

GOTC

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

OPEN SOURCE , OPEN WORLD

分布式数据库与存储

老树发新芽 —— 开箱即用的开源PostgreSQL发行版 Pigsty

冯若航 2021年08月01日

About Me

Ruohang Feng (Vonng)

<https://vonng.com>

rh@vonng.com

独立开源贡献者, Pigsty作者, 《Designing Data-Intensive Application》《PG Internal》译者
PostgreSQL中国社区开源技术委员会委员, PostgreSQL专家, 全栈开发者, 曾任职于Alibaba, Apple, TanTan

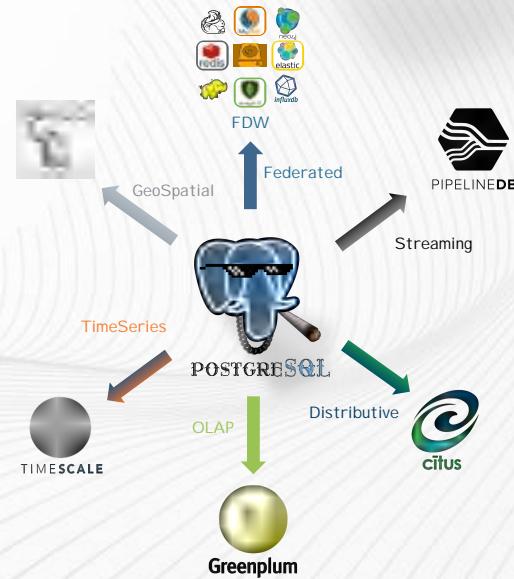
- <http://demo.pigsty.cc>

POSTGRESQL: THE WORLD'S MOST ADVANCED OPEN SOURCE RELATIONAL DATABASE

Get Started Now!

世界上最先进的开源关系型数据库

- 性能优异，功能强大
 - 扩展丰富，协议友善
 - 历史悠久，稳定可靠
 - 代码清晰，定制灵活
- 一招鲜，吃遍天



一专多长的全栈数据库



成熟的应用可能会用到许许多多的数据组件：
缓存，OLTP，OLAP/批处理/数据仓库，流处理/
消息队列，搜索索引，NoSQL/文档数据库，地理
数据库，空间数据库，时序数据库，图数据库。

传统的架构选型可能会组合使用多种组件，典型
如：Redis + MySQL + Greenplum/Hadoop
+ Kafka/Flink + Elasticsearch，
一套组合拳基本能应付大多数需求。

带来的新问题就是 **异构系统集成**

需编写大量重复繁琐搬砖代码，
把数据从A组件搬运到B组件。

“开源的Oracle”

