

基于隐私保护机制的辐射光源衍射图像筛选

许康^{1,2}, 祝永新^{2*}, 吴波², 郑小盈², 陈凌曜³¹中国科学院大学集成电路学院, 北京 100049;²中国科学院上海高等研究院, 上海 201210;³上海信息技术研究中心, 上海 201210

摘要 同步辐射光源产生超高速的衍射图像数据流, 需要通过数据筛选降低数据传输和存储的压力。但互相竞争的研究小组不愿意分享数据, 现有基于深度学习的筛选方法难以应对隐私保护下有效训练的挑战, 因此首次将联邦学习技术应用在辐射光源衍射图像筛选中, 通过数据和模型分离, 实现隐私保护下的训练数据增广。提出筛选方法 Federated Kullback-Leibler (FedKL), 基于改进的 KL 散度和数据量权重, 对全局模型更新进行改进, 在获得高准确率的同时降低算法的复杂度, 满足高速数据流高精度处理要求。针对异地光源多中心数据同步训练的困难, 又提出同步和异步相结合的混合训练方式, 在不降低模型识别准确率的同时, 显著提升了模型的训练速度。在光源 CXIDB-76 公开数据集上的实验结果表明, 相比 FedAvg, FedKL 能够提升准确率和 F1 分数, 分别提升了 25.2 个百分点和 0.419。

关键词 隐私保护; 图像筛选; 联邦学习; 布拉格斑点; 相对熵

中图分类号 O432

文献标志码 A

DOI: 10.3788/LOP220950

Diffraction Image Screening of Radiation Facilities Based on Privacy Protection Mechanism

Xu Kang^{1,2}, Zhu Yongxin^{2*}, Wu Bo², Zheng Xiaoying², Chen Lingyao³¹School of Integrated Circuits, University of Chinese Academy of Sciences, Beijing 100049, China;²Shanghai Advanced Research Institute, Chinese Academy of Sciences, Shanghai 201210, China;³Shanghai Information Technology Research Center, Shanghai 201210, China

Abstract Synchrotron radiation facilities generate ultra-high-speed diffraction image data streams, which require data screening to reduce the pressure on data transmission and storage. However, competing research groups are reluctant to share such data, and existing deep learning-based screening methods cannot easily achieve effective training under privacy protection. Therefore, for the first time, this study applies the federated learning technology to the screening of radiation source diffraction images, and training data augmentation under privacy protection is realized by separating the data and the model. The Federated Kullback-Leibler (FedKL) screening method is also proposed to improve the global model update based on Kullback-Leibler divergence and data volume weights, thus reducing the complexity of the algorithm while obtaining high accuracy; further, this satisfies the high-precision processing requirements for high-speed data streams. To address the difficulties encountered in data synchronization training for multiple centers of remote light sources, this paper also proposes a hybrid training method that combines the synchronous and asynchronous approaches; this significantly improves the training speed of the model without reducing the recognition accuracy. Experiments on the light source CXIDB-76 public dataset reveal that FedKL can improve the accuracy and F1 score by 25.2 percentage points and 0.419, respectively, compared with FedAvg.

Key words privacy protection; image screening; federated learning; Bragg spot; relative entropy

1 引言

同步辐射光源^[1-2]具有高强度、高亮度、高准直性

等特性, 可用于多学科的前沿基础研究, 其中, 使用硬 X 射线对蛋白质等大分子进行晶体衍射是常用的实验方法。但是, 硬 X 射线自由电子激光^[3-5]的平均数据带

收稿日期: 2022-03-10; 修回日期: 2022-04-01; 录用日期: 2022-04-09; 网络首发日期: 2022-04-20

基金项目: 国家自然科学基金(U2032125)、科技部 SKA 专项(2020SKA0120202)

通信作者: *zhuyongxin@sari.ac.cn

宽是 2~20 GB/s, 峰值为 100 GB/s, 高速的数据流给数据的存储和传输带来了巨大的压力。为了减轻后续数据传输和存储的压力, 可以采用深度学习对无效的图像数据进行筛选和抛弃。

神经网络模型的训练需要大量的图像数据, 这些图像数据源于不同的研究小组, 包含重大科学发现, 具有高度的科研价值, 但是存在数据版权保护的困难, 为了避免数据泄露造成的科研损失, 需要一种基于数据隐私保护的深度学习方法。为了解决数据隐私保护的问题, McMahan 等^[6]提出一种联邦学习方法 (FedAvg), FedAvg 采用均值化的方式聚合多个客户端的模型梯度信息, 从而对全局模型进行更新, 该模型的复杂度较低, 但是该方法对衍射图像的识别准确率较低, 无法应用于衍射图像筛选。Reisizadeh 等^[7]提出了一种周期平均和量化的处理方法 (FedPAQ), FedPAQ 对客户端和服务端通讯的数据进行了压缩。Hamer 等^[8]提出了 FedBoost 算法, 该算法通过学习一组预先训练好的基本预测因子 (base predictors) 实现联邦集成 (federated ensembles), 降低了下行通信成本。但是 FedPAQ 和 FedBoost 算法的复杂度较高, 时间开销较大, 对高速数据的处理效率比较低。

晶体衍射图像中包含的布拉格斑点^[9]是筛选图像的重要参考, 但是布拉格斑点是极小尺寸^[10-12]的斑点, 很难和噪声产生的斑点区分开, 此外晶体衍射图像的尺寸过大, 图像数据产生速度过快, 这些都给晶体衍射图像的筛选带来了巨大的困难。

因此, 本文提出了一种基于联邦学习的晶体衍射图像筛选方法 (Federated Kullback-Leibler, FedKL), 采用联邦学习方式对全局模型进行训练, 用户只需要上传训练参数而无需上传数据本身, 从而保护了归属于不同研究小组的数据的隐私。相比于 FedAvg, FedKL 的准确率大幅提高, 同时, 复杂度显著低于 FedPAQ 和 FedBoost, 能够满足对高速图像数据处理的要求。此外通过灰度转换、随机裁剪等预处理方法, 提高了对晶体衍射图像的识别准确率。并且, 采用同步和异步相结合的混合训练方式, 在不降低模型识别准确率的同时, 大大加快了模型的训练速度。FedKL 可以应用于辐射光源衍射图像的隐私保护和筛选, 也可以推广到工业互联网等低延迟、高带宽的隐私计算应用场景。

