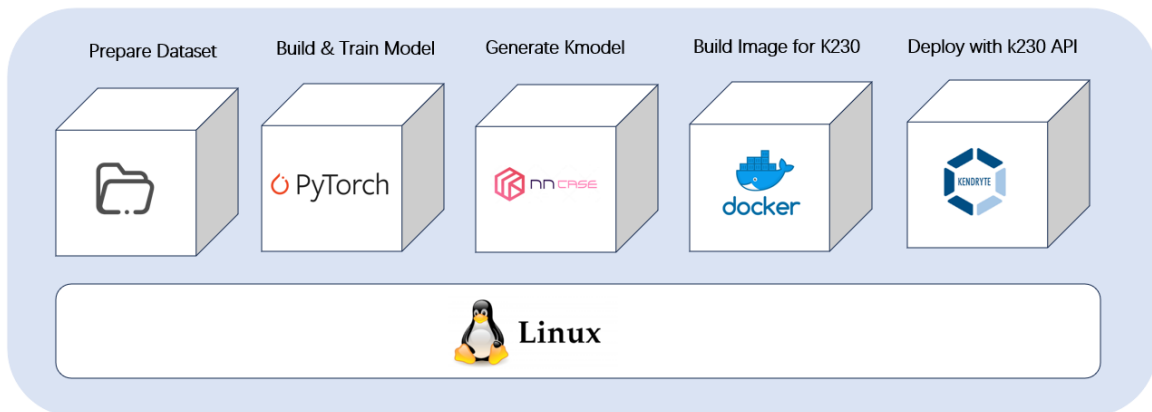


k230关键词唤醒教程

简介

K230芯片是嘉楠科技 Kendryte®系列AIoT芯片中的最新一代SoC产品。该芯片采用全新的多异构单元加速计算架构，集成了2个RISC-V高效计算核心，内置新一代KPU（Knowledge Process Unit）智能计算单元，具备多精度AI算力，广泛支持通用的AI计算框架，部分典型网络的利用率超过了70%。

该芯片同时具备丰富多样的外设接口，以及2D、2.5D等多个标量、向量、图形等专用硬件加速单元，可以对多种图像、视频、音频、AI等多样化计算任务进行全流程计算加速，具备低延迟、高性能、低功耗、快速启动、高安全性等多项特性。



本教程将介绍如何训练自己的**关键词唤醒**模型，并将模型转换为kmodel格式，部署在嘉楠Kendryte230芯片上。

实现该过程需要具备python和C++编程的基础知识，了解linux系统的简单操作，了解一定的深度学习知识，但并不是必须的。

本教程将实现从数据准备、模型训练和测试、k230镜像编译烧录、C++示例代码编译可执行文件、PC端和K230之间网络配置和文件传输、k230端部署的全部流程。操作系统为linux操作系统，深度学习框架选择PyTorch实现。

本教程选择训练自己的关键词——小楠小楠为例。

环境说明

显卡环境

本教程默认使用CUDA的用户已经安装好合适的显卡驱动，且已搭建好CUDA环境。

安装anaconda

如果已安装anaconda或miniconda，请忽略此步骤。

anaconda用于创建虚拟环境，将PyTorch模型训练环境和其他环境隔离。

```
1 apt-get install -y wget
2 wget https://repo.anaconda.com/archive/Anaconda3-5.3.0-Linux-x86_64.sh #可以选择
  合适的版本安装
3 chmod +x Anaconda3-5.3.0-Linux-x86_64.sh
4 ./Anaconda3-5.3.0-Linux-x86_64.sh
```

出现如下界面：

```
root@dev-wyf-react:~/Downloads# ls
Anaconda3-5.3.0-Linux-x86_64.sh
root@dev-wyf-react:~/Downloads# chmod +x Anaconda3-5.3.0-Linux-x86_64.sh
root@dev-wyf-react:~/Downloads# ./Anaconda3-5.3.0-Linux-x86_64.sh

Welcome to Anaconda3 5.3.0

In order to continue the installation process, please review the license
agreement.
Please, press ENTER to continue
>>> 
```

点击Enter(回车键)

此时显示Anaconda的信息，并且会出现More，继续按Enter，直到如下图所示：

```
Anaconda Distribution contains open source software packages from third parties. These are available on an "as is" basis and subject to their individual license agreements.
These licenses are available in Anaconda Distribution or at http://docs.anaconda.com/anaconda/pkg-docs. Any binary packages of these third party tools you obtain via Anaconda Distribution are subject to their individual licenses as well as the Anaconda license. Anaconda, Inc. reserves the right to change which third party tools are provided in Anaconda Distribution.

In particular, Anaconda Distribution contains re-distributable, run-time, shared-library files from the Intel(TM) Math Kernel Library ("MKL binaries"). You are specifically authorized to use the MKL binaries with your installation of Anaconda Distribution. You are also authorized to redistribute the MKL binaries with Anaconda Distribution or in the conda package that contains them. Use and redistribution of the MKL binaries are subject to the licensing terms located at https://software.intel.com/en-us/license/intel-simplified-software-license. If needed, instructions for removing the MKL binaries after installation of Anaconda Distribution are available at http://www.anaconda.com.

Anaconda Distribution also contains cuDNN software binaries from NVIDIA Corporation ("cuDNN binaries"). You are specifically authorized to use the cuDNN binaries with your installation of Anaconda Distribution. You are also authorized to redistribute the cuDNN binaries with an Anaconda Distribution package that contains them. If needed, instructions for removing the cuDNN binaries after installation of Anaconda Distribution are available at http://www.anaconda.com.

Anaconda Distribution also contains Visual Studio Code software binaries from Microsoft Corporation ("VS Code"). You are specifically authorized to use VS Code with your in
--More--
```

输入 yes

```
Do you accept the license terms? [yes|no]
[no] >>>
Please answer 'yes' or 'no':
>>>
Please answer 'yes' or 'no':
>>> yes

Anaconda3 will now be installed into this location:
/root/anaconda3

- Press ENTER to confirm the location
- Press CTRL-C to abort the installation
- Or specify a different location below

[/root/anaconda3] >>> 
```

继续点击 Enter

```
Anaconda3 will now be installed into this location:
/home/wangke/anaconda3

- Press ENTER to confirm the location
- Press CTRL-C to abort the installation
- Or specify a different location below

[/home/wangke/anaconda3] >>>
PREFIX=/home/wangke/anaconda3
installing: python-3.6.4-hc3d631a_1 ...
Python 3.6.4 :: Anaconda, Inc.
installing: ca-certificates-2017.08.26-h1d4fec5_0 ...
```

输入 yes，添加环境变量

检查是否安装成功：

```
1 | conda -V
```

若返回conda版本，表示安装成功。

安装docker

若已安装docker，请忽略此步骤。

Docker官方和国内daocloud都提供了一键安装的脚本，使得Docker的安装更加便捷。

官方的一键安装方式：

```
1 | curl -fSSL https://get.docker.com | bash -s docker --mirror Aliyun
```

