



(12) 发明专利

(10) 授权公告号 CN 114492778 B

(45) 授权公告日 2024. 09. 06

(21) 申请号 202210141274.7

G06F 16/22 (2019.01)

(22) 申请日 2022.02.16

G06F 7/483 (2006.01)

(65) 同一申请的已公布的文献号

G06F 7/499 (2006.01)

申请公布号 CN 114492778 A

G06N 3/082 (2023.01)

(43) 申请公布日 2022.05.13

(56) 对比文件

(73) 专利权人 安谋科技(中国)有限公司

Cliein, X.A Quantization Model Based on a Floating-point Computing-in-Memory Architecture.2022 IEEE ASIA PACIFIC CONFERENCE ON CIRCUITS AND SYSTEMS, APCCAS.2023,493-496.

地址 200233 上海市闵行区田林路1016号

科技绿洲三期11号楼

(72) 发明人 鲁若荻 韩冥生 余宗桥

审查员 雷冬

(74) 专利代理机构 上海华诚知识产权代理有限公司

公司 31300

专利代理师 肖华

(51) Int. Cl.

G06N 3/063 (2023.01)

G06N 3/0464 (2023.01)

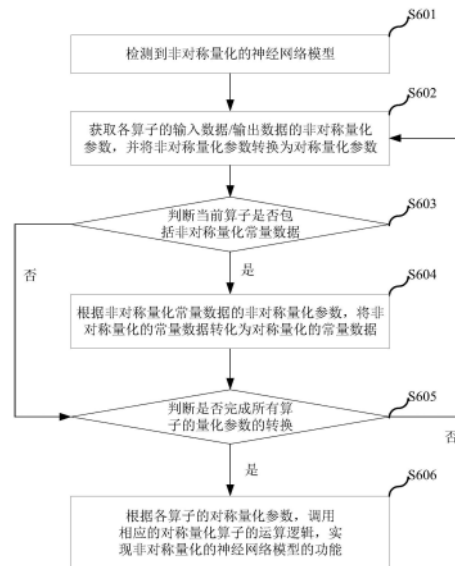
权利要求书3页 说明书19页 附图6页

(54) 发明名称

神经网络模型的运行方法、可读介质和电子设备

(57) 摘要

本申请涉及人工智能领域,公开了一种神经网络模型的运行方法、可读介质和电子设备。该方法应用于电子设备,包括:检测到第一神经网络模型,第一神经网络模型为非对称量化的神经网络模型,并且第一神经网络模型中包括第一神经网络模型的各算子的非对称量化参数;将各算子的非对称量化参数转换为对称量化参数;利用各算子的对称量化参数,调用预设的对称量化的算子的运算逻辑,得到第一神经网络模型的推理结果。如此,只能运行对称量化的神经网络模型的电子设备,可以调用预设的对称量化的算子的运算逻辑来实现非对称量化的神经网络模型的功能,增加了电子设备能够运行的神经网络模型的类型。



1. 一种神经网络模型的运行方法,应用于电子设备,其特征在于,包括:

将图像数据输入第一神经网络模型,所述第一神经网络模型为非对称量化的神经网络模型,并且所述第一神经网络模型中包括所述第一神经网络模型的各算子的非对称量化参数;

其中,所述非对称量化参数包括:常量数据的非对称量化参数,所述常量数据的非对称量化参数包括非对称量化常量数据、非对称量化常量数据的非对称量化缩放系数、非对称量化常量数据的非对称量化零点,所述非对称量化常量数据还包括非对称量化的查找表,所述非对称量化的查找表中包括非对称量化的查表索引和各非对称量化的查表索引对应的非对称量化的查表结果;

将各所述算子的非对称量化参数转换为对称量化参数;

其中,所述将各所述算子的非对称量化参数转换为对称量化参数,包括:

根据所述非对称量化常量数据的非对称量化缩放系数和所述非对称量化常量数据的非对称量化零点,确定出各所述非对称量化的查表索引对应的浮点数查表索引;

根据各所述非对称量化的查表索引对应的浮点数查表索引和各所述算子的运算逻辑,确定各所述浮点数查表索引对应的浮点数查表结果;

根据所述浮点数查表索引的对称量化缩放系数得到对称量化的查表索引、根据所述浮点数查表结果的对称量化缩放系数得到对称量化的查表结果,其中,所述浮点数查表索引的对称量化缩放系数基于所述非对称量化的查表索引的数据类型确定、所述浮点数查表结果的对称量化缩放系数基于所述非对称量化的查表结果的数据类型确定;

基于各所述对称量化的查表索引和相对应的对称量化的查表结果,得到对称量化的查找表;

利用各所述算子的对称量化参数,调用预设的对称量化的算子的运算逻辑,得到所述第一神经网络模型的推理结果;

基于所述推理结果确定所述图像数据的类别。

2. 根据权利要求1所述的方法,其特征在于,所述电子设备包括第一处理器,所述第一处理器能够对对称量化的神经网络模型进行推理,不能对非对称量化的神经网络模型进行推理;并且

由所述第一处理器运行所述第一神经网络模型。

3. 根据权利要求1或2所述的方法,其特征在于,所述非对称量化参数还包括以下参数中的至少一种:

输入数据的非对称量化参数,所述输入数据的非对称量化参数包括输入数据的非对称量化缩放系数、输入数据的非对称量化零点;

输出数据的非对称量化参数,所述输出数据的非对称量化参数包括输出数据的非对称量化缩放系数、输出数据的非对称量化零点。

4. 根据权利要求3所述的方法,其特征在于,所述将各所述算子的非对称量化参数转换为对称量化参数,包括:

根据所述输入数据或所述输出数据的数据类型、所述输入数据或所述输出数据的非对称量化参数,确定所述输入数据或所述输出数据对应的浮点数的最大值和最小值;

根据所述输入数据或所述输出数据对应的浮点数的最大值和最小值,确定所述输入数

据或所述输出数据的对称量化缩放系数。

5. 根据权利要求4所述的方法,其特征在于,根据所述输入数据或所述输出数据的数据类型、所述输入数据或所述输出数据的非对称量化参数,确定所述输入数据或所述输出数据对应的浮点数的最大值和最小值,包括:

根据各算子的所述输入数据或所述输出数据的数据类型,确定所述输入数据或所述输出数据的定点数的最大值和最小值;

根据所述输入数据或所述输出数据的非对称量化参数以及所述输入数据或所述输出数据的定点数的最大值和最小值,确定所述输入数据或所述输出数据对应的浮点数的最大值和最小值。

6. 根据权利要求3所述的方法,其特征在于,所述非对称量化常量数据包括非对称量化常数、非对称量化矩阵;并且,所述将各所述算子的非对称量化参数转换为对称量化参数,包括:

根据所述非对称量化常量数据的数据类型、所述非对称量化常量数据的非对称量化缩放系数和所述非对称量化常量数据的非对称量化零点,确定所述非对称量化常量数据对应的浮点数的最大值和最小值;

根据所述非对称量化常量数据对应的浮点数的最大值和最小值,确定所述非对称量化常量数据对应的浮点数的对称量化缩放系数;

根据确定出的所述非对称量化常量数据的对称量化缩放系数,将所述非对称量化常量数据对应的浮点数常量数据转换为对称量化常量数据,其中,所述非对称量化常量数据对应的浮点数常量数据,由所述常量数据的非对称量化参数确定。

7. 根据权利要求6所述的方法,其特征在于,所述根据所述非对称量化常量数据的定点数的数据类型、所述非对称量化常量数据的非对称量化缩放系数和所述非对称量化常量数据的非对称量化零点,确定所述非对称量化常量数据对应的浮点数的最大值和最小值,包括:

根据所述非对称量化常量数据的数据类型,确定所述非对称量化常量数据的定点数的最大值和最小值;

根据所述非对称量化常量数据的非对称量化缩放系数、所述非对称量化常量数据的非对称量化零点以及确定出的所述非对称量化常量数据的定点数的最大值和最小值,确定所述非对称量化常量数据对应的浮点数的最大值和最小值。

8. 根据权利要求1所述的方法,其特征在于,还包括:

根据所述非对称量化的查表索引或所述非对称量化的查表结果的数据类型,确定所述非对称量化的查表索引或所述非对称量化的查表结果对应的定点数的最大值和最小值,并基于确定出的最大值和最小值,根据所述非对称量化常量数据的非对称量化缩放系数、所述非对称量化常量数据的非对称量化零点,确定非对称量化的查表索引对应的符点数查表索引的最大值和最小值或所述非对称量化的查表结果对应的符点数查表结果的最大值和最小值;

根据确定出的非对称量化的查表索引对应的符点数查表索引的最大值和最小值或所述非对称量化的查表结果对应的符点数查表结果的最大值和最小值,确定所述浮点数查表索引或所述浮点数查表结果的对称量化缩放系数。

9. 一种可读介质,其特征在于,所述可读介质中包含有指令,当所述指令被电子设备的处理器执行时使电子设备实现权利要求1至8中任一项所述的神经网络模型的运行方法。

10. 一种电子设备,其特征在于,包括:

存储器,用于存储由电子设备的一个或多个处理器执行的指令;

以及处理器,是所述电子设备的处理器之一,用于运行所述指令以使所述电子设备实现权利要求1至8中任一项所述的神经网络模型的运行方法。

## 神经网络模型的运行方法、可读介质和电子设备

### 技术领域

[0001] 本申请涉及人工智能领域,特别涉及一种神经网络模型的运行方法、可读介质和电子设备。

### 背景技术

[0002] 随着人工智能(artificial intelligence, AI)的迅速发展,神经网络模型在人工智能领域的应用越来越广泛。由于运行神经网络模型的运算单元,例如神经网络处理器(Neural-Network Processing Unit, NPU),通常为定点运算单元,为提高神经网络模型的运行速度,通常将神经网络模型的各算子进行量化,得到定点运算的神经网络模型,再由电子设备来运行。对神经网络模型各算子的量化包括非对称量化或对称量化,但是,为了节省开发、制造成本,部分NPU中只预设有针对性对称量化的算子的运算逻辑,该类NPU只能运行对称量化的神经网络模型,而无法运行非对称量化的神经网络模型。

### 发明内容

[0003] 有鉴于此,本申请实施例提供了神经网络模型的运行方法、可读介质和电子设备。电子设备通过将非对称量化的神经网络模型的非对称量化参数转换为对称量化参数,即可根据得到的对称量化参数调用预设的对称量化的算子的运算逻辑来实现非对称量化的神经网络模型的功能,增加了电子设备能够运行的神经网络模型的类型,提高了电子设备的NPU的通用性。

[0004] 第一方面,本申请实施例提供了一种神经网络模型的运行方法,应用于电子设备,该方法包括:检测到第一神经网络模型,第一神经网络模型为非对称量化的神经网络模型,并且第一神经网络模型中包括第一神经网络模型各算子的非对称量化参数;将各算子的非对称量化参数转换为对称量化参数;利用各算子的对称量化参数,调用预设的对称量化的算子的运算逻辑,得到第一神经网络模型的推理结果。

[0005] 通过本申请实施例提供的方法,若电子设备中用于运行第一神经网络模型的处理器(例如NPU)只能调用对称量化的算子的运算逻辑来实现神经网络模型的推理,则可以通过将第一神经网络模型各算子的非对称量化参数转换为对称量化参数,并通过得到的对称量化参数来调用对称量化的算子的运算逻辑,来对第一神经网络模型进行推理,得到第一神经网络模型的推理结果。如此,增加了电子设备能够运行的神经网络模型的类型,提高了电子设备的NPU的通用性。此外,由于过程中无需先将非对称量化的神经网络模型转换为浮点型的神经网络模型,再转换为对称量化的神经网络模型,提高了神经网络模型的部署速度。

[0006] 在上述第一方面的一种可能实现中,上述电子设备包括第一处理器,第一处理器能够对对称量化的神经网络模型进行推理,不能对非对称量化的神经网络模型进行推理;并且由第一处理器运行第一神经网络模型。

[0007] 也即是说,电子设备的第一处理器只能够调用对称量化的算子的运算逻辑来对神

神经网络模型进行推理,通过本申请实施例提供的方法,第一处理器可以通过将第一神经网络模型的非对称量化参数转换为对称量化参数,并基于得到的对称量化参数调用对称量化的算子的运算逻辑来实现第一神经网络模型的功能,增加了第一处理器能够运行的神经网络模型的类型,提高了第一处理器的通用性。

