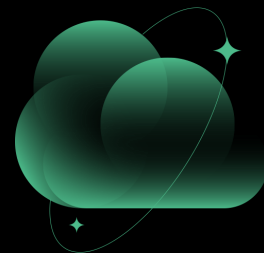
2023/04/22



讲师介绍



微软学生大使、中国发明协会会员、北信科在读20级本科生
热爱发明创新、开源、Coding、理化生实验&机器人仿生学综合研究
已获发明专利2项、实用新型专利多项



目录

CATALOGUE

0

经典深度学习云推理架构分析

1

向量化计算概念及应用场景

2

向量化计算优点

3

向量化计算编程注意事项

4

OpenVINO简介

5

Demo



经典深度学习云推理管道部署模式

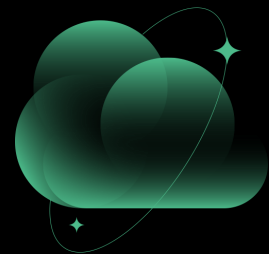
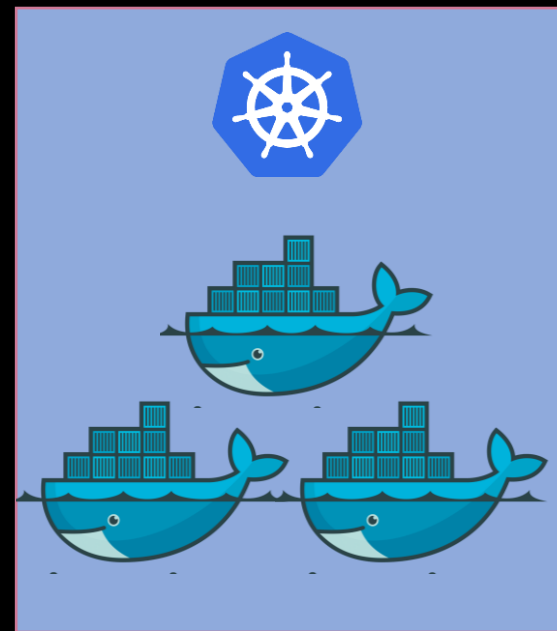
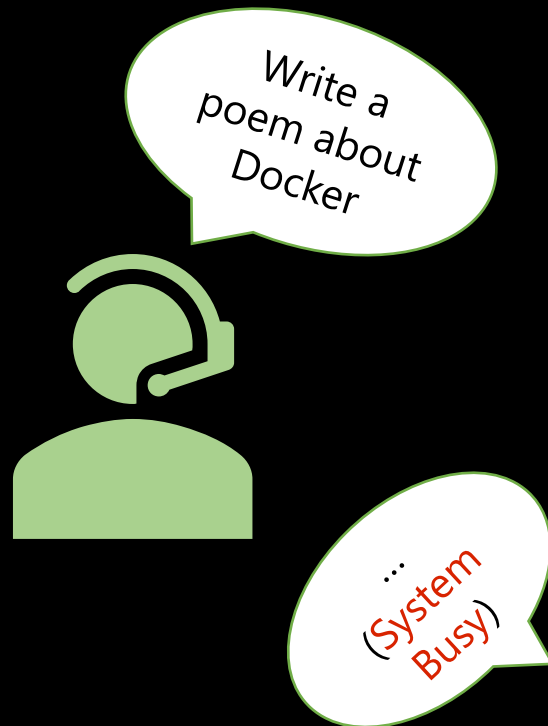
云原生部署：构建Docker镜像，部署到K8S，使用CPU推理