国内 daocloud一键安装命令：

```
1 | curl -SSL https://get.daocloud.io/docker | sh
```

执行上述任一条命令，耐心等待即可完成Docker的安装。

创建模型训练环境

```
1 | # 使用anaconda创建模型训练的虚拟环境
2 | conda create -n myenv python=3.8.0
3 | # 激活虚拟环境
4 | conda activate myenv
5 | # 按照项目内的requirements.txt安装训练所用的python库,等待安装
6 | pip install -r requirements.txt
```

在requirements.txt中会安装模型转换的包nncase和nncase-kpu，`nncase` 是一个为 AI 加速器设计的神经网络编译器，参考[nncase](#)。

安装dotnet

```
1 | wget https://packages.microsoft.com/config/ubuntu/20.04/packages-microsoft-prod.deb -O packages-microsoft-prod.deb
2 | dpkg -i packages-microsoft-prod.deb
3 | apt update
4 | apt install -y apt-transport-https
5 | apt install -y dotnet-sdk-7.0
```

添加nncase插件到环境变量

注意：此步骤需要根据个人机器中实际情况完成。如果使用anaconda虚拟环境，在anaconda安装位置下的envs目录下，选择为训练创建的虚拟环境myenv，在其下面选择lib/python3.9/site-packages/，也就是pip安装requirements.txt内nncase和nncase-kpu的安装位置。source后会退出当前虚拟环境，需要再次激活。如果直接使用机器上的python，则需要添加其下的lib/python3.9/site-packages/。具体python版本自己控制。

```
1 | # python安装路由自己机器实际情况修改
2 | export NNCASE_PLUGIN_PATH=$NNCASE_PLUGIN_PATH:/usr/local/lib/python3.9/site-packages/
3 | export PATH=$PATH:/usr/local/lib/python3.9/site-packages/
4 | source /etc/profile
5 | conda activate myenv
```

准备自己的唤醒词数据

组织数据

本方法借助开源数据集speech_commands来训练自己的唤醒词，开发者需要准备两份数据：

1. 下载speech_commands_v0.01.tar.gz作为负样本数据。

```
1 | # 该数据集下载地址：
   | http://download.tensorflow.org/data/speech_commands_v0.01.tar.gz
2 | wget http://download.tensorflow.org/data/speech_commands_v0.01.tar.gz
```

2.在K230开发板上录制wav_from_k230.zip作为正样本。

该数据集需要开发者自行在k230上录制，录制步骤示例：

```
1  cd /sharefs/app
2
3  # 创建音频存储文件夹
4  mkdir wav_from_k230
5  mkdir wav_from_k230/xiaonanxiaonan
6  # 如果训练多个唤醒词，需要您按照唤醒词拼音创建所有文件夹
7  # 比如：
8  #mkdir wav_from_k230/xiaonantongxue
9  #mkdir wav_from_k230/nannan
10
11 # 在开发板上采集音频，每个唤醒词的音频数据放到对应的文件夹下
12 ./sample_audio.elf -type 0 -samplerate 16000 -bitwidth 16 -channels 1 -
enablecodec 1 -filename ./wav_from_k230/xiaonanxiaonan/xiaonan_1.wav
13
14 # 将录制的音频从开发板传输到服务器或者主机用于训练
15 scp wav_from_k230/* your_user_name@your_IP:/path/to/wav_from_k230
```

上述步骤i， ii， iii在大核执行，步骤iv在小核执行，负责将板子上录制的音频发送到开发者自己的服务器以用于模型训练。其中iii重复执行多次并修改保存的wav名称则可以录制多条wav音频。每一个wav音频的长度为15s

，录制者可以在15s内重复多次唤醒词，后续可以根据项目脚本分割为多个含有单个唤醒词的样本。

录制的样本示例可以参考 `./resource/xiaonan.wav`

将准备的数据按照以下格式放在项目下的 `example/speech_commands_v1` 下：

```
1  `-- example
2      |-- speech_commands_v1
3          |-- s0
4          |-- speech_commands_v0.01.tar.gz # 负样本数据集
5      |-- wav_from_k230.zip # 正样本数据集
```

配置训练参数

模型训练用到的参数文件为 `example/speech_commands_v1/s0/conf/ds_tcn.yaml`：

```
1  dataset_conf:
2      filter_conf:
3          max_length: 2048
4          min_length: 0
5      resample_conf:
6          resample_rate: 16000
7      speed_perturb: false
8      feature_extraction_conf:
9          feature_type: 'fbank'
10         num_mel_bins: 40
11         frame_shift: 10
12         frame_length: 25
13         dither: 1.0
14     spec_aug: true
15     spec_aug_conf:
```

```

16         num_t_mask: 1
17         num_f_mask: 1
18         max_t: 20
19         max_f: 10
20         shuffle: true
21         shuffle_conf:
22             shuffle_size: 1500
23         batch_conf:
24             batch_size: 100
25
26 model:
27     output_dim: 2    # todo: 需要根据自己的训练的唤醒词数量修改, 2代表只有一个唤醒词, 负样本占一个维度
28     input_dim: 40    # 输入特征维度
29     hidden_dim: 256 # 隐藏层特征维度
30     preprocessing:
31         type: linear # 预处理层类型
32     backbone:
33         type: tcn # backbone的类型
34         ds: true # 是否depth-wise
35         num_layers: 4 # backbone的层数
36         kernel_size: 8 # 卷积的kernel size
37         dropout: 0.1 # dropout的概率
38
39 optim: adam # 优化器
40 optim_conf:
41     lr: 0.001 # 学习率
42     weight_decay: 0.0001
43
44 training_config:
45     grad_clip: 5
46     max_epoch: 100
47     log_interval: 10

```

模型训练&模型测试&导出kmodel

进入到工程的 `example/speech_commands_v1/s0` 目录, 执行训练代码:

```

1 ./run.sh stage stop_stage project_path my_keyword num_keyword gpu_index
2 # stage 为run.sh脚本的开始阶段
3 # stop_stage 为run.sh的结束阶段
4 # project_path 为项目地址的绝对路径, 例如我的项目为end2end_kws, 其在/mnt/end2end_kws, 则该参数为
5 # end2end_kws_doc, 注意: /end2end_kws/这种写法是错误的!
6 # my_keyword 想训练的关键词, 用字母表示, 例如小楠小楠=xiaonanxiaonan, 多个关键词用“-”连接, 比如: xiaonanxiaonan_xiaonantongxue
7 # num_keyword 关键词数量, 一个关键词则为2, 两个关键词为3, n个关键词加上负样本总共为n+1
8 # gpu_index 为GPU索引, '-1'/'0'/'1'/'1,2,3', -1为CPU, 其他为GPU索引
9
10 # run.sh中包含6个阶段:
11 # -1阶段为将数据划分为训练集, 验证集和测试集
12 # 0阶段为将数据处理为kaldi格式的文件
13 # 1阶段为计算CMVN并组织成模型可以使用的数据格式
14 # 2阶段为训练模型
15 # 3阶段平均模型和测试模型准确率

