

# TC609

## 全国数据标准化技术委员会技术文件

TC609-6-2025-03

### 全国一体化算力网 算力算效衡量技术要求

National integrated computing power network—Technical requirements for  
computing power and efficiency measurement

2025-08-29 发布

2025-08-29 实施

全国数据标准化技术委员会 发布



# 目 次

前 言 .....	II
1 范围 .....	1
2 规范性引用文件 .....	1
3 术语和定义 .....	1
4 缩略语 .....	1
5 算力网算力算效衡量功能概述 .....	2
5.1 总体要求 .....	2
5.2 算力算效衡量模块与资源层各模块的关系 .....	3
5.3 算力算效衡量指标 .....	3
6 算力算效衡量功能框架 .....	3
7 算力算效衡量功能技术要求 .....	4
7.1 资源采集技术要求 .....	4
7.2 计算资源衡量技术要求 .....	5
7.3 网络资源衡量技术要求 .....	5
7.4 存储资源衡量技术要求 .....	5
7.5 指标映射技术要求 .....	5
7.6 指标查询技术要求 .....	6
8 算力算效衡量指标技术要求 .....	6
8.1 应用场景 .....	6
8.2 每秒完成基本操作数定义 .....	7
8.3 每秒完成基本操作数计算方法 .....	8
8.4 每秒完成基本操作数和现有指标的映射 .....	9
参考文献 .....	10

# 前 言

请注意本文件的某些内容可能涉及专利。本文件的发布机构不承担识别专利的责任。

本文件由全国数据标准化技术委员会（SAC/TC609）提出并归口。

本文件起草单位：中国联合网络通信集团有限公司、中国科学院计算技术研究所、联通数字科技有限公司、飞腾信息技术有限公司、浪潮通信信息系统有限公司、紫金山实验室、鹏程实验室、中国移动通信有限公司研究院、中移（苏州）软件技术有限公司、联通智能制造科技产业（广东）有限公司、京东科技信息技术有限公司、曙光信息产业股份有限公司、曙光智算信息技术有限公司、湖北曙光三峡云大数据中心有限公司、深圳市尚数网科技有限公司、广东国腾量子科技有限公司、国家信息中心、中国信息通信研究院、国家数据发展研究院、中国电子技术标准化研究院、中国电信集团有限公司、天翼云科技有限公司、江苏未来网络集团有限公司。

# 全国一体化算力网 算力算效衡量技术要求

## 1 范围

本文件规定了全国一体化算力网的算力算效衡量功能框架、技术要求和算力算效衡量指标要求。  
本文件适用于全国一体化算力网的算力算效衡量功能开发和评估、算力算效的衡量工作。