本文的创新点和主要贡献总结如下: 1) 在光源晶体衍射图像筛选中应用联邦学习方法, 支持模型训练数据和全局模型的分离, 在不泄露数据隐私的前提下, 实现了训练数据的增广; 2) 提出筛选方法 FedKL, 基于改进的 KL 散度和数据量权重对全局模型更新进行改进, 并配合灰度处理、随机裁剪等预处理方法, 降低了算法的复杂度并提高了对晶体衍射图像的识别准确率; 3) 针对异地光源多中心数据的同步训练困难, 提出同步和异步相结合的混合训练方式, 在不降低识别准

准确率的同时显著加快了模型的训练速度; 4) 在数据集 CXIDB-76 的子数据集 L498 上完成了实验, 使用 FedKL 的模型的准确率、精准度、召回率及 F1 分数均得到很大提高。

2 相关工作

2.1 传统图像处理方法

传统的晶体衍射图像检测^[13-15]很难实现全自动化的晶体定心和晶体检测。Ito 等^[16]使用基于深度学习的方法构造了一个全自动化的晶体定心方法, 该方法在不使用 X 射线辐照晶体的情况下, 实现了全自动的精确晶体定心。在中子晶体学中, Sullivan 等^[17]使用基于 U-Net 的神经网络 (通常用于图像分割) 对布拉格衍射峰进行操作, 包括预测峰形和细化峰位。因此, 深度学习成为一种有效的图像处理方法, 一些经典的深度学习模型也可以用于检测布拉格斑点。相比于 Alex-Net^[18], VGG^[19-20]使用连续的 3 个 3×3 的卷积核来替代 Alex-Net 中较大的卷积核, 使用多个小的卷积核代替大的卷积核, 可以减少模型训练的参数量, 并且可以增加网络的深度, 保证了学习更多复杂模式的能力。ResNet^[21-22]是残差网络的简称, 它的出现解决了网络深度对训练精度的影响, 解决了梯度消失的问题。DenseNet^[23-24]是一种更为激进的密集连接机制, 所有的层都是互相连接的, 具体来说, 每一层都使用之前所有层的特征图作为输入, 同时当前层的特征图也是后续所有层的输入。DenseNet 解决了梯度消失的问题, 并且具有很好的抗过拟合的特性, 有效地减少了模型的参数量和计算量。

上述方法都可以应用于晶体衍射图像的识别中, 但是这些模型都需要依托于数据才能训练, 模型和数据无法分开, 对于研究小组的私有数据来说, 这样可能造成数据的泄露, 不能保证涉密信息的安全。

2.2 联邦学习实现隐私保护

联邦学习^[25-26]是近年来新兴的一种基础的人工智能技术, 模型结构如图 1 所示。McMahan 等提出一种联邦学习方法, 在训练数据没有共享的情况下, 通过该方法可以得到联合训练后的模型。具体来说, 各个数据的所有者, 无论是企业还是个人, 或者是机构, 他们的私有数据是保存在本地的, 在这些所有者的本地数据上单独进行训练, 可以得到包含各自特征的模型, 这些模型是属于客户端的模型。联邦学习的关键思想是将这些模型的梯度信息上传到服务器, 在服务器上对模型进行聚合, 聚合后的模型会被再次分发到各个客户端, 经过多次迭代, 使目标函数最小化, 最终达到模型收敛的效果。

$$\min_{\omega} f(\omega) = \sum_{i=1}^n p_i F_i(\omega) = E_i[F_i(\omega)], \quad (1)$$

式中: n 表示客户端的数量; p_i 表示每个客户端在聚合模型中所占的比例。

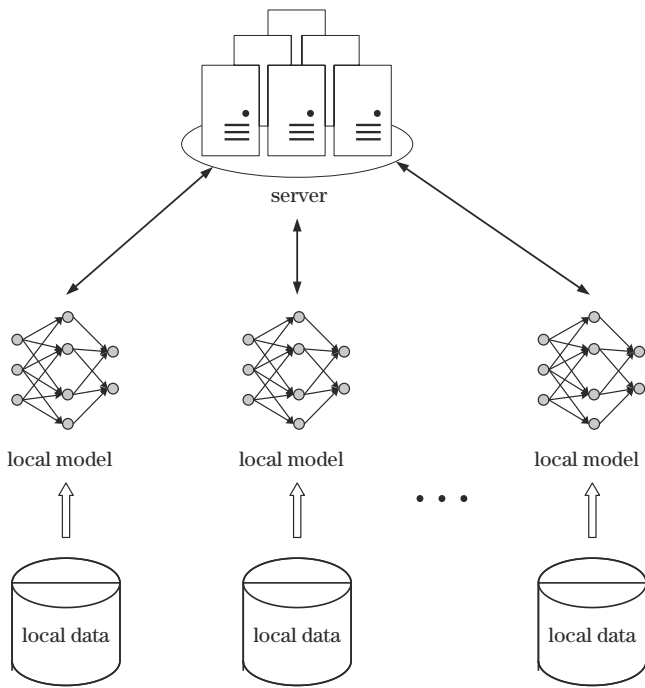


图 1 联邦学习模型

Fig. 1 Federated learning model

尽管与联邦学习相关的研究近年来受到广泛关注,联邦学习给数据隐私保护提供了有效的方式,但是针对不同的数据集、不同的应用场景、不同的网络结构、不同的优化器,联邦学习的应用效果也不尽相同。Reisizadeh等^[7]提出了一种周期平均和量化的处理方法(FedPAQ),FedPAQ允许网络中的客户端在与中央服务器同步之前进行本地训练,仅将活跃客户端的量化版本的更新发送回中央服务器,压缩了客户端和服务器通讯的数据。FedAvg^[6]采用均值化的方式对全局模型进行聚合和更新,能够使全局模型综合所有客户端模型的特点,缺点是遇到训练很差的客户端模型时,会大大降低全局模型的拟合效果。联邦学习的复杂度受到客户端模型复杂度的影响,文献^[27]表明,VGG、ResNet等网络的复杂度较高,因此必须选择复杂度较低的联邦学习框架,避免复杂度叠加。相比FedPAQ和FedBoost,FedAvg的复杂度较低,适合高数据流量、低延迟要求的光源场景,因此采用FedAvg作为对比方法。

2.3 联邦学习训练策略

联邦学习算法的调度方式分为同步训练和异步训练。FedAvg算法采用同步的方式对客户端进行调度,全局模型需要所有客户端训练完成后再进行模型更新。同步训练的优点是全局模型的更新是基于所有的客户端的梯度信息的,能够保证模型训练的稳定性;缺点是由于各个客户端的训练时间长短不一,通信开销各不相同,并且在实际运行中,网络环境的差异也导致中央服务器接收到客户端梯度参数的时间也不尽相同,耗时短的客户终端要等待耗时长客户终端本次训练