```

```

16 # 4阶段为导出ONNX模型
17 # 5阶段为导出kmodel模型
18
19 # 例如:
20 ./run.sh -1 -1 end2end_kws_doc xiaonanxiaonan 2 0
21 # 多唤醒词时
22 # ./run.sh -1 -1 end2end_kws_doc xiaonanxiaonan_xiaonantongxue 3 0
23 # 表示只执行-1阶段
24
25 ./run.sh -1 5 end2end_kws_doc xiaonanxiaonan 2 0
26 # 多唤醒词时
27 # ./run.sh -1 5 end2end_kws_doc xiaonanxiaonan_xiaonantongxue 3 0
28 # 上述命令表示使用 GPU 0 完成从数据处理到导出kmodel的整个流程,注意替换自己的项目路径
29
30 # 重训时需要删除example/speech_commands_v1目录下的my_data文件夹

```

如果训练成功，在 `example/speech_commands_v1/s0/exp/xiaonanxiaonan` 路径下可以找到

`avg_10.pt` ,

`avg_10.onnx` 和 `avg_10.kmodel` 。

使用k230部署模型

环境准备和镜像编译

注意：训练环境中nncase和nncase-kpu的版本和SDK的版本要对应，nncase和nncase-kpu版本为2.4.0，SDK版本为1.1。

K230 SDK需要在**Linux环境**下编译，推荐使用Ubuntu Liunx 20.04。

使用docker编译环境，下载[k230 sdk](#)。

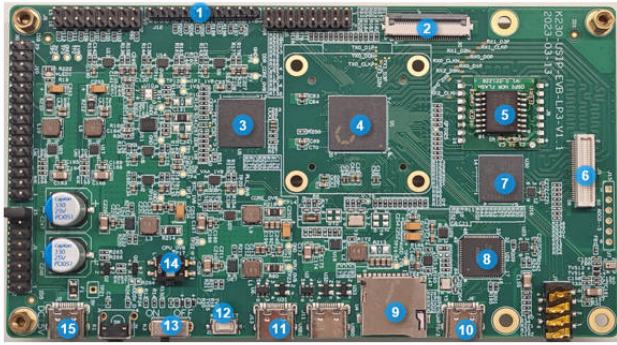
```

1 # 下载docker编译镜像
2 docker pull ghcr.io/kendryte/k230_sdk
3 # 可以使用以下命令确认docker镜像拉取成功
4 docker images | grep k230_sdk
5 # 下载sdk
6 git clone -b v1.1 --single-branch https://github.com/kendryte/k230_sdk.git
7 cd k230_sdk
8 # 下载工具链, make prepare_sourcecode 会自动下载Linux和RT-Smart toolchain,
  buildroot package, AI package等. 请确保该命令执行成功并没有Error产生, 下载时间和速度
  以实际网速为准。
9 make prepare_sourcecode
10 # 创建docker容器, $(pwd):$(pwd)表示系统当前目录映射到docker容器内部的相同目录下, 将系统
   下的工具链目录映射到docker容器内部的/opt/toolchain目录下
11 docker run -u root -it -v $(pwd):$(pwd) -v $(pwd)/toolchain:/opt/toolchain -w
   $(pwd) ghcr.io/kendryte/k230_sdk /bin/bash

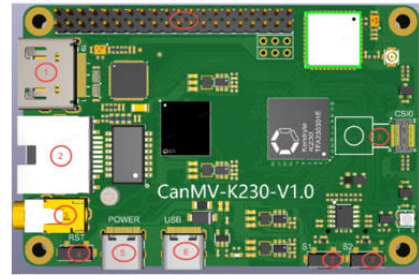
```

K230现有两种开发板，分别为CANMV-K230-V1.0（以下简称CANMV-K230）和K230-USIP-LP3-EVB（以下简称K230-EVB）。两种开发板区别如图：

K230-USIP-LP3-EVB



CANMV-K230



```

1 # 在docker中编译镜像，请耐心等待完成，不同类型开发板编译命令不同
2 # 如果是CANMV-K230开发板
3 make CONF=k230_canmv_defconfig
4 # 如果是K230-EVB开发板
5 make CONF=k230_evb_defconfig

```

SD卡镜像也可在嘉楠开发者社区下载：开发者社区-->资料下载-->K230-->Images。

镜像烧录

CANMV-K230开发板：

编译结束后在output/k230_canmv_defconfig/images目录下可以找到编译好的镜像文件：

```

1 k230_canmv_defconfig/images
2 |— big-core
3 |— little-core
4 |— sysimage-sdcard.img # SD卡镜像
5 |— sysimage-sdcard.img.gz # SD卡镜像压缩包

```

CANMV-K230开发板支持SD卡镜像启动。

K230-EVB开发板：

编译结束后在output/k230_evb_defconfig/images目录下可以找到编译好的镜像文件：

```

1 k230_evb_defconfig/images
2 |— big-core
3 |— little-core
4 |— sysimage-sdcard.img # SD和emmc非安全启动镜像
5 |— sysimage-sdcard.img.gz # SD和emmc的非安全启动镜像压缩包
6 |— sysimage-spinor32m.img # norflash非安全启动镜像
7 |— sysimage-spinor32m_jffs2.img # norflash jffs2非安全启动镜像

