

# Build A High Performance Remote Storage for Prometheus China 201 with Unlimited Time Series

Yang Xiang (向阳) Yunshan Networks (云杉网络)



**OPEN SOURCE SUMMIT** 

### **TSDB:** Time Series DataBase



#### **Prometheus Samples**

cpu\_total{instance=1.2.3.4:80, job=k8s-pod, cluster=prod, cpu=1, state=system} 83 1234567890

Metric Name + Lables

Metric Value & Timestamp



# **High Cardinality Problem**



Cardinality =~ N(metrics\_name) \* N(label\_1\_values) \* N(label\_2\_values) \* .....



N(Series) = 20M, Label Values:

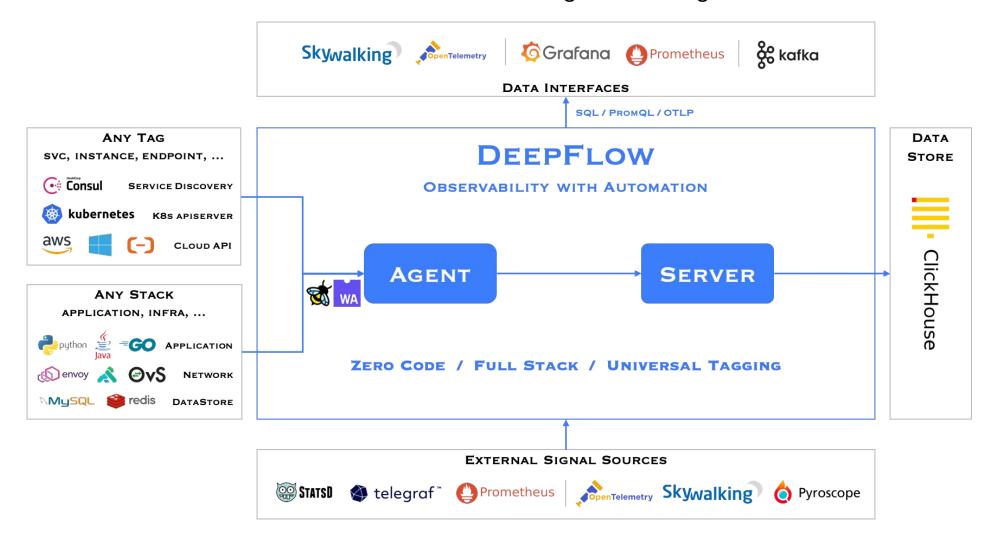
name	c I	max(value)
id	172878	ygc
container id	169890	docker://ffffdbc070dc1e00000db448390c605c19e7b9e625741f2490624b4c8c174185
name	144573	x01Vehicle
mountpoint	63361 I	/var/lib/kubelet/pods/fe02be86-35a9-4a6e-bd4d-e900000fbd33/volumes/kubernetes.io~projected/kube-api-access-x52sh
uid	46462	fffd733b-bbd4-421b-9afa-8885b282c846
pod	38254	yearlyreport-58d9849b-5rbcl
pod name	28501	yearlyreport-58d9849b-5rbcl
address	24578	fe:ff:f3:00:00:00
device	20703	xxxx nodev:/inf-hdmap
pod ip	20282	172.00.000.99
replicaset	20200	yearlyreport-c6b57fd94
uri	19837	root
path	19140 i	stats
route path	16372	获取认证字符串,使用正式私钥签名
routeId	14606	servicel
owner name	13380 i	yearlyreport-58d9849b
routeUri	13151	webjars/springfox-swagger-ui/*
image		xxxxxxxxcar.xxxxxxcloudcr.com/xxxxxxx/wp-weixin-provider@sha256:f30536baf1bcb0d47e000009088baef914d5a8b251b77192e66811d5ad5c660e
instance		prometheus.kube-system.svc:9093
created by name		yearlyreport-58d9849b
cluster ip	7690 i	
podip	7330 i	10.28.99.56
podname	7256	xxx-recommend-api-bffb7b7d5-ltrgx
label rollouts pod template hash		ffff9bbcc
resource		xxxxolumes.zfs.openebs.io
image id		docker-pullable://xxxxxxxxcar.xxxxxxcloudcr.com/xxxxxxx/wp-weixin-provider@sha256:f30536ba00000047eb97949088baef914d5a8b251b77192e66811d5ad5c6
image spec		xxxxxxxxcar.xxxxxxxcloudcr.com/xxxxxxx/wp-weixin-provider:1.0.0-230725.1059-32d4dfed
filename		vsp-recommend-api.log
host		yearlyreport.dev.k8s.chj.cloud
endpoint		yearlyreport

## Metrics Storage in DeepFlow



Metrics Lables: client\_ip, server\_ip, client\_port, server\_port, protocol, endpoint, url, process\_name, pod, ...