完成,才能进行下一次训练,这会导致大量客户端处于等待状态,浪费了设备资源,同时也会造成整体训练时间的增加。

异步训练的中央服务器采用“先到先更新”的原则,当服务器接收到客户端的梯度信息后,马上对全局模型进行更新,并将更新后的模型下发至该客户端。异步训练能充分调用客户端的资源,所有客户端都无需等待,在相同的时间内,异步训练的迭代次数远远大于同步训练的迭代次数,大大节约模型训练时间。但是,异步训练的服务器每次只采用一个客户端的梯度信息进行更新,并且可能会收到很早之前训练的客户端的模型参数,造成客户端之间版本差别过大和全局模型收敛的不稳定。

3 基于联邦学习的晶体衍射图像筛选算法 FedKL

提出了一种基于联邦学习的晶体衍射图像筛选方法 FedKL。针对晶体衍射图像数据的特征,在数据预处理阶段采用了灰度转换、中心裁剪、随机裁剪的方法对晶体衍射图像进行预处理,使各个研究小组训练模型的准确率更高,速度更快,并且根据调研分析,这应该是首次将联邦学习应用于光源晶体衍射图像的筛选中。

在中央服务器进行全局模型聚合时,由于各个研究小组的数据量和数据分布不同,根据各个研究小组的本地数据量和数据分布的权重对全局模型进行聚合和更新。具体来说,数据集的大小直接影响模型的训练效果,小样本数据集得到的模型有一定的局限性,大样本数据集得到的模型更具有泛化性,所以,样本量越多的数据集模型在全局模型中的占比就越大。除此之外,数据集的数据分布也对模型的准确率有很大的影响,例如,正负样本失衡的数据集得到的模型的鲁棒性更低,识别准确率也更低,这样的数据集训练出的研究小组模型在中央服务器模型聚合时会影响全局模型的泛化性能,造成全局模型的准确率下降。所以,FedKL策略采用研究小组数据的数据量和数据分布作为权重对全局模型进行聚合和更新。

3.1 数据预处理

本实验采用的晶体衍射图像是位深度为 16 的单通道灰度图像,相比于普通的 8 位深单通道灰度图像,16 位深灰度图像的像素值会超过 255。因此,在进行模型训练前,需要将 16 位深灰度图像转换为 8 位深灰度图像。

为了提高模型识别的准确率,降低训练的开销,并且尽可能地增强晶体衍射图像的特征,针对晶体衍射图像的特点,本文采用了数据增强的方式来增加标注数据的多样性。对于数据增强,通过翻转、裁剪、调整大小、旋转原始图像以及缩放图像的像素值来额外生成可以用于训练的数据。本实验的数据集中有图片 2000 张,属于小样本数据集,因此需要对数据进行数

据增强,这样的数据增强可以减少模型训练中的过拟合现象。布拉格衍射是电子均匀地撞击晶体内的原子时产生的布拉格衍射峰,在图像中反映出来的就是布拉格斑点。衍射图像的布拉格斑点多数位于图像的中心区域,相比于图 2(b)标记为 Miss 的图像,图 2(a)标记为 Hit 的图像的中心区域的布拉格斑点更为密集,因此可以根据图像布拉格斑点的分布特点对图像进行预处理。具体来说,原始的晶体衍射图像的分辨率为 960×960 ,由于布拉格斑点位于图像的中心区域,大概 $1/3$,为了降低计算成本,在图像中心区域裁剪出 320×320 大小的子图像,这种裁剪方式是中心裁剪,可以减小训练图像的尺寸,减少模型训练的参数量,大大加快模型的训练速度。

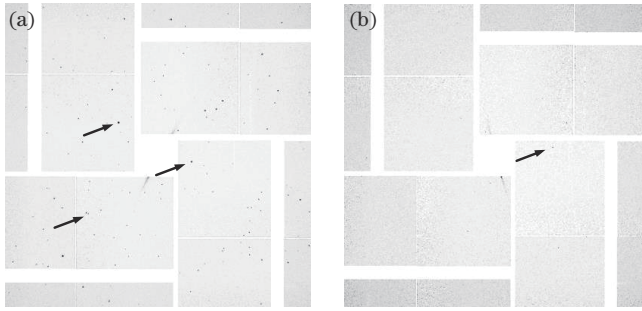


图 2 晶体衍射图像。(a)包含大量布拉格斑点;(b)包含少量布拉格斑点

Fig. 2 Crystal diffraction images. (a) Including a lot of Bragg spots; (b) including a few Bragg spots

除了对图像进行中心裁剪,实验还采取了随机裁剪的方式,在中心裁剪的基础上,将图像上下或者左右移动若干个像素点,然后进行裁剪。随机裁剪通过引入输入图像的不同区域来增加训练样本的数据量,这样可以降低模型过拟合的可能性,并且可以提高模型的性能。

3.2 全局模型聚合

3.2.1 Kullback-Leibler 散度作为模型权重

KL 散度^[28-30]也叫相对熵,用于度量两个概率分布之间的差异程度。设 $P(x)$ 和 $Q(x)$ 是随机变量 X 上的两个概率分布,则在离散和连续随机变量的情形下,相对熵的定义分别为

$$KL(P||Q) = \sum P(x) \log_2 \frac{P(x)}{Q(x)}, \quad (2)$$

$$KL(P||Q) = \int P(x) \log_2 \frac{P(x)}{Q(x)} dx. \quad (3)$$

实验中采用的数据集的正负样本属于离散型数据,所以计算公式为

$$D_{KL}(p||q) = \sum_{i=1}^N p(x_i) \times \log_2 \frac{p(x_i)}{q(x_i)} = \sum_{i=1}^N p(x_i) [\log_2 p(x_i) - \log_2 q(x_i)]. \quad (4)$$

在信息理论中,相对熵用来度量使用基于 P 的编码来编码来自 Q 的样本平均所需的额外比特个数。典型情况下, P 表示数据的真实分布, Q 表示数据的理论分布、模型分布,或 Q 的近似分布。在神经网络的训练中,正负样本的数量越均衡,对模型的训练效果的影响就越小,其他参数相同的情况下,模型的训练效果就越好。对于二分类的图像处理模型,理想数据的正负样本比例为 1:1,实际样本的比例越接近 1:1,模型分类结果越好,因此,以概率为 0.5 的 0-1 分布作为标准分布,研究小组数据分布与标准分布的距离作为研究小组的权重来衡量数据集的优劣,距离越小,权重越大。二分类的 KL 散度的公式为

$$D_{KL}(p||q) = \sum_{i=1}^2 p(x_i) \log_2 p(x_i) - \sum_{i=1}^2 p(x_i) \log_2 q(x_i) = \sum_{i=1}^2 p(x_i) \log_2 p(x_i) - \sum_{i=1}^2 p(x_i) \log_2 2^{-1} = \sum_{i=1}^2 p(x_i) \log_2 p(x_i) + 1, \quad (5)$$