使用CPU通常具有如下优势：

- CPU节点比GPU节点租金便宜很多
- 避免数据从RAM到VRAM的拷贝延迟
- 现代CPU指令集变化较GPU平缓
- 在多样化环境兼容性好

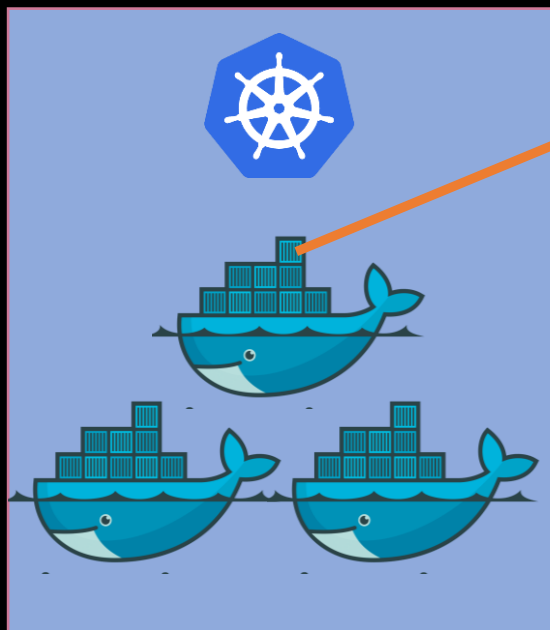
使用CPU通常具有如下缺点：

- 并行处理能力比GPU弱，模型参数量大时只能串行遍历



深度学习推理框架性能直接影响集群响应时间

...
(System
Busy)