[0008] 在上述第一方面的一种可能实现中,上述非对称量化参数包括以下参数中的至少一种:输入数据的非对称量化参数,输入数据的非对称量化参数包括输入数据的非对称量化缩放系数、输入数据的非对称量化零点;输出数据的非对称量化参数,输出数据的非对称量化参数包括输出数据的非对称量化缩放系数、输出数据的非对称量化零点;常量数据的非对称量化参数,常量数据的非对称量化参数包括非对称量化常量数据、非对称量化常量数据的非对称量化缩放系数、非对称量化常量数据的非对称量化零点。

[0009] 在上述第一方面的一种可能实现中,上述将各算子的非对称量化参数转换为对称量化参数,包括:根据输入数据或输出数据的数据类型、输入数据或输出数据的非对称量化参数,确定输入数据或输出数据对应的浮点数的最大值和最小值;根据输入数据或输出数据对应的浮点数的最大值和最小值,确定输入数据或输出数据的对称量化缩放系数。

[0010] 例如,在一个算子的输入数据或输出数据的数据类型为UINT8时,若输入数据或输出数据的非对称量化缩放系数为1、零点为0,则确定输入数据或输出数据的对应的浮点数的最大值为浮点数255,最小值为浮点数0,进而根据下文中的公式(2)得到输入数据或输出数据的对称量化缩放系数为0.498。

[0011] 在上述第一方面的一种可能实现中,上述根据输入数据或输出数据的数据类型、输入数据或输出数据的非对称量化参数,确定输入数据或输出数据对应的浮点数的最大值和最小值,包括:根据各算子的输入数据或输出数据的数据类型,确定输入数据或输出数据的定点数的最大值和最小值;根据输入数据或输出数据的非对称量化参数以及输入数据或输出数据的定点数的最大值和最小值,确定输入数据或输出数据对应的浮点数的最大值和最小值。

[0012] 例如,在输入数据或输出数据的数据类型为UINT8的情况下,则输入数据或输出数据的定点数的最大值为255,最小值为0,从而可以根据下文中的公式(11)确定输入数据或输出数据对应的浮点数的取值范围为浮点数 $[0, 255]$ ,即是输入数据或输出数据对应的浮点数的最大值为浮点数255,最小值为浮点数0。

[0013] 在上述第一方面的一种可能实现中,上述非对称量化常量数据包括非对称量化常数、非对称量化矩阵;并且,将各算子的非对称量化参数转换为对称量化参数,包括:根据非对称量化常量数据的定点数的数据类型、非对称量化常量数据的非对称量化缩放系数和非对称量化常量数据的非对称量化零点,确定非对称量化的常量数据对应的浮点数的最大值和最小值;根据非对称量化常量数据对应的浮点数的最大值和最小值,确定非对称量化常量数据对应的浮点数的对称量化缩放系数;根据确定出的非对称量化常量数据的对称量化缩放系数,将非对称量化常量数据对应的浮点数常量数据转换为对称量化常量数据,其中,非对称量化常量数据对应的浮点数常量数据,由常量数据的非对称量化参数确定。

[0014] 例如,某一算子的常量数据包括非对称量化的常数100,该非对称量化常量数据的非对称量化缩放系数为2、非对称量化零点为0,则根据公式(11)可以得到该常数对应的浮点数常数为50;基于该常量数据的数据类型,例如UINT8,可以得到常量数据对应的浮点数

的最大值为127.5,最小值为0,进而根据下文中的公式(2)可以得到该常量数据的对称量化缩放系数为0.9961。

[0015] 在上述第一方面的一种可能实现中,上述根据非对称量化常量数据的定点数的数据类型、非对称量化常量数据的非对称量化缩放系数和非对称量化常量数据的非对称量化零点,确定非对称量化常量数据对应的浮点数的最大值和最小值,包括:根据非对称量化常量数据的数据类型,确定非对称量化常量数据的定点数的最大值和最小值;根据非对称量化常量数据的非对称量化缩放系数、非对称量化常量数据的非对称量化零点以及确定出的非对称量化常量数据的定点数的最大值和最小值,确定非对称量化常量数据对应的浮点数的最大值和最小值。

[0016] 例如,在非对称量化常量数据的数据类型为UINT8的情况下,若该非对称量化常量数据的非对称量化缩放系数为2、非对称量化零点为0,则非对称量化常量数据的定点数的最大值为255,最小值为0,从而可以根据下文中的公式(11)确定非对称量化常量数据对应的浮点数的取值范围为浮点数[0,255],即是浮点数非对称量化常量数据最大值为浮点数255,最小值为浮点数0。

[0017] 在上述第一方面的一种可能实现中,上述非对称量化常量数据还包括非对称量化的查找表,非对称量化的查找表中包括非对称量化的查表索引和各非对称量化的查表索引对应的非对称量化的查表结果;并且,将各算子的非对称量化参数转换为对称量化参数,包括:根据非对称量化常量数据的非对称量化缩放系数和非对称量化常量数据的非对称量化零点,确定出各非对称量化的查表索引对应的浮点数查表索引;根据各非对称量化的查表索引对应的浮点数查表索引和各算子的运算逻辑,确定各浮点数查表索引对应的浮点数查表结果;根据浮点数查表索引的对称量化缩放系数得到对称量化的查表索引、根据浮点数数据查表结果的对称量化缩放系数得到对称量化的查表结果,其中,浮点数查表索引的对称量化缩放系数基于非对称量化的查表索引的数据类型确定、浮点数查表结果的对称量化缩放系数基于非对称量化的查表结果的数据类型确定;基于各对称量化的查表索引和相对应的对称量化的查表结果,得到对称量化的查找表。

[0018] 例如,对于下文中的Softmax算子,电子设备可以先根据查找表LUT中的查表索引的数据类型(UINT8),查表索引的非对称量化缩放系数( $2.2 \times 10^{-5}$ )、非对称量化零点(0),确定各查表索引对应的浮点数查表索引,例如对于查表索引[59,104,182],根据公式(11)可以得到浮点数查表索引为[2681818,4727273,8272727];再将浮点数查表索引代入到Softmax算子的运算逻辑(公式(12))中,得到浮点数查表索引对应的浮点数查表结果[0.0069,0.042,0.9511],再根据浮点数查表结果的对称量化缩放系数( $1.1 \times 10^{-5}$ ),将浮点数查表结果量化为定点数查表结果[1,5,121]。

[0019] 在上述第一方面的一种可能实现中,上述方法还包括:根据非对称量化的查表索引或非对称量化的查表结果的数据类型,确定非对称量化的查表索引或非对称量化的查表结果对应的定点数的最大值和最小值,并基于确定出的最大值和最小值,根据非对称量化常量数据的非对称量化缩放系数、非对称量化常量数据的非对称量化零点,确定非对称量化的查表索引对应的浮点数查表索引的最大值和最小值或非对称量化的查表结果对应的浮点数查表结果的最大值和最小值;根据确定出的非对称量化的查表索引对应的浮点数查表索引的最大值和最小值或非对称量化的查表结果对应的浮点数查表结果的最大值和

最小值,确定浮点数查表索引或浮点数查表结果的对称量化缩放系数。

[0020] 例如,在下文中的Softmax算子的非对称量化的查表结果、非对称量化的查表结果的数据类型为UINT8的情况下,非对称量化的查表结果的定数的最大值为255,最小值为1,根据非对称量化的查表索引的非对称量化缩放系数( $2.2 \times 10^{-5}$ )、非对称量化的查表索引的非对称量化零点(0),非对称量化的查表结果的非对称量化缩放系数(255)、非对称量化的查表结果非对称量化零点(0),基于下文公式(11)可以得到非对称量化查表索引对应的浮点数的最大值为11590909、最小值为0,非对称量化查表结果对应的浮点数的最大值为1、最小值为0,进而基于下文中的公式(2)可以得到浮点数查表索引的对称量化缩放系数为 $1.1 \times 10^{-5}$ 、浮点数查表结果的对称量化缩放系数为127。

[0021] 第二方面,本申请实施例提供了一种可读介质,该可读介质中包含有指令,当指令被电子设备的处理器执行时使电子设备实现上述第一方面及上述第一方面的各种可能实现提供的任意一种神经网络模型的运行方法。

[0022] 第三方面,本申请实施例提供了一种电子设备,该电子设备包括:存储器,用于存储由电子设备的一个或多个处理器执行的指令;以及处理器,是电子设备的处理器之一,用于运行指令以使电子设备实现上述第一方面及上述第一方面的各种可能实现提供的任意一种神经网络模型的运行方法。

## 附图说明

[0023] 图1A根据本申请的一些实施例,示出了一种8位对称量化的示意图;

[0024] 图1B根据本申请的一些实施例,示出了一种8位非对称量化的示意图;

[0025] 图2根据本申请的一些实施例,示出了一种将非对称量化的神经网络模型部署到电子设备100中的场景示意图;

[0026] 图3根据本申请的一些实施例,示出了一种电子设备100运行非对称量化的神经网络模型的场景图;

[0027] 图4根据本申请的一些实施例,示出了一种神经网络模型10的结构示意图;

[0028] 图5根据本申请的一些实施例,示出了一种电子设备利用非对称量化的神经网络模型10对图像20进行分类的过程示意图;

[0029] 图6根据本申请的一些实施例,示出了一种神经网络模型运行方法的流程示意图;

[0030] 图7根据本申请的一些实施例,示出了一种电子设备100调用对称量化的算子的运算逻辑对图像20进行分类的过程示意图;

[0031] 图8根据本申请的一些实施例,示出了一种电子设备100的结构示意图。

## 具体实施方式

[0032] 本申请的说明性实施例包括但不限于神经网络模型的运行方法、可读介质和电子设备。

[0033] 为了便于理解,首先介绍本申请实施例涉及的术语。

[0034] (1) 对称量化

[0035] 对称量化,即是将浮点数转换为取值范围为 $[-2^{n-1}, 2^{n-1}-1]$ 的有符号整形数(integral numeric types, INT),其中n为对称量化的位数。假设待量化的浮点数为 $x_f$ ,量



化目标为将 $x_f$ 进行 $n$ 位对称量化,即量化后的定点数的取值范围为,对称量化过程表示为如下公式(1)。

$$[0036] \quad x_q = \text{round}\left(x_f \cdot \frac{2^{n-1} - 1}{\max(\text{abs}(\max(x_f)), \text{abs}(\min(x_f)))}\right) \quad (1)$$

[0037] 公式(1)中, $\text{abs}()$ 为求绝对值函数, $\max()$ 为求最大值函数, $\min()$ 为求最小值函数, $\text{round}$ 为求四舍五入函数, $x_q$ 为定点数。此外,在公式(1)中,如下公式(2)所示 $S_c$ 项可以称为对称量化缩放系数。也即是说,对于一个对称量化的定数,可以根据该定点数对应的对称量化缩放系数确定该定点数对应的浮点数。

$$[0038] \quad S_c = \frac{2^{n-1} - 1}{\max(\text{abs}(\max(x_f)), \text{abs}(\min(x_f)))} \quad (2)$$

[0039] 具体地,图1A根据本申请的一些实施例,示出了一种对浮点数 $x_f$ 进行8位对称量化的示意图。参考图1A,量化的目标为将 $x_f$ 量化进行8位对称量化,即是将 $x_f$ 转换为INT8型(取值范围为 $[-128, 127]$ ),假设 $x_f$ 的绝对值的最大值为 $\max(|x_f|)$ ,则对 $x_f$ 进行8位对称量化的过程即是区间 $[-\max(|x_f|), \max(|x_f|)]$ 映射到区间 $[-128, 127]$ 中。

[0040] (2) 非对称量化

[0041] 非对称量化,即是将浮点数转换为取值范围为 $[0, 2^n - 1]$ 的无符号整形(unsigned integral numeric types, UIN)数,其中 $n$ 为非对称量化的位数。假设待量化的浮点数为 $x_f$ ,量化目标为将 $x_f$ 进行 $n$ 位非对称量化,即量化后的定点数的取值范围为 $[0, 2^n - 1]$ ,则非对称量化过程表示为如下公式(3)。

$$[0042] \quad x_q = \text{round}\left((x_f - \min(x_f)) \frac{2^n - 1}{\max(x_f) - \min(x_f)}\right) \quad (3)$$