去O扛旗者

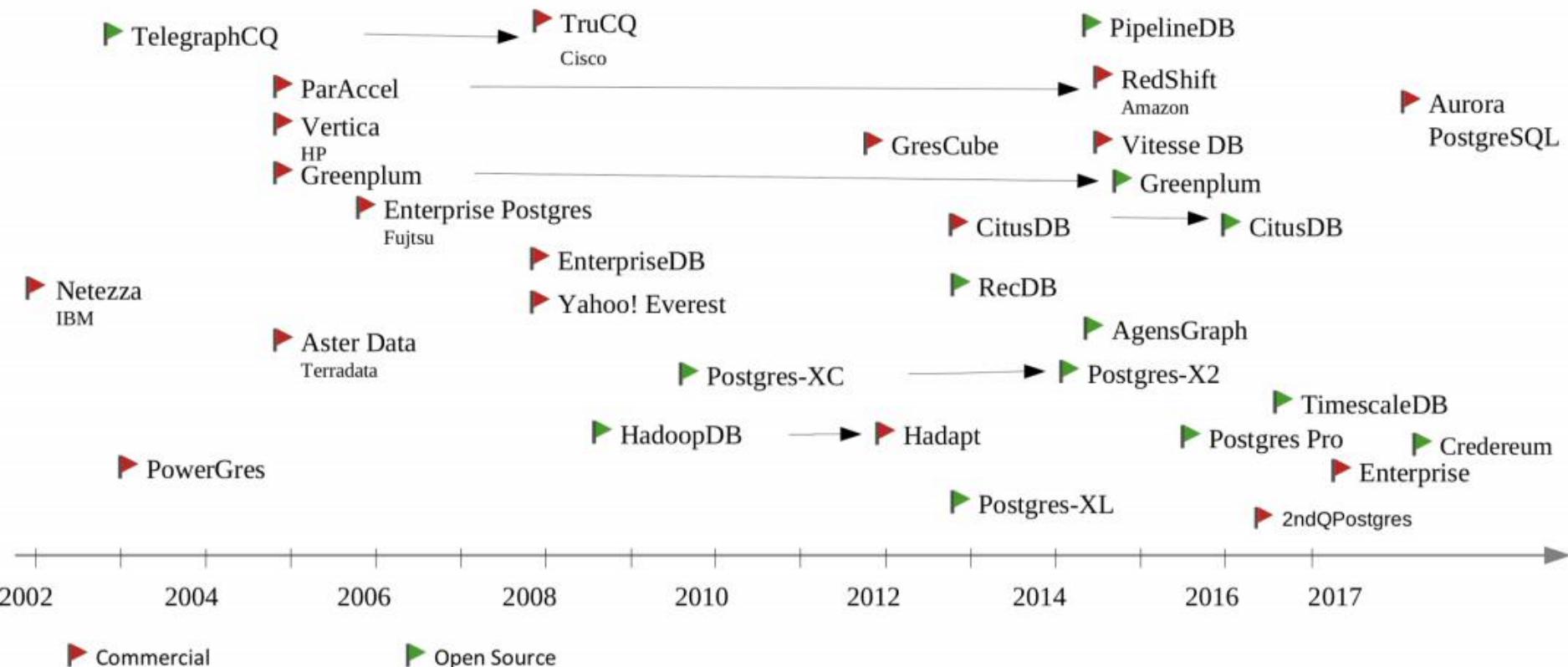
原生支持Oracle 八九成功能

唯一能对Oracle构成直接威胁的开源数据库



开放扩展的数据库

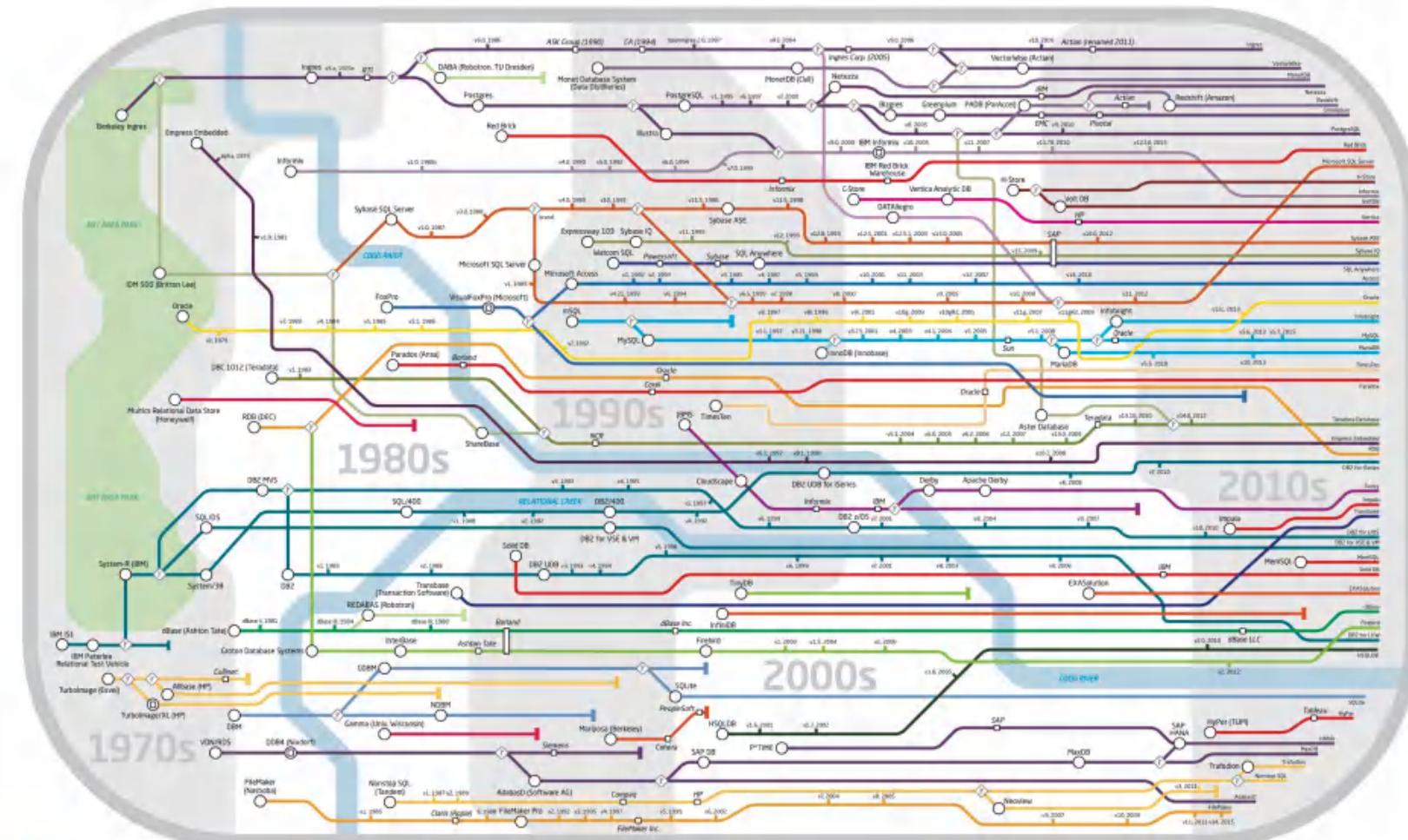
健康的社区，友善的协议，丰富的扩展，繁荣的生态
开支散叶，子孙满堂



老树

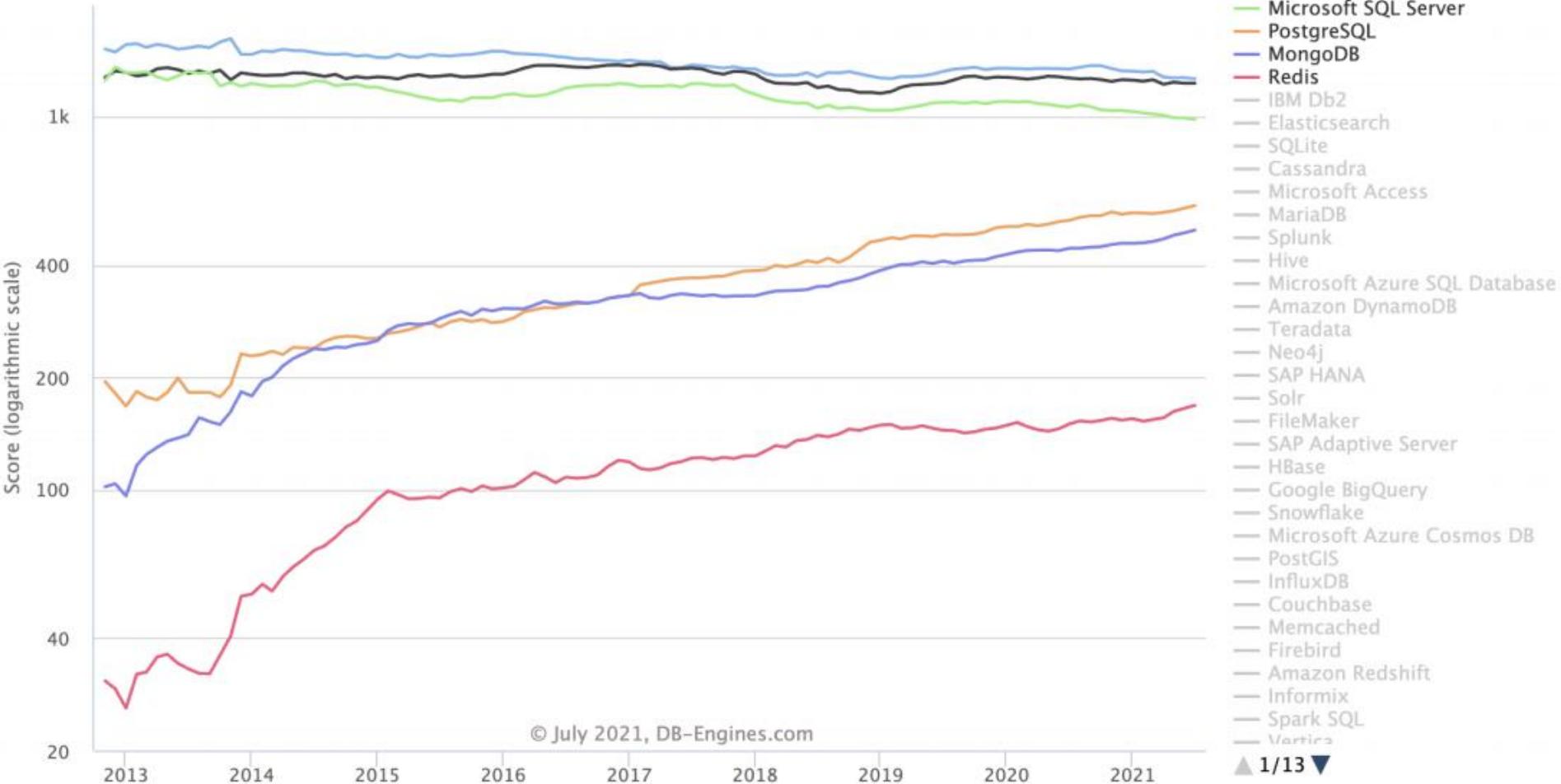
祖宗级开源项目

Genealogy of Relational Database Management Systems



现状 与 未来

DB-Engines Ranking



Pigsty

v1.0.0

Battery-Included Open Source PostgreSQL Distribution

Postgres in Graphic STYle

Ruohang Feng

Outline

What Where Who Why How

Monitoring

Providing

High Available

Sandbox

Analysis

Open Source



Lord of PostgreSQL

Fellowship of Pigsty

What is Pigsty

The Battery-Included Open Source PostgreSQL Distribution

PostgreSQL
数据库
发行版

开箱即用的开源PostgreSQL数据库发行版

打包最新PG内核：PostgreSQL 13.3/14-beta2，以及海量扩展插件（强力扩展：时序数据Ti mescal eDB 2.3，地理空间PostGIS 3.1，分布式Citus 10），全部开箱即用。

数据部署时带有可靠的中间件： Patroni 2.1, Pgbounce1.15, Haproxy 2.2, Consul 1.10

打包有生产级的运行时环境： Grafana 8.0, Prometheus 2.28, Nginx 1.20, Ansible

无需互联网访问，在全新CentOS 7机器上一行命令10分钟完成所有安装流程，所有组件已为生产使用做好准备（不计Pigsty下载，约1GB）

监控系统
解决方案

全面专业的开源PostgreSQL监控解决方案

1000+监控指标 描述了从数据库、中间件、机器节点以及负载均衡器的方方面面，从最顶层的全局概览到最小的单个数据库对象（表、索引、函数等）都一览无余。

带有了20余个精心设计的监控面板，将监控指标转化为多层次全方位的分析与洞察。所有图表都带有丰富的导航连接供下钻、上卷、横跳，可快速定位故障与性能问题。内置了基于Prometheus+AlertManager与Grafana的报警系统，两者等效且互为备份，预置常用数据库与机器报警规则，开箱即用。

可选装基于Logstash与promtail的实时日志收集方案，并提供用于分析PG CSV日志样本的应用 pglog

部署管控
解决方案

简单易用的开源PostgreSQL管控解决方案

声明式配置文件允许用户通过160+参数描述与定制数据库与基础设施运行时的方方面面，而无需操心具体的实现，整体提供类似Kubernetes的云原生使用体验。

自愈的高可用数据库集群 Pigsty部署的数据库集群可以自动从常见故障中恢复，集群成员对外表现等价，只要仍有任意实例存活，集群对外提供的各种服务就不会停止。

提供了一系列内置脚本 用于部署、扩容、销毁PostgreSQL数据库集群与实例，管理业务数据库与用户，使用简单方便。提供自带的CMDB供扩展与定制，提供可选的CLI与GUI内置服务发现机制。除用于提供高可用的DCS外，数据库与基础设施采用松耦合模式。数据库变动（扩缩容，上下线）自动适配，无需手工维护基础设施与周边系统。

便捷全能
沙箱环境

一键拉起的全能PostgreSQL本地沙箱环境

Pigsty针对大规模生产环境的PostgreSQL管理而设计，但亦设计为可完整运行于本地1核1GB虚拟机中，便于保持开发、测试、预发、生产环境的一致性。

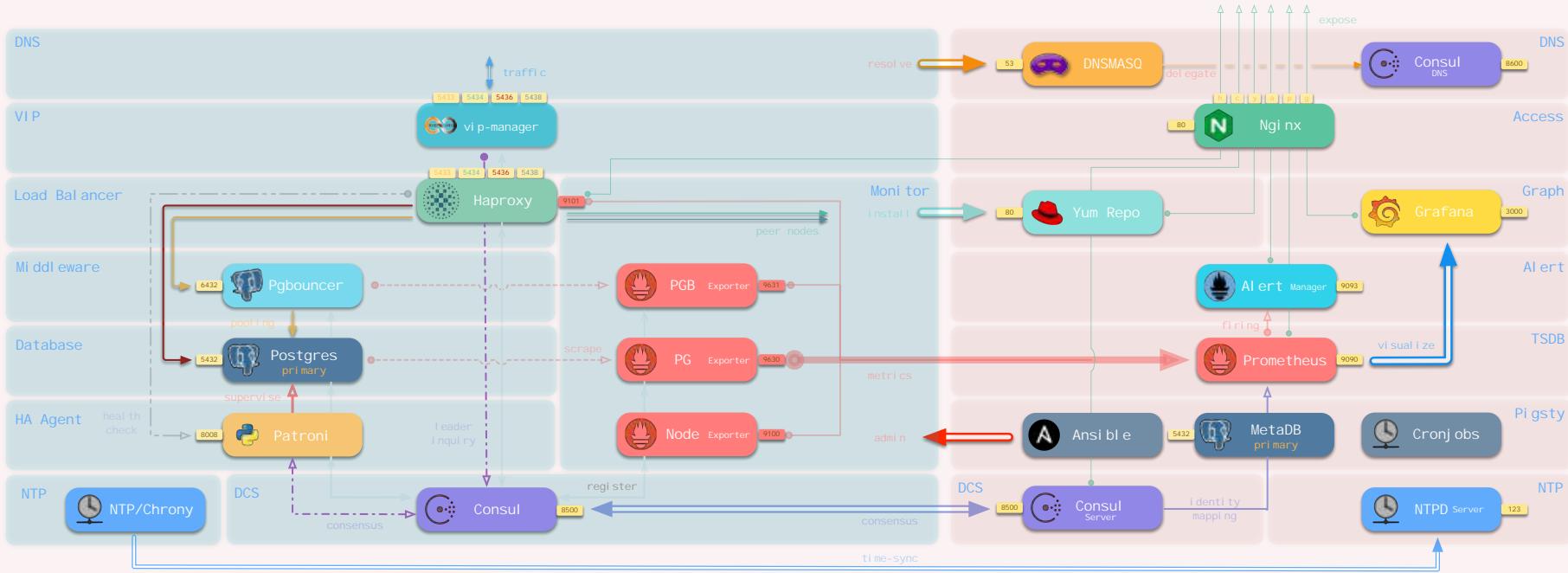
用户在自己的笔记本(Macbook)上可一键拉起功能完整的本地沙箱（由vagrant与virtual box驱动），或通过简单的配置在自己的虚拟机或云ECS/VPS上完成部署。

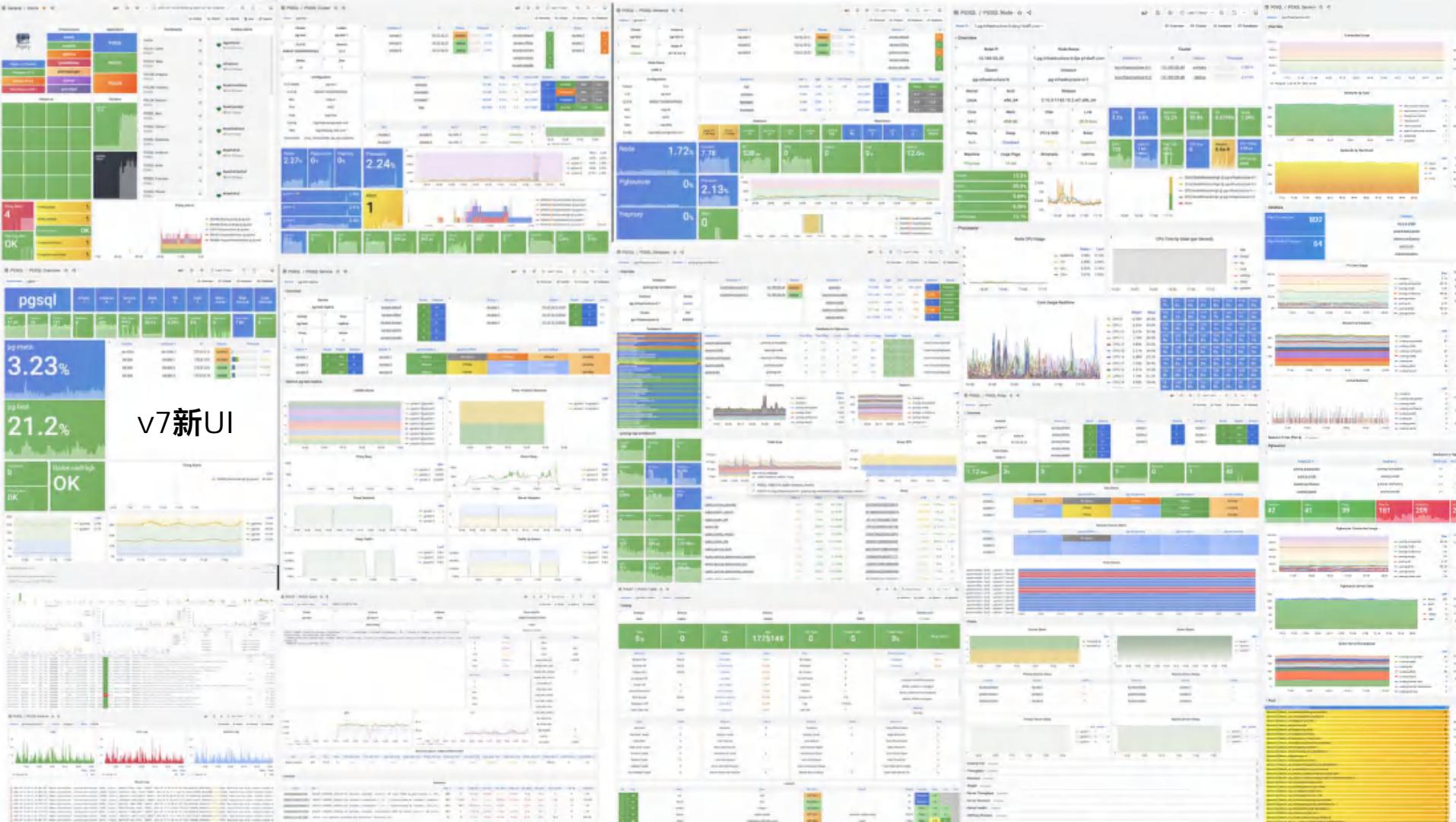
用户可以在与实际生产环境（除规格外）完全一致的沙箱环境中学习探索PostgreSQL，进行各种实验，仿真，测试，开发，并从生产级监控系统中获取全面且实时的数据反馈。