式中: $P(x)$ 表示研究小组的数据分布; $Q \sim B(1, 0.5)$ 。令研究小组的正类数据的概率为 θ , 负类数据的概率为 $1 - \theta$, KL 散度表示为

$$D_{KL}(p||q) = \theta \log_2 \theta + (1 - \theta) \log_2 (1 - \theta) + 1. \quad (6)$$

令 W 为研究小组数据的 KL 散度指标,研究小组数据越均衡,即 θ 越趋近于 0.5, W 越大,所以有

$$W = 1 - D_{KL}(p||q) = 1 - \theta \log_2 \theta - (1 - \theta) \log_2 (1 - \theta) - 1 = -\theta \log_2 \theta - (1 - \theta) \log_2 (1 - \theta). \quad (7)$$

3.2.2 数据量作为模型权重

各个研究小组的数据量不同,大样本数据集得到的模型的鲁棒性更强,泛化能力更好,因此在聚合全局模型时,研究小组数据量也是影响模型聚合的重要参数。考虑到中央服务器和研究小组之间的通信开销,在实验中,研究小组模型每进行 5 轮训练,和中央服务器通信 1 次,更新全局模型,中央服务器将更新后的模型下发到各个研究小组继续训练,经过一定轮次的迭代,当全局模型在中央服务器的验证集上的损失不再下降,各个研究小组的模型在自有的本地数据的验证集上的损失不再下降,就可以得到训练完成的全局模型。

图 3 为 FedKL 模型的示意图, FedKL 算法的具体实现如下。

1) 中央服务器 S 将初始化的模型下发至各个研究小组。

2) 研究小组 C_i 基于本地自有数据训练模型 γ , 计算得到模型的梯度参数 g_i :

$$g_i = \text{train}(C_i, \gamma). \quad (8)$$

3) 研究小组 C_i 将梯度参数 g_i 、数据集 KL 散度指标 W_i 、数据集样本数量 n_i 上传到中央服务器 S 。

4) 中央服务器 S 整合所有接入服务器的研究小

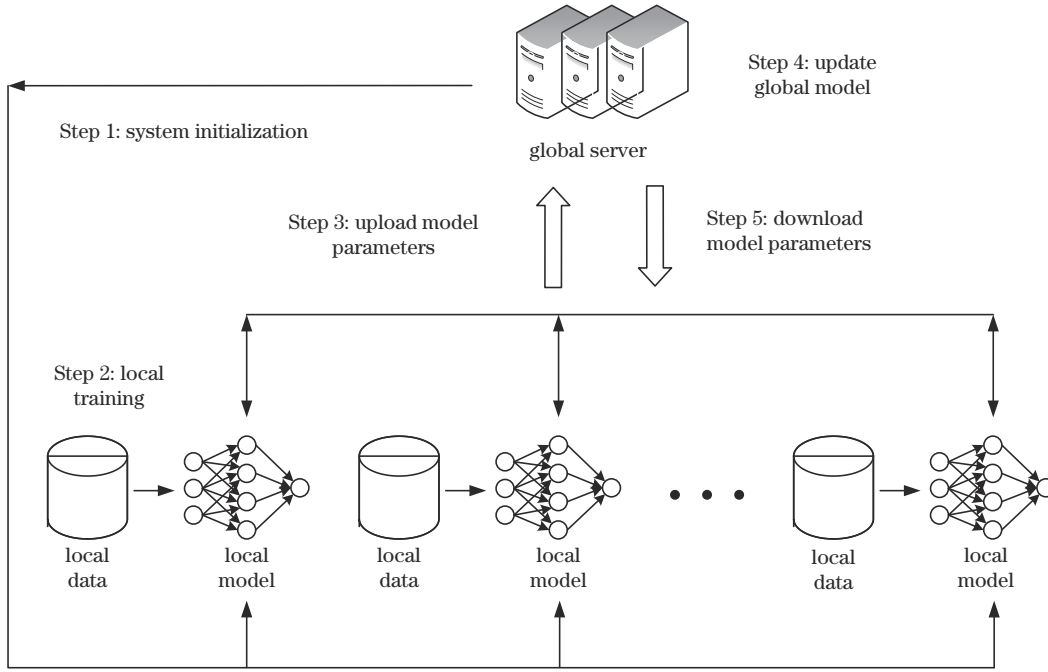


图 3 FedKL 模型
Fig. 3 FedKL model

组上传的梯度信息,计算出数据量权重 p_i 、数据分布权重 ω_i ,并更新全局模型的参数 g_s 。

$$\omega_i = \frac{W_i}{\sum W_i}, p_i = \frac{n_i}{\sum n_i}, \quad (9)$$

$$g_s = \frac{1}{2} \sum [(p_i + \omega_i) \times g_i], \quad (10)$$

$$\gamma = \text{update}(g_s, \gamma). \quad (11)$$

5)中央服务器S将全局模型的参数下发至研究小组,研究小组更新本地模型。

6)重复步骤2)~5),直到研究小组模型的损失不再降低。

3.3 基于混合训练的模型训练策略

为了加速模型训练,同时保证模型训练的稳定性 and 可靠性,FedKL算法采用同步训练和异步训练相结合的混合训练模式。

作为联邦学习模型的基础训练方式,同步训练的优点是稳定,中央服务器进行每次更新是基于所有客户端模型的,个别客户端训练的效果对全局模型的影响比较小,这样的代价是全局模型每次更新的周期比较长,模型整体的训练时间比较长。异步训练是基于单个客户端进行全局模型更新的,全局模型版本的迭代速度比较快,能够快速得到表现较好的模型。但是,当不同客户端训练时间的差异过大时,各个客户端之间的模型版本差别比较大。此外,异步训练很可能因为某个客户端的宕机造成全局模型的准确率突变式下降,直接影响到后续的训练,相比于同步训练,异步训练的稳定性比较差。因此,为了结合两者的优点,采用混合训练的方式,伪代码如图4所示。