## 2 规范性引用文件

本文件没有规范性引用文件。

## 3 术语和定义

下列术语和定义适用于本文件。

### 3.1

**算力** `computing power`

图形处理器（GPU）、中央处理器（CPU）等设备执行计算密集型任务的计算能力。

### 3.2

**算力资源** `computing power resources`

计算资源、存储资源以及节点内部网络资源，通过该节点的管控系统/运营平台进行抽象并对外提供算力资源服务，或称算力资源节点。

### 3.3

**算力网** `computing power network`

支撑数字经济高质量发展的关键基础设施，可通过网络连接多源异构、海量泛在算力，实现资源高效调度、设施绿色低碳、算力灵活供给、服务智能提供。

### 3.4

**通用计算** `general computing`

也称通算，是指由 CPU（中央处理器）芯片为主提供的计算能力，主要用于通用计算任务。

### 3.5

**智能计算** `intelligent computing`

也称智算，是指由 GPU（图形处理器）、FPGA（现场可编程门阵列）、ASIC（专用集成电路）等芯片为主提供的计算能力，主要用于人工智能相关的计算任务。

### 3.6

**超级计算** `supercomputing`

也称超算，是指基于大规模并行架构的超高性能计算能力，主要用于科学计算、工程仿真、气象预测等领域的超大规模、高复杂度计算任务。

## 4 缩略语

下列缩略语适用于本文件。

BOPS: 每秒完成基本操作数 (Basic Operations Per Second)  
 BOPs: 基本操作数 (Basic OPerations)  
 CE: 算效 (Computational Efficiency) 算力与功率的比值, BOPs/J.  
 CPU: 中央处理器 (Central Processing Unit)  
 FLOPS: 每秒浮点运算次数 (Floating-point Operations Per Second)  
 GRPC: 谷歌远程过程调用 (Google Remote Procedure Call)  
 GPU: 图形处理器 (Graphics Processing Unit)  
 IB: 无限带宽网络 (InfiniBand)  
 IOPS: 每秒输入/输出操作数 (Input/Output Operations Per Second)  
 JSON: JavaScript对象表示法 (JavaScript Object Notation)  
 NPU: 神经网络处理器 (Neural Processing Unit)  
 NVLink: 英伟达高速互联技术 (NVIDIA Link)  
 OPS: 每秒操作数 (Operations Per Second)  
 PCIe: 高速串行计算机扩展总线标准 (Peripheral Component Interconnect Express)  
 RESTful: 表征状态转移架构风格 (Representational State Transfer)  
 RoCE: 基于融合以太网的RDMA (RDMA over Converged Ethernet)  
 TPU: 张量处理器 (Tensor Processing Unit)  
 XML: 可扩展标记语言 (eXtensible Markup Language)

## 5 算力网算力算效衡量功能概述

### 5.1 总体要求

《全国一体化算力网 监测调度平台建设指南》给出了全国一体化算力网监测调度平台的建设指南, 提出了建设的参考架构和功能规范。全国一体化算力网监测调度平台由算力网资源层、调度层、运营层、监测层组成, 提供国家级、区域级、城市级的层次化算力网服务, 见图1, 资源层应通过资源并网将异属异构异地算力资源接入算力网, 调度层通过标准化接口对异属异构异地算力资源进行统一的算力资源管理, 运营层为供需多方提供运营服务能力, 监测层支持多维度的监测数据采集与分析。

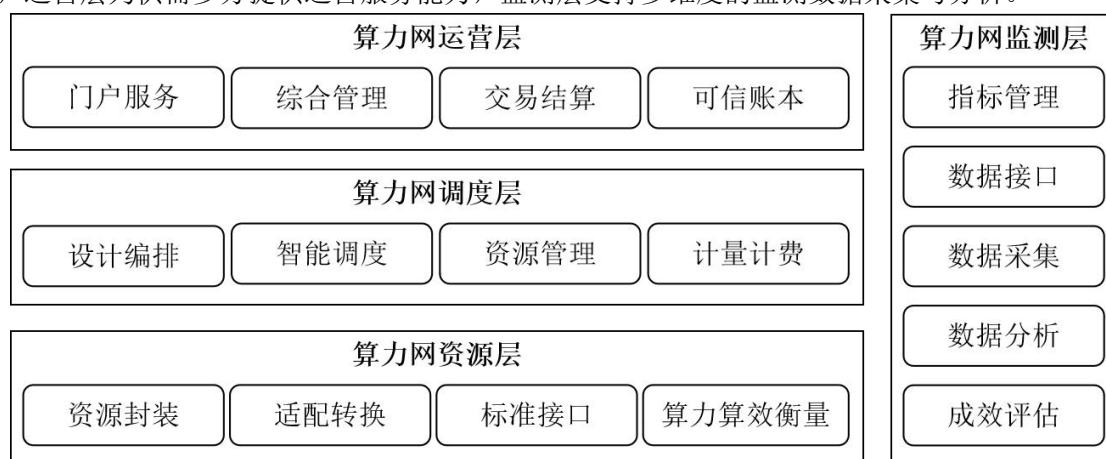


图1 全国一体化算力网监测调度平台总体框架示意图

作为算力网资源层的一个功能模块, 算力算效衡量模块对各类算力资源进行统一的衡量, 确保接入算力的质量。其主要功能有:

- 算力算效衡量支持计算资源的衡量、网络资源的衡量和存储资源的衡量；在计算资源衡量方面，综合考量其运算速度、核心数量、缓存大小等参数，精准评估计算能力；对于网络资源，聚焦网络带宽、延迟、丢包率等指标，确保数据传输的高效与稳定；在存储资源衡量上，依据存储容量、读写速度、可靠性等要素，确定存储服务的质量水平；
- 算力算效衡量支持对算力资源的量化描述，确认算力资源规格、功能、性能等信息，对符合接入条件的算力资源进行分类，根据分类确认对接方式。算力算效衡量应支持对算力资源的统一描述和建模，将各类算力资源抽象成统一描述，对算力资源的服务能力进行准确描述和评估。

算力算效衡量为通算、智算、超算、量子计算等各类算力资源解决异构算力资源的标准化描述问题，确保算力资源可量化、可比较、可交易。通过对算力算效衡量结果进行标准化，为算力网调度层提供统一的算力资源描述模型，支撑跨区域、跨行业、跨类型算力资源的协同调度与优化配置，推动算力资源在全国一体化算力网中实现高效流转、合理分配、有效交易，充分释放算力价值。

算力算效衡量功能基于算力算效衡量指标实现异构算力资源的标准化描述，支撑算力资源的可量化、可比较、可交易。

## 5.2 算力算效衡量模块与资源层各模块的关系

算力算效衡量作为资源层的核心功能模块，需与资源层其他模块（资源封装、适配转换、标准接口）协同，实现算力资源从接入到调度的标准化管理。算力算效衡量模块的衡量结果可同步至监测层，为算力资源运行状态监测提供数据支撑。

- 资源封装模块。资源封装模块在将算力资源信息进行抽象化、结构化的封装处理后提供给算力算效衡量模块，算力算效衡量模块基于封装后的资源类型，自动匹配衡量模型。衡量结果作为资源属性的核心字段，供其它模块进行查询与调用，实现资源的高效利用；
- 适配转换模块。通过适配转换模块的用户账户统一管理功能，通过模块内嵌的用户账户管理与身份认证机制，对算力算效衡量结果与具体用户或用户组进行绑定，确保衡量结果按需授权访问。适配转换模块对算力资源数据（如资源总量、余量、业务用量）进行标准化转换后，算力算效衡量模块基于统一格式的数据执行衡量计算；
- 标准接口模块。标准接口模块为算力算效衡量模块提供统一规范的接入通道，确保算力资源能顺畅接入衡量体系。算力算效衡量模块借助标准接口模块提供的规范接口，高效获取算力资源的基础信息，为后续准确执行算力算效衡量操作奠定基础，实现算力资源从接入到调度全流程的标准化、规范化管理。

算力算效衡量作为资源层关键功能模块，与资源封装、适配转换、标准接口模块紧密关联。通过各模块协同，实现算力资源从接入、衡量到调度的全流程标准化管理，提升算力网资源管理的规范性、准确性与高效性，为算力网的稳定运行和资源合理调配筑牢基础。

## 5.3 算力算效衡量指标

算力算效衡量功能对异构算力资源进行统一描述和评估，而算力算效衡量指标（BOPS）是衡量结果的核心输出，两者的关系体现在：算力算效衡量模块采集计算、网络、存储资源的静态/动态参数（如FLOPS、时延、IOPS），为BOPS等量化指标提供原始数据支撑；基于资源衡量结果，通过归一化方法将多维度指标转换为BOPS，实现用户侧与服务侧算力的统一量化表达；算力算效衡量指标通过衡量查询接口反馈至调度层，支撑算力资源的可量化、可比较、可交易。