[0043] 公式(3)中, $\max()$ 为求最大值函数, $\min()$ 为求最小值函数, $\text{round}$ 为求四舍五入函数, $x_q$ 为定点数。此外,在公式(3)中,如下公式(4)所示 $AS_c$ 项可以称为非对称量化缩放系数,如下公式(5)所示的 $Z_p$ 项可以称为非对称量化零点。也即是说,对于一个非对称量化的定点数,可以根据该定点数对应的非对称量化缩放系数和非对称量化零点,确定该定点数对应的浮点数。

$$[0044] \quad AS_c = \frac{2^n - 1}{\max(x_f) - \min(x_f)} \quad (4)$$

$$[0045] \quad Z_p = \text{round}(-AS_c \cdot \min(x_f)) \quad (5)$$

[0046] 基于公式(3)至公式(5),可以得到公式(6)所示的非对称量化的另一种表示方式:

$$[0047] \quad x_q = \text{round}(x_f \cdot AS_c + Z_p) \quad (6)$$

[0048] 具体地,图1B根据本申请的一些实施例,示出了一种对浮点数 $x_f$ 进行8位非对称量化的示意图。参考图1B,量化的目标为将 $x_f$ 量化进行8位非对称量化,即是将 $x_f$ 转换为UINT8型(取值范围为 $[0, 255]$ ),假设 $x_f$ 的最大值为 $\max(x_f)$ ,最小值为 $\min(x_f)$ ,则对 $x_f$ 进行8位非对称量化的过程即是区间 $[\min(x_f), \max(x_f)]$ 映射到区间 $[0, 255]$ 中。

[0049] 可以理解,在一些实施例中UINT型可以表示为UINT $n$ ,其中 $n$ 可以取4、8、16、32等,也可以取其他整数,并且UINT $n$ 型的无符号数据的取值范围为 $[0, 2^n - 1]$ 。即是对于一个给定的UINT型数据,电子设备可以根据 $n$ 的值确定该类数据的取值范围,并基于公式(6)计算出

该数据对应的浮点数的取值范围。

[0050] (3) 神经网络模型量化

[0051] 神经网络模型量化,即是将神经网络模型各算子中的输入数据、输出数据、常量数据从大数据类型的浮点数(例如,32位浮点数)转换为较小数据类型的定点数(例如,4/8/16位定点数),并且定点数的位数通常和运行神经网络模型的运算单元,例如NPU,所支持的定点数位数相匹配,以提高NPU运行神经网络模型的速度。

[0052] 一般地,神经网络模型的量化的过程,即是根据各算子浮点数的输入数据、浮点数的输出数据的取值范围,以及量化的定点数的类型(例如UINT8),确定各算子的浮点数输入数据、浮点数输出数据和浮点数常量数据的量化参数(例如,非对称量化的量化参数包括非对称量化缩放系数和非对称量化零点,对称量化的量化参数包括对称量化缩放系数)。也即经过量化后的算子的量化参数中,包括了该算子的输入数据的量化参数、输出数据的量化参数、量化后的常量数据以及常量数据的量化参数。而用于运行量化后的神经网络模型的电子设备中,预设有量化后算子的运算逻辑。该量化后的算子的运算逻辑,以定点数为输入、定点数为输出,并且该运算逻辑中的参数包括输入数据的量化参数、输出数据的量化参数、量化后的常量数据以及常量数据的量化参数,电子设备在运行一个量化后的算子时,根据该算子的量化参数,调用预设的运算逻辑,即可通过定点运算实现该算子的功能。

[0053] 下面结合附图介绍本申请实施例的技术方案。

[0054] 为便于理解,下面先以卷积算子为例,说明采用对称量化和非对称量化的算子的运算逻辑。

[0055] 假设 $B_f$ 为浮点数输入矩阵, $C_f$ 为浮点数卷积核, $D_f$ 为浮点数卷积结果,则 $B_f$ 、 $C_f$ 、 $D_f$ 间的关系可以表示为公式(7)。

[0056]  $D_f = B_f * C_f$  (7)

[0057] 公式(7)中,“\*”为卷积运算符号,卷积运算的具体过程将在后文进行介绍,在此不做赘述。

[0058] 根据公式(1)和公式(2)可以得到 $B_f = B_q / B_{Sc}$ ,  $C_f = C_q / C_{Sc}$ ,  $D_f = D_q / D_{Sc}$ ,其中: $B_q$ 为 $B_f$ 对应的对称量化的定点数矩阵, $B_{Sc}$ 为将 $B_f$ 量化为 $B_q$ 的对称量化缩放系数; $C_q$ 为 $C_f$ 对应的对称量化的定点数卷积核, $C_{Sc}$ 为将 $C_f$ 量化为 $C_q$ 的对称量化缩放系数; $D_q$ 为 $D_f$ 对应的对称量化的定点数矩阵, $D_{Sc}$ 为将 $D_f$ 对称量化为 $D_q$ 的对称量化缩放系数。进而公式(7)可以表示为如下公式(8)。

[0059] 
$$\frac{D_q}{D_{Sc}} = \frac{B_q}{B_{Sc}} * \frac{C_q}{C_{Sc}}$$
 (8)

[0060] 对公式(8)进行变形可以得到如下公式(9)所示的对称量化的卷积算子的运算逻辑。

[0061] 
$$D_q = B_q * C_q * \frac{D_{Sc}}{B_{Sc} * C_{Sc}}$$
 (9)

[0062] 也即是说电子设备中预设有公式(9)所示的运算逻辑,该运算逻辑的输入包括了定点数输入矩阵 $B_q$ ,定点数卷积核 $C_q$ ,以及定点数输入数据 $B_q$ 的对称量化参数 $B_{Sc}$ ,定点数输出数据 $D_q$ 的对称量化参数 $D_{Sc}$ 、卷积核的对称量化参数 $C_{Sc}$ ,输出为定点数卷积结果 $D_q$ 。

[0063] 类似地,非对称量化的卷积算子的运算逻辑可以表示表如下公式(10),具体推导过程可以参考对称量化的卷积算子的运算逻辑,在此不做赘述。

$$[0064] \quad D_q = \text{round} \left\{ (B_q - Z_p - B) * (C_q - Z_p - C) \times \frac{D\_ASc}{B\_ASc \times C\_ASc} + Z_p - D \right\} \quad (10)$$

[0065] 公式(10)中, $Z_p - B$ 为输入数据 $B_q$ 的非对称量化零点、 $Z_p - C$ 为卷积核 $C_q$ 的非对称量化零点、 $Z_p - D$ 为输出数据 $D_q$ 的非对称量化零点、 $B\_ASc$ 为输入数据 $B_q$ 的非对称量化缩放系数、 $C\_ASc$ 为卷积核 $C_q$ 的非对称量化缩放系数、 $D\_ASc$ 为输出数据 $D_q$ 的非对称量化缩放系数。

[0066] 从公式(9)和公式(10)中可以看出,对于同一算子,采用非对称量化和对称量化的运算逻辑和输入参数并不相同。如前所述,部分NPU只能运行对称量化的神经网络模型,即只能运行对称量化的算子的运算逻辑(例如公式(9)所示的卷积算子的运算逻辑),而无法运行非对称量化的神经网络模型。若将非对称量化的神经网络模型部署到包括该部分NPU的电子设备中,需要先将非对称量化的神经网络模型转换浮点型神经网络模型(即神经网络模型的各算子的输入数据、输出数据及常量数据都为浮点数),再将浮点型神经网络模型量化为对称量化的神经网络后,才能部署到该电子设备中。

[0067] 例如,参考图2,电子设备100的NPU只能运行对称量化的神经网络模型,而待运行的神经网络模型为非对称量化的神经网络模型。从而需要由电子设备200将非对称量化的神经网络模型转换为浮点型神经网络模型,并将浮点型神经网络模型量化为对称量化的神经网络模型后,再将对称量化后的神经网络模型部署到电子设备100中,由电子设备100的NPU来运行。由于将浮点数神经网络模型量化为定点数神经网络模型(即神经网络模型的各算子的输入数据、输出数据、常量数据等为定点数)的过程中占用大量计算资源,耗时较长,不利于神经网络模型的快速部署。

[0068] 为了解决上述问题,本申请实施例提供了一种神经网络模型的运行方法,电子设备100在检测到非对称量化的神经网络模型后,获取该非对称量化的神经网络模型的各算子的非对称量化参数,并将各算子非对称量化参数转换为对应的对称量化参数,电子设备100的NPU再根据各算子的对称量化参数调用预设的对称量化的算子的运算逻辑,来实现该非对称量化的神经网络模型的相关功能。如此,参考图3,虽然电子设备100的NPU并不能运行非对称量化的神经网络模型,但电子设备100通过将非对称量化的神经网络模型中的各算子的非对称量化参数转换为对称量化参数,电子设备100的NPU即可调用预设的对称量化的算子的运算逻辑,来实现该非对称量化神经网络模型的相关功能,而无需由其他电子设备将非对称量化的神经网络模型转换为浮点型神经网络模型,并将浮点型神经网络模型转化为对称量化的神经网络模型后,再由电子设备100的NPU来运行,增加了电子设备100能够运行的神经网络模型的类型,提高了NPU的通用性,提升了神经网络模型的部署速度。

[0069] 可以理解,各算子的非对称量化参数包括以下参数中的至少一项:输入数据的非对称量化缩放系数、输入数据的非对称量化零点、输出数据的非对称量化缩放系数、输出数据的非对称量化零点、非对称量化常量数据、非对称量化常量数据的非对称量化零点、非对称量化常量数据的非对称量化缩放系数。各算子的对称量化参数包括以下参数中的至少一项:输入数据的对称量化缩放系数、输出数据的对称量化缩放系数、对称量化常量数据、对称量化常量数据的对称量化缩放系数。

[0070] 可以理解,在另一些实施例中,各算子的非对称量化参数/对称量化参数也可以包括更多或更少的参数,在此不做限定。

[0071] 可以理解,预设的对称量化的算子的运算逻辑可以由电子设备100的NPU的开发商预设于NPU中,也可以由NPU的开发商提供给电子设备100的开发商,由电子设备100的开发商预设于电子设备100的存储器中。电子设备100的NPU可以根据基于算子的对称量化的输入数据及对称量化的输入数据的对称量化参数、输出数据的对称量化参数、对称量化常量数据调用对称量化的算子的运算逻辑,得到算子的对称量化的定点数输出数据。

[0072] 具体地,在一些实施例中,对上述公式(6)进行变形可以得到如公式(11)所示的、将非对称量化的定点数转化为浮点数的计算公式。

$$[0073] \quad x_f = \frac{x_q - Z_p}{AS_c} \quad (11)$$

[0074] 电子设备100可以根据非对称量化的算子的定点数输入数据的数据类型确定定点数输入数据的取值范围(例如若输入数据为UINT8,说明取值范围为定点数[0,255],最大值为255,最小值为0),基于非对称量化的输入数据的最大值和最小值、非对称量化的输入数据的非对称量化缩放系数和非对称量化零点,基于公式(11)确定输入数据对应的浮点数的最大值和最小值;再基于输入数据对应的浮点数的最大值和最小值,根据公式(2)确定输入数据的对称量化缩放系数。

[0075] 类似地,电子设备100可以根据非对称量化的算子的输出数据的类型确定定点数输出数据的取值范围(例如若输入数据为UINT8,说明取值范围为定点数[0,255],最大值为255,最小值为0),并根据非对称量化的输出数据的最大值和最小值、输出数据的非对称量化缩放系数和非对称量化零点,基于公式(11)确定输出数据对应的浮点数的最大值和最小值;再基于确定的浮点数的最大值和最小值,根据公式(2)确定输出数据的对称量化缩放系数。

[0076] 类似地,电子设备100也可以根据非对称量化的算子中非对称量化常量数据的数据类型,确定非对称量化常量数据的定点数的最大值和最小值,再基于确定的最大值和最小值、非对称量化常量数据的非对称量化缩放系数和非对称量化零点,根据公式(11)确定非对称量化常量数据对应的浮点数的最大值和最小值;再基于确定的浮点数的最大值和最小值,根据公式(2)确定常量数据的对称量化缩放系数;然后,根据公式(11)将非对称量化常量数据转换对应的浮点数,再基于确定的对称量化缩放系数和非对称量化常量数据对应的浮点数常量数据,根据公式(1)和公式(2)将非对称量化常量数据对应的浮点数常量数据转换为对称量化常量数据。

[0077] 例如,假设某一量化为UINT8的非对称量化的卷积算子的输入数据矩阵B<sub>f</sub>中的各

元素的取值范围为浮点数[0,1],浮点数卷积核C<sub>f</sub> =  $\begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.1 \end{bmatrix}$ ,则可以得到浮点数卷