用户可以使用内置的Postgres & Grafana & Echarts开发交互性数据应用与数据可视化作品。快速产出作品原型，并以标准的方式、在确定性的环境中分享，演示与交付。

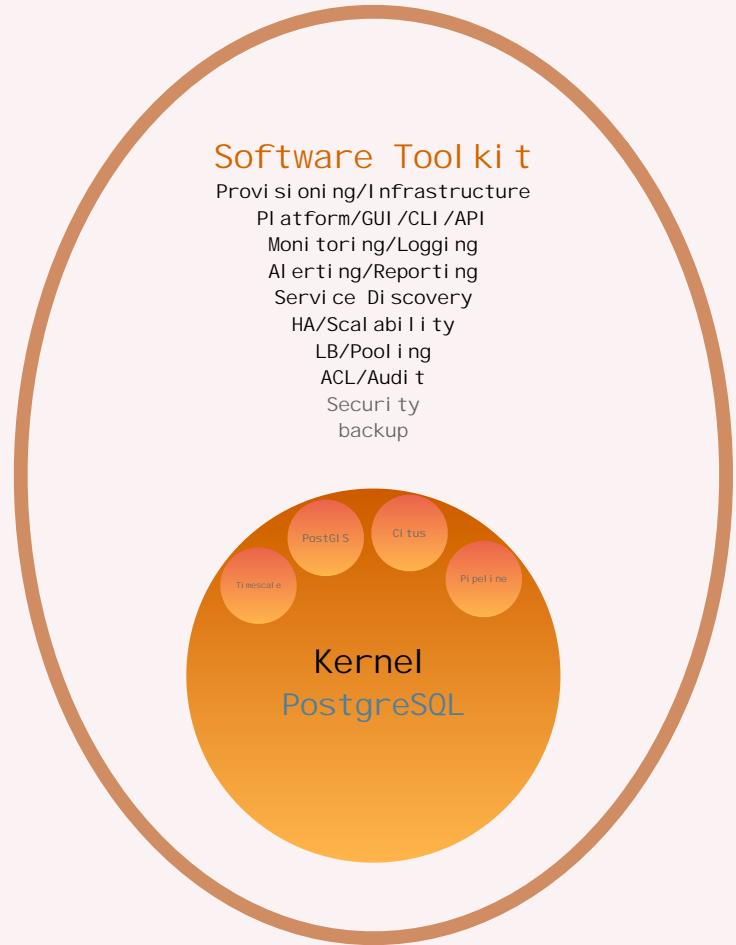
开箱即用

Pigsty将集群部署，扩容缩容，主从复制，故障切换，流量代理，连接池，服务发现，访问控制，监控系统，告警系统，日志采集等生产级成熟解决方案封装为发行版，一次性解决生产环境、个人使用PostgreSQL数据库的各类问题，真正做到 **开箱即用**。





发行版



Who wants Pigsty

互联网企业

中小型互联网企业 (日活千万量级, QPS百万量级, 数据库集群数百量级)

- 体量可观: 实例动辄数百个, QPS动辄数百万。
- 追求敏捷: 集群实例变动频繁, 对敏捷交付与快速伸缩有需求。
- 追求可用: 对可用性有高要求, 故障自愈, 自动切换。
- 追求技术, 希望通过数据驱动提高治理水平, 偏爱开源。

传统企业

有信息化需求的各类传统企业

- 不想折腾, 简单好用快速上手, 易维护。
- 追求功能, 地理, 空间, 时序, JSON, 存储过程, 物化视图, 倒排索引等。
- 追求可靠, 数据不丢不错不乱, 备份容灾, 有兜底。
- 自主可控, 不想上云, 但又想要类似云的管理体验

软件厂商

软件集成商, 软件厂商

- 需要人才: 用到PostgreSQL, 却没有相关人才储备。
- 需求方案: 需要成熟、完整的解决方案, 比 `$ yum install && systemctl start` 更靠谱
- 追求成本: 用最少的成本办最多的事

个人用户

DBA, OPS, DEV, 数据研发, 数据分析师, 爱好者, 学生

- 上手快: 一键拉起, 配套齐全, 开箱即用
- 门槛低: 1核1GB的底配虚拟机也可以跑起来。
- 一致性: 本地开发测试环境与生产环境保持一致
- 全功能: 希望有数据处理与分析的整合环境: Python + SQL + 可视化

Where to run Pigsty

生产环境

管理节点：3-5 64核 400GB PCI-E SSD x n

使用高规格物理机/虚拟机部署，使用3-5个管理节点控制监管几百套数据库集群，实际服务于线上业务。
一套Pigsty部署支持几百套数据库集群，几百万QPS，上千万日活不在话下。

测试/开发
仿真环境

管理节点：1-3 8核 16B SSD x n

与生产环境保持一致的Staging/UAT/预发环境
搭建公用的开发、测试环境，仿真线上环境

云虚拟机
演示环境

管理节点：1 2核 4GB x 1~4

使用云服务器部署，分享公开演示程序，快速服务小微数据业务。

个人电脑
沙箱环境

管理节点：1 1核 2GB x 1

运行于个人笔记本&台式机，Mac & PC 虚拟机中。
用于搭建Devbox，测试，实验，学习PostgreSQL
制作数据应用，交互式数据可视化应用

痛点！痒点！爽点！

					易用性?	生态	是否开源
					复杂度?	费用	社区活力
					机器成本?	厂商支持	协议
功能?	实施?	可靠性	伸缩性	SQL变更?			
ACID事务	环境?	高可用?	分区分片	数据迁移?			
地理空间?	安装?	监控系统?	高并发	扩容缩容			
时间序列?	平台?	告警系统	性能?	用户管理			
全文检索?	演示	可观测性	调度?	软件升级			
存储过程?	部署?	主从延迟	弹性伸缩?	插件管理			
JSON/XML	管理?	负载均衡	参数调优	配置管理	审计		
JOIN	维护?	中间件	HTAP	恢复方案	日志		
可扩展性	自动化?	接入	读写分离	备份方案	访问控制		
扩展性			ETL	故障预案	权限管理	审美?	谈资?
			CDC				

痛点?

从无到有

发行版, 部署, 解决方案

痒点?

从有到优

监控系统, 高可用架构

爽点?

从优到易

开箱即用, 全能沙箱

不可替代性?

从易到廉

开源免费

安装

Installation Wizard



Someone told me that each equation I included in the book would halve the sales.

— Stephen Hawking

有人告诉我，安装脚本里每多一行命令，用户数量就会减半

How to get Pigsty

Get pigsty with three lines of code

准备

内核: Li nux
架构: x86_64
版本: CentOS 7. 8. 2003 (或其他RHEL7等效版本)
用户: root (或免密码sudo用户)
规格: 1核2G (推荐使用2核4G以上规格)

下载

源代码 (必选项)
pigsty.tgz ---下载解压至--> ~/pigsty

软件包 (可选项)
pkg.tgz -----> /tmp/pkg.tgz

配置

使用何种配置样板? (通过-m参数指定)
demo (单节点), demo4 (4节点沙箱), pg14, tiny (常规默认), oltp, ...
是否下载软件包? (如果~/tmp/pkg.tgz不存在则会询问)
Y: 现在从GitHub上下载pkg.tgz
N: 后续安装时从原始上游源下载 (如果操作系统非CentOS 7. 8请使用此项)
主要的IP地址? (如果发现多块网卡与多个可用IP地址会询问)
输入: 当前节点 (管理节点) 使用的IP地址
(访问该节点时使用的IP地址, 但不要使用公网IP)

安装

访问该节点的3000端口即可浏览Pigsty监控首页
<http://10.10.10.10:3000> # 替换为您自己的IP地址
默认用户名与密码: admin:pigsty
使用域名将带来更好的使用体验, 例如: http://g.pigsty
执行 sudo make dns 会将Pigsty相关域名写入/etc/hosts
10.10.10.10 meta.pigsty.g.pigsty

```
$ echo "$(uname -s) $(uname -m) $(cat /etc/redhat-release)"
```

```
Li nux x86_64 CentOS Li nux release 7. 8. 2003 (Core)
```

确认无误后即可开始下载配置安装

\$ cd ~

```
$ git clone https://github.com/Vonng/pigsty
```

或使用curl 下载

```
$ curl -SL https://github.com/Vonng/pigsty/releases/download/v1.0.0-beta1/pigsty.tgz | gzip -d | tar -xC ~
```

下载可选离线软件安装包, 可加速安装, 目标节点无互联网访问时为必选。离线软件包也可以在configure过程中下载

```
$ curl -SL https://github.com/Vonng/pigsty/releases/download/v1.0.0-beta1/pkg.tgz -o /tmp/pkg.tgz
```

\$ cd ~/pigsty

```
$ ./configure
```

```
[-n] [--non-interactive]  
[-m] --mode <configuration-mode>  
[-d] --download  
[-i] --ip <primary_ip>
```

```
# 如果您已经确定了三个参数选项, 可以使用 -n 跳过交互式向导  
# 指定需要使用的配置文件模板: demo, demo4, pg14, tiny, oltp等  
# 如果~/tmp/pkg.tgz不存在, 就从GitHub下载离线安装包  
# 存在多个IP地址时, 指定其中一个作为首要IP地址
```

```
$ make install
```

Ansiible剧本infra.yml 将在meta分组上执行以完成安装

ansible将在configure过程中从离线软件包或本地可用的其他源中安装

Installation

download → configure → install



A WIZARD!

```
[vagrant@meta ~]$ curl -fsSL https://pigsty.cc/pigsty.tgz | gzip -d | tar -xC ~; cd ~/pigsty # DOWNLOAD
[vagrant@meta pigsty]$ make config
./configure
configure pigsty begin
[ OK ] kernel = Linux
[ OK ] machine = x86_64
[ OK ] release = 7.8.2003 , perfect
[ OK ] sudo = vagrant ok
[ OK ] ssh = vagrant@127.0.0.1 ok
[WARN] Multiple IP address candidates found:
        (1) 10.0.2.15          inet 10.0.2.15/24 brd 10.0.2.255 scope global noprefixroute dynamic eth0
        (2) 10.10.10.10         inet 10.10.10.10/24 brd 10.10.10.255 scope global noprefixroute eth1
[ OK ] primary_ip = 10.10.10.10 (from demo)
[ OK ] admin = vagrant@10.10.10.10 ok
[ OK ] mode = demo (vagrant demo)
[ OK ] config = demo@10.10.10.10
[ IN ] Cache /tmp/pkg.tgz not exists, download? (y/n):
=> y
[ OK ] cache = download from internet
      % Total    % Received % Xferd  Average Speed   Time     Time     Time  Current
                                     Dload  Upload   Total   Spent   Left  Speed
100 1031M  100 1031M    0     0  3338k      0  0:05:16  0:05:16  --:--:-- 7061k
[ OK ] repo = extract from /tmp/pkg.tgz
[ OK ] repo file = /etc/yum.repos.d/pigsty-local.repo
[ OK ] utils = install from local file repo
[ OK ] ansible : ansible 2.9.18
configure pigsty end. Use 'make install' to proceed
[vagrant@meta pigsty]$ make install
./infra.yml -l meta
```