## 6 算力算效衡量功能框架

算力算效衡量结果通过标准业务接口向上层调度层、运营层和监测层传输。

输出算力算效衡量指标至调度/运营层/监测层

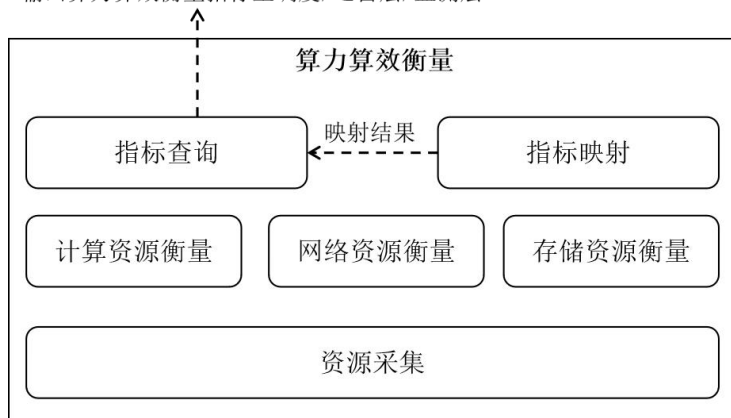


图2 算力算效衡量功能架构示意图

算力算效衡量的功能模块包括资源采集、计算资源衡量、网络资源衡量、存储资源衡量、指标映射和指标查询，见图2。各功能模块通过标准化数据接口实现协同联动：资源采集模块向衡量类模块（计算/网络/存储资源衡量）输出原始数据，衡量类模块向指标映射模块输出单维度结果，指标映射模块结合校准验证结果生成最终指标，通过指标查询模块向上层输出，安全与权限管理模块贯穿全流程。各个功能模块的具体含义如下：

- 资源采集，通过标准化接口，实时或周期性从资源封装模块采集资源信息。采集的资源信息包括通算、智算、超算、量子计算等静态/动态计算信息、网络资源的静态/动态资源信息和存储资源信息；
- 计算资源衡量，构建覆盖通算、智算、超算、量子计算等类型的计算资源衡量模型，包括 FLOPS、TOPS、BOPS 等算力单位及规格参数，实现算力资源的量化描述；
- 网络资源衡量，对网络资源进行标准化衡量，包括带宽、时延、丢包率、吞吐率等关键性能指标，评估网络连接稳定性与传输效率；
- 存储资源衡量，针对存储资源的容量、性能与可靠性进行衡量，包括总容量、可用容量、数据冗余率等容量指标，支持存储资源的动态分配与回收；
- 指标映射，实现不同量纲、多维度资源衡量的值和指标进行转换和统一；
- 指标查询，通过统一数据接口，为调度层、运营层、监测层提供多维度的资源衡量数据查询服务。

## 7 算力算效衡量功能技术要求

### 7.1 资源采集技术要求

规定从资源封装模块采集算力资源信息的接口标准、采集频率及数据格式要求，保障数据完整性与时效性。

- 计算资源采集应明确计算资源的类型，包括通用算力（CPU）、智能算力（GPU/TPU/NPU）及超算算力、量子计算等类型；
- 计算资源采集支持多源异构采集能力，从多种来源和异构平台采集计算资源信息；
- 计算资源的采集指标应包括算力节点的计算负载率、处理器主频、核心数、峰值计算能力；
- 网络资源采集应采集网络的带宽利用率、时延、丢包率、吞吐率、抖动等性能指标；



- 网络资源采集支持多协议适配，在高并发场景下保证采集数据的完整性；
- 存储资源采集包括总容量、可用容量、数据冗余率等容量指标，采集输入输出操作数、吞吐量、访问延迟、故障恢复时间及存储硬件状态等指标；
- 采集接口应支持通用监控协议，兼容主流云服务商、智算、超算服务商的 API 接口，支持资源信息的自动化上报与更新；
- 采集数据需遵循统一的元数据规范，支持数据输出至时序数据库、数据仓库及对象存储；
- 采集频率应支持实时采集与周期性采集两种模式，建议周期性采集的默认时间间隔为 5 分钟，可根据实际业务需求动态调整，最小采集间隔不低于 1 分钟，最大不超过 1 小时；
- 采集方法宜支持资源配置信息读取、运行时状态监控、主动性能测试等方法；
- 宜支持实时和定时批量采集方式；
- 宜支持全面的数据采集安全管理，至少包括权限、审计、加密、脱敏。

## 7.2 计算资源衡量技术要求

针对通算、智算、超算、量子计算等不同类型计算资源，规范计算衡量模型，具体要求如下：

- 支持对通算、智算、超算、量子计算的资源计算衡量，各类型计算资源衡量需支持包括 fp64、fp32、fp16、fp8/int8、fp4/int4 等精度的计算能力；
- 支持用户侧和供给侧的统一计算衡量；
- 计算资源的计算能力衡量单位应包括 FLOPS、OPS、BOPS 等；
- 各类型计算资源衡量可选支持动态指标、静态指标，可选支持静态动态指标结合的混合指标衡量。