积结果D<sub>f</sub>中的各元素的取值范围为[0,0.7]。

[0078] 根据公式(4)和公式(5),可以得到输入数据的非对称量化零点Z<sub>p</sub>B=0、非对称量化缩放系数B<sub>ASc</sub>=255,输出数据的非对称量化零点Z<sub>p</sub>D=0、非对称量化缩放系数D<sub>ASc</sub>

=364.29。假设卷积核C<sub>f</sub>的量化参数与输入数据B<sub>f</sub>的非对称量化参数相同(即C<sub>ASc</sub>=

255, Z<sub>p\_C</sub>=0), 根据公式(6)可以得到C<sub>q</sub>=
$$\begin{bmatrix} 26 & 0 & 0 \\ 0 & 128 & 0 \\ 0 & 0 & 26 \end{bmatrix}$$
。也即是该非对称量化的卷积算

子的非对称量化参数包括:输入数据的非对称量化零点Z<sub>p\_B</sub>、非对称量化缩放系数B<sub>ASc</sub>, 输出数据的非对称量化零点Z<sub>p\_D</sub>、非对称量化缩放系数D<sub>ASc</sub>,非对称量化常量数据C<sub>q</sub>,以及非对称量化常量数据的非对称量化缩放系数C<sub>ASc</sub>、非对称量化零点Z<sub>p\_C</sub>。

[0079] 电子设备在检测到上述非对称量化的卷积算子时,先根据输入数据的数据类型UINT8,确定定点数输入数据D<sub>q</sub>的取值范围为[0,255],最大值为255,最小值为0;将x<sub>q</sub>=255、Z<sub>p\_B</sub>=0、B<sub>ASc</sub>=255代入公式(11)得到输入数据对应的浮点数最大值为1;将x<sub>q</sub>=0、Z<sub>p\_B</sub>=0、B<sub>ASc</sub>=255代入公式(11)得到输入数据对应的浮点数最小值为0,进而根据公式(2)可以得到输入数据的对称量化缩放系数B<sub>Sc</sub>=(2<sup>8-1</sup>-1)/1=127。

[0080] 类似地,电子设备可以根据输出数据的数据类型UINT8,确定定点数输出数据D<sub>q</sub>的取值范围为[0,255],最大值为255,最小值为0;将x<sub>q</sub>=255、Z<sub>p\_D</sub>=0、D<sub>ASc</sub>=364.29代入公式(11)得到输出数据对应的浮点数最大值为(255-0)/364.29=0.7;将x<sub>q</sub>=0、Z<sub>p\_D</sub>=0、D<sub>ASc</sub>=364.29代入公式(11)得到输出数据对应的浮点数最小值为(0-0)/364.29=0,进而根据公式(2)可以得到输出数据的对称量化缩放系数D<sub>Sc</sub>=(2<sup>8-1</sup>-1)/0.7=181.43。

[0081] 由于卷积核C<sub>q</sub>的非对称量化参数与输入数据B<sub>q</sub>的非对称量化参数相同,卷积核C<sub>q</sub>的对称量化参数也应当与输入数据B<sub>q</sub>的对称量化参数相同,即是卷积核C<sub>q</sub>的对称量

化缩放系数C<sub>Sc</sub>=127。将C<sub>ASc</sub>=255、Z<sub>p\_C</sub>=0、C<sub>q</sub>代入公式(11)可以得到C<sub>f</sub>=
$$\begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.1 \end{bmatrix}$$
,

再根据C<sub>f</sub>和C<sub>Sc</sub>=127、公式(1)和公式(2)可以得到对称量化的卷积核C<sub>q'</sub>=
$$\begin{bmatrix} 13 & 0 & 0 \\ 0 & 64 & 0 \\ 0 & 0 & 13 \end{bmatrix}$$
。

[0082] 进而电子设备100的NPU可以根据B<sub>Sc</sub>=127、D<sub>Sc</sub>=181.43、C<sub>Sc</sub>=127以及输入到该卷积算子的定点数输入数据B<sub>q</sub>,即可调用预设的公式(9)所示的对称量化的卷积算子的运算逻辑来实现前述非对称量化的卷积算子的功能,增加了电子设备100的NPU能够运行的神经网络模型的类型,提高了NPU的通用性。

[0083] 下面结合具体的神经网络模型介绍本申请实施例的技术方案。

[0084] 图4根据本申请的一些实施例,示出了一种神经网络模型10的结构示意图。如图4所示,神经网络模型10为非对称量化的神经网络模型,包括输入层11、卷积算子12、全连接算子13、Softmax算子14和输出层15,用于对输入的图像数据进行分类。其中,输入层11用于对输入的图像进行预处理,将输入的图像数据转换为非对称量化的输入数据,例如UINT8型的数据;卷积算子12用于对非对称量化的输入数据进行卷积运算,得到输入的图像对应的非对称量化的特征矩阵;全连接算子13用于对非对称量化的特征矩阵进行全连接运算,得到输入的图像数据属于各预设类别的得分;Softmax算子14,用于根据输入的图像数据属于各预设类别的得分,得到输入的图像数据属于各预设类别的概率;输出层15用于根据输入的图像数据属于各预设类别的概率,确定出输入的图像数据的类别,例如以输入的图像数

据为各预设类别的概率最大的类别作为该输入的图像数据的类别。

[0085] 进一步,图5根据本申请实施例,示出了电子设备利用神经网络模型10对图像20进行分类的过程示意图。

[0086] 参考图5,电子设备先利用输入层11对图像20进行预处理,得到UINT8的图像矩阵H,输入层11的非对称量化缩放系数为 $ASc\_out1=1$ ,非对称量化零点为 $Z_{p\_out1}=0$ 。

[0087] 其次,利用卷积算子12将图像矩阵H分别与卷积核 $K_i$  ( $i=1,2,3$ )进行卷积运算,得到三个特征矩阵 $A_i$  ( $i=1,2,3$ ),卷积算子12的输入数据的非对称量化缩放系数为 $ASc\_in2=1$ 、输入数据的非对称量化零点 $Z_{p\_in2}=0$ 、输出数据的非对称量化缩放系数为 $ASc\_out2=0.0833$ 、输出数据的非对称量化零点为 $Z_{p\_out2}=0$ 、非对称量化常量数据包括卷积核 $K_i$  ( $i=1,2,3$ ) (卷积核的非对称量化参数与卷积算子12的输入数据的非对称量化参数相同)。

[0088] 再利用全连接算子13对特征矩阵 $A_i$  ( $i=1,2,3$ )进行全连接运算,例如分别将特征矩阵 $A_i$  ( $i=1,2,3$ )与权重矩阵W做内积运算,得到图像20为预设类别(兔子/狗/猫)的得分,全连接算子13的输入数据的非对称量化缩放系数 $ASc\_in3=0.08333$ 、输入数据的非对称量化零点 $Z_{p\_in3}=0$ 、输出数据的非对称量化缩放系数为 $ASc\_out3=2.2 \times 10^{-5}$ 、输出数据的非对称量化零点为 $Z_{p\_out3}=0$ 、非对称量化常量数据包括权重矩阵W(权重矩阵的非对称量化参数与全连接算子13的输入数据的非对称量化参数相同)。

[0089] 然后,利用Softmax算子14根据图像20为预设类别的得分,从查找表LUT中获取图像20为预设类别的概率,Softmax算子14的输入数据的非对称量化缩放系数 $ASc\_in4=2.2 \times 10^{-5}$ 、输入数据的非对称量化零点 $Z_{p\_in4}=0$ 、输出数据的非对称量化缩放系数为 $ASc\_out4=255$ 、输出数据的非对称量化零点为 $Z_{p\_out4}=0$ 、非对称量化常量数据包括查找表LUT。

[0090] 最后,利用输出层15根据图像20为各预设类别的概率,确定图像20的类别,例如比较图像20为各预设类别的概率,将图像20为预设类别的概率中最大的概率确定为图像20的类别(猫)。

[0091] 下面结合图4所示的神经网络模型10的结构和图5所示的神经网络模型10对图像20进行分类的过程,介绍本申请实施例的技术方案。

[0092] 具体地,图6根据本申请的一些实施例,示出了一种神经网络模型的运行方法的流程示意图。该方法的执行主体为电子设备100,如图6所示,该流程包括如下步骤。

[0093] S601:检测到非对称量化的神经网络模型。

[0094] 电子设备100在检测待运行的神经网络模型为非对称量化的神经网络模型的情况下,触发本申请实施例提供的神经网络模型的运行方法。

[0095] 在一些实施例中,电子设备100可以根据待运行神经网络模型中的数据的数据类型,来确定该神经网络是否是非对称量化的神经网络模型。具体地,电子设备100在检测到待运行神经网络模型中的数据的数据类型为UINT(UINT包括但不限于UINT4、UINT8、UINT16、UINT32等)时,例如检测到神经网络模型10中的数据的数据类型为UINT8时,确定待运行的神经网络模型为非对称量化的神经网络模型。

[0096] 在一些实施例,电子设备100也可以根据待运行神经网络模型各算子的量化参数来确定待运行神经网络模型是否是非对称量化的神经网络模型。例如,电子设备100可以在检测到待运行神经网络模型的量化参数包括缩放系数和零点时,确定该神经网络模型为非

对称量化的神经网络模型。

[0097] 可以理解,在另一些实施例中,电子设备100也可以通过其他方式确定待运行的神经网络模型是否为非对称量化的神经网络模型,并在检测到待运行的神经网络模型是非对称量化的神经网络模型的情况下,触发本申请实施例提供的神经网络模型的运行方法。

[0098] S602:获取各算子的输入数据/输出数据的非对称量化参数,并将非对称量化参数转换为对称量化参数。

[0099] 即是电子设备100依次获取各算子的输入数据和输出数据的非对称量化参数,并将输入数据的非对称量化参数(输入数据的非对称量化缩放系数、输入数据的非对称量化零点)转换为输入数据的对称量化参数(输入数据的对称量化缩放系数),将输出数据的非对称量化参数(输出数据的非对称量化缩放系数、输出数据的非对称量化零点)转换为输出数据的对称量化参数(输出数据的对称量化缩放系数)。

[0100] 例如,对于图4所示的神经网络模型10,输入层11的非对称量化参数包括:输出数据的非对称量化缩放系数 $ASc\_out1=1$ 、输出数据的非对称量化零点 $Zp\_out1=0$ 。由于UINT8对应的取值范围为 $[0,255]$ ,从而根据公式(11)可以得到输出数据对应的浮点数的最大值为浮点数 $(255-0)/1=255$ ,最小值为浮点数 $(0-0)/1=0$ ,电子设备100可以基于公式(2),确定输入层11的输出数据的对称量化缩放系数 $Sc\_out1=(2^7-1)/255=0.498$ 。

[0101] 又例如,对于图4所示的神经网络模型10,卷积算子12的非对称量化参数包括:输入数据的非对称量化缩放系数 $ASc\_in2=2$ 、输入数据的非对称量化零点 $Zp\_in2=0$ 、输出数据的非对称量化缩放系数 $ASc\_out2=0.0833$ 、输入数据的非对称量化零点 $Zp\_out2=0$ 。由于UINT8对应的取值范围为 $[0,255]$ ,从而根据公式(11)可以得到输入数据对应的浮点数的最大值为浮点数 $(255-0)/1=255$ 、最小值为浮点数 $(0-0)/1=0$ ,输出数据对应的浮点数的最大值为浮点数 $(255-0)/0.0833=3061$ ,最小值为浮点数 $(0-0)/1=0$ 。电子设备100可以基于公式(2),确定卷积算子12的输入数据的对称量化缩放系数 $Sc\_in2=(2^7-1)/255=0.498$ ,输出数据的对称量化缩放系数 $Sc\_out2=(2^7-1)/3061=0.0417$ 。类似地,可以得到全连接算子13的输入数据的对称量化缩放系数 $Sc\_in3=0.0417$ 、输出数据的对称量化缩放系数 $Sc\_out3=1.1\times 10^{-5}$ 。

