

Prometheus 监控服务

产品简介



腾讯云

【 版权声明 】

©2013–2025 腾讯云版权所有

本文档（含所有文字、数据、图片等内容）完整的著作权归腾讯云计算（北京）有限责任公司单独所有，未经腾讯云事先明确书面许可，任何主体不得以任何形式复制、修改、使用、抄袭、传播本文档全部或部分内容。前述行为构成对腾讯云著作权的侵犯，腾讯云将依法采取措施追究法律责任。

【 商标声明 】



及其它腾讯云服务相关的商标均为腾讯云计算（北京）有限责任公司及其关联公司所有。本文档涉及的第三方主体的商标，依法由权利人所有。未经腾讯云及有关权利人书面许可，任何主体不得以任何方式对前述商标进行使用、复制、修改、传播、抄录等行为，否则将构成对腾讯云及有关权利人商标权的侵犯，腾讯云将依法采取措施追究法律责任。

【 服务声明 】

本文档意在向您介绍腾讯云全部或部分产品、服务的当时的相关概况，部分产品、服务的内容可能不时有所调整。您所购买的腾讯云产品、服务的种类、服务标准等应由您与腾讯云之间的商业合同约定，除非双方另有约定，否则，腾讯云对本文档内容不做任何明示或默示的承诺或保证。

【 联系我们 】

我们致力于为您提供个性化的售前购买咨询服务，及相应的技术售后服务，任何问题请联系 4009100100或 95716。

文档目录

产品简介

Prometheus 监控概述

Prometheus 监控优势

Prometheus 应用场景

基本概念

相关限制

功能特性

开服地域

产品简介

Prometheus 监控概述

最近更新时间：2025-01-21 15:33:22

Prometheus 监控服务（TencentCloud Managed Service for Prometheus，TMP）是基于开源 Prometheus 构建的高可用、全托管的服务，与腾讯云容器服务（TKE）高度集成，兼容开源生态丰富多样的应用组件，结合腾讯云可观测平台的告警功能和 Prometheus Alertmanager 能力，为您提供免搭建的高效运维能力，减少开发及运维成本。

开源 Prometheus 简介

Prometheus 是一个开源监控系统。与 Kubernetes 相似，Prometheus 受启发于 Google 的 Borgmon 监控系统，而 Kubernetes 也是从 Google 的 Borg 演变而来的。Prometheus 始于2012年，并由 SoundCloud 内部工程师开发，于2015年1月发布。2016年5月，其成为继 Kubernetes 之后第二个正式加入 [Cloud Native Computing Foundation \(CNCF\)](#) 基金会的项目。现在最常见的 Kubernetes 容器管理系统中，通常会搭配 Prometheus 进行监控。

Prometheus 具有如下特性：

- 自定义多维数据模型（时间序列数据由 Metric 和一组 Key/Value Label 组成）。
- 灵活而强大的查询语言 PromQL，可利用多维数据完成复杂的监控查询。
- 不依赖分布式存储，支持单主节点工作。
- 通过基于 HTTP 的 Pull 方式采集时序数据。
- 可通过 PushGateway 的方式来实现数据 Push 模式。
- 可通过动态的服务发现或者静态配置去获取要采集的目标服务器。
- 结合 Grafana 可方便地支持多种可视化图表及仪表盘。

产品功能

根据监控分层，Prometheus 监控服务覆盖了业务监控、应用层监控、中间件监控、系统层监控，结合腾讯云可观测平台告警和开源 Grafana 可以提供一站式全方位的监控体系，帮助业务快速发现和定位问题，减轻故障给业务带来的影响。

- 系统层监控：例如 CPU、Memory、Disk 和 Network 等。
- 中间组件层监控：例如 Kafka、MySQL 和 Redis 等。
- 应用层监控：应用服务监控，例如 JVM、HTTP 和 RPC 等。
- 业务监控：业务黄金指标，例如登录数和订单量等。

Prometheus 监控优势

最近更新时间：2024-12-05 16:07:33

基于开源自建 Prometheus，会遇到哪些问题？

使用开源 Prometheus 需要自行购买相关资源并部署系统（简称自建）。由于开源 Prometheus 自身的短板，自建 Prometheus 也给企业带来了不少困扰。

1. 对于中小企业，使用成本高

自建 Prometheus 的使用成本，包括机器资源成本和人力成本，最主要的是人力成本。其中人力成本又包括：

- 前期调研成本
- 中期搭建成本
- 后期维护成本

由于中小企业的运维团队规模较小，一般不多于5人，有的甚至只有一两个人，要自建和维护一套 Prometheus 监控服务，显然非常吃力。

2. 对于大企业，可扩展性差，容易出现性能瓶颈

大企业或快速发展的中型企业，在业务发展初期自建 Prometheus 监控，但随着业务量高速增长，意味着更多资源的投入，对监控也有了更高的要求。自建 Prometheus 就会开始暴露出可扩展性差、性能瓶颈的问题，使企业运维面临巨大的挑战。