## 7.3 网络资源衡量技术要求

网络资源衡量模块规范网络衡量模型，应满足以下要求：

- 支持对入算网络、算间网络和算内网络等各类网络连接进行全面、准确的衡量；
- 支持多种网络互联结构下对网络协议的衡量，包括但不限于 PCIe、NVLink、InfiniBand/RoCE、UB-Mesh 等；
- 支持通用网络衡量指标，包括带宽、时延、抖动、丢包率等；
- 支持分层次分区域的网络资源衡量，涵盖国家级至接入网各层次、跨区域端到端、核心至接入各层级、物理与虚拟网络协同、网络域边界及资源自动发现，全面反映全国一体化算力网络状况；
- 可选支持静态网络资源参数和动态网络性能指标的衡量，能够实时反映网络资源的可用状态、性能水平及负载情况。

## 7.4 存储资源衡量技术要求

存储资源衡量模块规范存储资源衡量模型，应满足以下要求：

- 支持对块存储、文件存储、对象存储等各类存储类型的衡量；
- 支持对存储容量、性能、可靠性等多维度指标的全面衡量，能够实时反映存储资源的使用状态、剩余容量及性能水平；
- 支持通用存储衡量指标，包括：物理总容量、可用容量、已用容量、容量利用率、输入输出操作数、访问延迟等。

## 7.5 指标映射技术要求

规范不同类型资源的指标关联和映射关系模型，具体要求如下：

- 支持能效指标的计算，包括计算资源的算效衡量指标、单位功耗的算力指标、存储能效的指标等；
- 宜支持对算力中心整体算力资源总量的度量，包括计算、存储、网络等资源的功能与性能指标，反映算力中心的总体服务能力与规模水平；
- 可选支持异构计算的指标统一建模和指标映射，包括通用算力、智能算力、超算能力、算效能力等异构计算能力；
- 可选支持静态指标、动态指标和静动态结合的混合指标的统一映射；
- 可选支持多维度资源统一映射，如计算维度（FLOPS/IOPS）、存储维度（容量/带宽）、网络维度（时延/吞吐量）等不同量纲的指标的融合和统一映射；
- 可选支持指标映射的验证与纠错机制。

## 7.6 指标查询技术要求

指标查询模块通过统一数据接口，为调度层、运营层、监测层提供多维度的资源衡量数据查询服务，满足以下总体要求：

- 支持调度层、运营层、监测层对算力算效衡量结果的实时查询与历史查询；
- 提供标准化的查询接口，支持多种查询方式和查询协议；
- 支持细粒度的权限控制，确保衡量数据的安全访问；
- 支持多维度查询，包括：按资源类型(计算、网络、存储)查询衡量结果、按地理位置(区域、数据中心)查询衡量结果、按时间范围(实时、历史、自定义时段)查询衡量结果、支持按资源标识(资源 ID、资源组)查询衡量结果、按指标类型(性能、容量、可用性)查询衡量结果；
- 采用 JSON、XML 等标准数据格式返回查询结果；
- 支持记录所有查询请求的日志，包括查询者、查询时间、查询条件等，用于安全审计和故障排查；
- 权限控制应支持基于角色的访问控制（RBAC），细粒度至指标类型、资源范围、时间维度，权限变更需经审批并记录审计日志；
- 可选支持查询结果可视化输出，包括趋势曲线图、分布柱状图、对比热力图等格式，可视化数据更新频率与指标采集频率一致；
- 可选支持查询结果的压缩传输，减少网络带宽占用；
- 可选支持高并发查询请求的处理，保证查询响应时间；
- 可选支持对关键指标的简易特征分析，包括但不限于当前时间范围内的最大最小值、平均、方差、趋势、和历史周期变化程度等。