[0102] 再例如,对于图4所示的神经网络模型10,Softmax算子14的非对称量化参数包括:输入数据的非对称量化缩放系数 $ASc\_in4=2.2\times 10^{-5}$ 、输入数据的非对称量化零点 $Zp\_in4=0$ 、输出数据的非对称量化缩放系数 $ASc\_out4=255$ 、输出数据的非对称量化零点 $Zp\_out4=0$ 。由于UINT8对应的取值范围为 $[0,255]$ ,从而根据公式(11)可以得到输入数据对应的浮点数的最大值为浮点数 $(255-0)/(2.2\times 10^{-5})=11590909$ 、最小值为浮点数 $(0-0)/(2.2\times 10^{-5})=0$ ,输出数据对应的浮点数的最大值为浮点数 $(255-0)/255=1$ ,最小值为浮点数 $(0-0)/255=0$ 。电子设备100可以基于公式(2),确定Softmax算子14的输入数据的对称量化缩放系数 $Sc\_in4=(2^7-1)/11590909=1.1\times 10^{-5}$ ,输出数据的对称量化缩放系数 $Sc\_out4=(2^7-1)/1=127$ 。

[0103] S603:判断当前算子是否包括非对称量化常量数据。

[0104] 电子设备100判断当前算子是否包括非对称量化常量数据,如果有,说明需要将非对称量化常量数据转换为对称量化常量数据,转至步骤S604;否则,说明不需要将非对称量化常量数据转换为对称量化常量数据,转至步骤S605。

[0105] 例如,对于前述神经网络模型10,在当前算子为输入层11或输出层15时,当前算子不存在非对称量化常量数据,转至步骤S605;对于前述神经网络模型10,卷积算子12存在非对称量化常量数据卷积核 $K_i$  ( $i=1,2,3$ )、全连接算子13存在非对称量化常量数据权重矩阵 $W$ ,Softmax算子14存在非对称量化常量数据查找表LUT,在当前算子为卷积算子12、全连接算子13或Softmax算子14时,电子设备100可以确定出当前算子包括非对称量化常量数据,转至步骤S604。

[0106] S604:根据非对称量化常量数据的非对称量化参数,将非对称量化常量数据转化为对称量化常量数据。

[0107] 电子设备100在当前算子存在非对称量化常量数据的情况下,根据当前算子的非对称量化常量数据的非对称量化参数,先根据非对称量化常量数据的定点数的数据类型,确定非对称量化常量数据的定点数的最大值和最小值,从而根据公式(11)确定非对称量化常量数据对应的浮点数的最大值和最小值;再根据公式(2)确定将浮点数的常量数据转化为对称量化常量数据的对称量化缩放系数;然后根据公式(11)将非对称量化常量数据转化为对应的浮点数的常量数据,再根据公式(1)将浮点数的常量数据转换为对称量化常量数据。

[0108] 例如,在当前算子为前述卷积算子12时,非对称量化常量数据包括卷积核 $K_i$  ( $i=1,2,3$ )。由于卷积核 $K_i$ 的非对称量化参数与卷积算子12的输入数据的非对称量化参数相同,则卷积核 $K_i$ 的对称量化缩放系数与卷积算子12的输入数据的对称量化缩放系数相同(均为 $Sc\_in2$ )。电子设备100可以先根据公式(11),将 $K_i$ 转换为对应的浮点数 $K_{i\_f} = (K_i - Zp\_in2) / ASc\_in2$ ;再根据公式(1),将 $K_{i\_f}$ 转换为对称量化的卷积核 $K_i' = \text{round}(K_{i\_f} \times Sc\_in2)$ 。

具体地,参考图7,假设 $K_1 = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 5 \end{bmatrix}$ ,则 $K_1' = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$ 。

[0109] 又例如,在当前算子为前述全连接算子13时,非对称量化常量数据包括权重矩阵 $W$ 。由于权重矩阵 $W$ 的非对称量化参数与全连接算子13的输入数据的非对称量化参数相同,则权重矩阵 $W$ 的对称量化缩放系数与全连接算子13的输入数据的对称量化缩放系数相同(均为 $Sc\_in3$ )。电子设备100可以先根据公式(11),将 $W$ 转换为对应的浮点数 $W\_f = (W - Zp\_in3) / ASc\_in3$ ;再根据公式(1)和公式(2),将 $W\_f$ 转换为对称量化的权重矩阵 $W' = \text{round}(W\_f \times Sc\_in3)$ 。

[0110] 具体地,参考图7,假设

[0111]  $W = \begin{bmatrix} 10 & 10 & 20 & 18 & 10 & 10 \\ 0 & 1 & 5 & 19 & 20 & 1 \\ 0 & 0 & 2 & 8 & 10 & 10 \\ 5 & 20 & 19 & 0 & 12 & 18 \\ 5 & 12 & 0 & 0 & 0 & 18 \\ 1 & 11 & 15 & 10 & 11 & 0 \end{bmatrix}$ ,

[0112] 则



$$[0113] \quad W' = \text{round}\left(\frac{W-0}{0.0833} \times 0.0415\right) = \begin{bmatrix} 5 & 5 & 10 & 9 & 5 & 5 \\ 0 & 0 & 2 & 9 & 10 & 0 \\ 0 & 0 & 2 & 4 & 5 & 5 \\ 2 & 10 & 9 & 0 & 6 & 9 \\ 2 & 6 & 0 & 0 & 0 & 9 \\ 1 & 5 & 7 & 5 & 5 & 0 \end{bmatrix}。$$

[0114] 再例如,在当前算子为前述Softmax算子14时,非对称量化常量数据包括查找表LUT。Softmax运算可以表示为如下公式(12)。

$$[0115] \quad P_{k\_f} = \frac{e^{\frac{in_{k\_f}}{1159091}}}{\sum_{j=1}^3 e^{\frac{in_{j\_f}}{1159091}}} \quad (k=1, 2, 3) \quad (12)$$

[0116] 公式(12)中, $in_{k\_f}$ 表示Softmax算子的浮点数输入数据,即是输入的图像数据在第k类别的浮点数得分; $P_{k\_f}$ 是Softmax算子的浮点数输出结果,表示输入的图像数据为第k类的概率,其中, $k=1$ 表示类别为兔子, $k=2$ 表示类别为狗, $k=3$ 表示类别为猫。从公式(12)可知, $P_{k\_f}$ 的取值范围为[0,1]。

[0117] 查找表LUT用于存储不同的定点数输入数据 $in_{k\_q}$ 对应的Softmax运算的定数结果。例如,假设全连接算子13的输出数据的非对称量化缩放系数为 $2.2 \times 10^{-5}$ ,非对称量化零点为0, $in_{1\_q}=59$ , $in_{2\_q}=104$ , $in_{3\_q}=182$ ,则可以得到 $in_{1\_f}=2681818$ , $in_{2\_f}=4727273$ , $in_{3\_f}=8272727$ ,将 $in_{1\_f}=2681818$ , $in_{2\_f}=4727273$ , $in_{3\_f}=8272727$ 代入前述公式(12)可以得到 $P_{1\_f}=0.0069$ , $P_{2\_f}=0.042$ , $P_{3\_f}=0.9511$ ,再将 $P_{1\_f}=0.0069$ , $P_{2\_f}=0.042$ , $P_{3\_f}=0.9511$ 进行8位非对称量化得到: $P_{1\_q}=2$ , $P_{2\_q}=11$ , $P_{3\_q}=243$ ,也即说,在查找表LUT中,存储有查表索引[59,104,182](对应[ $in_1, in_2, in_3$ ])对应的查表结果为[2,11,243](对应 $P_1, P_2, P_3$ )。其他的查表索引对应的查表结果可以通过类似的方法得到,在此不做赘述。

[0118] 电子设备100在检测到上述查找表LUT时,先根据公式(11)将查找表LUT的查表索引转换为浮点数查表索引,将浮点数查表索引转化为对称量化的定点数查表索引,并将浮点数查表索引代入前述公式(12)得到浮点数 $P_{k\_f}$ ,再将浮点数 $P_{k\_f}$ 进行对称量化的对称量化查表结果作为前述对称量化的定点数查表索引在新的查找表LUT'的查表结果。例如,将查表索引[59,104,182]转换为浮点数得到 $in_{1\_f}=2681818$ , $in_{2\_f}=4727273$ , $in_{3\_f}=8272727$ ,将 $in_{1\_f}=2681818$ , $in_{2\_f}=4727273$ , $in_{3\_f}=8272727$ 进行8位对称量化得到新的查表索引[30,52,91],并将 $in_{1\_f}=2681818$ , $in_{2\_f}=4727273$ , $in_{3\_f}=8272727$ 代入前述公式(12)可以得到 $P_{1\_f}=0.0069$ , $P_{2\_f}=0.042$ , $P_{3\_f}=0.9511$ ,再将 $P_{1\_f}=0.0069$ , $P_{2\_f}=0.042$ , $P_{3\_f}=0.9511$ 进行8位对称量化得到 $P_{1\_q}=1$ , $P_{2\_q}=5$ , $P_{3\_q}=121$ ,即是在图7所示的查找表LUT'中,查表索引[30,52,91](对应[ $in_1', in_2', in_3'$ ])对应的查表结果为[1,5,121](对应 $P_1', P_2', P_3'$ )。其他的查表索引对应的结果可以通过类似的方法得到,在此不做赘述。

[0119] 可以理解,以上将卷积算子、全连接算子和Softmax算子中非对称量化常量数据转换为对称量化常量数据只是一种示例,对于其他算子(包括但不限于池化算子、激活算子、排序算子、归一化算子等)中的非对称量化常量数据,可以使用类似的方法转换为对称量化

常量数据,在此不做赘述。

[0120] S605:判断是否完成所有算子的量化参数的转换。

[0121] 电子设备100判断是否完成所有算子的转换,如果完成,则转至步骤S606;否则转至步骤S602进行下一个算子的量化参数的转换。

[0122] S606:根据各算子的对称量化参数,调用相应的对称量化算子的运算逻辑,实现非对称量化的神经网络模型的功能。

[0123] 电子设备100在完成所有算子的量化参数的转换后,根据各算子的对称量化参数,通过NPU调用相对应的对称量化算子的运算逻辑,实现非对称量化的神经网络模型的功能。

[0124] 具体地,参考图7,电子设备100将神经网络模型10的各算子的非对称量化参数转换为对称量化参数后,各算子可以表示为图7所示的输入层11',卷积算子12',全连接算子13',Softmax算子14',和输出层15'。

[0125] 电子设备100的NPU可以先在对称量化缩放系数 $Sc\_out1=0.0498$ 的情况下,基于公式(1)将图像20量化为图像矩阵 $H'$ 。

[0126] 其次,电子设备100的NPU调用对称量化的卷积算子的运算逻辑,例如前述公式(9)所示的运算逻辑,在卷积核为 $K_i'$  ( $i=1,2,3$ )的情况下,将图像矩阵 $H'$ 分别与 $K_i'$ 进行卷积,得到定点数特征矩阵 $A_i'$  ( $i=1,2,3$ )。即是NPU获取对称量化的输入数据 $B\_q$ (例如前述 $H'$ )、输入数据的对称量化缩放系数 $B\_Sc$ (例如由前述步骤S602得到的 $Sc\_in2$ )、对称量化的卷积核 $C\_q$ (例如由前述步骤S603得到的卷积核 $K_i'$ )、卷积核的对称量化缩放系数 $C\_Sc$ (例如前述 $Sc\_in2$ )、输出数据的对称量化缩放系数 $D\_Sc$ (例如由前述步骤S602得到的 $Sc\_out2$ ),再根据前述公式(9)得到 $D\_q$ 。例如,在 $B\_q=H'$ , $C\_q=K_1'$ 的情况下,可以得到前述特征矩阵 $A_1'$ 。

[0127] 可以理解,由于NPU中并没有可以直接实现除法运算的电路,在一些实施例中,公式(9)中的除法运算可以通过乘法移位来实现,以提高NPU运行卷积算子的速度。例如,假设 $B\_Sc \times C\_Sc = 0.498^2 = 0.248$ ,可以将0.248表示为 $1/1 \times 2^{-2}$ ,从而将 $B\_q * C\_q \times D\_Sc$ 的结果对应的二进制数向右移-2位再乘以1,即可得到 $B\_q * C\_q \times D\_Sc / (B\_Sc \times C\_Sc)$ 的运算结果。

[0128] 然后,电子设备100的NPU调用对称量化的全连接算子的运算逻辑(例如下方公式(17)所示的运算逻辑),分别将特征矩阵 $A_i'$ 与权重矩阵 $W'$ 作全连接运算,例如做内积,得到图像20为各预设类别的得分 $in1'$ 、 $in2'$ 和 $in3'$ ;再调用对称量化的Softmax算子的运算逻辑,即是以 $[in1', in2', in3']$ 为查表索引从查找表 $LUT'$ 中查找得到图像20属于各预设类别的概率;最后再调用对称量化的输出层的运算逻辑,将图像20属于各预设类别的概率中最大的概率对应的预设类别作为图像20的类别,例如将图像20的类别确定为猫。对称量化全连接算子的运算逻辑的推导过程将在下文进行介绍,在此不做赘述。