```

K230 支持SDCard、eMMC、norflash等多种启动方式。

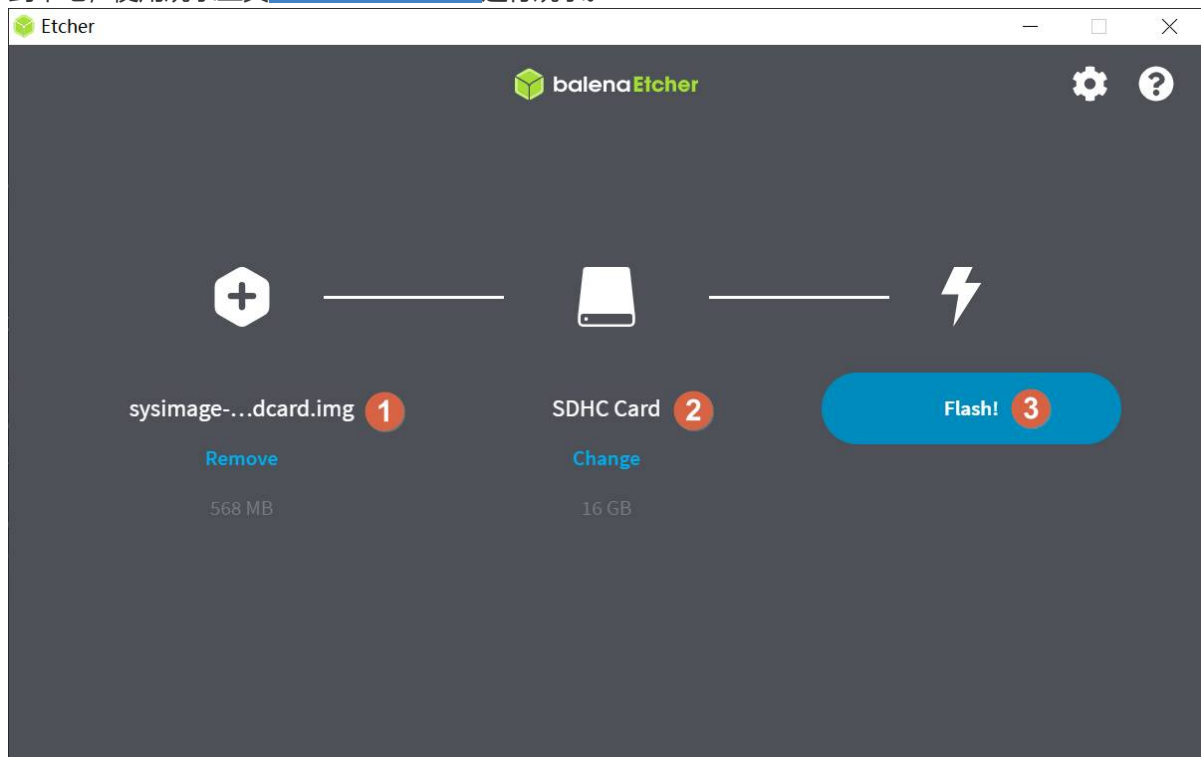
烧录TF卡

详细烧录步骤参考[K230 SDK 使用说明](#)。

Linux: 如使用Linux烧录TF卡,需要先确认SD卡在系统中的名称/dev/sdx, 并替换如下命令中的/dev/sdx

```
1 | sudo dd if=sysimage-sdcard.img of=/dev/sdx bs=1M oflag=sync
```

Windows: 如使用Windows烧录, 建议使用[the balena Etcher](#)工具。将生成的sysimage-sdcard.img下载到本地, 使用烧录工具[the balena Etcher](#)进行烧录。



其它更详细的烧录方法, 请参考[K230 SDK 使用说明](#)。

上电启动K230 EVB开发板

K230-EVB开发板上电启动

K230 EVB支持SDCard、eMMC、norflash等多种启动方式, 用户可以通过改变开板上启动拨码开关的设置, 来切换不同启动模式。为方便开发, 建议您准备一张TF卡, 并将**拨码开关切换至SD卡启动模式**, 后续可考虑将镜像文件固化至emmc中。

1. 请先**确认启动开关SW1选择在SD卡启动模式**下 (详情可参考[开机上电方式](#))
2. 将烧录完成的TF卡插入开发板TF卡槽中
3. 开发板接上电源
4. **将电源开关K1拨到ON位置**, 系统可上电启动
5. 如果您有接好串口, 可在串口中看到启动日志输出。

CanMV-K230开发板上电启动

K230 CanMV-K230开发板支持SDCard启动方式、HDMI输出显示, 因此, 需要准备一张TF卡, 此外建议准备一个HDMI显示器。

1. 将烧录完成的TF卡插入开发板TF卡槽中
2. 开发板上电, 此时, 系统可上电启动

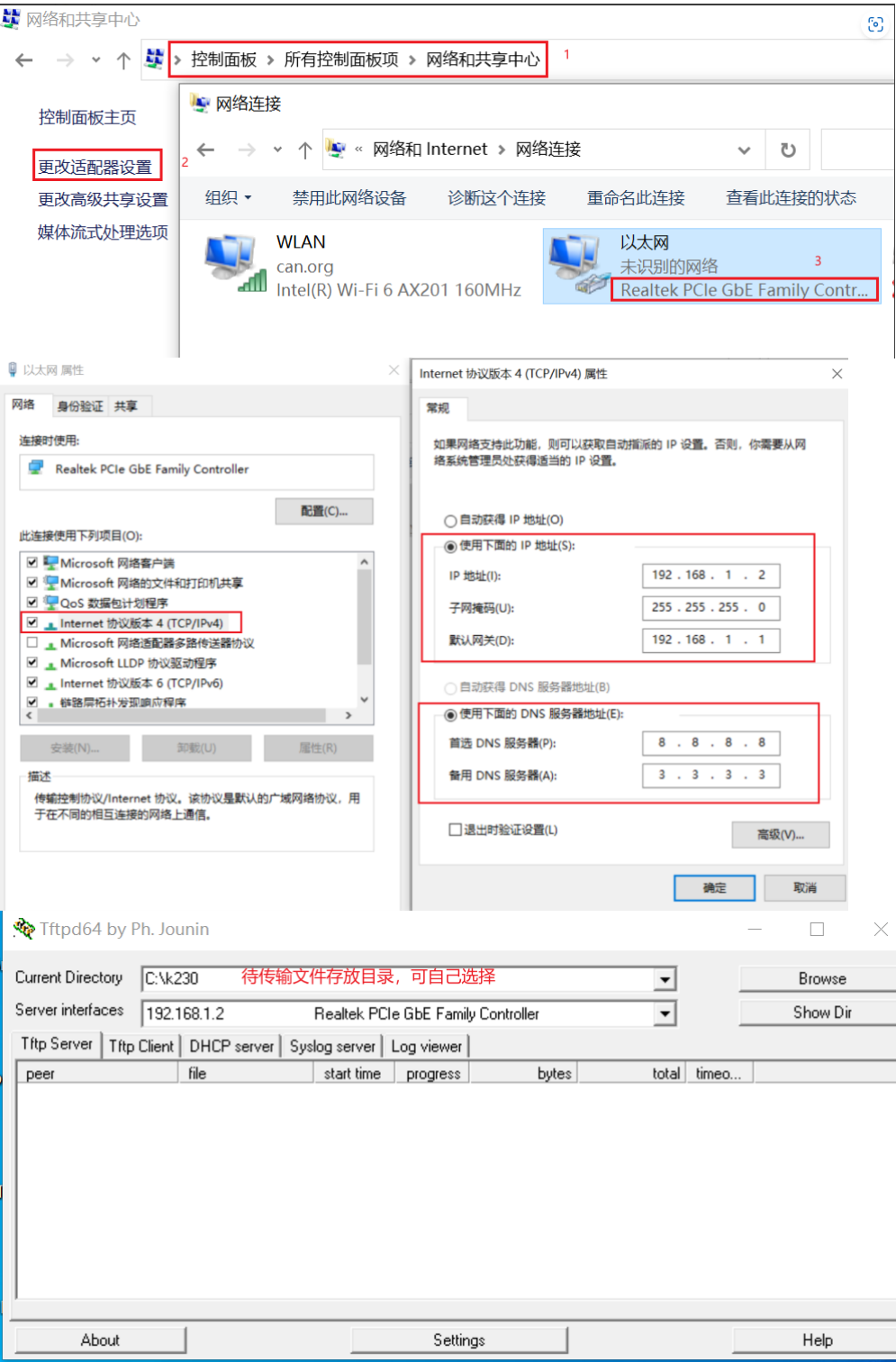
系统上电后, 默认会有**两个串口设备**, 可分别用于访问小核Linux和大核RTSmart