阻塞: 前端Wait:
推理线程x

当前Stack:
嵌入层1->矩阵乘1-
>For Loop

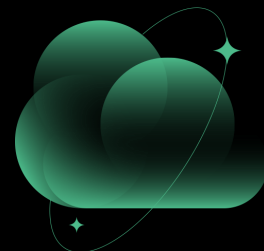
\$\$\$: 按CPU时间计费

如果使用传统循环遍历神经网络参数,
将是非常低效的做法

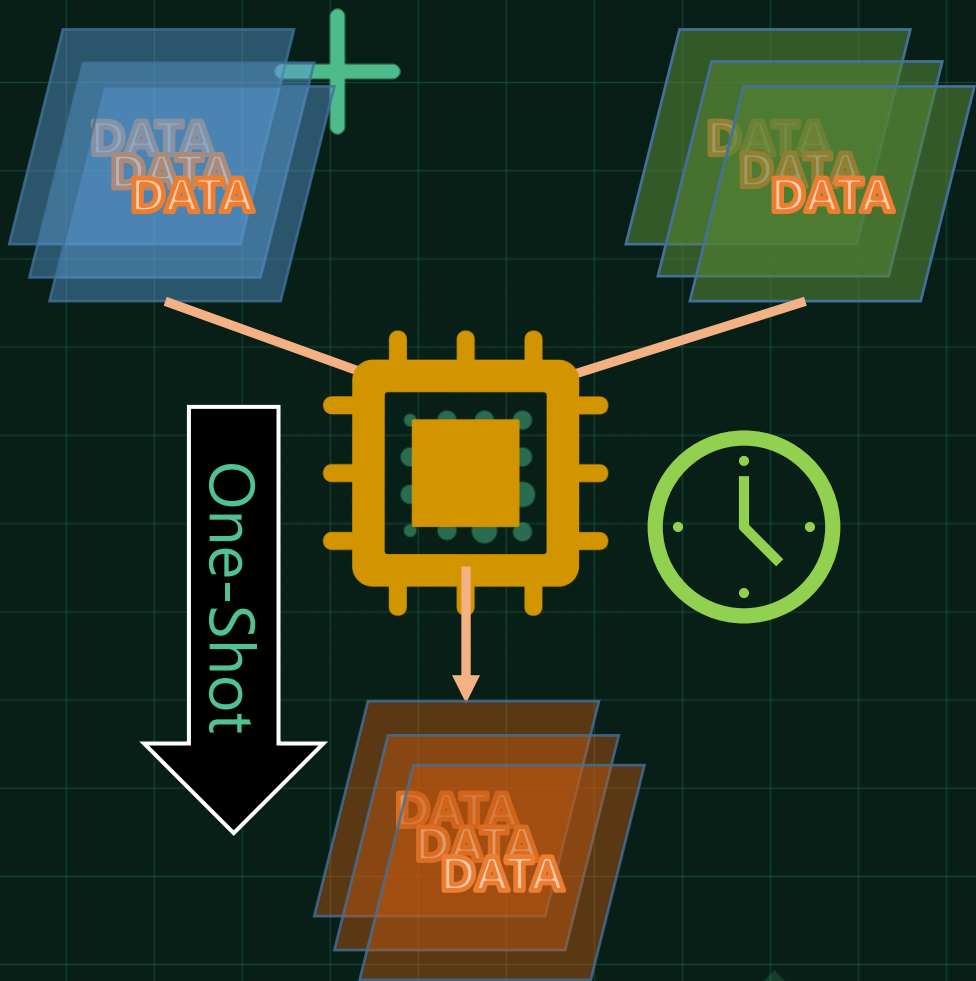
编译器自动优化
经常不可靠

改善编程
习惯/使用
优质框架

CPU其实也有针对大向量计算的优化:
向量化计算指令集-单指令多数据处理



向量化计算概念



向量化计算是一种利用处理器的SIMD（单指令多数据）能力来加速数学运算的技术。在向量化计算中，数值数据被存储在连续的内存地址中，处理器可以一次性地对整个数据向量执行相同的操作。

这种方法能够提高运算效率和吞吐量，因为它可以减少内存访问和操作指令的数量，从而降低了处理器的负载和延迟。向量化计算在许多科学计算、数据分析和机器学习应用中得到广泛应用。

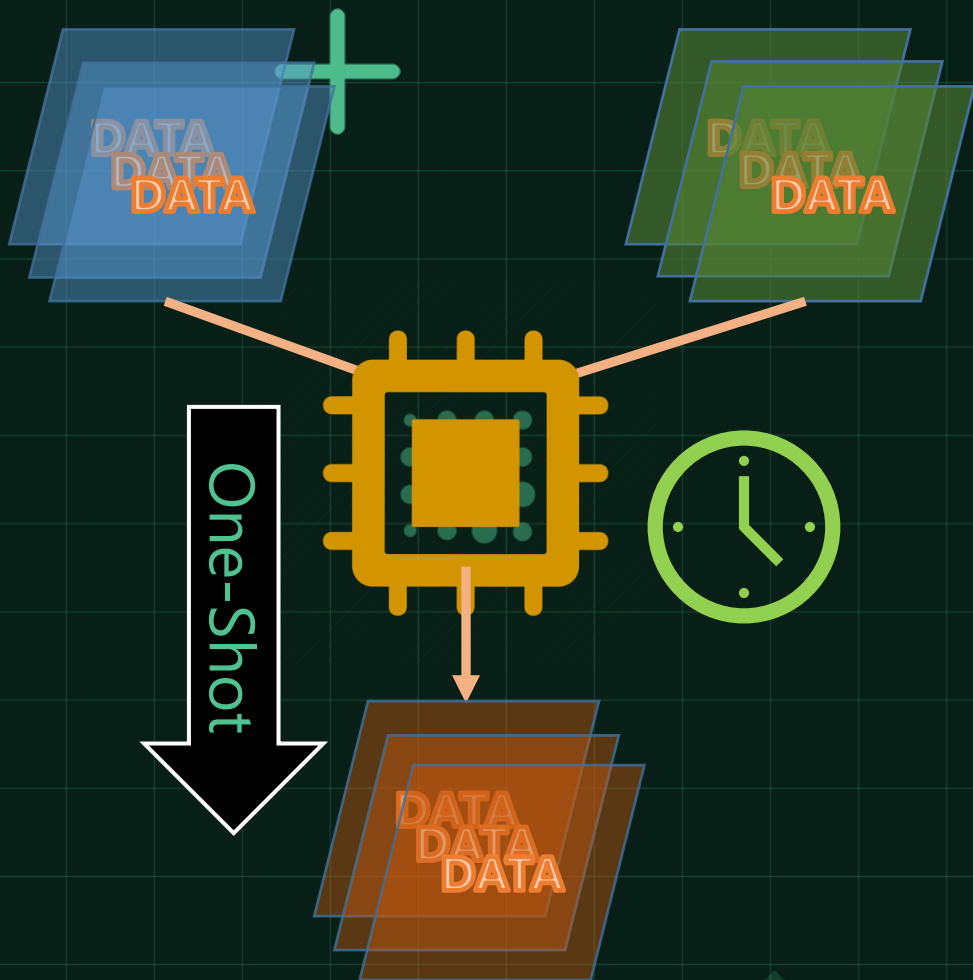
向量化计算应用场景

需求实例：两个向量求和

```
__declspec(align(16)) float vecA[] = {1.1, 2.2, 3.3, 4.4};
```

```
__declspec(align(16)) float vecB[] = {5.5, 6.6, 3.3, 4.4};
```