We use ClickHouse as a long-term storage.



#### An Intuitive Idea



```
CREATE TABLE prometheus samples
    `time` DateTime('Asia/Shanghai') CODEC(DoubleDelta),
    `metric name` String,
    `tag names` Array(String),
    `tag values` Array(String),
    `value` Float64
SELECT tag_values[indexOf(tag names, 'host')] AS `host`
FROM prometheus samples
WHERE (tag values[indexOf(tag names, 'host')]) = 'some-host'
```

# Too Many Lables, Very Slow Query



Query id: 78d4852c-b61c-4acc-a732-ba7dd8904e3c

Row 1: metric\_name: prometheus.container\_memory\_failures\_total any(tag\_names): ['container', 'failure\_type', 'id', 'image', 'name', 'namespace', 'pod', 'scope', 'beta\_kubernetes\_io\_arch', 'beta\_kubernetes\_io\_instance\_gpu', 'beta\_kub ernetes\_io\_instance\_type','beta\_kubernetes\_io\_os','cce\_baidubce\_com\_kubelet\_dir','cluster','cluster\_id','cluster\_role','failure\_domain\_beta\_kubernetes\_io\_region ','failure\_domain\_beta\_kubernetes\_io\_zone','feature\_node\_kubernetes\_io\_cpu\_cpuid\_ADX','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AESNI','feature\_node\_kubernetes\_io\_c pu\_cpuid\_AVX','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX2','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512BITALG','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512BW','f eature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512CD','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512DQ','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512F','feature\_node\_kuber netes\_io\_cpu\_cpuid\_AVX512IFMA','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512VBMI','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512VBMI2','feature\_node\_kubernetes\_io\_ cpu\_cpuid\_AVX512VL','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512VNNI','feature\_node\_kubernetes\_io\_cpu\_cpuid\_AVX512VPOPCNTDQ','feature\_node\_kubernetes\_io\_cpu\_cpu id\_CMPXCHG8','feature\_node\_kubernetes\_io\_cpu\_cpuid\_FMA3','feature\_node\_kubernetes\_io\_cpuid\_FXSR','feature\_node\_kubernetes\_io\_cpu\_cpuid\_FXSROPT','feature\_node e\_kubernetes\_io\_cpu\_cpuid\_GFNI','feature\_node\_kubernetes\_io\_cpu\_cpuid\_HYPERVISOR','feature\_node\_kubernetes\_io\_cpu\_cpuid\_IBPB','feature\_node\_kubernetes\_io\_cpu\_cp uid\_LAHF','feature\_node\_kubernetes\_io\_cpu\_cpuid\_MOVBE','feature\_node\_kubernetes\_io\_cpu\_cpuid\_OSXSAVE','feature\_node\_kubernetes\_io\_cpu\_cpuid\_SCE','feature\_node\_k ubernetes\_io\_cpu\_cpuid\_SHA','feature\_node\_kubernetes\_io\_cpu\_cpuid\_VAES','feature\_node\_kubernetes\_io\_cpu\_cpuid\_VMX','feature\_node\_kubernetes\_io\_cpu\_cpuid\_VPCLMUL QDQ','feature\_node\_kubernetes\_io\_cpu\_cpuid\_WBNOINVD','feature\_node\_kubernetes\_io\_cpu\_cpuid\_X87','feature\_node\_kubernetes\_io\_cpu\_cpuid\_XSAVE','feature\_node\_kuber netes\_io\_cpu\_hardware\_multithreading','feature\_node\_kubernetes\_io\_cpu\_model\_family','feature\_node\_kubernetes\_io\_cpu\_model\_id','feature\_node\_kubernetes\_io\_cpu\_mo del\_vendor\_id','feature\_node\_kubernetes\_io\_kernel\_config\_NO\_HZ','feature\_node\_kubernetes\_io\_kernel\_config\_NO\_HZ\_FULL','feature\_node\_kubernetes\_io\_kernel\_version \_full','feature\_node\_kubernetes\_io\_kernel\_version\_major','feature\_node\_kubernetes\_io\_kernel\_version\_minor','feature\_node\_kubernetes\_io\_kernel\_version\_revision', 'feature\_node\_kubernetes\_io\_pci\_0300\_1013\_present','feature\_node\_kubernetes\_io\_system\_os\_release\_ID','feature\_node\_kubernetes\_io\_system\_os\_release\_VERSION\_ID', feature\_node\_kubernetes\_io\_system\_os\_release\_VERSION\_ID\_major','instance','instance\_group\_id','job','kubernetes\_io\_arch','kubernetes\_io\_hostname','kubernetes\_io \_os','kubernetes\_io\_tor','node\_kubernetes\_io\_instance\_type','topology\_kubernetes\_io\_region','topology\_kubernetes\_io\_zone','container\_name','pod\_name','receive\_r eplica', 'role'] 6aî 🗆