小核Linux默认用户名root, 密码为空。大核RTSmart系统中开机会自动启动一个应用程序, 可按 **q** 键退出至命令提示符终端。

PC和k230文件传输配置与实现

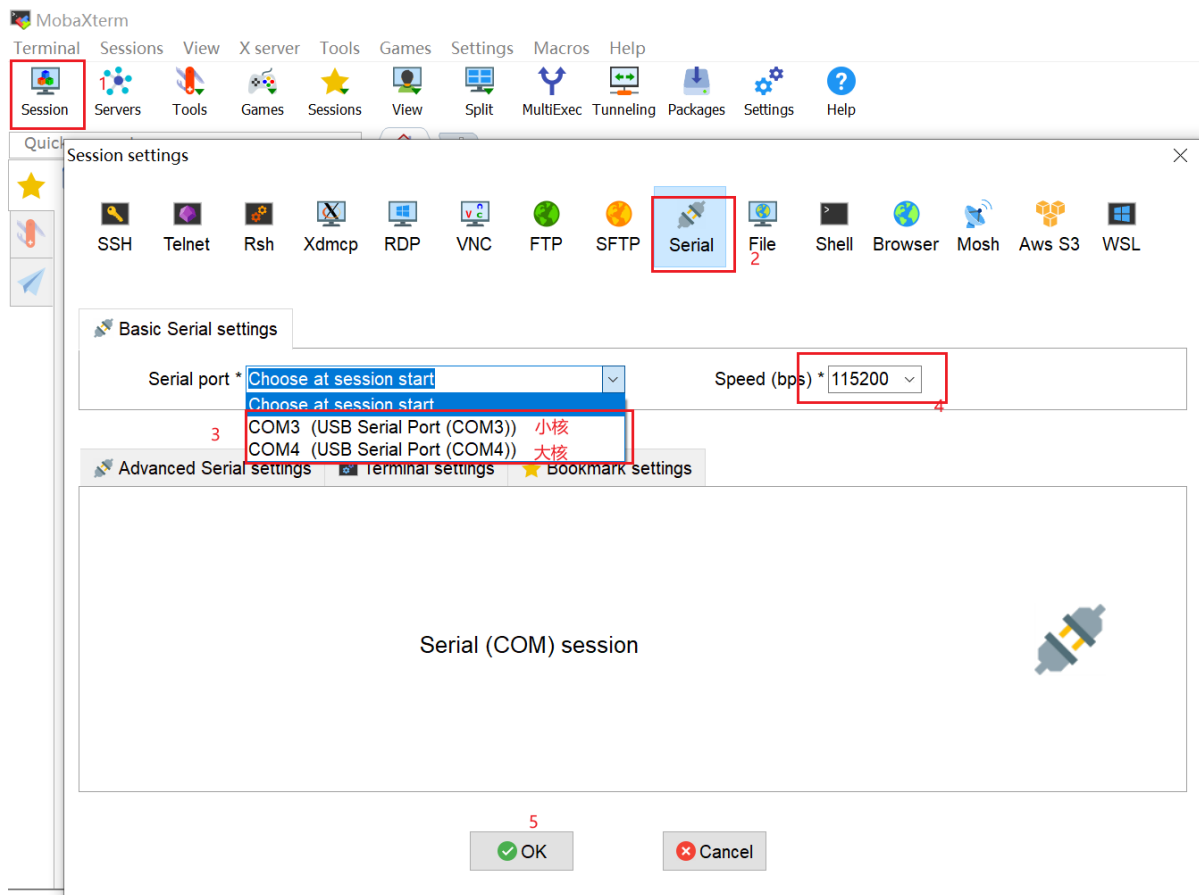
windows系统

- (1) Tftpd64安装, 在<https://bitbucket.org/phjounin/tftpd64/downloads/>下载。
- (2) MobaXterm安装: 在<https://mobaxterm.mobatek.net/download.html>下载安装。
- (2) 配置PC网络:



- (4) 开发板网络配置:

开发板上电, 电源线、网线、COM口连接线配置见文档: [K230_SDK_使用说明](#)。打开MobaXterm, 通过两路COM串口连接开发板, COM编号不固定, 较小为小核串口, 较大为大核串口。



小核进入后回车，进入如下界面，使用root登录：

```

Welcome to Buildroot
canaan login: root
[root@canaan ~]#
  
```

大核进入后回车，进入如下界面：

```

msh />
  
```

在小核配置网络：

```

Welcome to Buildroot
canaan login: root
[root@canaan ~]#ls
[root@canaan ~]#cd /
[root@canaan /]#ls
app      init      linuxrc   opt        sbin      usr
bin      lib       lost+found proc        sharefs   var
dev      lib64     media     root       sys
etc      lib64xthead mnt       run        tmp
[root@canaan /]#ifconfig eth0 192.168.1.22
  
```

在小核端口配置网络，和PC在同一网段即可

大小核共享存储区域：/sharefs

```
2. COM3 (USB Serial Port (COM3)) x 3. COM4 (USB Serial Port (COM4))
Welcome to Buildroot
canaan login: root
[root@canaan ~]#ls
[root@canaan ~]#cd /
[root@canaan /]#ls
app      init      linuxrc   opt        sbin      usr
bin      lib       lost+found  proc       sharefs   var
dev      lib64     media     root       sys
etc      lib64xthead  mnt       run        tmp
[root@canaan /]#ifconfig eth0 192.168.1.22
[root@canaan /]#cd sharefs/
[root@canaan /sharefs]#ls
```

当要从Tftpd64配置的文件中拷贝数据时，在小核界面使用如下命令：

```
1 # 192.168.1.2 为PC的局域网IP
2 tftp -g -r your_file_name 192.168.1.2
```

当将开发板文件拷贝到PC端Tftpd64配置的文件夹下时，在小核使用如下命令：

```
1 # 192.168.1.2 为PC的局域网IP
2 tftp -p -r your_file_name 192.168.1.2
```

Linux系统

在Linux系统中，PC正常连接网络，开发板可以通过网线连接PC所在网关下其他网口，通过scp命令实现文件传输。

开发板上电，进入大小核COM界面，在小核执行scp传输命令：

```
1 # 从PC拷贝文件至开发板
2 scp 用户名@域名或IP:文件所在目录 开发板目的目录
3 # 从开发板拷贝文件至PC
4 scp 开发板待拷贝目录 用户名@域名或IP:PC目的目录
```