<https://pigsty.cc/zh/docs/quick-start/>

Monitoring



You can't manage what you don't MEASURE.
— Peter F. Drucker

监控系统

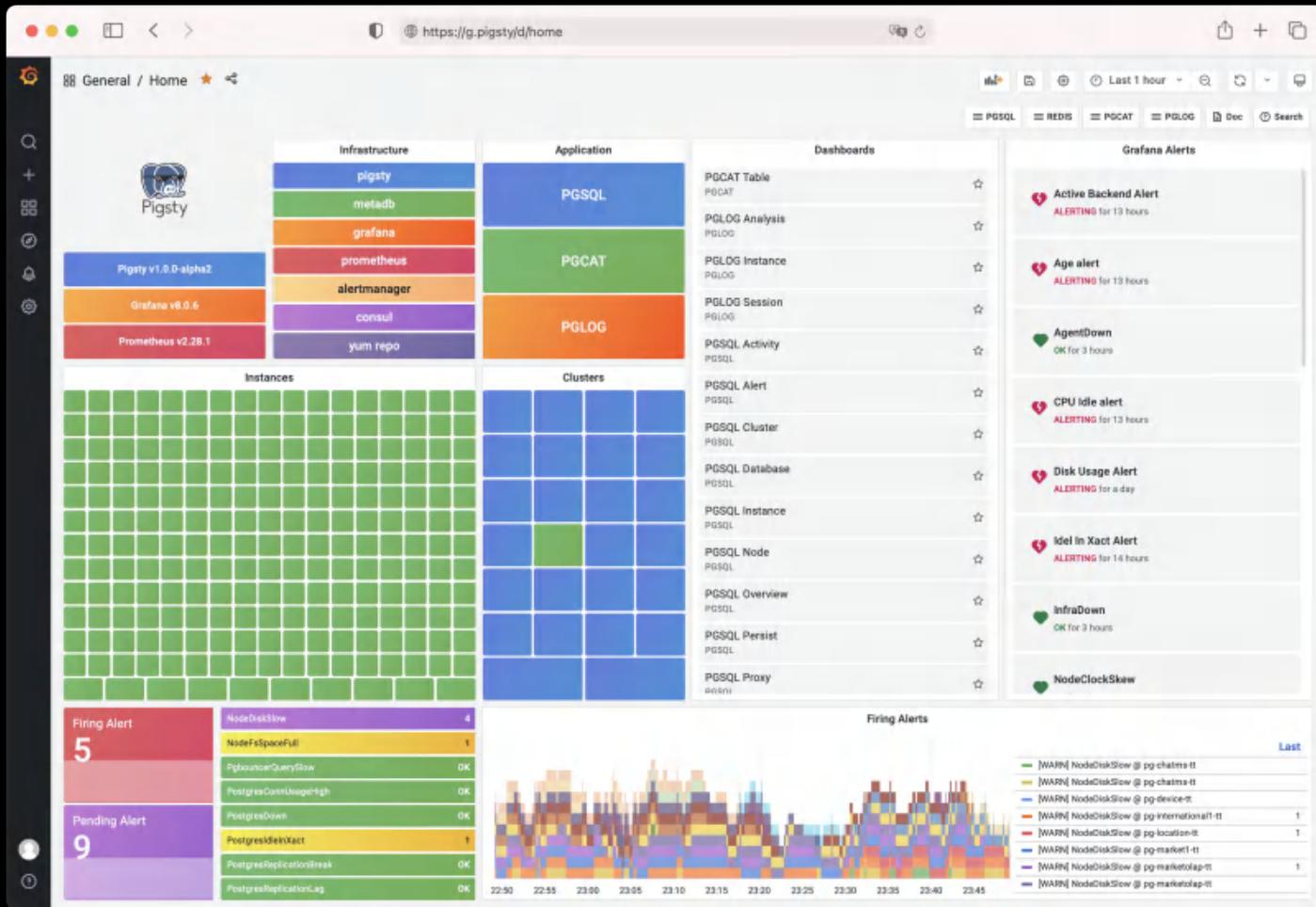
开源监控最佳实践

大规模集群管理

多层次，全覆盖

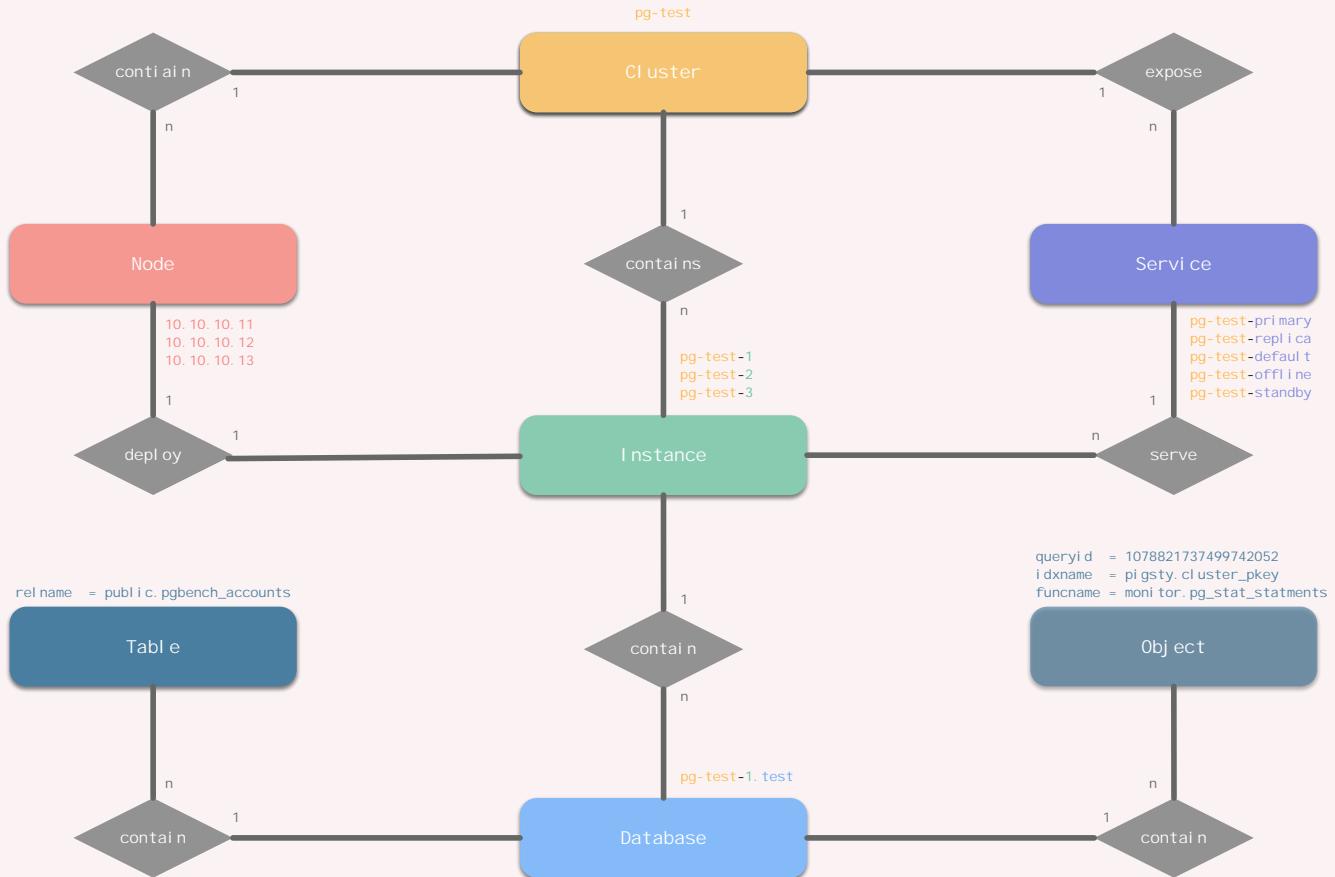
生产验证方案

最专业，最全面



Domains

Entity
Relation
Identifier



Entities & Identifiers

Entity

Identifier

Label

Overview

pgsql
redis
mysql

```
{ job = "pgsql" }
```

Shard

pg-test

```
{ job = "pgsql", cls =~ 'pg-test\d+' }
```

Cluster

pg-meta

pg-test
pg-test1
pg-test2

```
{ job = "pgsql", cls = 'pg-test' }
```

Service

pg-meta-primary
pg-meta-replica
pg-meta-default
pg-meta-offline

pg-test-primary
pg-test-replica
pg-test-default
pg-test-offline
pg-test-standby

```
{ job = "pgsql", cls = 'pg-test', svc = 'pg-test-primary' }
```

Instance

pg-meta-1

pg-test-1
pg-test-2
pg-test-3

```
{ job = "pgsql", cls = 'pg-test', ins = 'pg-test-1' }
```

Database

pg-meta-1.meta

pg-test-1.test
pg-test-2.test
pg-test-3.test

```
{ job = "pgsql", cls = 'pg-test', ins = 'pg-test-1' , datname = 'test' }
```

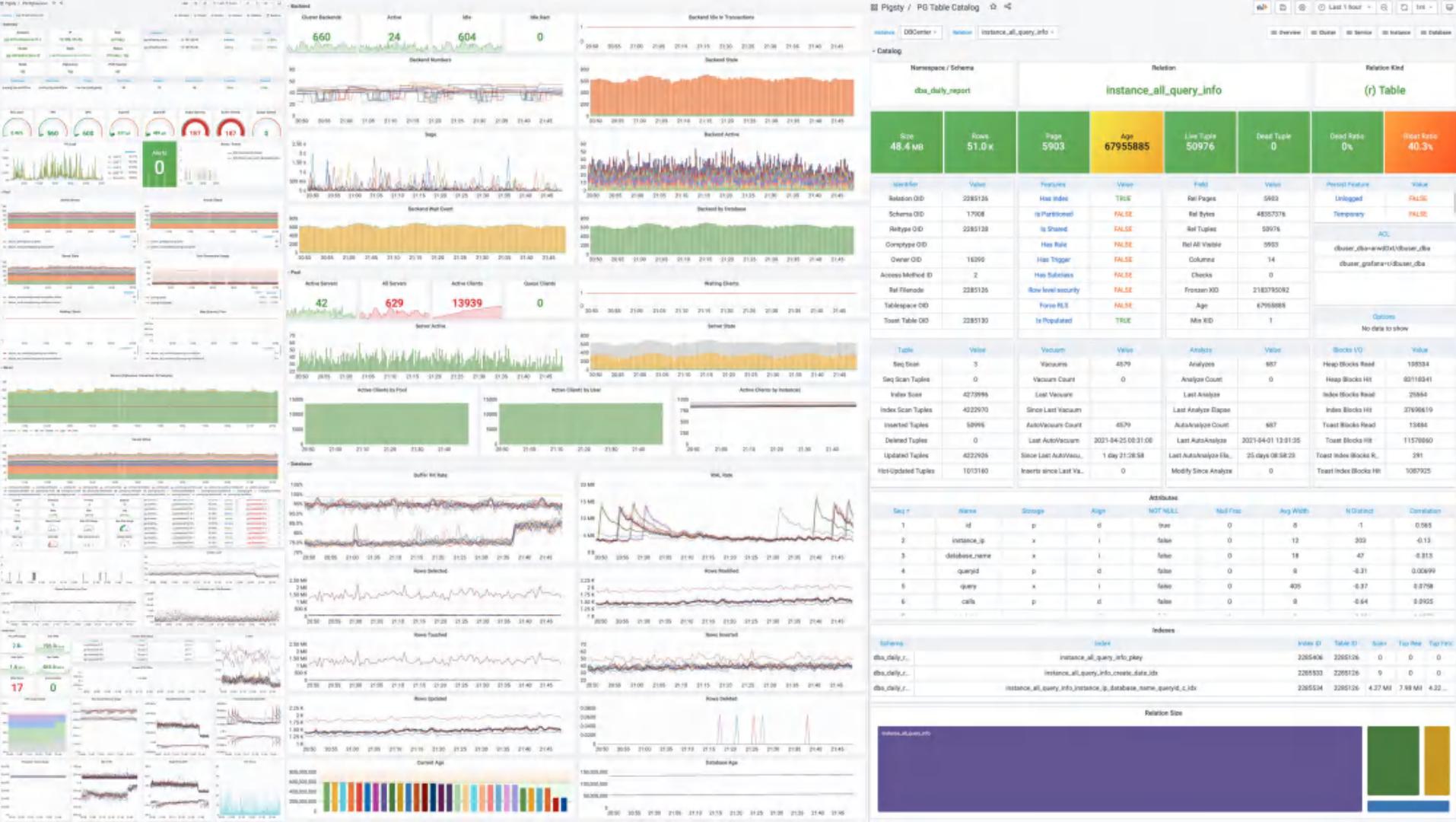
Object

```
queryid = 1078821737499742052  
relname = public.pgbench_accounts  
idxname = pgsty_cluster_pkey  
funcname = monitor.pg_stat_statements
```