tag\_count: 79

# **Too Many Columns?**



Number of columns is not limited explicitly. 1000 columns will work Ok.

Jun 28, 2016



Are there any limit on number of columns in ClickHouse

But actually we don't need so many columns, this will be a very sparse table.

## Align Columns by Metric Name

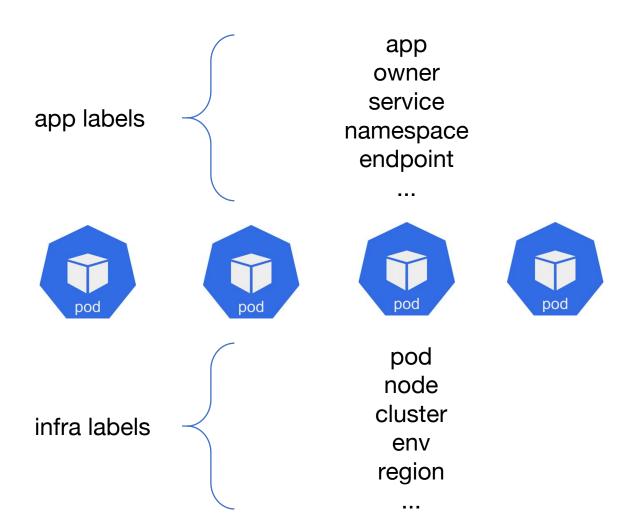


#### One metric usually has no more than a few hundred labels

```
CREATE TABLE prometheus samples
    `time` DateTime('Asia/Shanghai') CODEC(DoubleDelta),
    `metric name` String,
    `label name 1` String,
    `label name 2` String,
    // ...
    `label name 256` String, // Assume that metric will not have
                              // more than 256 labels
    `value` Float.64
SELECT host, value
FROM prometheus samples
WHERE host = 'some-host'
```

### Can Columns be Reduced Further?

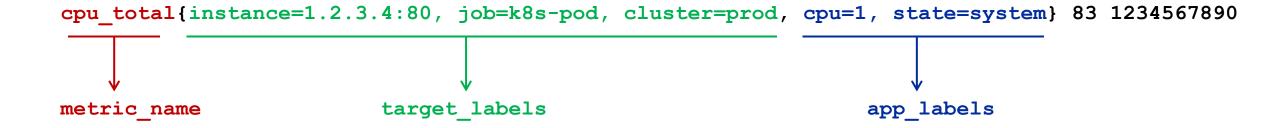




Pods are often rebuilt

# Distinguish Target Labels an App Lables SOPEN SOUR





Cardinality(all series) = 20M

V.S.

N(targets) = 14K Cardinality(app labels) = 1M

# Distinguish Target Labels an App Lables SOPEN SOUR



### Can Columns be Reduced Further? / 2





meta label pod container

```
apiVersion: v1
items:
- apiVersion: v1
  kind: Pod
  metadata:
    annotations:
      cni.projectcalico.org/containerID: b85f57e1d0a98d274293cefd4b9...
      cni.projectcalico.org/podIP: 100.89.13.129/32
      cni.projectcalico.org/podIPs: 100.89.13.129/32
    labels:
      k8s-app: calico-kube-controllers
      pod-template-hash: 57fbd7bd59
  spec:
    containers:
      - name: KUBE CONTROLLERS_CONFIG_NAME
       value: default
      - name: DATASTORE TYPE
       value: kubernetes
      - name: ENABLED CONTROLLERS
        value: node
      - name: KUBERNETES_SERVICE_HOST
        value: 10.96.0.1
      - name: KUBERNETES SERVICE PORT
        value: "443"
```