腾讯云 Prometheus 与自建 Prometheus 功能对比

对比类型	具体功能	腾讯云 Prometheus 监控	自建 Prometheus 监控
数据集成	集成腾讯云容器服务	一键自动集成	手动接入，配置复杂
	跨 VPC/跨地域集成容器	自动支持	需自行做网络的打通
	集成基础云产品数据	一键安装	需自行安装 Exporter
	常用监控组件集成	一键安装	需自行安装 Exporter
	标签自动识别资源变化	支持	不支持
可视化	关联 Grafana 可视化	支持快速关联托管 Grafana 服务	需自行搭建 Grafana

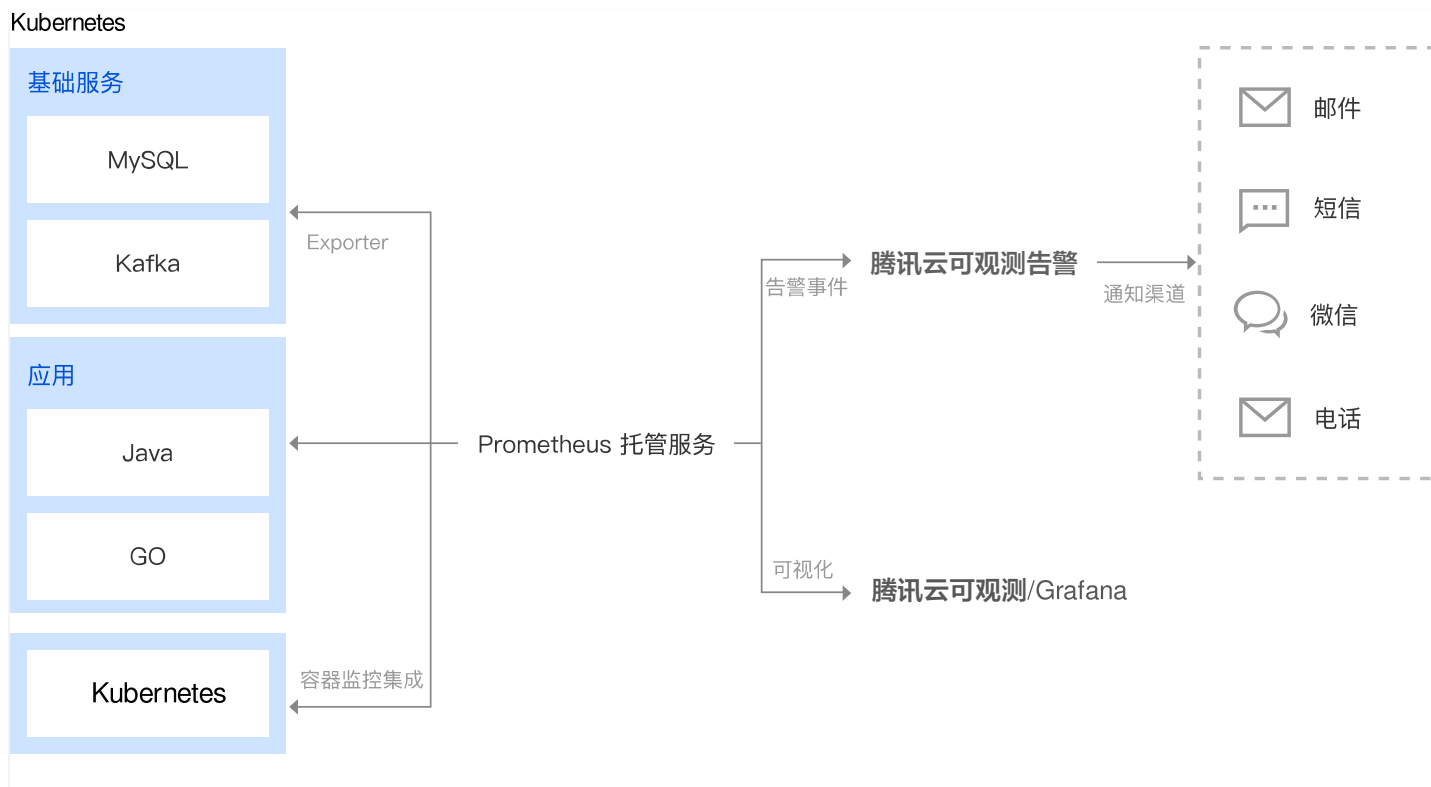
	预设 Dashboard 模板	支持	不支持
告警	告警通知渠道	可复用腾讯云可观测平台-告警管理的渠道	需自行搭建
	告警通知模板	支持	不支持
其他能力	健康巡检	支持	不支持
	预聚合	支持	不支持
高可用性	多副本	支持	不支持
	水平拓展	结合腾讯云自研的分片和调度技术，实现动态扩缩，满足用户的弹性需求，同时支持负载均衡	无法水平扩展
	数据存储	数据存储能力无上限	数据存储受限于本地磁盘大小
安全管理	数据安全	基于腾讯云安全体系，支持鉴权管理	不支持
成本	人力成本	一次性配置，免运维	<ul style="list-style-type: none"> ● 前期调研 ● 中期搭建 ● 后期维护
	资源成本	<ul style="list-style-type: none"> ● 按需使用 ● 按量计费 ● 容器核心基础指标免费 	<ul style="list-style-type: none"> ● 固定费用 ● 存在资源浪费的可能