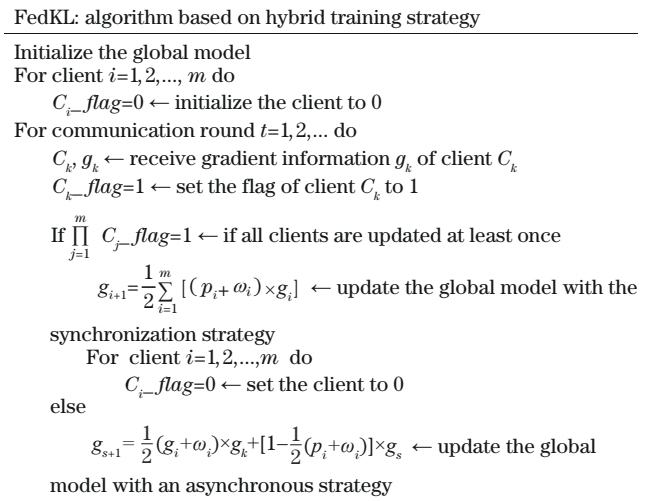


图 4 混合训练方式的流程
Fig.4 Process of mixed training method

在混合训练中,当所有的客户端更新至少一次时,采取同步更新策略,所有的客户端都会接收到最新的模型参数,客户端较早的参数就会被覆盖掉,不参与后续的参数混合。相比仅考虑客户端的数据量权重,结合了客户端数据量和分布情况的FedKL策略的权重更能准确反映客户端的重要程度。不同数据量的客户端分布情况造成的权重不同,源于数据量权重和数据分布权重(即改进的KL散度),FedKL策略得到的综合权重更多面化、更准确。

4 实验及分析

4.1 数据集描述

本实验采用开源 coherent X-ray imaging data bank

(CXIDB)数据集^[31],它是一个相干X射线成像数据库,其中CXIDB-76数据集包含了嗜热菌蛋白酶、氢化酶、亲环蛋白A等蛋白质分子的晶体衍射图像,L498是CXIDB-76数据集的一部分,包含了2000张晶体衍射图像。L498的图像分为三类,具体分布如表1所示,Hit表示衍射图像中包含布拉格斑点,Miss表示衍射图像中不包含布拉格斑点,Maybe表示衍射图像中可能包含布拉格斑点。所有标记为Miss的图像在进行数据传输前会被丢弃,这会大大减少数据传输和处理的工作量。

表1 L498数据分布
Table 1 L498 data distribution

Label	Number of images	Proportion / %
Hit	148	7.4
Miss	1355	67.76
Maybe	498	24.9

为了模拟联邦学习研究小组的数据分布,将L498的数据按照8:2分为总训练集和总测试集。实验中模拟了5个客户端的数据分布,每个客户端随机地从总训练集中抽取晶体衍射图像作为本地数据。然后将每个客户端的数据按照8:2分为训练集和验证集,并且保证训练集和验证集中各个标签类别的比例相同。实验中各个客户端的数据分布如图5所示。

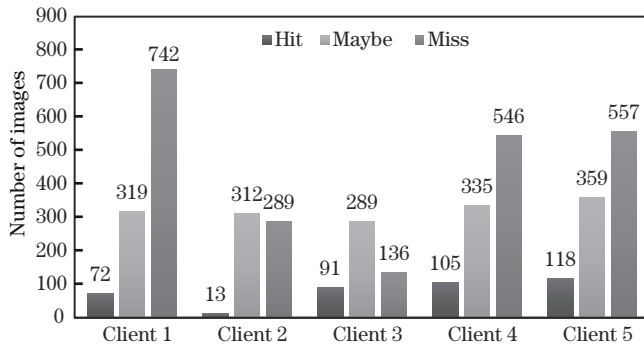


图5 客户端数据分布

Fig. 5 Client data distribution

4.2 实验设置

本文的实验是在包含4块RTX2080 GPU的服务器上进行的,在Docker中模拟了1个中央服务器和5个客户端。中央服务器包含400张图片,用于测试,各个客户端的晶体衍射图像的数据量随机分布。实验中使用了机器学习框架PyTorch作为底层的机器学习训练库,并且利用Python 3编程软件实现了各个客户端的模型训练和中央服务器的模型聚合。

在客户端的模型选择中,选取ResNet、VGG、DenseNet网络分别对中央服务器的模型进行聚合。由于ResNet、VGG、DenseNet模型的复杂度比较高,为了避免复杂度叠加导致的算法复杂度剧增,所以需要

要选择复杂度较低的联邦学习框架作为对比。相比FedPAQ和FedBoost,FedAvg的复杂度更低,适合高数据流量、低延迟要求的光源场景,因此,选择FedAvg作为对比方法。在实际的训练中,客户端根据自有数据进行15个epoch的预训练,然后将训练后的模型的梯度信息上传至中央服务器,服务器判断是否所有的客户端已经更新至少一次,根据判断结果选择同步或者异步更新全局模型,中央服务器聚合全局模型后下发至各个客户端,客户端继续训练5个epoch后,重复上传-聚合模型(同步/异步)-分发-训练的步骤,直到损失函数不再下降。损失函数不再下降的具体阈值是根据选择的模型不同有所不同的,在本实验中,以DenseNet121为例,客户端的损失函数的阈值为0.05。相比传统同步训练,混合训练的策略速度更快,同时准确率也没有下降。

4.3 实验评价指标

模型分类准确率是模型测试中一个非常重要的指标,采用准确率、精准度、召回率及F1分数对模型识别的准确率进行评价。各指标的计算公式分别为

$$A_c = \frac{N_{TP} + N_{TN}}{N_{TP} + N_{FP} + N_{TN} + N_{FN}}, \quad (12)$$

$$P = \frac{N_{TP}}{N_{TP} + N_{FP}}, \quad (13)$$

$$R = \frac{N_{TP}}{N_{TP} + N_{FN}}, \quad (14)$$

$$F_1 = \frac{2PR}{P + R}, \quad (15)$$

式中:TP表示标签为正类,预测结果为正类;FP表示标签为负类,预测结果为正类;FN表示标签为正类,预测结果为负类。精准度针对模型预测结果,表示预测为正类的样本中有多少是真正的正类样本;召回率针对被测试样本,表示样本中有多少样本被正确地预测了;F1分数是精准度和召回率的加权调和平均,在尽可能提高精准度和召回率的同时,希望两者之间的差异尽可能小。