### Meta Labels & Unlimited Custom Labels



```
CREATE TABLE prometheus samples
    `time` DateTime('Asia/Shanghai') CODEC(DoubleDelta),
    `metric name` String,
    `target id` Uint32,
     `<mark>pod</mark>` String,
    `container` String,
    `label name 1` String,
    // ...
    `label name 16` String, // Assume that metric will not have
                              // more than 16 labels injected by instrumentation
    `value` Float.64
```

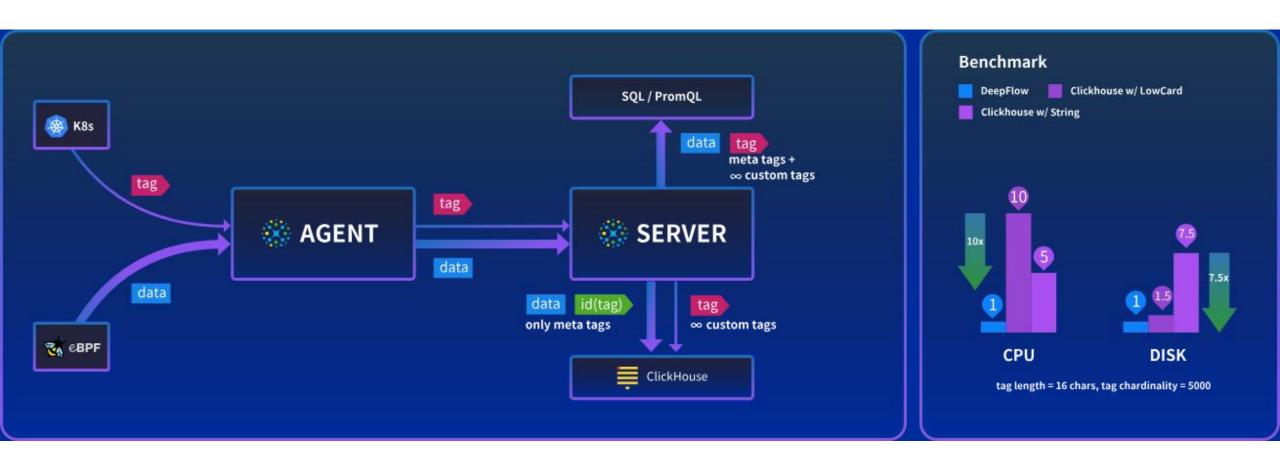
# **Query Custom Lables**



```
CREATE TABLE prometheus samples
    `time` DateTime('Asia/Shanghai') CODEC(DoubleDelta),
    `metric name` String,
    `target id` Uint32,
    `pod` String,
    `container` String,
    `label name_1` String,
    // ...
    `label name 16` String, // Assume that metric will not have
                            // more than 16 custom labels (by instrumentation)
    `value` Float64
CREATE DICTIONARY custom labels.k8s_label
                                                 SELECT dictGet(
                                                     custom labels.k8s label,
   `pod id` Uint64,
                                                     'value',
    `key` String,
                                                     (pod id, 'app')
    `value` String,
                                                 ) AS `app`
                                                 FROM `cpu total`
PRIMARY KEY pod, key
                                                WHERE `app` = 'some-app'
LIFETIME (MIN 0 MAX 60)
                                                 LIMIT 1
LAYOUT (COMPLEX KEY HASHED())
```

# **SmartEncoding in DeepFlow**





# Benchmark (v.s. VictoriaMetrics)



Similar write resource consumption to VictoriaMetrics, but supports unlimited labels and has an order of magnitude better query performance when querying high-cardinality series.

	总资源 [1]	压测参数 [2]	读写压力	CPU 均值	CPU 峰值	内存均值	内存峰值
DeepFlow (DF-Server 及 CK)	<b>24</b> U <b>96</b> GB	targetCount = <b>4200</b> scrapeInterval = 15s queryInterval = <b>60</b> s	367k samples/s 5.8M series 1.5 QPS	12.0 U	15.3 U	12.1 GB	18.2 GB
VictoriaMetrics (不含 VMInsert)	44 U 208 GB 8642	targetCount = 6000 scrapeInterval = 15s queryInterval = 15s	367k samples/s 5.8M series 1.5 QPS	12.1 U	13.8 U	14.8 GB	18.6 GB
Mimir 8642	43.2 U 283 GB	targetCount = 6000 scrapeInterval = 15s queryInterval = 15s	367k samples/s 5.8M series 1.5 QPS	20.3 U 向阳 8642	24.1 U	102 GB	120 GB