Prometheus 应用场景

最近更新时间：2024-12-05 16:07:33

一体化监控场景

Prometheus 监控服务提供一站式开箱即用的 Prometheus 全托管服务，天然集成开源 Grafana 大盘和腾讯云可观测平台告警。支持基础服务、应用层、容器服务等监控场景。



应用服务监控场景

场景一

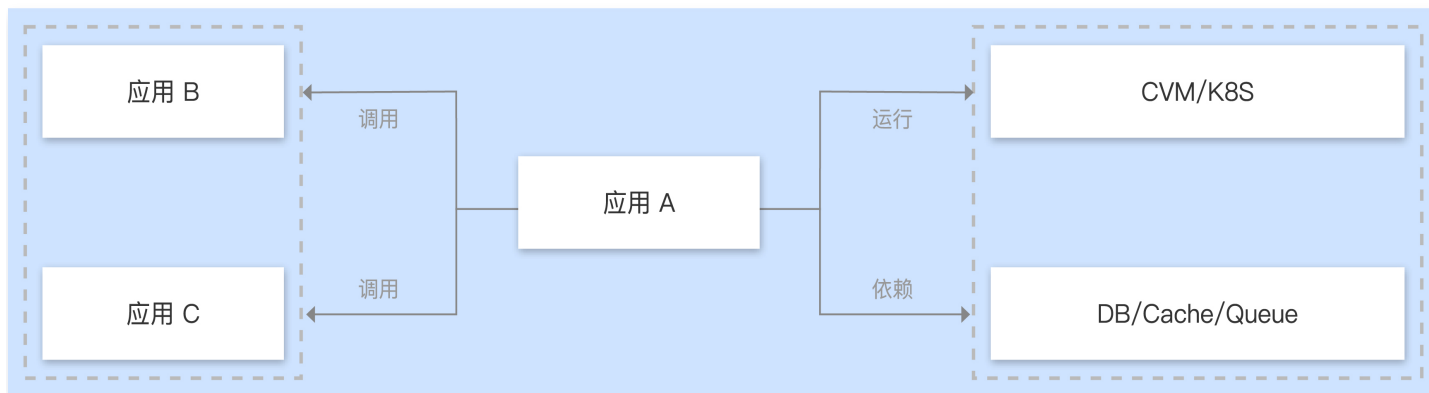
某应用提供了对外的接口服务，但无法了解该接口服务质量。Prometheus 监控服务可对开发语言进行集成，实时对接口的访问量/延时/成功率进行监控。

场景二

Prometheus 监控服务同时也会对服务进行异常检测，可了解该异常影响了哪些接口、发生在哪些主机，或者了解该异常是单机问题还是整个集群的共性问题。

场景三

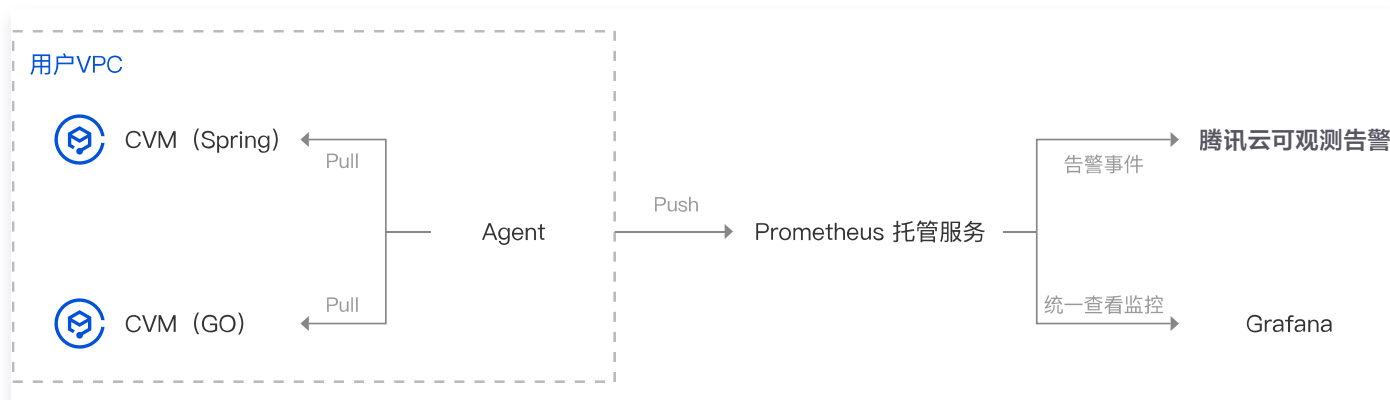
对于 Java 应用来说，可进行单机的 GC/内存/线程状态等监控，全方面地了解 JVM 内部的状态。