[0129] 可以理解,对称量化的算子的运算逻辑可以由电子设备100的NPU的开发商预先设置于NPU中,也可以由NPU的开发商提供给电子设备100的开发商,由电子设备100的开发商预设于电子设备100的存储器中。

[0130] 可以理解,上述步骤S601至步骤S605中的各步骤可以全部由电子设备100的CPU来完成,也可以全部由电子设备100的NPU来完成,还可以由电子设备100的CPU和NPU分别完成部分步骤,在此不做限定。

[0131] 可以理解,上述步骤S601至步骤S605的运行顺序只是一种示例,在另一些实施例中,可以调整部分步骤的运行顺序,也可以合并或拆分部分步骤,本申请实施例不做限定。

[0132] 通过本申请实施例提供的方法,电子设备100的NPU可以通过调用预设的对称量化的算子来实现该非对称量化神经网络模型的相关功能,而无需由其他电子设备将非对称量化的神经网络模型转换为浮点型神经网络模型,并将浮点型神经网络模型转化为对称量化的神经网络模型后,再由电子设备100的NPU来运行,增加了电子设备100能够运行的神经网络模型的类型,提高了NPU的通用性,提升了神经网络模型的部署速度。

[0133] 下面介绍卷积计算的具体过程和对称量化的全连接算子的运算逻辑。

[0134] 首先介绍卷积运算的计算过程。

[0135] 假设输入数据矩阵B的大小为 $M \times M$ ,卷积核C的大小为 $N \times N$ ,卷积步长为 $k$ ,则矩阵B与卷积核C的卷积结果D可以表示为:

$$[0136] \quad D(m, n) = \sum_{i=1}^N \sum_{j=1}^N B(k(m-1)+i, k(n-1)+j) \times C(i, j) \quad (13)$$

[0137] 在公式(13)中, $D(m, n)$ 为矩阵D第 $m$ 行第 $n$ 列的元素; $m, n$ 满足以下关系式:

$$[0138] \quad 1 \leq m \leq \lfloor \frac{(M-N)}{k} + 1 \rfloor, \quad 1 \leq n \leq \lfloor \frac{(M-N)}{k} + 1 \rfloor。$$

[0139] 其中 $\lfloor X \rfloor$ 为向下取整运算,即 $\lfloor X \rfloor$ 为小于 $X$ 的最大整数。由于 $M-N < M$ 且 $k$ 为正整数,可见 $\lfloor \frac{(M-N)}{k} + 1 \rfloor \leq M$ ,也即是说卷积结果D的大小总是小于或等于矩阵B的大小。

[0140] 为确保卷积结果对应的矩阵的大小与输入数据的大小相同,避免丢失数据矩阵边缘的数据特征,通常在卷积计算的过程中在输入矩阵的第一行前和最后一行后填充值为0的行以及在输入矩阵第一列之前和最后一列之后填充值为0的行或列,即在输入矩阵的四周填充值为0的行或列。设在输入矩阵B的四周各填充数量为 $P$ 的值为0的行或列,此时,输入矩阵B的大小变为 $(M+2P) \times (M+2P)$ 。此时,公式(1)中的 $m, n$ 满足以下关系式:

$$[0141] \quad 1 \leq m \leq \lfloor \frac{(M-N+2 \times P)}{k} + 1 \rfloor, \quad 1 \leq n \leq \lfloor \frac{(M-N+2 \times P)}{k} + 1 \rfloor。$$

[0142] 令 $\frac{(M-N+2 \times P)}{k} + 1 = M$ ,即可计算得到 $P$ 的值,例如,在卷积核大小为 $N=3$ ,步长 $k=1$ ,则 $P=1$ 。

[0143] 下面介绍对称量化的全连接算子的运算逻辑。

[0144] 全连接算子为对输入数据进行加权计算的算子,输入矩阵E和权重矩阵W的全连接计算结果F可以表示为如下公式(14)。

$$[0145] \quad F = \sum_{i=1}^M \sum_{j=1}^N E(i, j) \times W(i, j) \quad (14)$$

[0146] 其中, $E(i, j)$ 为输入矩阵的第 $i$ 行第 $j$ 列的元素, $W(i, j)$ 为权重矩阵的第 $i$ 行第 $j$ 列的元素,输入矩阵E和权重矩阵W的大小均为 $M \times N$ 。

[0147] 假设 $E_f$ 为浮点数输入矩阵, $W_f$ 为浮点数权重矩阵, $F_f$ 为浮点数全连接计算结果,基于公式(14), $E_f, W_f, F_f$ 间的关系可以表示为如下公式(15)。

$$[0148] \quad F_f = \sum_{i=1}^M \sum_{j=1}^N E(i, j)_f \times W(i, j)_f \quad (15)$$

[0149] 根据公式(1)和公式(2)可以得到 $E_f = E_q / E_{Sc}$ ,  $W_f = W_q / W_{Sc}$ ,  $F_f = F_q / F_{Sc}$ , 其中: $E_q$ 为 $E_f$ 对应的对称量化的定点数矩阵, $E_{Sc}$ 为将 $E_f$ 量化为 $E_q$ 的对称量化缩放系数; $W_q$ 为 $W_f$ 对应的对称量化的定点数权重矩阵, $W_{Sc}$ 为将 $W_f$ 量化为 $W_q$ 的对称量化缩放系数; $F_q$ 为 $F_f$ 对应的对称量化的定点数, $F_{Sc}$ 为将 $F_f$ 量化为 $F_q$ 的对称量化缩放系数。进而公式(15)可以表示为如下公式(16)。

$$[0150] \quad \frac{F_q}{F_{Sc}} = \sum_{i=1}^M \sum_{j=1}^N \frac{E_q(i, j)}{E_{Sc}} \times \frac{W_q(i, j)}{W_{Sc}} \quad (16)$$

[0151] 对公式(16)进行变形可以得到如下公式(17)所示的对称量化的全连接算子的运算逻辑。

$$[0152] \quad F_q = \text{round} \left\{ \left[ \sum_{i=1}^M \sum_{j=1}^N E_q(i, j) \times W_q(i, j) \right] \times \frac{F_{Sc}}{E_{Sc} \times W_{Sc}} \right\} \quad (17)$$

[0153] NPU在执行公式(17)所示的全连接算子的运算逻辑时,获取对称量化的输入数据 $E_q$ (例如前述 $A_i'$  ( $i=1, 2, 3$ ))及输入数据的对称量化缩放系数 $E_{Sc}$ (例如由前述步骤S602得到的 $S_{cin3}$ )、对称量化的权重矩阵 $W_q$ (例如由前述步骤S603得到的权重矩阵 $W'$ )、权重矩阵的对称量化缩放系数 $W_{Sc}$ (例如前述 $S_{cin3}$ )、输出数据的对称量化缩放系数 $F_{Sc}$ (例如由前述步骤S602得到的 $S_{cout3}$ ),再根据前述公式(17)得到 $A_i'$ 和 $W'$ 的全连接计算结果,例如在 $E_q = A_1'$ 的情况下,可以得到前述 $in1'$ 。

[0154] 可以理解,由于NPU中并没有可以直接实现除法运算的电路,在一些实施例中,公式(17)中的除法运算可以通过移位和乘法来实现,以提高NPU运行全连接算子的速度。例如,假设 $E_{Sc} \times W_{Sc} = 0.0417^2 = 0.00174$ ,可以将0.00174表示为 $1/9 \times 2^{-6}$ ,从而将 $B_q * C_q \times D_{Sc}$ 的结果对应的二进制数向右移-6位(即向左移6位)再乘以9,即可得到 $B_q * C_q \times D_{Sc} / (E_{Sc} \times W_{Sc})$ 的运算结果。

[0155] 可以理解,在另一些实施例中全连接算子也可以采用其他运算逻辑,本申请实施例不做限定。

[0156] 可以理解,对于其他的对称量化的算子,运算逻辑可以通过相似的方法得到,在此不做赘述。

[0157] 进一步,图8根据本申请的一些实施例,示出了一种电子设备100的结构示意图。如图8所示,电子设备100包括一个或多个处理器101A、NPU 101B、系统内存102、非易失性存储器(Non-Volatile Memory, NVM) 103、通信接口104、输入/输出(I/O)设备105、以及用于耦接处理器101A、系统内存102、非易失性存储器103、通信接口104和输入/输出(I/O)设备105的系统控制逻辑106。其中:

[0158] 处理器101A可以包括一个或多个处理单元,例如,可以包括中央处理器CPU(Central Processing Unit)、图像处理器GPU(Graphics Processing Unit)、数字信号处理器DSP(Digital Signal Processor)、微处理器MCU(Micro-programmed Control Unit)、AI(Artificial Intelligence,人工智能)处理器或可编程逻辑器件FPGA(Field Programmable Gate Array)的处理模块或处理电路可以包括一个或多个单核或多核处理器。

[0159] 神经网络处理器101B可以用于调用预设的对称量化的算子的运算逻辑,实现神经网络模型的推理。神经网络处理器101B可以是独立的处理器,也可以集成于处理器101A内

部。在一些实施例中,NPU可以用于运行本申请实施例提供的神经网络模型的运行方法对应的指令。

[0160] 系统内存102是易失性存储器,例如随机存取存储器(Random-Access Memory, RAM),双倍数据率同步动态随机存取存储器(Double Data Rate Synchronous Dynamic Random Access Memory,DDR SDRAM)等。系统内存用于临时存储数据和/或指令,例如,在一些实施例中,系统内存102可以用于存储上述神经网络模型10的相关指令、非对称/对称量化参数、非对称/对称量化常量数据等,也可以用于存储预设的对称量化的算子的运算逻辑。

[0161] 非易失性存储器103可以包括用于存储数据和/或指令的一个或多个有形的、非暂时性的计算机可读介质。在一些实施例中,非易失性存储器103可以包括闪存等任意合适的非易失性存储器和/或任意合适的非易失性存储设备,例如硬盘驱动器(Hard Disk Drive, HDD)、光盘(Compact Disc,CD)、数字通用光盘(Digital Versatile Disc,DVD)、固态硬盘(Solid-State Drive,SSD)等。在一些实施例中,非易失性存储器103也可以是可移动存储介质,例如安全数字(Secure Digital,SD)存储卡等。在另一些实施例中,非易失性存储器103可以用于存储上述神经网络模型10的相关指令、非对称/对称量化参数、非对称/对称量化常量数据等,也可以用于存储预设的对称量化的算子的运算逻辑。

[0162] 特别地,系统内存102和非易失性存储器103可以分别包括:指令107的临时副本和永久副本。指令107可以包括:由处理器101A和/或神经网络处理器101B中的至少一个执行时使电子设备100实现本申请各实施例提供的神经网络模型的运行方法。

[0163] 通信接口104可以包括收发器,用于为电子设备100提供有线或无线通信接口,进而通过一个或多个网络与任意其他合适的设备进行通信。在一些实施例中,通信接口104可以集成于电子设备100的其他组件,例如通信接口104可以集成于处理器101A中。在一些实施例中,电子设备100可以通过通信接口104和其他设备通信,例如,电子设备100可以通过通信接口104从其他电子设备获取待运行的神经网络模型。

[0164] 输入/输出(I/O)设备105可以包括输入设备如键盘、鼠标等,输出设备如显示器等,用户可以通过输入/输出(I/O)设备105与电子设备100进行交互。

[0165] 系统控制逻辑106可以包括任意合适的接口控制器,以电子设备100的其他模块提供任意合适的接口。例如在一些实施例中,系统控制逻辑106可以包括一个或多个存储器控制器,以提供连接到系统内存102和非易失性存储器103的接口。

[0166] 在一些实施例中,处理器101A中的至少一个可以与用于系统控制逻辑106的一个或多个控制器的逻辑封装在一起,以形成系统封装(System in Package,SiP)。在另一些实施例中,处理器101A中的至少一个还可以与用于系统控制逻辑106的一个或多个控制器的逻辑集成在同一芯片上,以形成片上系统(System-on-Chip,SoC)。