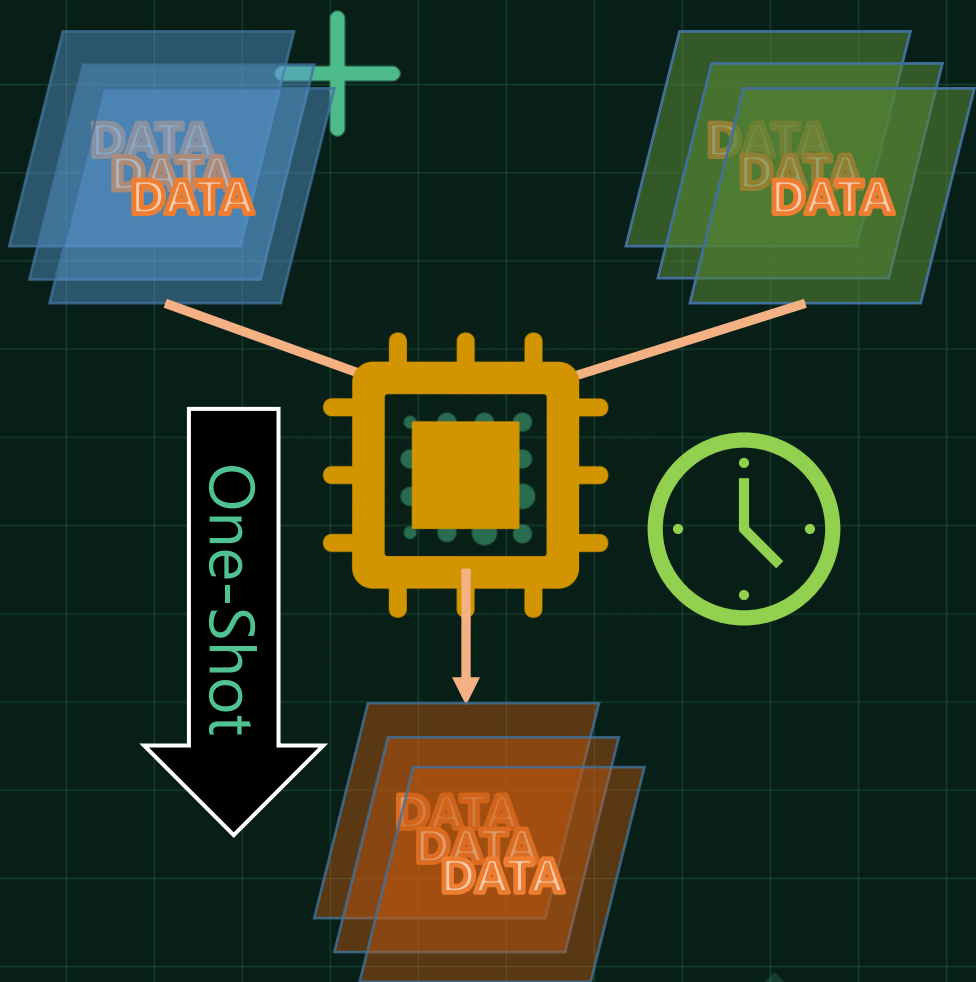
```
__declspec(align(16)) float res[4] = {0.0}; //res = vecA + vecB
```



```
;Code Block A  
LEA RAX, vecA;  
LEA RBX, vecB;  
LEA, RCX, vecC;  
MOV RDX, 0;  
LOOP:  
MOVSS XMM0, [RAX];  
MOVSS XMM1, [RBX];  
ADDSS XMM0, XMM1;  
MOVSS [RCX], XMM0  
  
ADD RAX, 4;  
ADD RBX, 4;  
ADD RCX, 4;  
ADD DL, 1;  
CMP DL, 4;  
JB LOOP;  
...;END  
SLOW
```