上板code解析

完成上述开发板的准备工作后，我们可以使用C++编写自己的代码，下面就关键词唤醒任务的示例代码进行解析。本教程给出相关关键词唤醒任务的示例代码，并进行简单解析。

代码结构

```
1 k230_code
2 |— kws_stream
3 |— ai_base.cc
4 |— ai_base.h # ai base基类
5 |— audio_stream.h # 封装了实时音频流的一些功能
6 |— blocking_queue.h # 封装了音频预处理用到的blocking queue
7 |— build
8 |— CMakeLists.txt
9 |— fbank.h # 封装了将音频转换为fbank特征的方法
10 |— feature_pipeline.cc
11 |— feature_pipeline.h # 封装了整个音频预处理流程
12 |— fft.cc
13 |— fft.h # 封装了音频预处理傅里叶变换的一些计算方法
14 |— kws.cc
```

```
15 | └─ kws.h # 封装了KWS类，主要实现了预处理，推理和后处理过程
16 | └─ log.h # 封装了log类
17 | └─ main.cc # 主函数文件，封装了实时检测音频流的功能
18 | └─ pcm_data.cc # 封装了解析实时音频流数据的功能
19 | └─ scoped_timing.hpp # 封装了计算某一scope运行所耗时间的功能
20 | └─ wav_ctrl.cc # 封装了音频预处理相关函数
21 | └─ CMakeLists.txt
22 | └─ build_app.sh
```

核心代码

```
1  /**
2   * @brief AI基类，封装nncase相关操作
3   * 主要封装了nncase的加载、设置输入、运行、获取输出操作，后续开发demo只需要关注模型的前处
  理、后处理即可
4   */
5  class AIBase
6  {
7  public:
8      /**
9       * @brief AI基类构造函数，加载kmodel,并初始化kmodel输入、输出
10      * @param kmodel_file kmodel文件路径
11      * @param debug_mode 0（不调试）、1（只显示时间）、2（显示所有打印信息）
12      * @return None
13      */
14  AIBase(const char *kmodel_file,const string model_name, const int debug_mode
    = 1);
15
16  /**
17   * @brief AI基类析构函数
18   * @return None
19   */
20  ~AIBase();
21
22  /**
23   * @brief 设置kmodel输入
24   * @param buf 输入数据指针
25   * @param size 输入数据大小
26   * @return None
27   */
28  void set_input(const unsigned char *buf, size_t size);
29
30  /**
31   * @brief 根据索引获取kmodel输入tensor
32   * @param idx 输入数据指针
33   * @return None
34   */
35  runtime_tensor get_input_tensor(size_t idx);
36
37  void set_input_tensor(size_t idx, runtime_tensor &tensor);
38
39  /**
40   * @brief 初始化kmodel输出
41   * @return None
42   */
```

```

43 void set_output();
44
45 /**
46  * @brief 推理kmodel
47  * @return None
48  */
49 void run();
50
51 /**
52  * @brief 获取kmodel输出，结果保存在对应的类属性中
53  * @return None
54  */
55 void get_output();
56
57
58
59 protected:
60 string model_name_;           // 模型名字
61 int debug_mode_;             // 调试模型，0（不打印），1（打印时间），2
    （打印所有）
62 vector<float*> p_outputs_;     // kmodel输出对应的指针列表
63 vector<vector<int>> input_shapes_; // {{N,C,H,W},{N,C,H,W}...}
64 vector<vector<int>> output_shapes_; // {{N,C,H,W},{N,C,H,W}...} 或 {{N,C},
    {N,C}...}等
65 vector<int> each_input_size_by_byte_;
    // {0,layer1_length,layer1_length+layer2_length,...}
66 vector<int> each_output_size_by_byte_;
    // {0,layer1_length,layer1_length+layer2_length,...}
67 private:
68 /**
69  * @brief 首次初始化kmodel输入，并获取输入shape
70  * @return None
71  */
72 void set_input_init();
73
74 /**
75  * @brief 首次初始化kmodel输出，并获取输出shape
76  * @return None
77  */
78 void set_output_init();
79
80 // kmodel解释器，从kmodel文件构建，负责模型的加载、输入输出设置和推理
81 vector<unsigned char> kmodel_vec_; // 通过读取kmodel文件得到整个kmodel数据，用于传
    给kmodel解释器加载kmodel
82 interpreter kmodel_interp_;
83 };

```

上述代码是ai_base.h文件中AIBase类的定义代码。主要定义了kmodel解释器，kmodel的相关信息，以及输入输出设置、推理过程的接口定义。具体实现在ai_base.cc中。

```

1 /**
2  * @brief KWS---关键词唤醒
3  * 封装了利用麦克风实时关键词唤醒的过程
4  */
5 class KWS : public AIBase

```

```

6 {
7     public:
8         /**
9          * @brief KWS构造函数，加载kmodel
10         * @param kmodel_file kmodel文件路径
11         * @param task_name 任务名称
12         * @param num_keyword 关键词数量
13         * @param debug_mode 0（不调试）、1（只显示时间）、2（显示所有打印信息）
14         * @return None
15         */
16         KWS(const char *kmodel_file, const std::string task_name, const int
num_keyword, const float spot_thresh, const int debug_mode);
17
18         ~KWS();
19
20         bool pre_process(std::vector<float> wav);
21
22         /**
23         * @brief kmodel推理
24         * @return None
25         */
26         void inference();
27
28         /**
29         * @brief kmodel推理结果后处理
30         * @return None
31         */
32         std::string post_process();
33
34         static wenet::FeaturePipelineConfig feature_config; // 音频预处理类
35         static wenet::FeaturePipeline feature_pipeline; // 音频预处理类
36
37
38     private:
39         int num_bin = 40; // 音频特征维度
40         const int chunk_size = 30; // 推理块大小
41         const std::string task_name; // 音频位置
42         const int batch_size = 1; // 音频预处理默认参数，不可更改！
43         const int hidden_dim = 256; // 隐藏层维度
44         const int cache_dim = 105; // in_cache最后一维数量，与模型参
数绑定，不可更改！
45         const int num_keyword; // 唤醒词数量，包含Deactivated!
46         const float spot_thresh; // 唤醒阈值
47         std::vector<string> idx2res; // ID与唤醒词的映射
48         runtime_tensor wav_feature; // Kmodel输入
49         runtime_tensor in_cache; // Kmodel输入
50         std::vector<std::vector<float>> cache; // 用于保存模型每次推理输出的cache，
cache用于流式推理
51     };
52 #endif

```

上述代码是实现实时关键词唤醒任务的类定义，主要定义关键词唤醒模型推理的前处理、推理、后处理接口。

```

1 void print_usage(const char *name)