云服务器监控场景

当您的服务部署在 CVM 上时，几乎每次服务的扩缩容都要修改 Prometheus 的抓取配置。针对这类场景，结合腾讯云平台提供的标签能力和 Prometheus Agent 对腾讯云标签的发现能力，用户只需合理地对 CVM 关联标签即可方便地管理监控目标对象，免去了需要不断手动更新配置的维护成本，例如：

1. 服务 A 同时部署在 2 台 CVM 上，并对其所在的 CVM 关联标签（服务名：A）。
2. 由于需要进行业务活动，原有 CVM 数量不满足业务活动需求，需再扩容 3 台 CVM，这时只需要对这 3 台 CVM 关联标签（服务名：A）。成功关联后，Agent 就会自动发现新增的3台 CVM，主动抓取监控指标。
3. 活动过后缩容下线 3 台 CVM，服务发现功能会自动感知服务下线，停止抓取监控指标。



自定义监控场景

您可以通过 Prometheus 监控服务自定义上报指标监控数据，对应用或者服务内部的一些状态进行监控，如请求处理数、下单数等，也可以对一些核心逻辑的处理耗时进行监控，如请求外部服务的耗时情况等。

基本概念

最近更新时间：2024-12-05 16:07:33

本文汇总使用 Prometheus 监控服务过程中涉及的基本概念，方便您查询和了解相关概念。

概念	说明
Exporter	Exporter 是一个采集监控数据并通过 Prometheus 监控规范对外提供数据的组件。目前有上百个官方或者第三方 Exporter 可供使用，请参见 Exporter详情 。
Job	一组 Target 的配置集合。定义了抓取间隔，访问限制等作用于一组 Target 的抓取行为。
Prometheus 实例	Prometheus 监控服务提供的管理 Prometheus 监控数据采集和数据存储分析的逻辑单元。
Prometheus 探针	部署在用户侧或者云产品侧 Kubernetes 集群。负责自动发现采集目标、采集指标和远程写到其他库。
PromQL	Prometheus 监控服务的查询语言。支持瞬时查询和时间跨度查询，内置多种函数和操作符。可以对原始数据进行聚合、切片、预测和联合。
Target	Prometheus Agent 要抓取的采集目标。采集目标暴露自身运行、业务指标，或者代理暴露监控对象的运行、业务指标。
告警规则	Prometheus 监控 Alerting Rule 格式的告警配置。可以通过 PromQL 描述。
标签	描述指标的一组 Key-Value 值。
服务发现	Prometheus 监控服务的功能特点之一，无需静态配置，可以自动发现采集目标。支持 Kubernetes SD、Consul、Eureka 等多种服务发现方式，支持通过 Service Monitor、Pod Monitor 的方式暴露采集目标。
预聚合	Prometheus 监控服务 Recording Rule 能力。可以通过 PromQL 将原始数据加工成新的指标，提升查询效率。
集成中心	集成了 Prometheus 监控服务支持的所有服务，您可以根据页面指引安装对应的服务，成功安装后即可在监控面板查看监控数据。
告警策略	用于定义告警如何触发，如何发送。
云产品监控	Prometheus 监控服务集成了腾讯云云产品的监控数据。可一键安装 Agent 即可查看监控数据。
指标	采集目标暴露的、可以完整反映监控对象运行或者业务状态的一系列标签化数据。

TPS	每秒数据点的上报总数。它是衡量系统处理能力的重要指标。
Series 上限	指标个数上限，Series 上限 = (单个指标 × 该指标的维度组合) × 指标个数。

相关限制

最近更新时间：2024-12-05 16:07:33

实例限制

每个付费实例的 Series 限制最大为450万，免费试用实例为200万。如需调整付费实例的限制可以 [联系我们](#)，在资源充足的情况下，我们将会为您合理地调整相关限制。