[0167] 可以理解,电子设备100可以是能够运行神经网络模型的任意电子设备,包括但不限于手机、可穿戴设备(如智能手表等)、平板电脑、桌面型、膝上型、手持计算机、笔记本电脑、超级移动个人计算机(ultra-mobile personal computer,UMPC)、上网本,以及蜂窝电话、个人数字助理(personal digital assistant,PDA)、增强现实(augmented reality,AR)/虚拟现实(virtual reality,VR)设备等,本申请实施例不做限定。

[0168] 可以理解,图8所示的电子设备100的结构只是一种示例,在另一些实施例中,电子

设备100可以包括比图示更多或更少的部件,或者组合某些部件,或者拆分某些部件,或者不同的部件布置。图示的部件可以以硬件,软件或软件和硬件的组合实现。

[0169] 本申请公开的机制的各实施例可以被实现在硬件、软件、固件或这些实现方法的组合中。本申请的实施例可实现为在可编程系统上执行的计算机程序或程序代码,该可编程系统包括至少一个处理器、存储系统(包括易失性和非易失性存储器和/或存储元件)、至少一个输入设备以及至少一个输出设备。

[0170] 可将程序代码应用于输入指令,以执行本申请描述的各功能并生成输出信息。可以按已知方式将输出信息应用于一个或多个输出设备。为了本申请的目的,处理系统包括具有诸如例如数字信号处理器(Digital Signal Processor,DSP)、微控制器、专用集成电路(Application Specific Integrated Circuit,ASIC)或微处理器之类的处理器的任何系统。

[0171] 程序代码可以用高级程序化语言或面向对象的编程语言来实现,以便与处理系统通信。在需要时,也可用汇编语言或机器语言来实现程序代码。事实上,本申请中描述的机制不限于任何特定编程语言的范围。在任一情形下,该语言可以是编译语言或解释语言。

[0172] 在一些情况下,所公开的实施例可以以硬件、固件、软件或其任何组合来实现。所公开的实施例还可以被实现为由一个或多个暂时或非暂时性机器可读(例如,计算机可读)存储介质承载或存储在其上的指令,其可以由一个或多个处理器读取和执行。例如,指令可以通过网络或通过其他计算机可读介质分发。因此,机器可读介质可以包括用于以机器(例如,计算机)可读的形式存储或传输信息的任何机制,包括但不限于,软盘、光盘、光碟、只读存储器(CD-ROMs)、磁光盘、只读存储器(Read Only Memory,ROM)、随机存取存储器(Random Access Memory,RAM)、可擦除可编程只读存储器(Erasable Programmable Read Only Memory,EPRM)、电可擦除可编程只读存储器(Electrically Erasable Programmable Read-Only Memory,EEPROM)、磁卡或光卡、闪存、或用于利用因特网以电、光、声或其他形式的传播信号来传输信息(例如,载波、红外信号数字信号等)的有形的机器可读存储器。因此,机器可读介质包括适合于以机器(例如计算机)可读的形式存储或传输电子指令或信息的任何类型的机器可读介质。

[0173] 在附图中,可以以特定布置和/或顺序示出一些结构或方法特征。然而,应该理解,可能不需要这样的特定布置和/或排序。而是,在一些实施例中,这些特征可以以不同于说明性附图中所示的方式和/或顺序来布置。另外,在特定图中包括结构或方法特征并不意味着暗示在所有实施例中都需要这样的特征,并且在一些实施例中,可以不包括这些特征或者可以与其他特征组合。

[0174] 需要说明的是,本申请各设备实施例中提到的各单元/模块都是逻辑单元/模块,在物理上,一个逻辑单元/模块可以是一个物理单元/模块,也可以是一个物理单元/模块的一部分,还可以以多个物理单元/模块的组合实现,这些逻辑单元/模块本身的物理实现方式并不是最重要的,这些逻辑单元/模块所实现的功能的组合才是解决本申请所提出的技术问题的关键。此外,为了突出本申请的创新部分,本申请上述各设备实施例并没有将与解决本申请所提出的技术问题关系不太密切的单元/模块引入,这并不表明上述设备实施例并不存在其它的单元/模块。

[0175] 需要说明的是,在本专利的示例和说明书中,术语“包括”、“包含”或者其任何其他

变体意在涵盖非排他性的包含,从而使得包括一系列要素的过程、方法、物品或者设备不仅包括那些要素,而且还包括没有明确列出的其他要素,或者是还包括为这种过程、方法、物品或者设备所固有的要素。在没有更多限制的情况下,由语句“包括一个”限定的要素,并不排除在包括所述要素的过程、方法、物品或者设备中还存在另外的相同要素。

[0176] 虽然通过参照本申请的某些优选实施例,已经对本申请进行了图示和描述,但本领域的普通技术人员应该明白,可以在形式上和细节上对其作各种改变,而不偏离本申请的精神和范围。