## 8 算力算效衡量指标技术要求

### 8.1 应用场景

算力基础设施化实现了算力使用从租用转变为按需使用模式，算力算效衡量的核心也从服务侧转变为用户侧，而这个过程必然需要统一的算力使用指标。存储和网络领域，延迟、带宽、容量等指标已成为业界公认的统一评价指标，而面向算力网尚缺乏统一的算力使用指标。算力使用指标需要抽象用户侧算力使用，而用户侧使用算力过程可以抽象为运行用户应用程序的过程，应用程序又可以抽象为计算操作和数据移动操作的集合。通过获取程序的计算操作和数据移动操作可以获取用户的算力消耗。用户应用程序包括 64 位浮点为主的科学计算、32 位/16 位浮点为主的深度学习模型训练、4 位/8 位/16 位整

型为主的深度学习模型推理、整型计算为主的大数据/数据库处理等不同类型的计算，需要将不同计算类型进行有效的归一化。算力算效衡量以每秒完成基本操作数作为基本指标。具体应用场景如下：

- 为算力容量规划提供参考依据。每秒完成基本操作数（BOPS）可以实现异构算力提供商的峰值算力统一表示，通过服务侧 BOPS 指标的计算，可以将 FP64 算力、FP16/32 算力、INT8/16/32 算力等不同类型的算力进行归一化统一表示。
- 为算力计量计费提供依据，在算力网场景下需要细粒度算力计量，通过基本操作数指标可以反映用户实际的算力使用量，为计量计费提供参考。即按照实际消费的基本操作数进行计费。
- 为算力调度提供参考依据。定义算力占用率为  $E$ ， $E = \text{BOPS（实际）} / \text{BOPS（峰值）}$ ；即当前实际 BOPS 除以其运行算力设备的峰值 BOPS 的比值为算力占用率  $E$ 。通过  $E$  可以评估不同的算力提供商的当前算力使用量，继而为调度的负载均衡提供参考。如 A、B 两个算力提供商，当前的算力占用率分别为 0.2 和 0.5，则任务调度可以优先选择 A 进行任务分配。

## 8.2 每秒完成基本操作数定义

每秒完成基本操作数通过基本操作数除以使用时长来获取。如表 1 所示，包括计算操作（C-operations）、比较操作（P-operations）、地址操作（A-operations）。每次操作都计为 1。例如，执行一次地址访问，执行一次加法计算。同时为了进行归一化，将所有操作归一为 64 位操作。如果是一个 32 位的操作则记为 0.5。在地址计算上，采用的方法是计算寻址操作次数，访问指针 P 的第 i 个位置，需要进行的操作是： $p+i$ ，即在 P 地址位置加上 i 的偏移，因此地址操作可以按照此方法计算。对于数组，同样计算它们寻址需要的操作。如：访问  $P[i]$ ，对应的操作是： $p+i$ ，即在 P 地址位置加上 i 的偏移，因此地址操作可以按照此方法计算，即进行一次加法运算。对于多维数组也可按此方法计算。

表 1 基本操作数基本操作定义

操作种类	具体操作
计算操作	四则计算（加、减、乘、除），位运算，逻辑运算
比较操作	大于、小于、大于等于、小于等于、等于、不等于
地址操作	数组，指针的寻址操作

表 2 展示了基本操作数的归一化，归一化计算主要考虑的对象是整型计算、浮点计算、比较操作、地址操作，并将相关指标归一化为 64 位操作。在操作的计数方面，基本操作数将所有操作都计为 1。

表 2 基本操作数归一化计算表（以 C 语言为例）

基本操作类型	描述	计数
整型计算	加法	1
	减法	1
	乘法	1
	除法	1
	移位计算（左移、右移）	1
	位运算（与、或、异或）	1
	逻辑运算（与、或、非）	1
浮点计算	加法	1
	减法	1
	乘法	1

	除法	1
	位运算（与、或、异或）	1
	逻辑运算（与、或）	1
比较操作	大于、大于等于	1
	小于、小于等于	1
	等于	1
	不等于	1
地址操作	一维数组、指针的寻址操作	1
	多维数组、指针的寻址操作	N（N 为数组的维度）

结合表 2 的定义，基本操作数定义的基本操作可以映射到用户程序（用户侧）和计算机系统（服务侧）的对应操作并适合不同计算类型，从而保证普适性。每秒可完成的基本操作数可以通过基本操作数除以使用时长来获取。