4.4 实验结果及分析

采用ResNet^[22]、VGG^[20]、DenseNet^[24]网络结构结合联邦学习架构,对晶体衍射图像进行分类测试,测试集的数据分布如表2所示。

表2 测试集数据分布

Table 2 Test set data distribution

Label	Number of images	Proportion / %
Hit	66	16.5
Miss	190	47.5
Maybe	144	36.0

在实验中,采用DenseNet121、ResNet18、ResNet50和VGG19分别在FedAvg^[6]和FedKL两种方式下对全局模型进行聚合,模型的输入为单通道8

位 960×960 大小的灰度图像,各个 Client 数据集中训练集和验证集的比例为 8:2。所有模型通过 PyTorch 导入标准网络模型架构中,未导入预训练模型参数。图 6 为在 FedAvg 和 FedKL 两种方式下,不同模型的准确率对比。从图 6 可以看出,相比于 FedAvg,采用 FedKL 的模型的准确率有了大幅度提高, DenseNet121、ResNet18、ResNet50 和 VGG19 四种模型的准确率分别提高了 25.2 个百分点、17.0 个百分点、14.0 个百分点、21.0 个百分点。对于晶体衍射图像的筛选来说,高准确率意味着能够减少图像分类中的失误,降低丢失重要观测数据的风险。

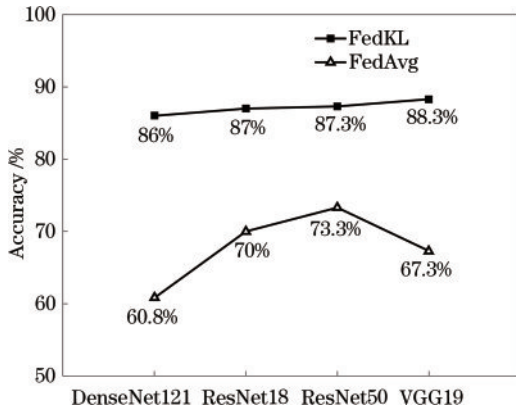


图 6 在 FedAvg 和 FedKL 两种方式下,不同模型的准确率对比
Fig. 6 Comparison of accuracy of different models under FedAvg and FedKL

图 7 为在 FedAvg 和 FedKL 两种方式下,不同模型运行时间和大小对比。从图 7 可以看出, FedAvg 和 FedKL 两种方式得到的模型大小一致, FedAvg 和 FedKL 的不同之处在于聚合全局模型时选择的策略不同,得到的全局模型的网络结构是一致的。FedAvg-DenseNet121 和 FedKL-DenseNet121 得到的全局模型都是 DenseNet121,区别在于参数值的不同。运行时间是各个模型测试 400 张图片所花费的总时

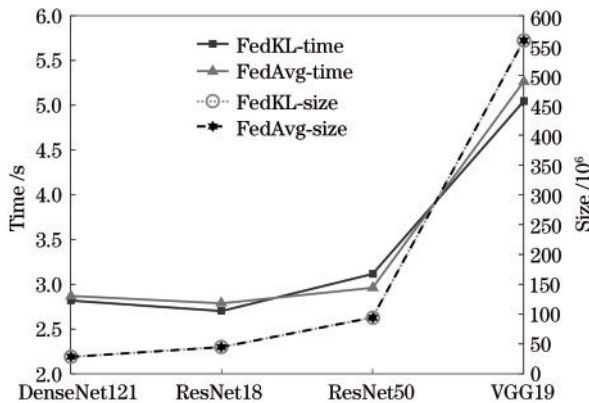


图 7 在 FedAvg 和 FedKL 两种方式下,不同模型运行时间和大小对比
Fig. 7 Comparison of running time and size of different models under FedAvg and FedKL

间,采用不同策略的模型运行时间没有明显的差异,耗时的区别主要来源于模型的大小,模型越大,加载模型的时间越长,模型的计算量越大,从而耗费的时间也就越多。从图 7 可以看出,随着模型大小的增加,模型的运行时间也在不断增加。

表 3 为采用 FedAvg 和 FedKL 算法的 DenseNet121 模型在测试集上分类结果的对比。从表 3 可以看出: FedAvg-DenseNet121 将测试集中 70% 的正类数据判别为负类,只有将 30% 的正类数据判断正确; FedKL-DenseNet121 将测试集中 84.8% 的正类数据判断正确,并且将 87.4% 的负类数据判断正确。在晶体衍射图像的筛选策略中,判定为负类的数据表示该数据不具有科研价值,需要被丢弃,所以 FedAvg-DenseNet121 将丢弃 70% 的正类数据和 94.7% 的负类数据,仅保留 30% 的正类数据; FedKL-DenseNet121 将丢弃 15.2% 的正类数据和 87.4% 的负类数据,保留 84.8% 的正类数据。对于晶体衍射图像筛选而言,丢弃正类数据过多会大大影响科研, FedAvg 会丢弃大量的正类数据; FedKL 在保留大量正类数据的情况下还能够丢弃大量的负类数据,在不影响科研的情况下大大减少负类数据占用的存储空间。

表 3 测试集的分类结果(以 DenseNet121 为例)
Table 3 Classification results on the test set (e.g. DenseNet121)

Model	Predict	Hit or Maybe Miss		
		Hit or Maybe	Miss	
FedAvg-DenseNet121	Label	Hit or Maybe	30.0%	70.0%
	Label	Miss	5.3%	94.7%
FedKL-DenseNet121	Label	Hit or Maybe	84.8%	15.2%
	Label	Miss	12.6%	87.4%

图 8~10 分别为在 FedAvg 和 FedKL 两种方式下,不同模型的精准度、召回率、F1 分数对比。可以看出: FedKL-DenseNet121 模型的精准度高于 FedAvg-DenseNet121; 相比于 FedAvg-DenseNet121, 召回率提

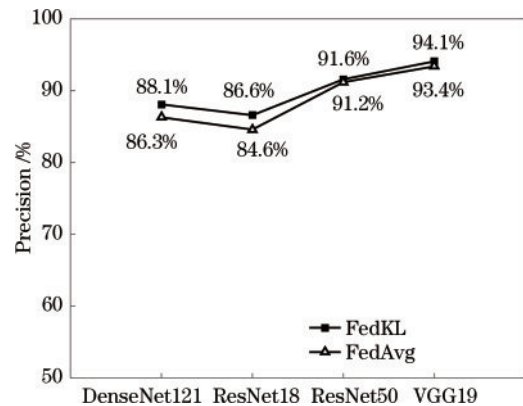


图 8 在 FedAvg 和 FedKL 两种方式下,不同模型的精准度对比
Fig. 8 Comparison of precision of different models under FedAvg and FedKL