说明

时间序列（时间线，Series）含义：由指标名和标签组成。相同的指标名和标签在时间序列中构成唯一的一条时间线。

数据上报限制

若您使用 Prometheus 监控服务上报监控数据，将会有下列指标（`__name__` 唯一的 Series）限制。

- 上报时必须要有指标名，即 `__name__` 标签，指标名必须符合规范，只支持英文字母开头，可包含 ASCII 字符、数字、下划线以及冒号，并必须符合正则表达式 `[a-zA-Z_:][a-zA-Z0-9_]*`，参见 [Prometheus 官方文档](#)。
- 每个指标最多32个标签。
- 标签名必须符合规范，可包含 ASCII 字符、数字、下划线，并必须符合正则表达式 `[a-zA-Z_][a-zA-Z0-9_]*`。标签名称开头为 `__` 仅供内部使用。
- 标签长度：标签名为1024字符，标签值为 2048 字符。
- 同一个指标下，和标签的维度组合不能超过10万（在 histogram 有较多 bucket 的情况下，histogram 类型的指标不支持调整）。
- 每秒上报的数据点总量限制：付费实例不能超过300000，免费试用实例不能超过100000。

说明

标签的作用：Prometheus 中存储的数据为时间序列，是由指标名和一系列的标签（键值对）唯一标识的，不同的标签代表不同的时间序列，即通过指定标签查询指定数据。添加的标签越多，查询的维度越细。

Prometheus 查询限制

为了保证查询效率及更好的用户体验，Prometheus 查询有如下限制：

- 单个查询涉及的 time series 不能超过100000。
- 单个查询涉及的数据量不能超过100MB。
- 对于查询频次没有限制，但是如果并发量超过 15 可能会有一定的排队延时（较多较慢的大查询可能会出现，一般情况不会出现。时间跨度超过两周的大查询出现延时的概率会相应升高）。

以上限制对告警规则和预聚合规则同样有效，建议根据业务场景来限制查询范围或者其他方式对大查询进行适当拆分，也可以采用先拆分后聚合的方式，例如对预聚合后的结果再次进行聚合。

告警与预聚合限制

告警相关的限制：

- 单实例告警规则上限：150。
- 单实例总告警数量上限：2000（超出则直接丢弃）。
- 单实例所有告警的标签、Annotation 等字段总的大小上限：20MiB（超出则直接丢弃）。

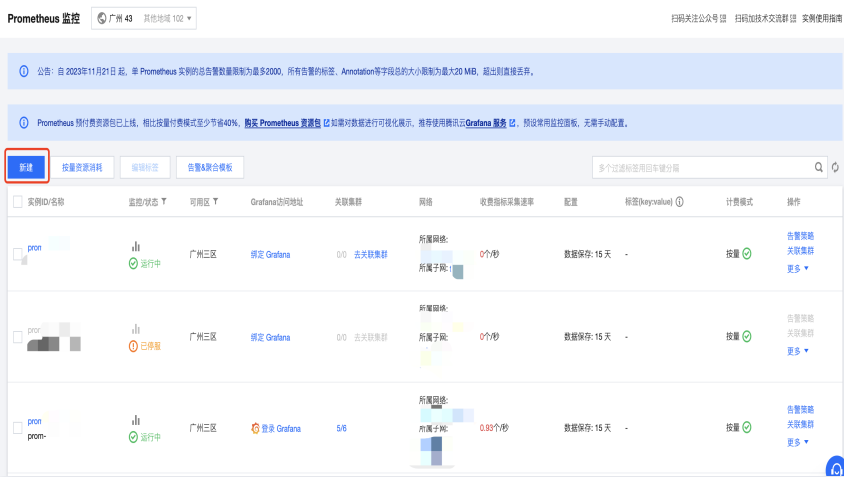
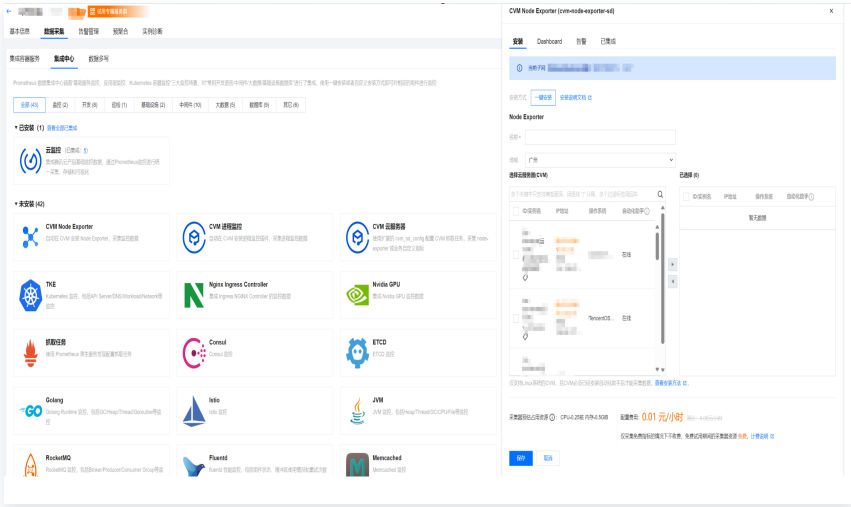
预聚合相关的限制：