```

```

2  {
3      cout << "Usage: " << name << "<kmodel_kws> <num_keyword> <task_name>
<debug_mode>" << endl
4          << "Options:" << endl
5          << "  kmodel_kws      语音唤醒 kmodel路径\n"
6          << "  task_name       任务名称\n"
7          << "  num_keyword     关键词数量\n"
8          << "  spot_thresh     唤醒词激活阈值\n"
9          << "  debug_mode      是否需要调试, 0、1、2分别表示不调试、简单调试、详细调
试\n"
10         << "\n"
11         << endl;
12 }
13
14 int main(int argc, char *argv[])
15 {
16     std::cout << "case " << argv[0] << " built at " << __DATE__ << " " <<
__TIME__ << std::endl;
17     if (argc != 6)
18     {
19         print_usage(argv[0]);
20         return -1;
21     }
22
23     else
24     {
25         std::thread thread_iss(video_proc, argv);
26         while (getchar() != 'q')
27         {
28             usleep(10000);
29         }
30
31         audio_stop = true;
32         thread_iss.join();
33     }
34     return 0;
35 }
36

```

上述代码是main.cc的一部分，主要实现参数解析功能。

```

1      bool enable_audio_input = true;
2      bool enable_audio_output = true;
3      bool disable_audio_input = true; //disable必须要和enable对应
4      bool disable_audio_output = true; //disable必须要和enable对应
5
6      if (K_SUCCESS != start_aio_stream(enable_audio_input,
enable_audio_output))
7      {
8          std::cout << "start_aio_stream failed\n" << std::endl;
9      }

```

上述代码是main.cc中开启音频流的代码，k230录音的必备动作。


```

1 while (!audio_stop)
2 {
3     ScopedTiming st("total time", 1);
4
5     // 获取音频流数据
6     auto pcm_chunk = get_audio_chunk();
7     k_u16 *pcm_data = pcm_chunk.first;
8     k_u32 pcm_size = pcm_chunk.second;
9
10    // 将PCM数据由uint16 -> int16_t -> float, 以满足PCM数据处理接口的输入要求
11    float* data_;
12    data_ = new float[pcm_size];
13    for (int i = 0; i < pcm_size; i++){
14        int16_t sample;
15        sample = pcm_data[i];
16        data_[i] = static_cast<float>(sample);
17    }
18    std::vector<float> wav(data_, data_+pcm_size);
19
20    std::cout << "===== " <<
std::endl;
21
22    // 将音频流PCM数据处理为音频Fbank特征
23    bool ok = kws.pre_process(wav);
24
25    // 模型推理
26    kws.inference();
27
28    // 获得检测结果
29    std::string results = kws.post_process();
30
31    if (results == "Deactivated!"){
32        std::cout << "Deactivated!" << std::endl;
33    }
34    else{
35        std::cout << results << std::endl;
36        std::srand(static_cast<unsigned int>(std::time(nullptr)));
37        std::default_random_engine randomEngine(std::rand());
38        std::uniform_int_distribution<int> distribution(0,
reply_wav_list.size() - 1);
39        int index = distribution(randomEngine);
40        const char * reply_wav_file = reply_wav_list[index];
41        const char *reply_wav = reply_wav_file;
42        if (K_SUCCESS != play_wav(reply_wav))
43        {
44            std::cout << "play_wav failed\n" << std::endl;
45        }
46    }
47
48 }

```

上述代码是main.cc中对音频流进行实时关键词检测的过程。

k230_code/kws_stream/CMakeLists.txt 脚本说明

```
1  set(src main.cc kws.cc fft.cc feature_pipeline.cc ai_base.cc utils.cc
    pcm_data.cc wav_ctrl.cc)
2  set(bin kws_stream.elf)
3
4  include_directories(${PROJECT_SOURCE_DIR})
5  include_directories(${nncase_sdk_root}/riscv64/rvvlib/include)
6  include_directories(${k230_sdk}/src/big/mpp/userapps/api/)
7  include_directories(${k230_sdk}/src/big/mpp/include)
8  include_directories(${k230_sdk}/src/big/mpp/include/comm)
9  include_directories(${k230_sdk}/src/big/mpp/userapps/sample/sample_vo)
10 include_directories(${k230_sdk}/src/big/mpp/userapps/sample/sample_audio)
11 link_directories(${nncase_sdk_root}/riscv64/rvvlib/)
12
13 add_executable(${bin} ${src})
14 target_link_libraries(${bin} -Wl,--start-group rvv Nncase.Runtime.Native
    nncase.rt_modules.k230 functional_k230 sys vicap vb cam_device cam_engine
15 hal oslayer ebase fpga isp_drv binder auto_ctrl common cam_caldb isi 3a
    buffer_management cameric_drv video_in virtual_hal start_engine cmd_buffer
16 switch cameric_reg_drv t_database_c t_mxml_c t_json_c t_common_c vo
    connector sensor atomic dma -Wl,--end-group ai ao)
17
18 target_link_libraries(${bin} opencv_imgcodecs opencv_imgproc opencv_core zlib
    libjpeg-turbo libopenjp2 libpng libtiff libwebp csi_cv)
19 install(TARGETS ${bin} DESTINATION bin)
```

这是 k230_code/kws_stream 目录下的CMakeLists.txt脚本，设置编译的C++文件和生成的elf可执行文件名，由下面：

```
1  set(src main.cc kws.cc fft.cc feature_pipeline.cc ai_base.cc utils.cc
    pcm_data.cc wav_ctrl.cc)
2  set(bin kws_stream.elf)
```

- 1 • include_directories(\${PROJECT_SOURCE_DIR}): 添加项目的根目录到头文件搜索路径中。
- 2 • include_directories(\${nncase_sdk_root}/riscv64/rvvlib/include): 添加 nncase RISC-V 向量库的头文件目录。
- 3 • include_directories(\${k230_sdk}/src/big/mpp/userapps/api/): 添加 k230 SDK 中的用户应用程序 API 头文件目录。
- 4 • include_directories(\${k230_sdk}/src/big/mpp/include): 添加 k230 SDK 的一般头文件目录。
- 5 • include_directories(\${k230_sdk}/src/big/mpp/include/comm): 添加与通信相关的头文件目录。
- 6 • include_directories(\${k230_sdk}/src/big/mpp/userapps/sample/sample_audio): 添加示例 aio (音频输入输出) 应用程序头文件目录。
- 7 • include_directories(\${k230_sdk}/src/big/mpp/userapps/sample/sample_vo): 添加示例 vo (视频输出) 应用程序头文件目录。

```

1  link_directories(${nncase_sdk_root}/riscv64/rvvlib/): 添加链接器搜索路径, 指向
   nncase RISC-V 向量库的目录。
2  add_executable(${bin} ${src}): 创建一个可执行文件, 将之前设置的源文件列表作为输入。
3  target_link_libraries(${bin} ...): 设置可执行文件需要链接的库。列表中列出了各种库, 包括
   nncase 相关的库、k230 SDK 的库, 以及其他一些库。
4  target_link_libraries(${bin} opencv_imgproc opencv_imgcodecs opencv_core zlib
   libjpeg-turbo libopenjp2 libpng libtiff libwebp csi_cv): 将一些 OpenCV 相关的库
   和其他一些库链接到可执行文件中。
5  install(TARGETS ${bin} DESTINATION bin): 安装生成的可执行文件到指定的目标路径 (bin
   目录) 中。