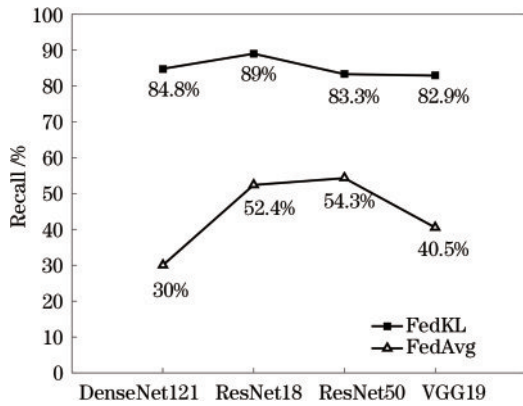


图 9 在 FedAvg 和 FedKL 两种方式下,不同模型的召回率对比
Fig. 9 Comparison of recall of different models under FedAvg and FedKL

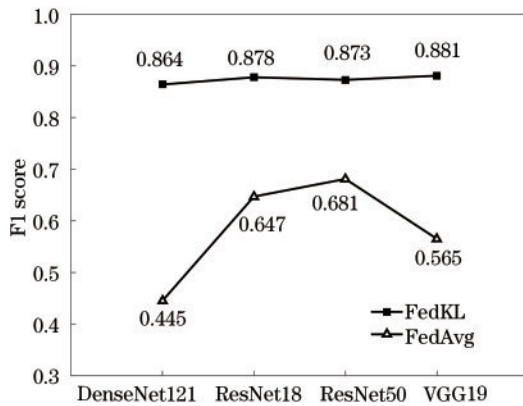


图 10 在 FedAvg 和 FedKL 两种方式下,不同模型的 F1 分数对比
Fig. 10 Comparison of F1 score of different models under FedAvg and FedKL

高了 54.8 个百分点,说明数据量的多少和数据分布的好坏确实会对全局模型分类效果造成影响;FedKL-DenseNet121 的 F1 分数提高了 0.419,说明 FedKL 策略的鲁棒性更好。

图 11 表示随着训练时间的增加,采用同步训练和混合训练两种方式进行训练的全局模型的准确率对

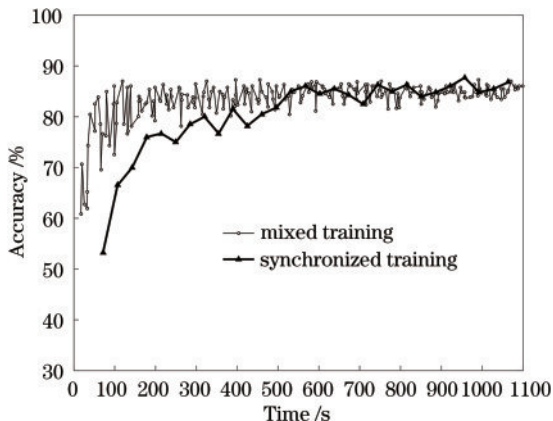


图 11 在 FedAvg 方式下,同步训练和混合训练的时间对比(以 DenseNet121 为例)

Fig. 11 Time comparison between synchronized training and mixed training under FedAvg (e.g. DenseNet121)

比。可以看出:同步训练每次更新全局模型需要的时间更长;在相同的时间内,混合训练更新的次数远远高于同步训练,并且混合训练的准确率在前期也大大高于同步训练;相比于同步训练,混合训练能够在更短的时间内达到模型收敛。

综上,对于晶体衍射图像分类的任务,FedKL 策略的分类效果更好,在保护研究小组数据安全的情况下得到了高准确度、高精度、高召回率及高 F1 分数的分类模型,并且在常用的网络模型上均取得了较好的效果。

5 结 论

为了应对光源单个线站有效数据不足、多中心数据源不同步、筛选准确率不足等困难,提出了一种基于联邦学习的晶体衍射图像筛选方法 FedKL。通过灰度转换、中心裁剪、随机裁剪的方法对晶体衍射图像进行预处理。在全局模型聚合阶段,创新性地使用改进的 KL 散度表示数据集的数据分布,并按照数据量和数据分布的权重对各个研究小组的模型参数进行加权计算,实现了研究小组数据和全局模型的分离。同时创新性地提出同步和异步相结合的混合训练方式,在不降低识别准确率的同时显著提升了模型的训练速度。在 CXIDB-76 数据集上进行了实验,相比于基线模型 FedAvg-DenseNet121, FedKL-DenseNet121 的准确率和 F1 分数分别提高了 25.2 个百分点和 0.419,实现了对晶体衍射图像的准确分类。此外联邦学习机制也能够有效地保护不同研究小组的数据隐私。

目前 FedKL 策略的优化器是 SGD,并且对所有的研究小组均采用同一模型。实际上不同的研究小组数据可以采用不同的网络模型进行训练,不同的网络模型可以采用不同的优化方法。如何对不同网络模型的梯度信息在中央服务器进行聚合、中央服务器要选择何种模型作为全局模型、采用何种优化方法寻找模型的最优解,这将是下一步工作研究的重点。

参 考 文 献

- [1] Liang C, Zhang W H, Wei Z S, et al. Transition-metal redox evolution and its effect on thermal stability of $\text{LiNi}_x\text{Co}_y\text{Mn}_z\text{O}_2$ based on synchrotron soft X-ray absorption spectroscopy[J]. Journal of Energy Chemistry, 2021, 59: 446-454.
- [2] 尚雷, 尚风雷, 孙振彪, 等. 先进同步辐射光源特种电源概述[J]. 强激光与粒子束, 2019, 31(4): 12-17.
Shang L, Shang F L, Sun Z B, et al. Overview of special power supplies for advanced synchrotron radiation source[J]. High Power Laser and Particle Beams, 2019, 31(4): 12-17.
- [3] 张豪, 郭海涛, 许彦涛, 等. 用于红外激光传输的硫系玻璃光纤研究进展[J]. 中国激光, 2022, 49(1): 0101007.
Zhang H, Guo H T, Xu Y T, et al. Research progress in chalcogenide glass fibers for infrared laser delivery[J]. Chinese Journal of Lasers, 2022, 49(1): 0101007.
- [4] 孙伟义, 黄家鹏, 陈丽明, 等. 10 W 量级高功率中红外