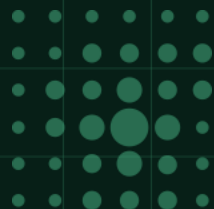
```
;Code Block B  
LEA RAX, vecA;  
LEA RBX, vecB;  
LEA, RCX, vecC;  
MOVAPS XMM0, [RAX]  
MOVAPS XMM1, [RBX]  
ADDPS XMM0, XMM1;  
MOVAPS [RCX], XMM0  
...;END  
FAST
```

使用向量化计算的好处



- 1 编译后的程序Bin体积小
- 2 程序加载快、运行快
- 4 操作系统无关性（兼容性好）
- 5 对深度学习类应用效果显著

向量化计算编程注意事项



0

不要指望编译器能优化的多好

1

处理好内存对齐、数据格式

2

手写汇编封装一个函数供上层调用
运行效率最高

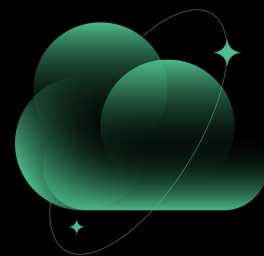
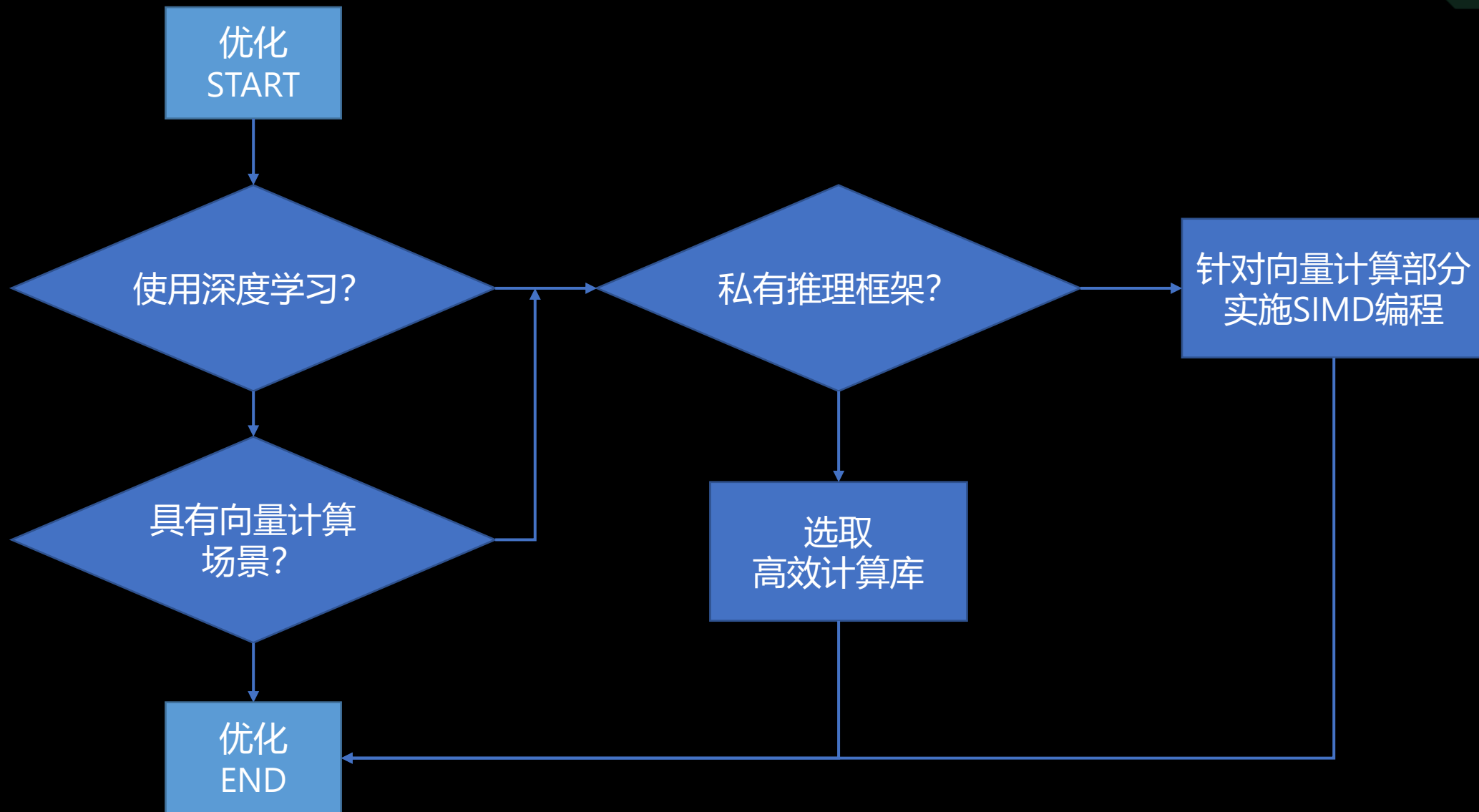
3