# Benchmark (one of our prod. env)



**Targets Count** 

28674

		Targets Metric Count	
Instance	Job	scrape_url	metrics_count 4
10.28.51.130:8000	kubernetes-pod		564
10.28.51.52:8000	kubernetes-pod		564
10.28.51.52:8000	kubernetes-pod		564
10.28.51.130:8000	kubernetes-pod		564
10.28.51.130:8000	kubernetes-pod	http://10.28.51.130:8000/metrics	564
10.28.51.52:8000	kubernetes-pod	http://10.28.51.52:8000/metrics	564
10 28 51 181-8000	kuhernetes-nod		5na

Metrics Count

Label Values Count

6036 790212

name	values_count
id	176287
container_id	173962
name	147769
mountpoint	63825
uid	46962
pod	38751
pod_name	28698
address	24763
device	20888
pod_ip	20421
replicaset	20222
uri	19840
path	19209

**Label Values Count** 

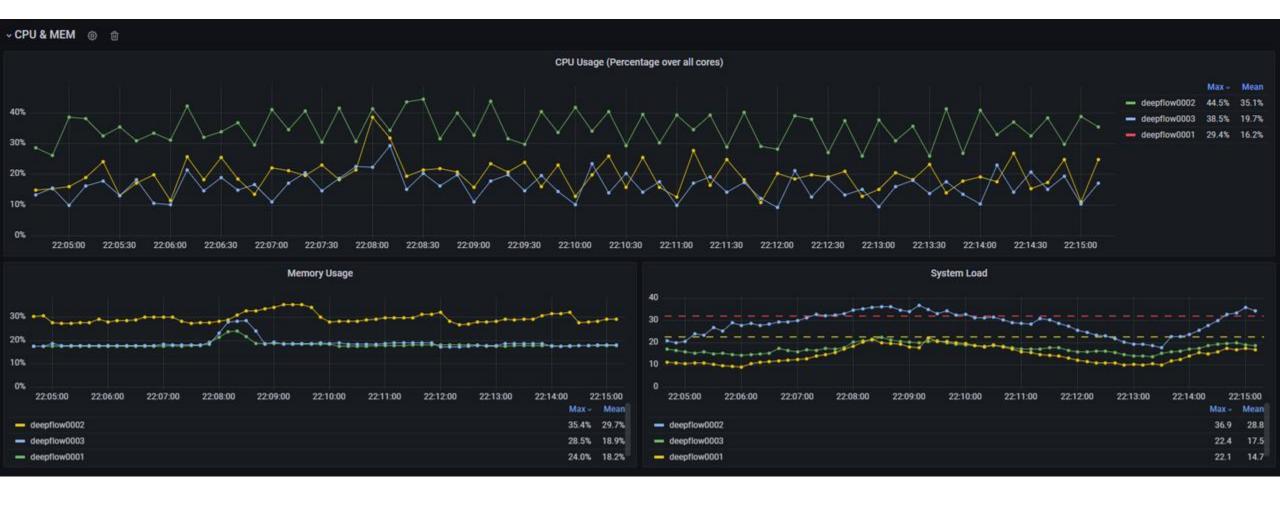
	Lables Count	pod
metric_name	labels_count	pod.
container_memory_failures_total	325744	
container_cpu_user_seconds_total	301511	addr
container_cpu_system_seconds_total	299934	devi
container_cpu_load_average_10s	297590	pod.
container_last_seen	297402	repli
container_memory_cache	295290	
container_tasks_state	294973	path

# Benchmark (one of our prod. env)



400K samples/s ingestion

CPU: 740%~1100% (3200% Total), MEM: 29G~37G (128GB Total) // only deepflow-server



# Benchmark (one of our prod. env)



+1500K QPS

CPU: +600% // only deepflow-server



### Thanks!





OpenSourced under the Apache 2.0 License https://github.com/deepflowio/deepflow https://deepflow.io

Our WeChat Group
Discord? Check our GitHub repo!

Consul

**ANY TAG** 

SVC, INSTANCE, ENDPOINT, ...

ANY STACK
APPLICATION, INFRA, ...

MySQL eredis DATASTORE

kubernetes K8S APISERVER

SERVICE DISCOVERY

CLOUD API

Skywalking Skywalking Grafana Prometheus kafka **DATA INTERFACES** SQL / PROMQL / OTLP DATA **DEEPFLOW** STORE **OBSERVABILITY WITH AUTOMATION** ClickHous **AGENT** SERVER ZERO CODE / FULL STACK / UNIVERSAL TAGGING **EXTERNAL SIGNAL SOURCES** telegraf Prometheus Skywalking Pyroscope