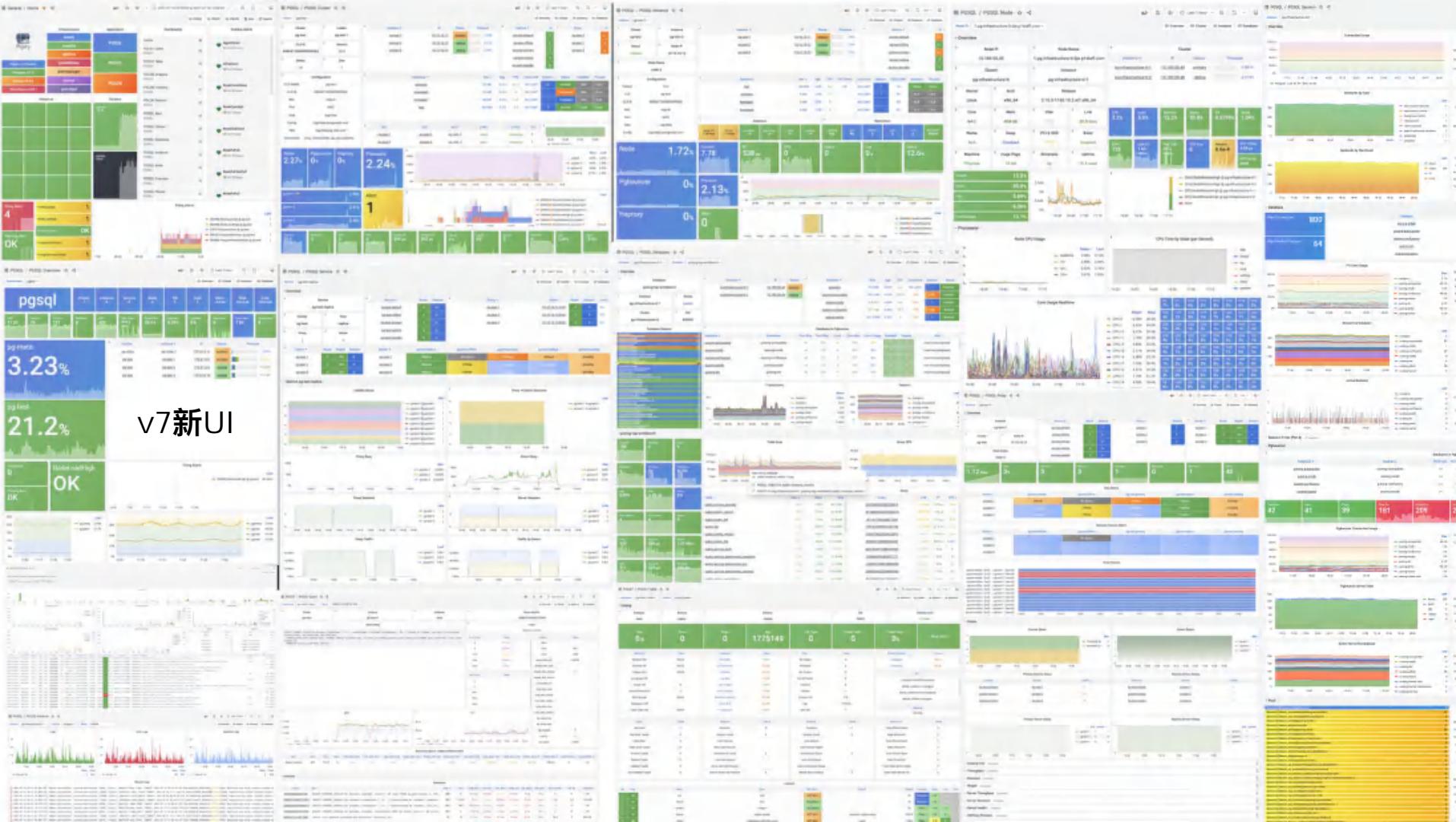
```
{ job = "pgsql", cls = 'pg-test', ins = 'pg-test-1' , datname = 'test' , relname = 'public.pgbench_accounts' }
```





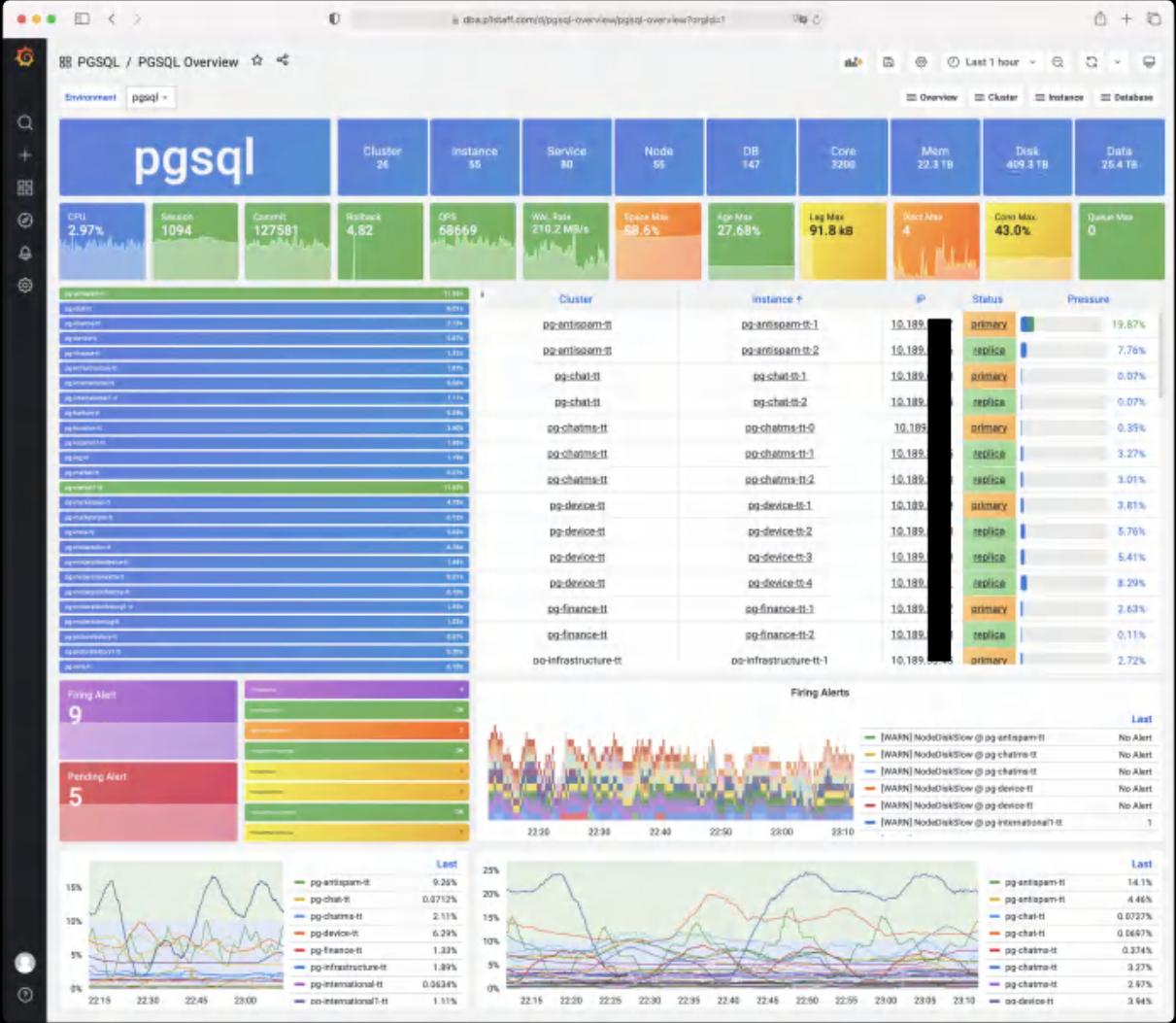
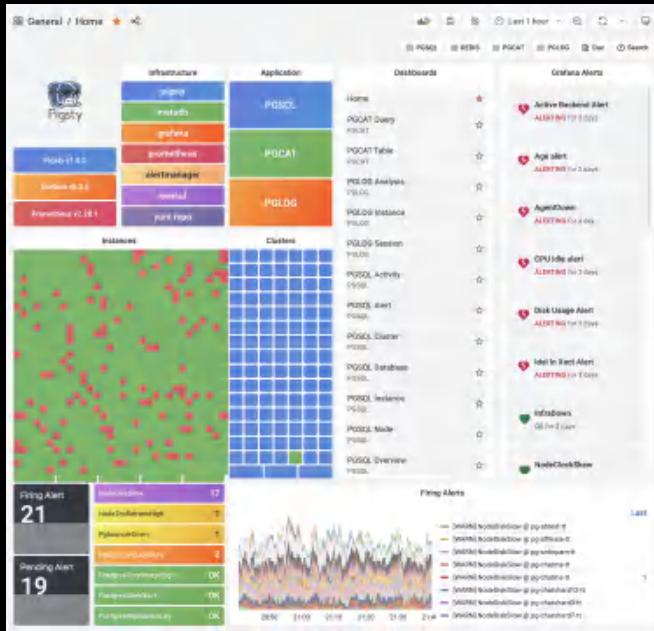