```

上述是 k230_code/k230_deploy 目录下的CMakeLists.txt脚本说明。

k230_code/CMakeLists.txt脚本说明

```

1  cmake_minimum_required(VERSION 3.2)
2  project(nncase_sdk C CXX)
3
4  set(nncase_sdk_root "${PROJECT_SOURCE_DIR}/../../nncase/")
5  set(k230_sdk ${nncase_sdk_root}/../../..)
6  set(CMAKE_EXE_LINKER_FLAGS "-T ${PROJECT_SOURCE_DIR}/cmake/link.lds --
   static")
7
8  # set opencv
9  set(k230_opencv ${k230_sdk}/src/big/utils/lib/opencv)
10 include_directories(${k230_opencv}/include/opencv4/)
11 link_directories(${k230_opencv}/lib ${k230_opencv}/lib/opencv4/3rdparty)
12
13 # set mmz
14 link_directories(${k230_sdk}/src/big/mpp/userapps/lib)
15
16 # set nncase
17 include_directories(${nncase_sdk_root}/riscv64)
18 include_directories(${nncase_sdk_root}/riscv64/nncase/include)
19 include_directories(${nncase_sdk_root}/riscv64/nncase/include/nncase/runtime)
20 link_directories(${nncase_sdk_root}/riscv64/nncase/lib/)
21
22 add_subdirectory(k230_deploy)

```

这是k230_code目录下的CMakeLists.txt脚本。该脚本重点关注如下部分

```

1  set(nncase_sdk_root "${PROJECT_SOURCE_DIR}/../../nncase/"):设置nncase目录
2  set(k230_sdk ${nncase_sdk_root}/../../..):设置k230_sdk的目录, 当前是从nncase目录
   的三级父目录得到
3  set(CMAKE_EXE_LINKER_FLAGS "-T ${PROJECT_SOURCE_DIR}/cmake/link.lds --
   static"):设置链接脚本路径, 链接脚本放于k230_code/cmake下
4  ...
5  add_subdirectory(k230_deploy): 添加待编译的工程子目录, 如您要编译自己的工程, 可以更换该
   行

```

以上是k230_code目录下的CMakeLists.txt脚本说明。

k230_code/build_app.sh脚本说明

```
1  #!/bin/bash
2  set -x
3
4  # set cross build toolchain
5  # 将交叉编译工具链的路径添加到系统的 PATH 环境变量中，以便在后续的命令中使用。该工具链是使用的
   大核编译工具链。
6  export PATH=$PATH:/opt/toolchain/riscv64-linux-musleabi_for_x86_64-pc-linux-
   gnu/bin/
7
8  clear
9  rm -rf out
10 mkdir out
11 pushd out
12 cmake -DCMAKE_BUILD_TYPE=Release \
13       -DCMAKE_INSTALL_PREFIX=`pwd` \
14       -DCMAKE_TOOLCHAIN_FILE=cmake/Riscv64.cmake \
15       ..
16
17 make -j && make install
18 popd
19
20 # 生成的main.elf可以在k230_code目录下的k230_bin文件夹下找到
21 k230_bin=`pwd`/k230_bin
22 rm -rf ${k230_bin}
23 mkdir -p ${k230_bin}
24
25 if [ -f out/bin/kws_stream.elf ]; then
26     cp out/bin/kws_stream.elf ${k230_bin}
27 fi
```

AI代码编译

将项目中的k230_code文件夹拷贝到k230_sdk目录下的src/big/nncase下，执行编译脚本，将C++代码编译成main.elf可执行文件。K230现有两种开发板，两种开发板编译命令不同，请注意区分。

CANMV-K230开发板

如果编译可以在CANMV-K230开发板执行的elf文件：

```
1  # 在k230_SDK根目录下执行
2  make CONF=k230_canmv_defconfig prepare_memory
3  # 回到当前项目目录下
4  ./build_app.sh
```

K230-EVB开发板

如果编译可以在K230-EVB开发板执行的elf文件：

```
1  # 在k230_SDK根目录下执行
2  make CONF=k230_evb_defconfig prepare_memory
3  # 回到当前项目目录下
4  ./build_app.sh
```

若权限不够，可使用如下代码赋予相关权限：

```
1 | chmod +x build_app.sh
2 | ./build_app.sh
```

文件拷贝

按照上文第4节配置好文件传输，在MobaXterm上的小核界面进入/sharefs，将训练得到的 avg_10.kmodel 和编译得到的 kws_stream.elf 传输到 /sharefs/kws 目录下，如果没有需要提前创建目录。

```
1 | kws
2 | └─reply_wav # 当被唤醒时，响应用户的音频文件，比如：“小楠小楠”，“我在！”，“我在”就是响应音频。
3 |   └─wozai.wav # “我在！”
4 |   └─wolaile.wav # “我来了！”
5 | └─avg_10.kmodel # 训练得到的kmodel
6 | └─kws_stream.elf # 编译的可执行文件
```

模型板上运行

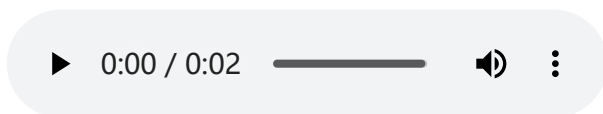
在大核COM口界面执行main.elf实现关键词唤醒：

```
1 | # "模型推理时传参说明："
2 | # Usage: kws_stream.elf<kmodel_kws> <num_keyword> <task_name> <debug_mode>
3 | #Options:
4 | # kmodel_kws      语音唤醒 kmodel路径
5 | # task_name       任务名称
6 | # num_keyword     关键词数量
7 | # spot_thresh     唤醒词激活阈值
8 | # debug_mode      是否需要调试，0、1、2分别表示不调试、简单调试、详细调试
9 | kws_stream.elf avg_10.kmodel xiaonan 2 0.5 0
```

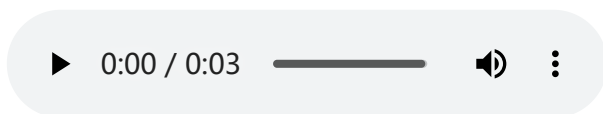
上板部署效果

我们多次唤醒小楠小楠，它随机选择 reply_wav 文件夹下的响应音频给出响应：

小楠响应“叫我干嘛”：



小楠响应“我在”：



工具

烧录工具the balena Etcher: <https://etcher.balena.io/>

局域网文件传输工具Tftpd64: <https://bitbucket.org/phjounin/tftpd64/downloads/>

MobaXterm下载地址: <https://mobaxterm.mobatek.net/download.html>

参考

k230_sdk github: https://github.com/kendryte/k230_sdk

k230_sdk_doc github: [k230_sdk 使用说明](#)

k230_sdk gitee: https://gitee.com/kendryte/k230_sdk

nncase github: [kendryte/nncase: Open deep learning compiler stack for Kendryte AI accelerator \(github.com\)](#)