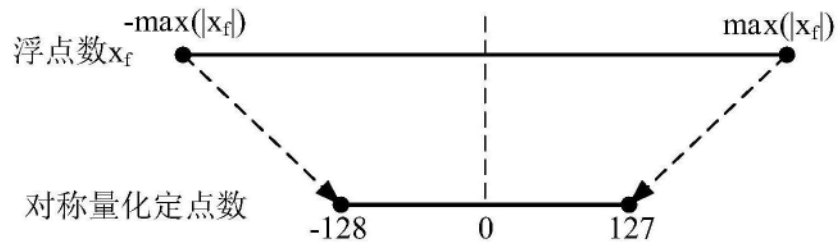


图1A

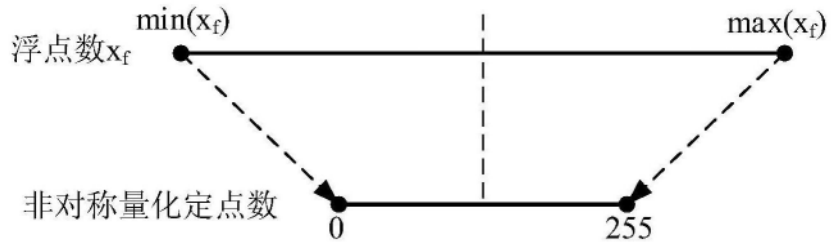


图1B

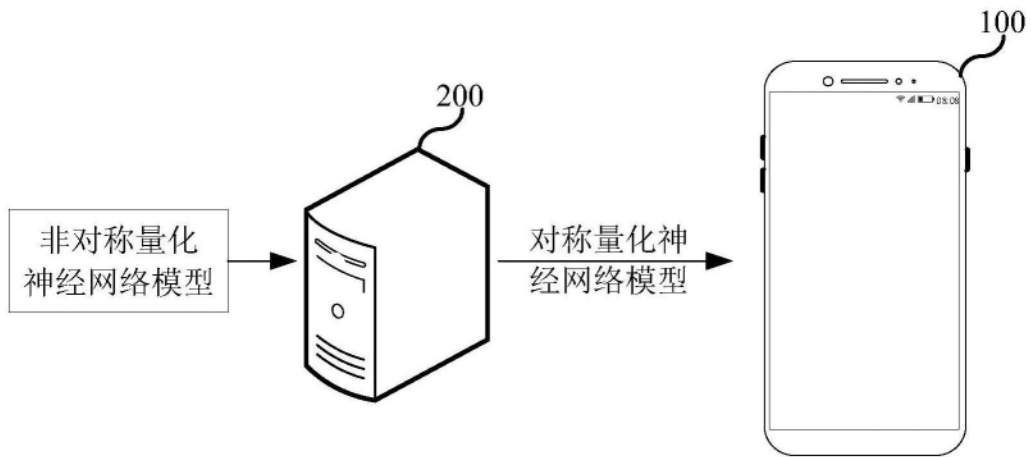


图2



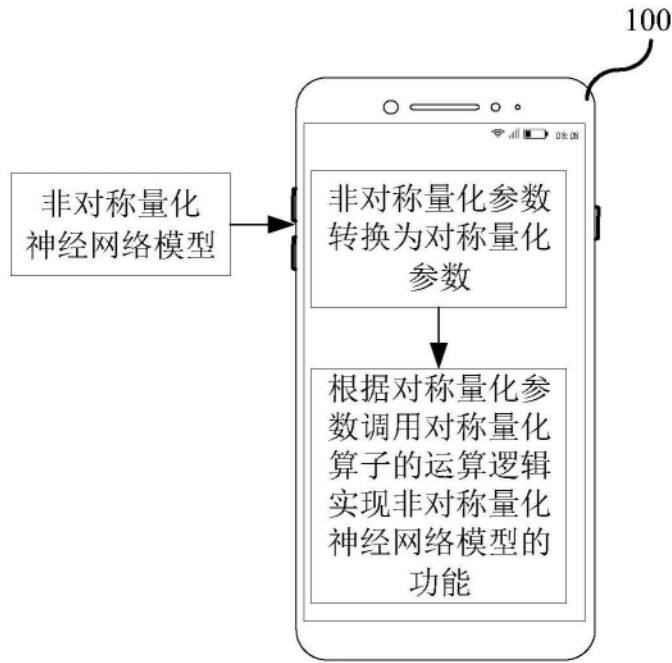


图3

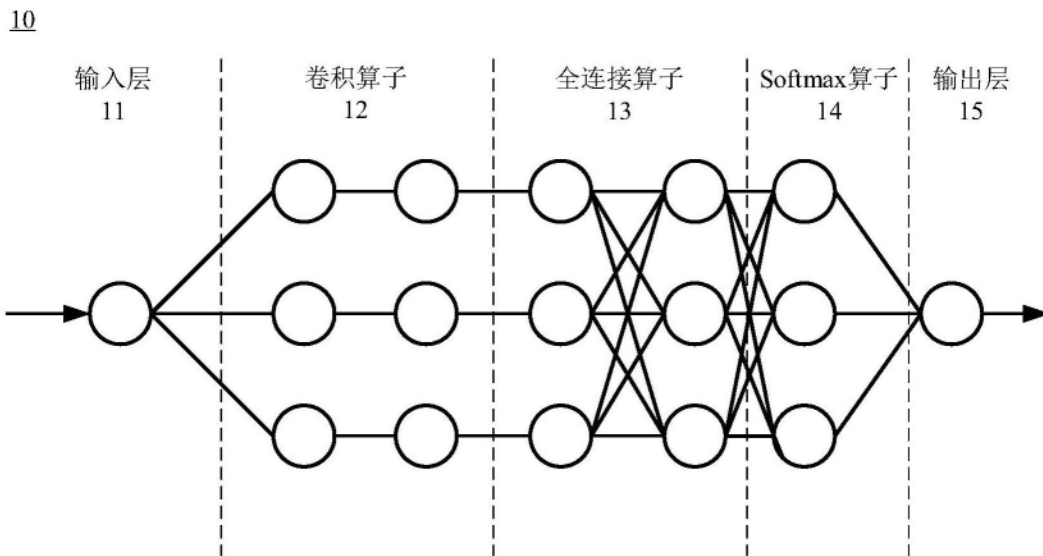


图4

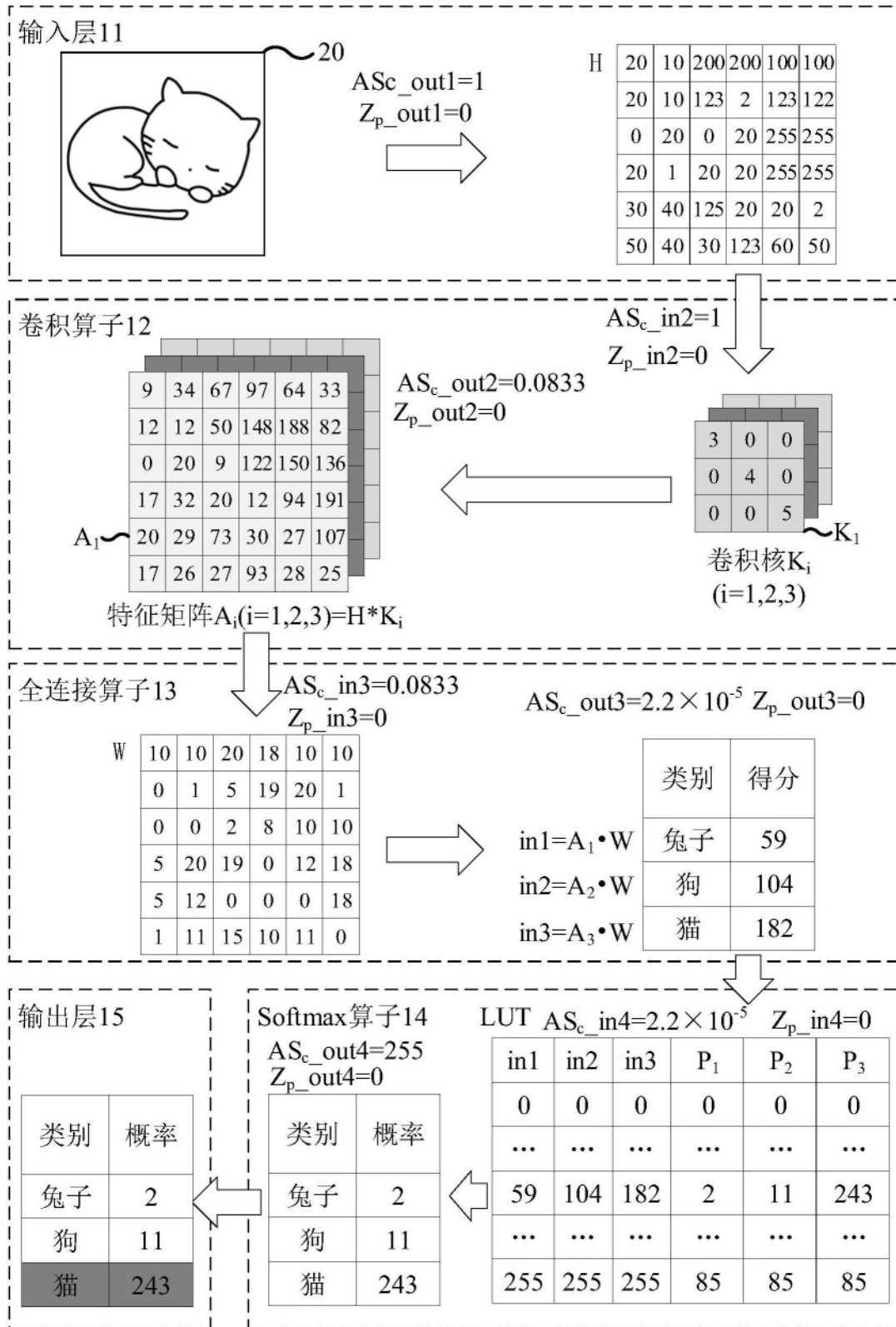


图5

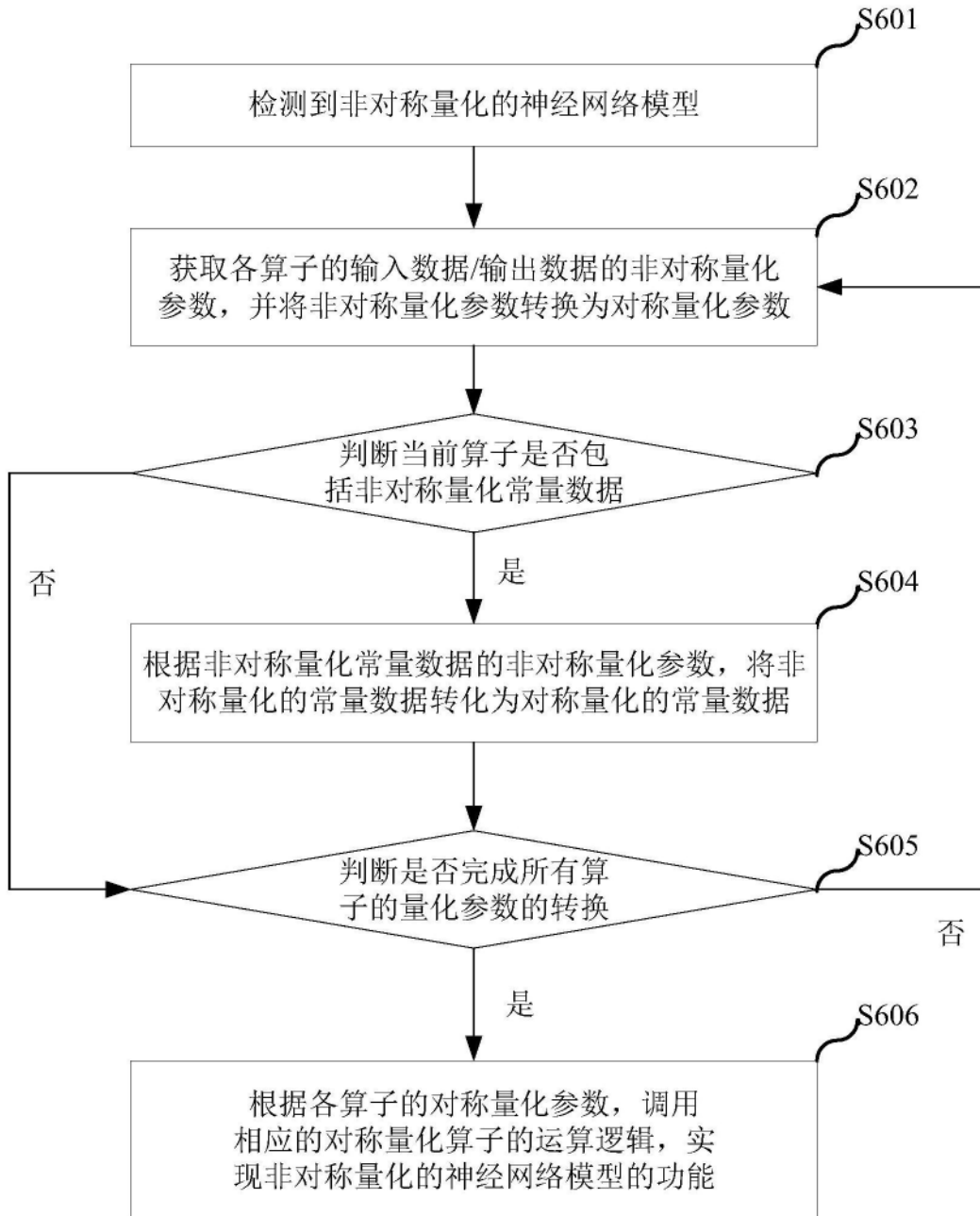


图6

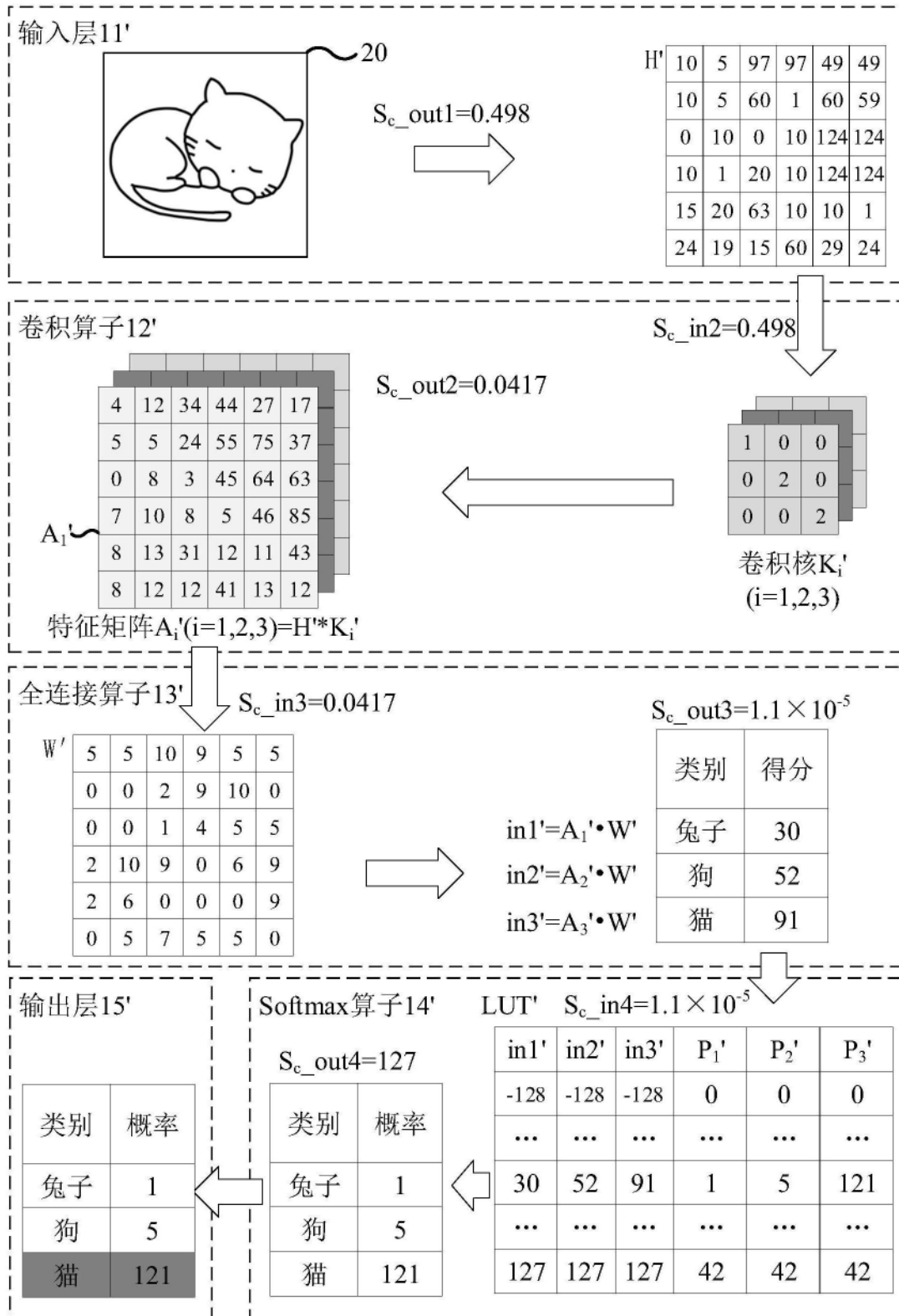


图7

100

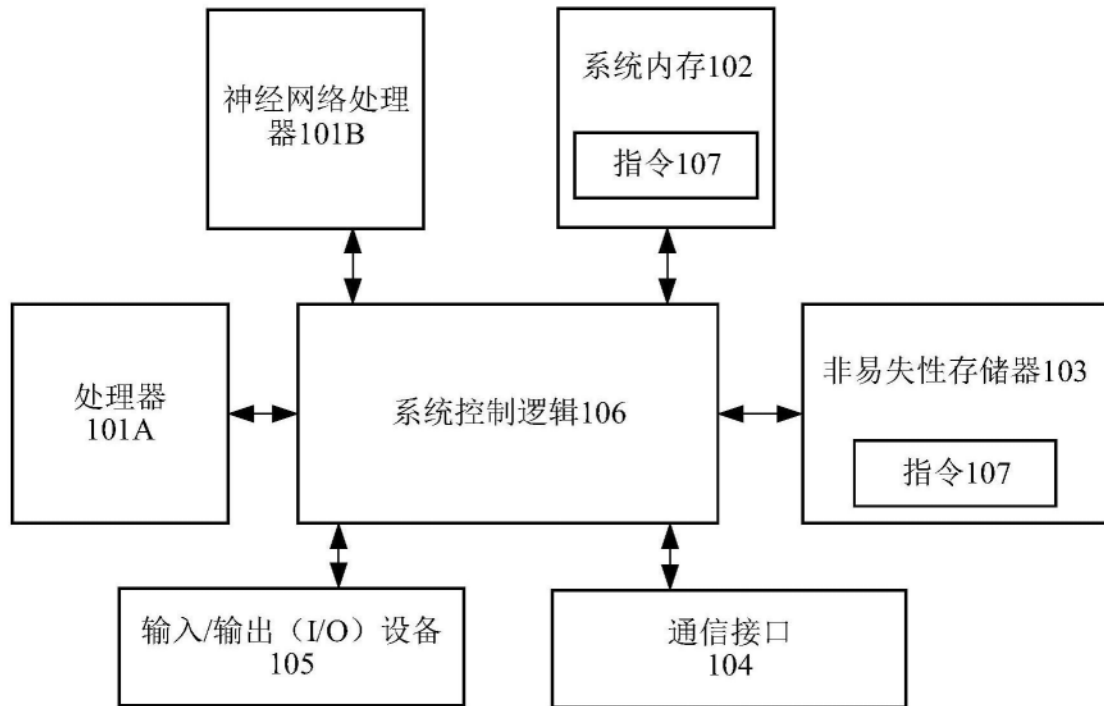


图8