v7新UI

PGSQL Overview

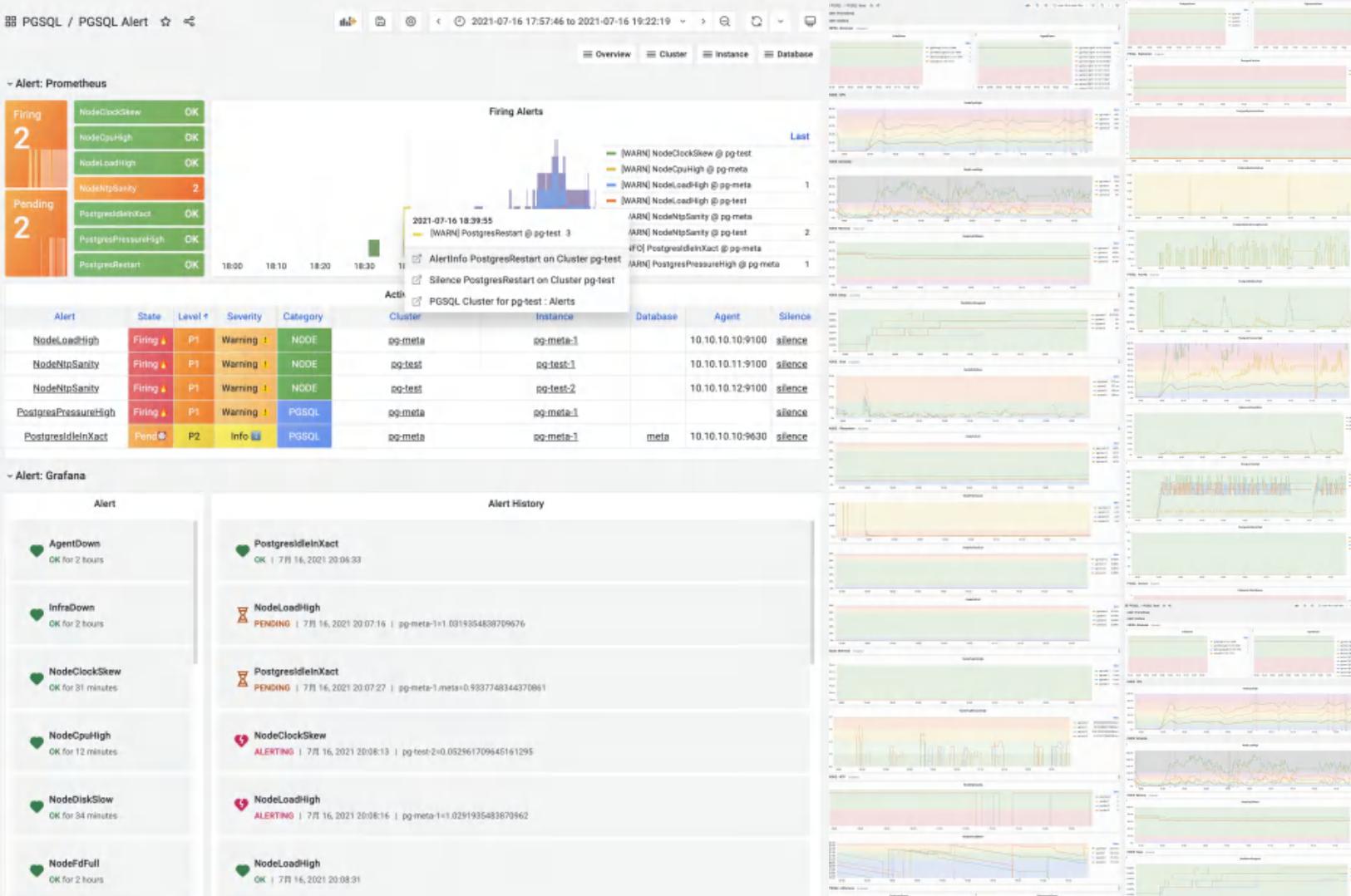


PGSQL Alert

Prometheus AlertManager

Grafana

等效实现
互为补充



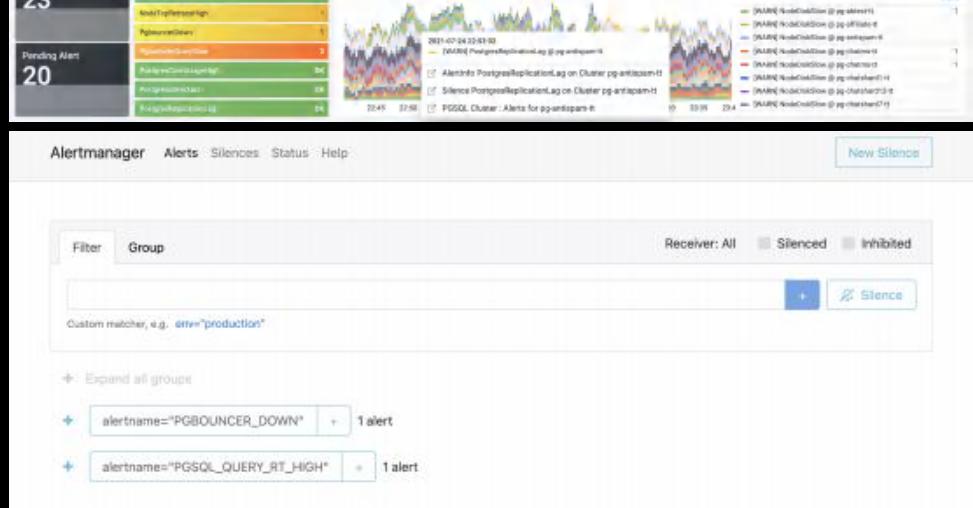
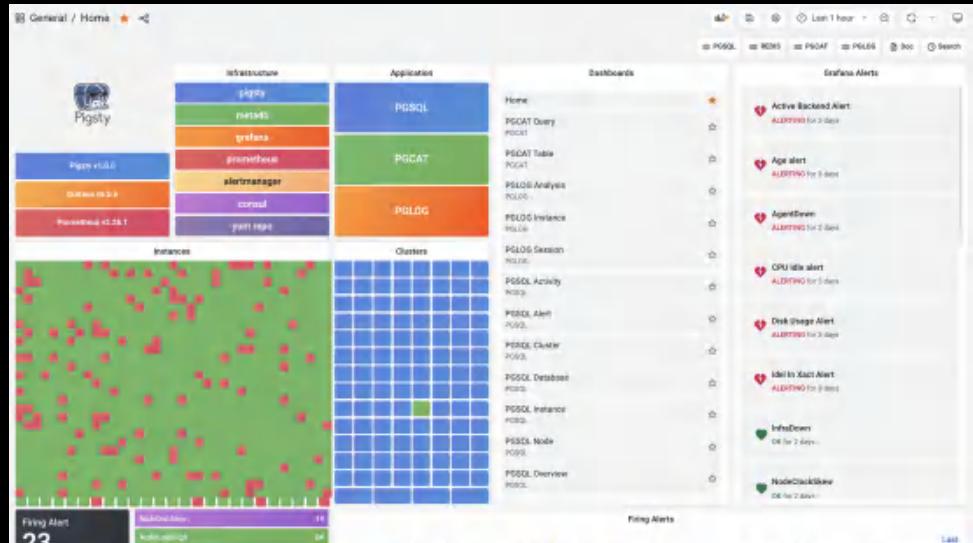
交互式图表