- 超快光纤激光系统中色散管理的仿真设计[J]. 中国激光, 2022, 49(1): 0101012.
- Sun W Y, Huang J P, Chen L M, et al. Design of a 10 W level dispersion-managed high-power ultrafast mid-infrared fiber laser system[J]. Chinese Journal of Lasers, 2022, 49(1): 0101012.
- [5] 徐昌骏, 张集权, 刘墨, 等. 基于钛掺杂 ZBYA 玻璃光纤的中红外激光研究[J]. 中国激光, 2022, 49(1): 0101016.
- Xu C J, Zhang J Q, Liu M, et al. Midinfrared laser in Ho³⁺-doped ZBYA glass fiber[J]. Chinese Journal of Lasers, 2022, 49(1): 0101016.
- [6] McMahan H B, Moore E, Ramage D, et al. Communication-efficient learning of deep networks from decentralized data[EB/OL]. (2016-02-17) [2022-02-04]. <https://arxiv.org/abs/1602.05629>.
- [7] Reiszadeh A, Mokhtari A, Hassani H, et al. FedPAQ: a communication-efficient federated learning method with periodic averaging and quantization[EB/OL]. (2019-09-28)[2022-02-04]. <https://arxiv.org/abs/1909.13014>.
- [8] Hamer J, Mohri M, Suresh A T. Fedboost: a communication-efficient algorithm for federated learning [C]//Proceedings of the 37th International Conference on Machine Learning, July 13-18, 2020, Virtual Event. Cambridge: PMLR, 2020: 3973-3983.
- [9] Nakajima H, Kotani A, Harada K, et al. Electron diffraction covering a wide angular range from Bragg diffraction to small-angle diffraction[J]. Microscopy, 2018, 67(4): 207-213.
- [10] 邓世杰, 王海晏, 徐安, 等. 基于对抗生长的目标检测方法[J]. 光学学报, 2022, 42(2): 0210002.
- Deng S J, Wang H Y, Xu A, et al. Target detection method based on antigrowth[J]. Acta Optica Sinica, 2022, 42(2): 0210002.
- [11] 张宇, 张焱, 石志广, 等. 基于图像衍生的红外无人机图像仿真方法研究[J]. 光学学报, 2022, 42(2): 0210003.
- Zhang Y, Zhang Y, Shi Z G, et al. Image simulation method of infrared UAV based on image derivation[J]. Acta Optica Sinica, 2022, 42(2): 0210003.
- [12] 朱江平, 王睿珂, 段智涓, 等. 基于多尺度注意力机制相位展开的三维人脸建模[J]. 光学学报, 2022, 42(1): 0112005.
- Zhu J P, Wang R K, Duan Z J, et al. Three-dimensional face modeling based on multi-scale attention phase unwrapping[J]. Acta Optica Sinica, 2022, 42(1): 0112005.
- [13] Frank M, Carlson D B, Hunter M S, et al. Femtosecond X-ray diffraction from two-dimensional protein crystals[J]. IUCrJ, 2014, 1(2): 95-100.
- [14] Harder R. Deep neural networks in real-time coherent diffraction imaging[J]. IUCrJ, 2021, 8(1): 1-3.
- [15] Rivenson Y, Wu Y C, Ozcan A. Deep learning in holography and coherent imaging[J]. Light: Science & Applications, 2019, 8: 85.
- [16] Ito S, Ueno G, Yamamoto M. DeepCentering: fully automated crystal centering using deep learning for macromolecular crystallography[J]. Journal of Synchrotron Radiation, 2019, 26(4): 1361-1366.
- [17] Sullivan B, Archibald R, Azadmanesh J, et al. BraggNet: integrating Bragg peaks using neural networks[J]. Journal of Applied Crystallography, 2019, 52(4): 854-863.
- [18] Krizhevsky A, Sutskever I, Hinton G E. ImageNet classification with deep convolutional neural networks [C]//Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 1, December 3-6, 2012, Lake Tahoe, Nevada, USA. New York: ACM Press, 2012: 1097-1105.
- [19] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition[EB/OL]. (2014-09-04)[2022-02-04]. <https://arxiv.org/abs/1409.1556>.
- [20] Zhou Y P, Chang H Y, Lu Y H, et al. Improving the performance of VGG through different granularity feature combinations[J]. IEEE Access, 2020, 9: 26208-26220.
- [21] He K M, Zhang X Y, Ren S Q, et al. Deep residual learning for image recognition[C]//2016 IEEE Conference on Computer Vision and Pattern Recognition, June 27-30, 2016, Las Vegas, NV, USA. New York: IEEE Press, 2016: 770-778.
- [22] Zhao T M, Liu J, Wang Y, et al. Towards low-cost sign language gesture recognition leveraging wearables[J]. IEEE Transactions on Mobile Computing, 2021, 20(4): 1685-1701.
- [23] Huang G, Liu Z, van der Maaten L, et al. Densely connected convolutional networks[C]//2017 IEEE Conference on Computer Vision and Pattern Recognition, July 21-26, 2017, Honolulu, HI, USA. New York: IEEE Press, 2017: 2261-2269.
- [24] Anandhi V, Vinod P, Menon V G. Malware visualization and detection using DenseNets[J]. Personal and Ubiquitous Computing, 2021: 1-17.
- [25] Deng Y H, Lyu F, Ren J, et al. AUCTION: automated and quality-aware client selection framework for efficient federated learning[J]. IEEE Transactions on Parallel and Distributed Systems, 2022, 33(8): 1996-2009.
- [26] Zhou Y H, Ye Q, Lü J C. Communication-efficient federated learning with compensated overlap-FedAvg[J]. IEEE Transactions on Parallel and Distributed Systems, 2022, 33(1): 192-205.
- [27] 张顺, 龚怡宏, 王进军. 深度卷积神经网络的发展及其在计算机视觉领域的应用[J]. 计算机学报, 2019, 42(3): 453-482.
- Zhang S, Gong Y H, Wang J J. The development of deep convolution neural network and its applications on computer vision[J]. Chinese Journal of Computers, 2019, 42(3): 453-482.
- [28] Bulinski A, Dimitrov D. Statistical estimation of the Kullback-Leibler divergence[J]. Mathematics, 2021, 9(5): 544.
- [29] Ji S Y, Zhang Z Z, Ying S H, et al. Kullback - Leibler divergence metric learning[J]. IEEE Transactions on Cybernetics, 2022, 52(4): 2047-2058.
- [30] Alexopoulos A. The fractional Kullback-Leibler divergence[J]. Journal of Physics A: Mathematical and Theoretical, 2021, 54(7): 075001.
- [31] Ke T W, Brewster A S, Yu S X, et al. A convolutional neural network-based screening tool for X-ray serial crystallography[J]. Journal of Synchrotron Radiation, 2018, 25(3): 655-670.