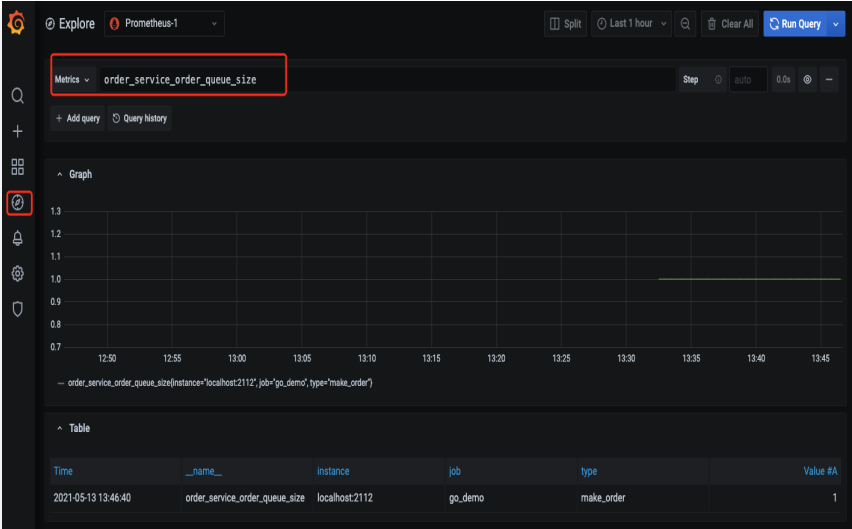
- 单实例预聚合规则上限：150。
- 单预聚合规则分组中的预聚合规则上限：35。

功能特性

最近更新时间：2024-12-05 16:07:33


监控对象接入

功能	功能说明	控制台示例图
<p>创建 Prometheus 实例</p>	<p>支持创建多地域 Prometheus 实例。</p>	
<p>集成中心</p>	<p>支持多种组件一键安装和自定义接入，成功安装后即可在 Grafana 中查看监控数据。</p>	

<p>健康巡检</p>	<p>通过定期探测对应服务的连通性，来检测其服务的健康情况，帮助您实时地掌握服务的健康状况，及时发现异常，来提升服务的SLA。</p>	
<p>自定义监控</p>	<p>支持自定义上报指标监控数据，对业务内部核心指标进行监控，如请求数，下单数，请求外部服务的耗时情况等。</p>	

监控指标采集

功能	功能说明	控制台示例图
<p>数据采集</p>	<p>在数据采集页面，可以进行数据采集配置，并直观查看采集目标。</p>	

<p>targ ets</p>	<p>支持通过 targets 可以直观查看正在被抓取的目标，以及抓取状态是否正常。</p>	
---------------------	--	--

监控数据处理

功能	功能说明	控制台示例图
<p>获取 Remote Write 地址</p>	<p>Remote Write 功能支持作为远程数据库存储 Prometheus 监控服务的数据。您可以使用 Remote Write 地址，将自建 Prometheus 的监控数据存储到 Prometheus 监控服务的实例中，实现远程存储，并可视化展示在同一 Grafana。</p>	