点击跳转

查看详情

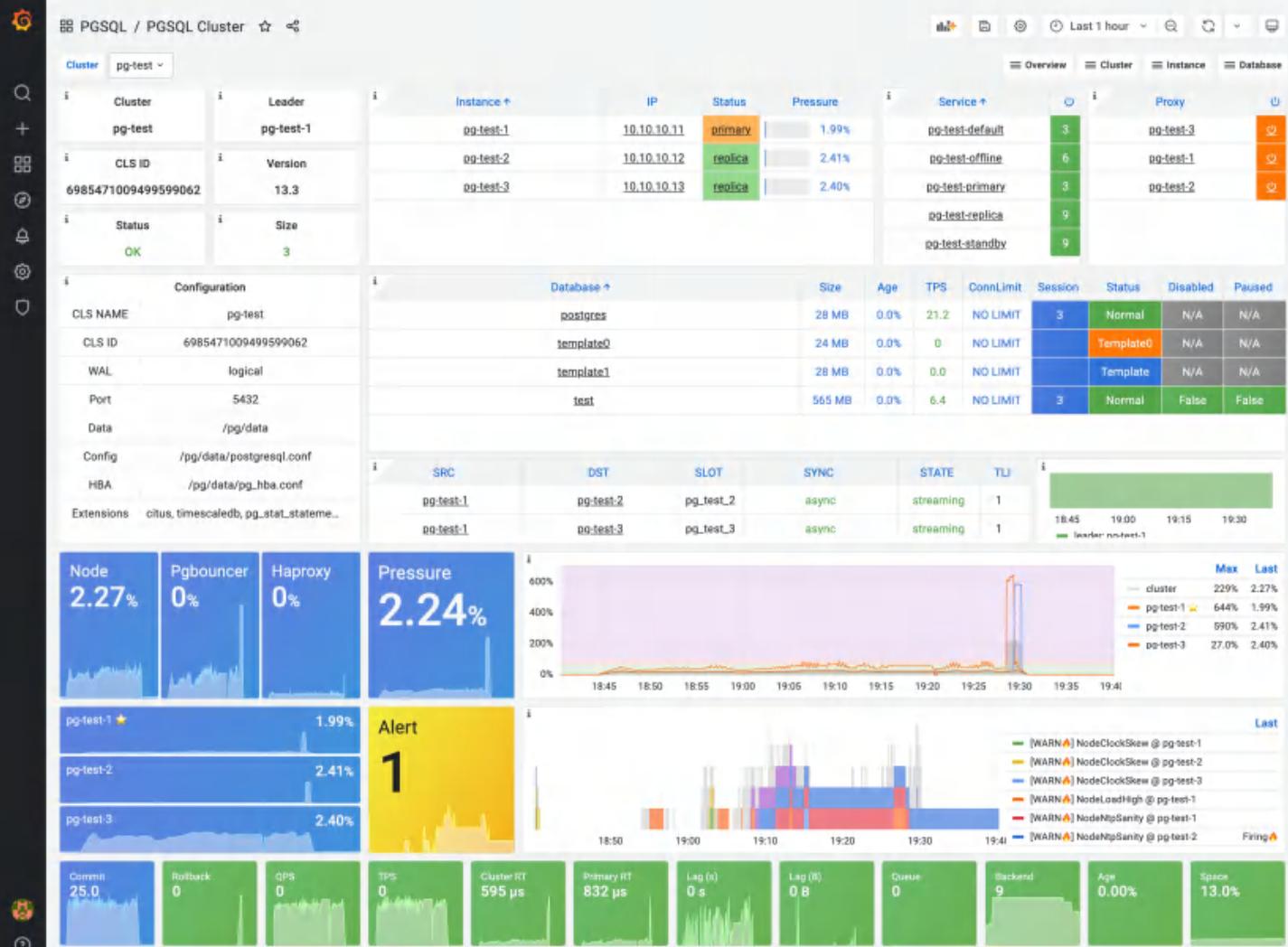
屏蔽告警

下钻至具体集群/实例

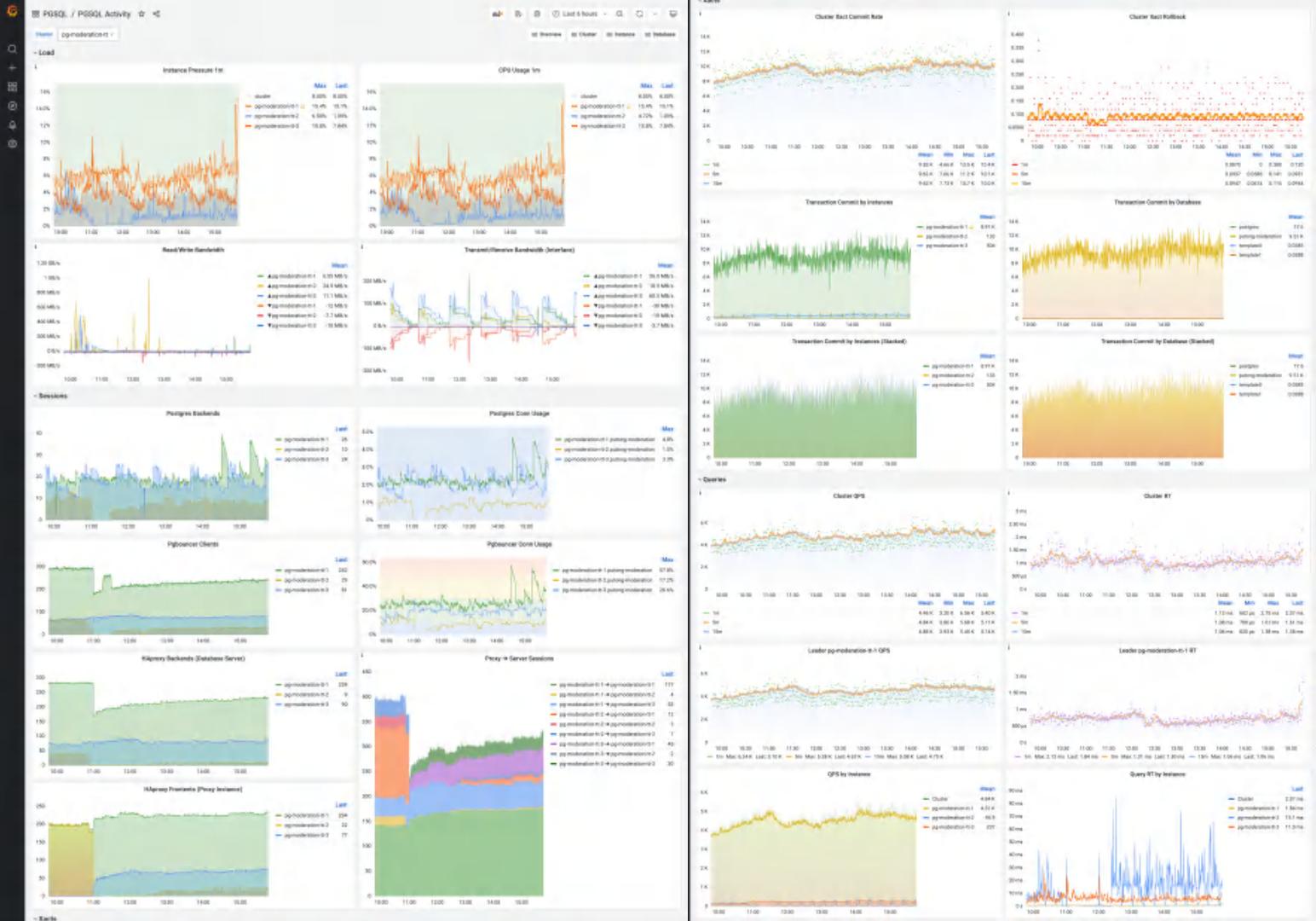


PGSQL Cluster

cluster-wide insights



PGSQL Activity



Replication

跟踪复制进度

监控复制延迟

定位复制延迟问题

逻辑复制进度

同步复制状态

Pigsty / PG Cluster Replication

Last 1 hour 1m

Cluster: pg-payment-tt

Identity (0 panels)

Cluster	Instance #	IP	Role	Load	Info	Value
pg-payment-tt	pg-payment-tt-0	10.189.11.21	primary	0.19%	Cluster ID	65948363840666798
Leader	pg-payment-tt-1	10.189.11.22	replica	0.23%	Data Dir	/export/postgres
pg-payment-tt-0	pg-payment-tt-2	10.189.11.23	replica	0.00%	Port	5432
Members	pg-payment-tt-3	10.189.11.24	replica	0%	Version	18.5
4					Extension	pg_stat_statement
Synchronous					WAL Level	logical
True						

Replication Topo (0 panels)

Sender	Receiver #	Dst IP	State	Mode	TLS	Lag
pg.payment-tt-0	pg_receiveonly	10.189.37.65	streaming	async	0	Nan
pg.payment-tt-0	slave1	10.189.11.22	streaming	quorum	1	0 B
pg.payment-tt-0	slave2	10.189.11.26	streaming	quorum	1	0 B

Primary Connection Info

Instance	Upstream
pg.payment-tt-0	N/A
pg.payment-tt-1	N/A
pg.payment-tt-2	N/A

Instances Primary Replicas Max Senders

Instances	Primary	Replicas	Max Senders
4	1	3	10

WAL Source (Primary/Bridge)

Instance	IP	Role	Downstreams
pg.payment-tt-0	10.189.11.21	primary	3

Instances WAL Senders WAL Receivers Cascading

Instances	WAL Senders	WAL Receivers	Cascading
4	3	2	∅

Sync Commit Mode Enabled Sync Standby Names ANY 1 (slave1, slave2)

Lag Bytes 0 B Lag Secs 147 µs NTP Offset

Leadership / Timeline

Downstream Count

Replication Metrics

Replication Slots

Monitoring dashboard showing various metrics for PostgreSQL replication, including replication slots, backlog, and timeline.

Service

PGSQL / PGSQL Service

Last 6 hours 1m

Overview

Service	Route	Session	Proxy	Admin	Route	Session	Load
pg-moderation-tt-replica	3	0	pg-moderation-tt-1	10.189.58.52:9101	3	41	0%
Cluster	6	10	pg-moderation-tt-2	10.189.58.51:9101	3	38	0%
pg-moderation-tt	3	258	pg-moderation-tt-3	10.189.58.36:9101	3	33	4%
Role	9	112					

Service +

Service	Route	Session
pg-moderation-tt-replica	3	0
pg-moderation-tt-default	3	0
pg-moderation-tt-offline	6	10
pg-moderation-tt-primary	3	258
pg-moderation-tt-replica	9	112

Cluster +

Role	Server	Route	Weight	Session
replica	pg-moderation-tt-1	3	100	0
replica	pg-moderation-tt-2	3	100	55
replica	pg-moderation-tt-3	3	100	57

Proxy +

Proxy	Server	Route	Session
pg-moderation-tt-1	pg-moderation-tt-1	Default	No Route
pg-moderation-tt-2	pg-moderation-tt-2	Candidate	Offline
pg-moderation-tt-3	pg-moderation-tt-3	Candidate	Offline

Admin

Route	Session	Load
10.189.58.52:9101	3	41
10.189.58.51:9101	3	38
10.189.58.36:9101	3	33

Server +

Server	Route	Weight	Session
pg-moderation-tt-1	3	100	0
pg-moderation-tt-2	3	100	55
pg-moderation-tt-3	3	100	57

Proxy → Server Sessions

Proxy Busy

Server Busy

Proxy Sessions

Server Sessions

Proxy Traffic

Traffic by Server

Legend for Servers:

- pg-moderation-tt-1
- pg-moderation-tt-2
- pg-moderation-tt-3

Load Balancer
Proxy Instance
Traffic Control

...

Service Registry

With

Consul



Help ▾ Settings

- Services
- Nodes
- KeyValue
- Intentions
- ACCESS CONTROLS •
 - Tokens
 - Policies
 - Roles
- Auth Methods

Nodes 416 total

Search

Search Across ▾

Health Status ▾

Unhealthy to Healthy ▾

✓ pg-moment shard10-tt-1

5 Services 10.189.4.51

✓ pg-moment shard11-tt-0

5 Services 10.189.4.56

✓ pg-moment shard11-tt-2

5 Services 10.189.55.64

✓ pg-moment shard12-tt-0

5 Services 10.189.4.57

✓ pg-moment shard12-tt-2

5 Services 10.189.31.64

✓ pg-moment shard13-tt-0

5 Services 10.189.4.58

✓ pg-moment shard13-tt-1

5 Services 10.189.4.34

✓ pg-moment shard14-tt-0

5 Services 10.189.4.33

✓ pg-moment shard14-tt-2

5 Services 10.189.42.58

✓ pg-moment shard15-tt-0

5 Services 10.189.4.48

✓ pg-moment shard15-tt-1

5 Services 10.189.4.32



Services

Nodes

KeyValue

Intentions

ACCESS CONTROLS •

- Tokens
- Policies
- Roles

Auth Methods

Services 168 total

Search

Health Status ▾ Service Type ▾

✓ redis-tm-tt

3 instances 8GB, C6, P1R2, redis, redis-tm-tt, sentinel

✓ redis-tsp-tt

3 instances 8GB, C6, P1R2, redis, redis-tsp-tt, sentinel

✓ redis-wallet-tt

3 instances 8GB, C6, P1R2, redis, redis-wallet-tt, sentinel

✓ pg-chat-tt

2 instances master, pg-chat-tt, replica

✓ pg-device-tt

4 instances master, pg-device-tt, replica

✓ pg-gift-tt

2 instances master, pg-gift-tt, replica

✓ pg-international-tt

2 instances master, pg-international-tt, replica

✓ pg-location-tt

2 instances master, pg-location-tt, replica

✓ pg-market-tt

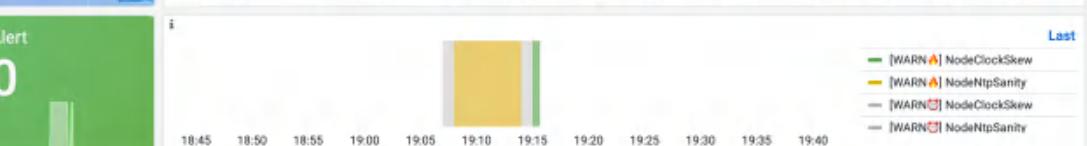
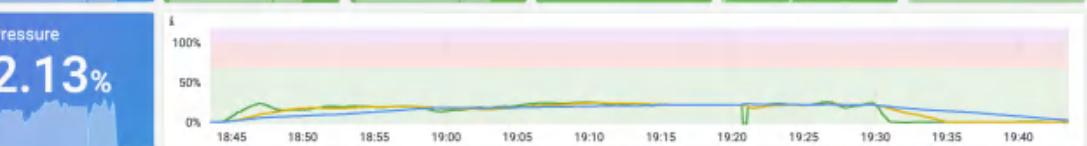
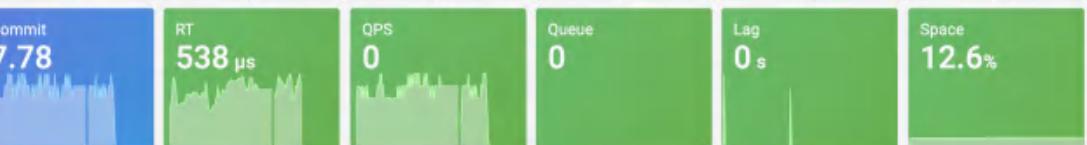
2 instances master, pg-market-tt, replica