注意处理器硬件架构，尤其是使用较新的指令集时，有必要检查处理器FLAG



TAKE ACTION

2023





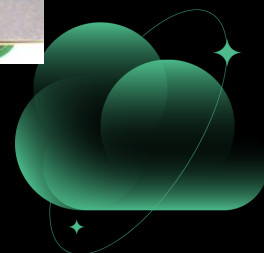
OpenVINO介绍

如果你不想自己编框架，
可用现成的，
SIMD Optimized



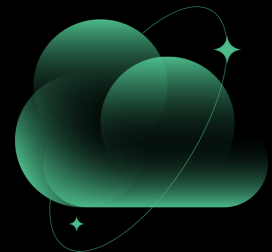
OpenVINO介绍

- OpenVINO是一个高性能的深度学习推理框架（你可用C/Cpp/Python调API）
- 诞生之初，OpenVINO主要对Intel自家产品优化，以发挥X86 CPU、核心显卡、VPU（加速棒）、FPGA性能为主。
- 现在已经支持ARM CPU，也能通过OpenCL使用NVIDIA GPU





Demo-Repo



Demo-VINO

2023

OpenVINO

Install Blog Forum Training GitHub

latest English

OpenVINO 2022.1 introduces a new version of OpenVINO API (API 2.0). For more information on the changes and transition steps, see the transition guide

Search the doc

Sample

Image Classification Sample

Hello Classification Sample

Hello Classification Sample

Hello Classification Sample

Model Creation C++ Sample

Model Creation Python* Sample

Automatic Speech Recognition C++ Sample

Automatic Speech Recognition Python* Sample

Administrator: C:\Windows\system32\cmd.exe

```
classification\public\alexnet\FP32\alexnet.xml
[ SUCCESS ] BIN file: D:\openvino_demo\w_openvino_toolkit_windows_2022.3.0.9052.9752fafe8eb_x86_64\samples\python\hello_classification\public\alexnet\FP32\alexnet.bin

D:\openvino_demo\w_openvino_toolkit_windows_2022.3.0.9052.9752fafe8eb_x86_64\samples\python\hello_classification>python hello_classification.py public\alexnet\FP16\alexnet.xml banana.jpg CPU
[ INFO ] Creating OpenVINO Runtime Core
[ INFO ] Reading the model: public\alexnet\FP16\alexnet.xml
[ INFO ] Loading the model to the plugin
[ INFO ] Starting inference in synchronous mode
[ INFO ] Image path: banana.jpg
[ INFO ] Top 10 results:
[ INFO ] class_id probability
[ INFO ] -----
[ INFO ] 954      0.9935312
[ INFO ] 317      0.0024962
[ INFO ] 934      0.0014471
[ INFO ] 110      0.0008863
[ INFO ] 939      0.0002137
[ INFO ] 940      0.0001649
[ INFO ] 941      0.0001622
[ INFO ] 112      0.0001263
[ INFO ] 998      0.0001044
[ INFO ] 953      0.0001029
[ INFO ] This sample is an API example, for any performance measurements please use the dedicated benchmark_app tool

D:\openvino_demo\w_openvino_toolkit_windows_2022.3.0.9052.9752fafe8eb_x86_64\samples\python\hello_classification>
```

On this page

How It Works

Running

Example

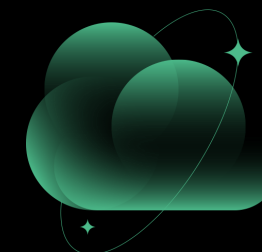
Sample Output

See Also

Download Docs



VINO-Video





Thanks