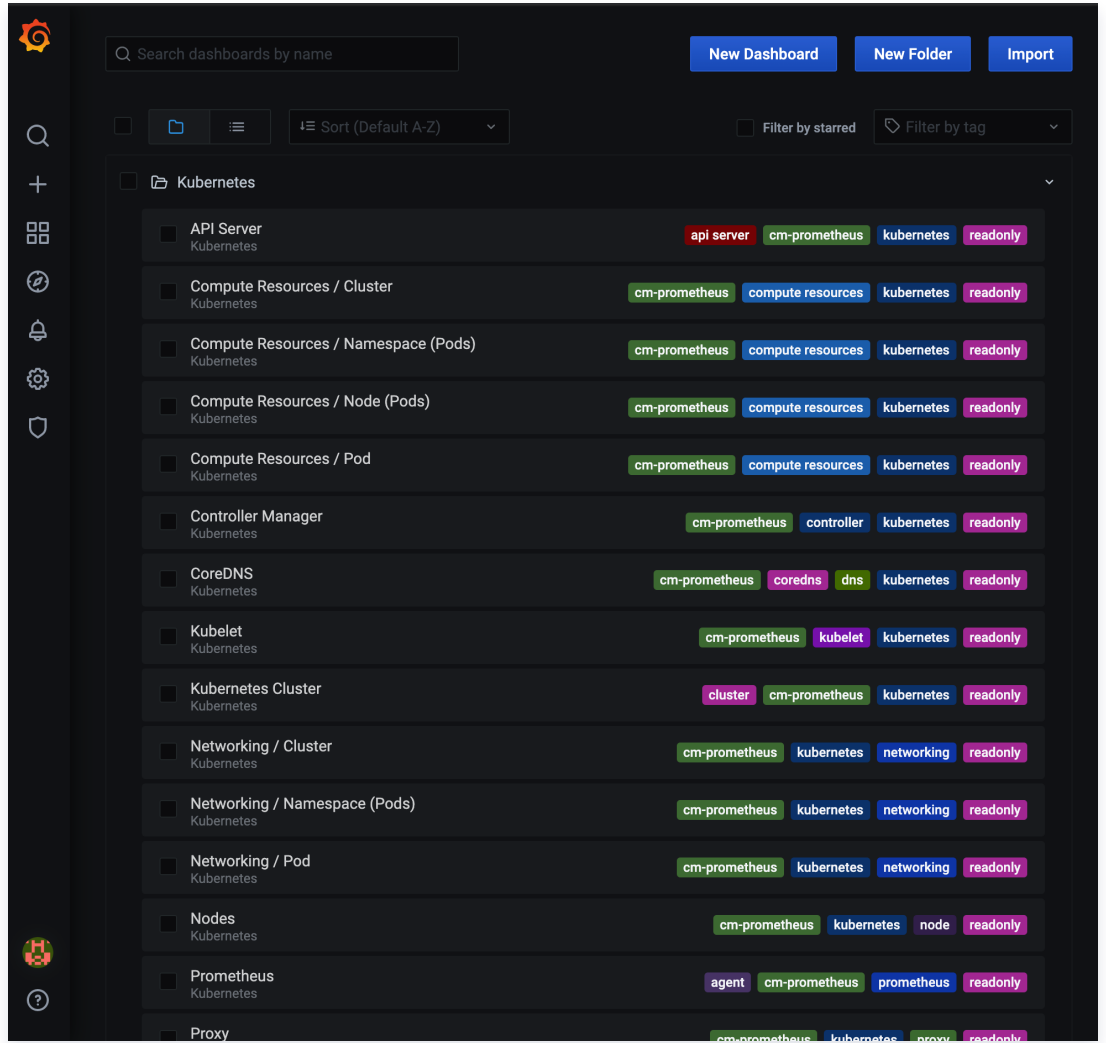
<p>预聚合</p>	<p>对一些常用的指标或者计算相对复杂的指标进行提前计算，然后将这些数据存储在新的数据指标中，提前计算好的指标查询速度更快，可以解决用户配置以及查询慢的问题。</p>	
------------	---	--

监控数据展示

功能	功能说明	控制台示例图
----	------	--------

Grafana

预置丰富的 Grafana 大盘，同时也支持自定义大盘。还预设了 Grafana 官网常用的插件，您可以在控制台一键安装。



实例监控

支持 Prometheus 实例状态和使用情况监控，包括 Prometheus 实例存储情况、告警发送情况、Grafana 请求和仪表盘数



	<p>量等，方便您实时了解 Prometheus 实例使用情况。</p>	
<p>HTTP API</p>	<p>获取 Prometheus 监控服务数据地址。您可以通过该地址将 Prometheus 实例的监控数据接入自建的 Grafana 大盘展示数据，也可以将 Prometheus 监控数据进行二次开发。</p>	<div data-bbox="384 804 1482 1456" style="border: 1px solid #ccc; padding: 10px;"> <p>服务地址</p> <p>Token ***** </p> <p>Remote Write 地址 </p> <p>Remote Read 地址 </p> <p>HTTP API </p> <p>Pushgateway 地址 </p> </div>

告警

功能	功能说明	控制台示例图
----	------	--------

创建告警

预置多种报警规则，也支持针对特定监控对象自定义告警规则。TMP 集成了腾讯云可观测平台告警通知模板，在某些指标发生异常时及时通知您采取措施。

创建告警策略
✕

创建方式

创建方式

选择模板
页面编辑
YAML编辑

基本信息

实例ID/名称 [实例ID]

策略名称 请输入策略名称

告警规则

规则

状态

✕

规则名称 请输入规则名称

PromQL rate(metrics0{ [2m]} > 1

[点击预览规则](#)

告警对象(Summary)

告警消息(Description)

Labels = ✕

[添加](#)

Annotations = ✕

[添加](#)

持续时间 - 0 + 分钟 ▼

触发规则的最小持续时间，若设置为1分钟，则告警在满足规则1分钟后被触发，1分钟内被视为正常波动不作告警。

[添加](#)

收敛时间 5分钟 ▼

告警收敛时间对alertmanager渠道不生效。在收敛时间周期内若多次满足告警条件，仅会发送一次通知，若设置为1小时，则1小时内该策略被触发后仅会发送1次告警通知。