### 8.3 每秒完成基本操作数计算方法

每秒完成操作数是用户侧出发的衡量指标，既能反映用户侧算力使用情况，也能反映服务侧算力的供给情况，在用户侧和服务侧可以被统一表达：

- 对于服务侧，服务算力的性能可以通过分析 CPU/计算设备的微架构实现，可以获取理论峰值 BOPS 峰值的基本计算公式为： $BOPS_{Peak} = (\text{设备核数}) * (\text{每时钟周期完成的 BOPs}) * (\text{主频})$ ；
- 对于用户侧，可以使用 CPU、GPU 等设备的性能分析工具获取。对于 CPU 而言，其性能分析可由其上的性能计数器所支持，性能计数器是特殊的硬件寄存器，在大多数现代 CPU/加速设备上都可以使用且不会降低应用程序的速度，实现非入侵分析。对于 GPU 而言，需要使用 Nsight Compute、nvprof、CUPTI（CUDA 性能分析接口）等性能分析工具对特定性能事件进行监测，这些性能事件可能由硬件或软件所支持。一般而言，此类性能分析工具支持对整数运算事件、浮点数运算事件发生次数予以监测，每次发生的事件对应一次运算过程；
- 对于 CPU 平台的计算，通过硬件计数器获取总指令数（ $ins_{total}$ ），分支指令（ $ins_{branch}$ ），加载指令（ $ins_{load}$ ），存储指令（ $ins_{store}$ ），使用计算公式： $BOPs \cong ins_{total} - ins_{load} - ins_{store} - ins_{branch}$  即总指令减去分支指令、加载指令和存储指令可得到 BOPs 近似数量。将获取的 BOPs 值除以对应的算力使用时长即可得到用户侧 BOPS；
- 对于 GPU 平台，为了满足归一化的需要，需要对性能事件按计算位宽分类讨论，设整数和浮点数性能事件的集合为 E，按位数划分，每种浮点数或整数运算事件会对应 8、16、32、64 等位数不等的运算，如 64 位浮点计算有其专属的性能事件进行计数。约定 k 为归一化系数，其含义为运算数的“标准位宽”，取 64。对于每种 E 中的事件 event，其事件计数记作  $Count_{event}$ ，事件的位宽为  $Width_{event}$ ，对应公式为  $BOPs \cong \frac{1}{k} \times \sum_{event \in E} Count_{event} \times Width_{event}$ ，即所有浮点和整数运算按运算位数的加权和除以归一化系数 k。例如，若事件属于 64 位数浮点运算，则其对于 BOPs 的贡献为其次数乘以位宽，即 64。其他位宽的运算事件依次类推并累加。例如，若事件属于 64 位数浮点运算，则其对于 BOPs 的贡献为其次数乘以位宽，即 64。其他位宽的

运算事件依次类推并累加。

#### 8.4 每秒完成基本操作数和现有指标的映射

每秒完成基本操作数和现有主流衡量指标的映射关系如下：

- FLOPS 与 BOPS 的映射：以 FLOPS 为主的计算设备的  $BOPS = (FLOPs + A\_operations) / Time$ ；其中 FLOPs 为浮点相关的计算量，A\_operations 为地址操作计算量，Time 为运行时间。如忽略数据移动量 BOPS 可近似为 FLOPS 的 64 位归一；
- OPS 与 BOPS 的映射：以 OPS 为主的计算设备的  $BOPS = (OPs + A\_operations) / Time$ ；其中 OPs 为特定计算的计算量，A\_operations 为地址操作计算量，Time 为运行时间。如忽略数据移动量 BOPS 可近似为 OPS 的 64 位归一。

基于上述描述，BOPS 指标可以适用于通算（FP64 计算为主）、超算（FP64 计算为主）、智算（FP16/FP32、OPS）等不同类平台。

#### 8.5 算效衡量方法

算效衡量的方法，依据计算公式采用和上述相同的 BOPs 获取方法，将获取的 BOPs 值除以对应的算力使用期间的能量消耗即可得到算效 CE。

## 参考文献

- [1] 《全国一体化算力网 监测调度平台建设指南》标准草案
-