告警渠道 腾讯云 Webhook 自建alertmanager

告警通知 选择模板 新建

管理告警

支持对告警策略执行开启、关闭、编辑、删除等操作，同时还支持其他实例告警一键导入。

编辑告警策略

创建方式

创建方式

选择模板
页面编辑
YAML编辑

基本信息

实例ID/名称

策略名称 *

告警规则

规则

状态 *

规则名称 *

PromQL *

test

点击预览规则 [↗](#)

告警对象(Summary) *

告警消息(Description) *

```

{{ printf
"node_filesystem_size(instance=%s",device=~"/dev/nvme.+|/dev/xvd.+]"
$labels.instance | query | first | value }}
                    
```

Labels [添加](#)

Annotations [添加](#)

持续时间

—

+

▼

触发规则的最小持续时间，若设置为1分钟，则告警在满足规则1分钟后被触发，1分钟内被视为正常波动不作告警。

添加

收敛时间 *

告警收敛时间对alertmanager渠道不生效。在收敛时间周期内若多次满足告警条件，仅会发送一次通知，若设置为1小时，则1小时内该策略被触发后仅会发送1次告警通知。

告警渠道 * 腾讯云 Webhook 自建alertmanager

告警通知 * 选择模板 [新建](#) [↗](#)

开服地域

最近更新时间：2025-03-18 15:06:02

说明：

开服地域指的是物理的数据中心或服务器位置，可简单理解为一个服务或资源在哪个地区或国家开放，即在开放的地域为用户提供服务和资源。当资源在特定地域的数据中心创建成功后，通常就无法更换地域。

Prometheus 支持的地域和可用区如下：

地域	可用区
华南地区（广州） ap-guangzhou	广州三区 ap-guangzhou-3
	广州四区 ap-guangzhou-4
	广州六区 ap-guangzhou-6
	广州七区 ap-guangzhou-7
华南地区（深圳金融） ap-shenzhen-fsi	深圳金融一区 ap-shenzhen-fsi-1
	深圳金融二区 ap-shenzhen-fsi-2
	深圳金融三区 ap-shenzhen-fsi-3
华东地区（上海） ap-shanghai	上海二区 ap-shanghai-2
	上海三区 ap-shanghai-3
	上海四区 ap-shanghai-4
	上海五区 ap-shanghai-5
	上海八区 ap-shanghai-8
华东地区（上海金融） ap-shanghai-fsi	上海金融一区 ap-shanghai-fsi-1
	上海金融三区 ap-shanghai-fsi-3
华东地区（南京） ap-nanjing	南京一区 ap-nanjing-1
	南京二区 ap-nanjing-2
	南京三区 ap-nanjing-3
华东地区（上海自动驾驶云） ap-shanghai-adc	上海自动驾驶云-区 ap-shanghai-adc-1

华北地区（北京） ap-beijing	北京三区 ap-beijing-3
	北京四区 ap-beijing-4
	北京五区 ap-beijing-5
	北京六区 ap-beijing-6
	北京七区 ap-beijing-7
华北地区（北京金融） ap-beijing-fsi	北京金融一区 ap-beijing-fsi-1
	北京金融二区 ap-beijing-fsi-2
西南地区（成都） ap-chengdu	成都一区 ap-chengdu-1
	成都二区 ap-chengdu-2
西南地区（重庆） ap-chongqing	重庆一区 ap-chongqing-1
港澳台地区（中国香港） ap-hongkong	中国香港一区 ap-hongkong-1
	中国香港二区 ap-hongkong-2
	中国香港三区 ap-hongkong-3
港澳台地区（中国台北） ap-taipei	中国台北一区 ap-taipei-1
东南亚地区（新加坡） ap-singapore	新加坡一区 ap-singapore-1
	新加坡二区 ap-singapore-2
	新加坡三区 ap-singapore-3
	新加坡四区 ap-singapore-4
美国西部地区（硅谷） na-siliconvalley	硅谷一区 na-siliconvalley-1
	硅谷二区 na-siliconvalley-2
欧洲地区（法兰克福） eu-frankfurt	法兰克福一区 eu-frankfurt-1
亚太地区（首尔） ap-seoul	首尔一区 ap-seoul-1
	首尔二区 ap-seoul-2
亚太地区（曼谷） ap-bangkok	曼谷一区 ap-bangkok-1
	曼谷二区 ap-bangkok-2

亚太地区（东京）ap-tokyo	东京一区 ap-tokyo-1
	东京二区 ap-tokyo-2
美东地区（弗吉尼亚）na-ashburn	弗吉尼亚一区 na-ashburn-1
	弗吉尼亚二区 na-ashburn-2
南美地区（圣保罗）sa-saopaulo	圣保罗一区 sa-saopaulo-1
亚太东南地区（雅加达）ap-jakarta	雅加达一区 ap-jakarta-1