PG Instance

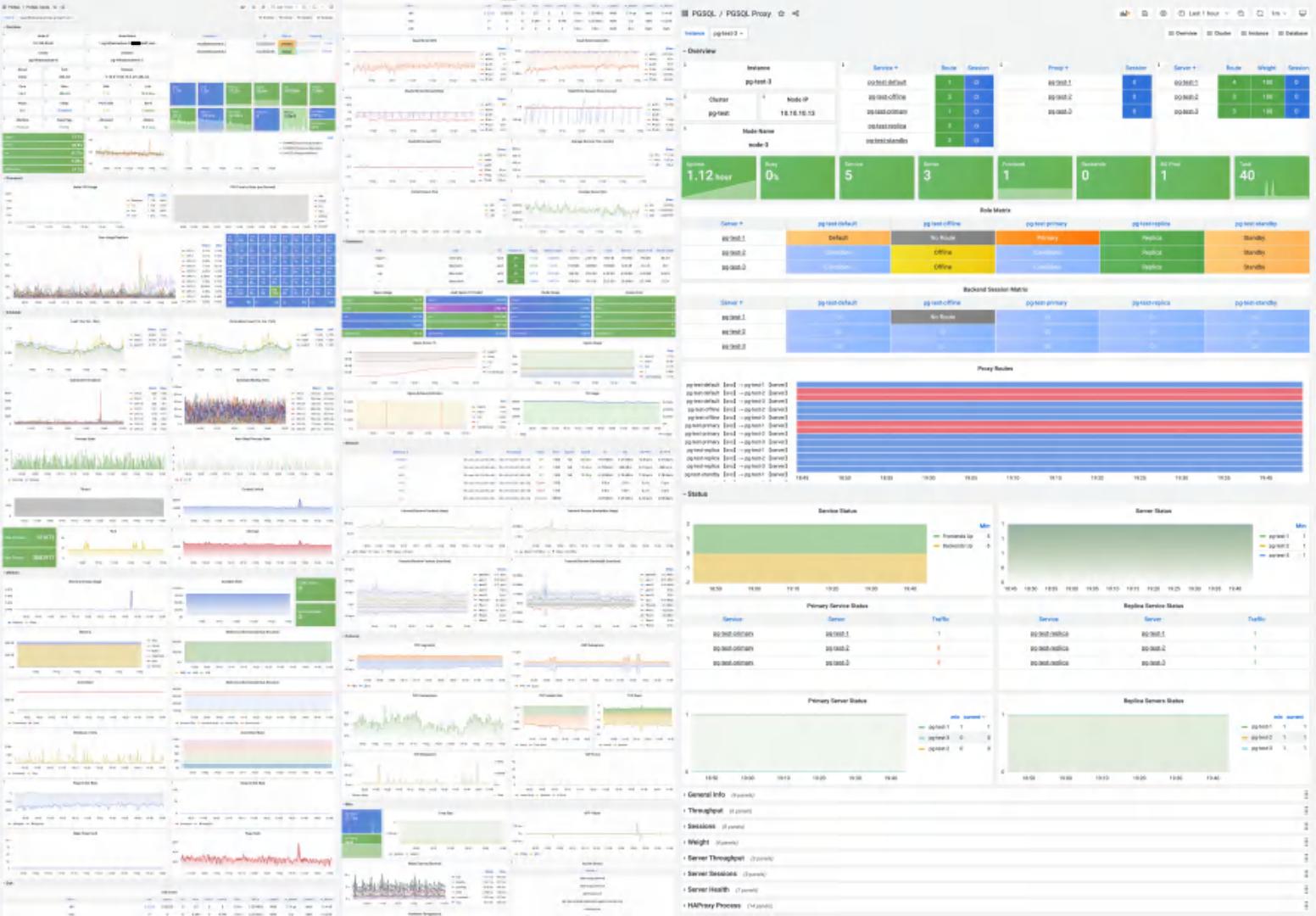
Instance	
Cluster	Instance
pg-test	pg-test-3
Status	Node IP
Follower	10.10.10.13
Node Name	
node-3	
Configuration	
Version	13.3
CLS	pg-test
CLS ID	6985471009499599062
WAL	logical
Port	5432
Data	/pg/data
Config	/pg/data/postgresql.conf

Instance	IP	Status	Pressure	Service
pg-test-1	10.10.10.11	primary	1.99%	pg-test-default
pg-test-2	10.10.10.12	replica	2.19%	pg-test-offline
pg-test-3	10.10.10.13	replica	2.13%	pg-test-primary
				pg-test-replica
				pg-test-standby

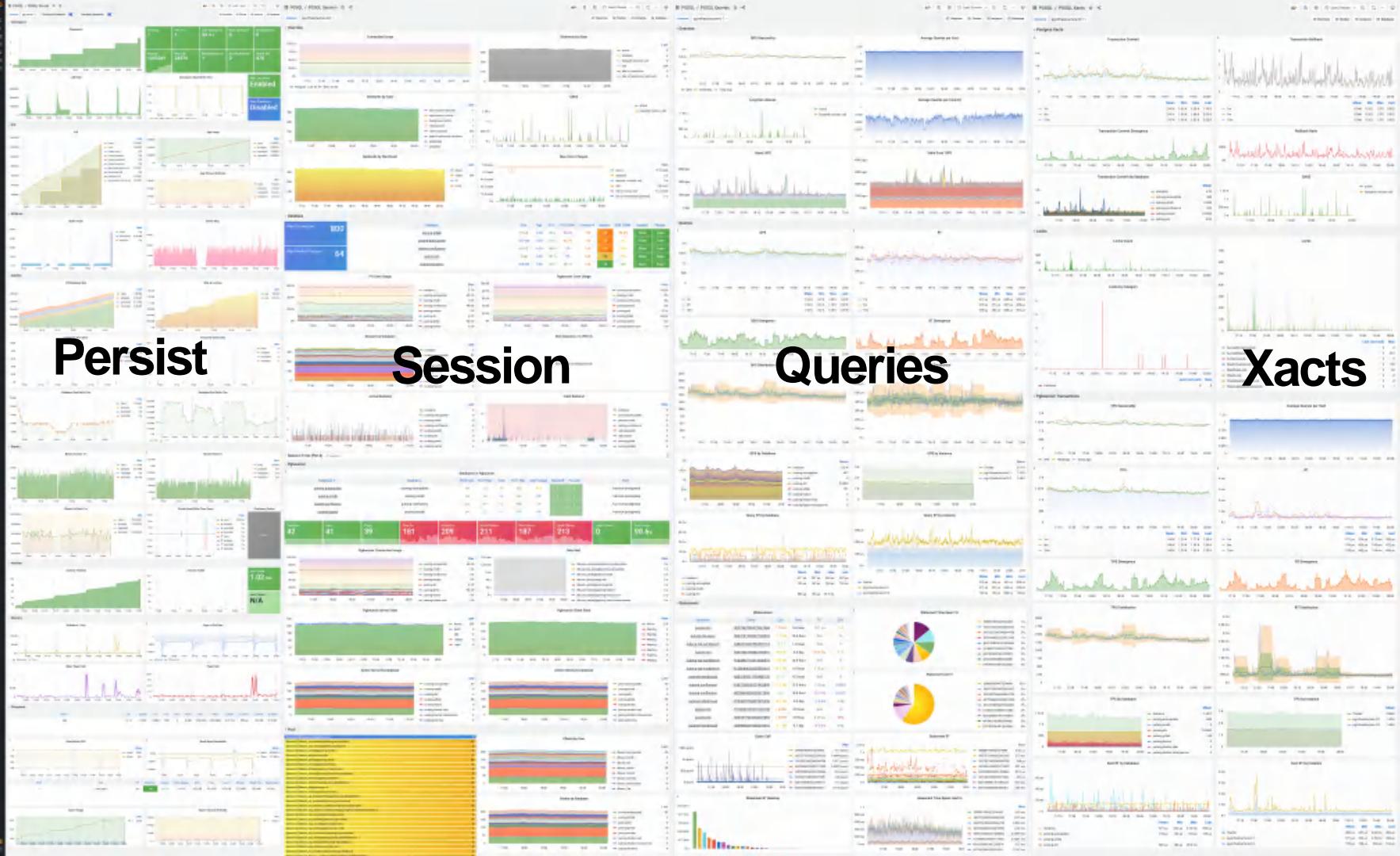
Database	Size	Age	TPS	PG CONN	ConnLimit	Session	PGB CONN	Disabled	Paused
test	188 MB	0.0%	0.7	0%	NO LIMIT	1	1%	False	False
postgres	9 MB	0.0%	7.1		NO LIMIT	1	2%	N/A	N/A
template1	9 MB	0.0%	0		NO LIMIT		0%	N/A	N/A
template0	8 MB	0.0%	0		NO LIMIT		0%	N/A	N/A



Node Proxy



Persist Session Queries Xacts



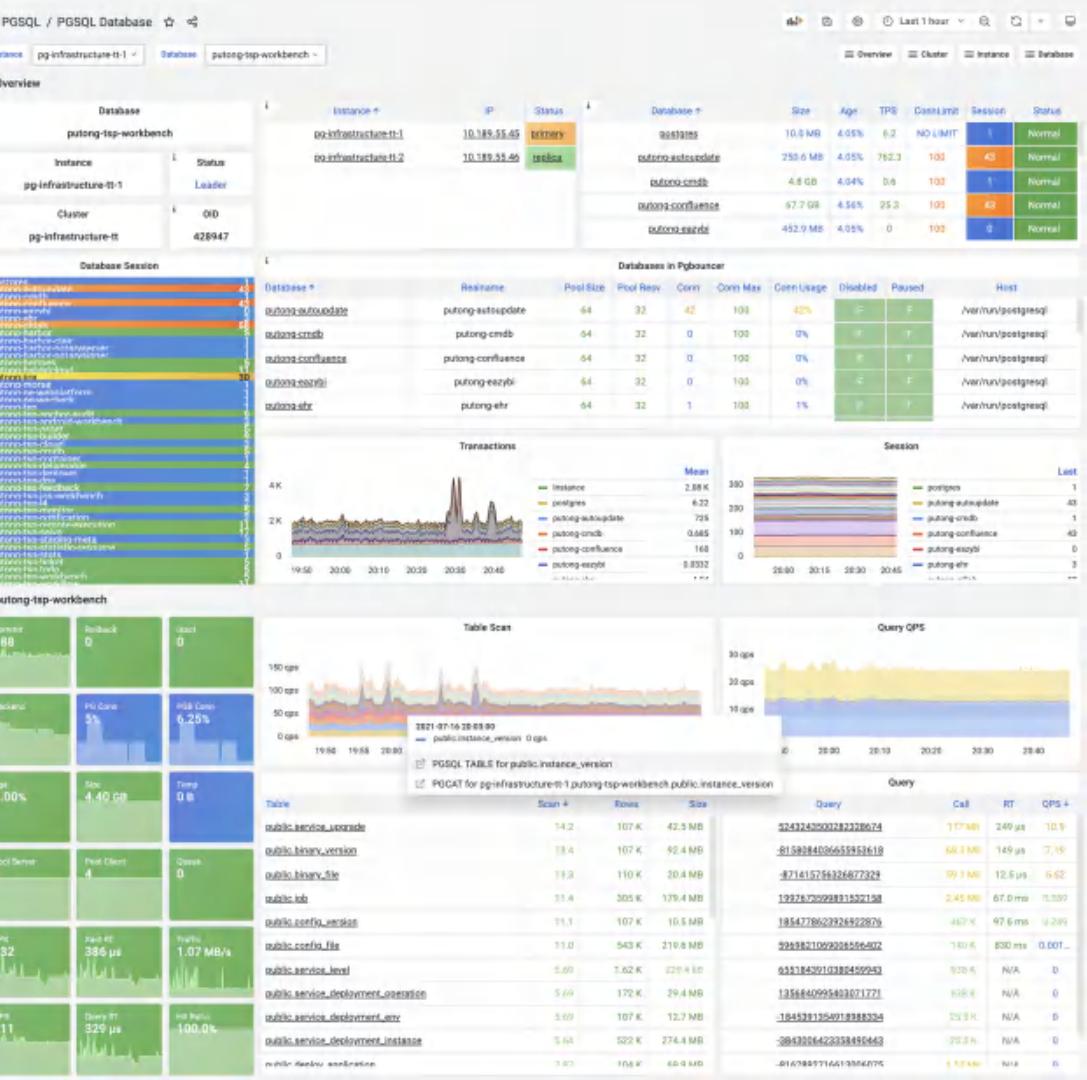
PGSQL Database

数据库级别指标

快速访问

每一张表

每一个查询



PGSQL Table

PGCAT / PGCAT Table

Data Source: pg-table-1-mta - Relayer (pgcgt.cluster)

- Catalog

Database	Scheme	Relation	Size	Rows	Page	Age	Live Tuple	Dead Tuple	Dead Ratio	Bloat Ratio	Retention Kind
meta	piggy	cluster	0 B	0	0	1775149	0	0	0%		(j) Table
Identifier	Value	Feature	Value	Field	Value	Persist Feature	Value				
Relation OID	28429	Has Index	TRUE	Rel Pages	0	Unlogged	FALSE				
Schema OID	2872	Is Partitioned	FALSE	Rel Tuples	0	Temporary	FALSE				
Reltype OID	28425	Is Shared	FALSE	Rel All Visible	0	AD					
Comptype OID		Has Rule	FALSE	Columns	6	postgres-readonly/postgres					
Owner OID	10	Has Trigger	TRUE	Checks	1	dbrole,readwriter-read/postgres					
Access Method ID	2	Has Subclass	FALSE	Frontend ID	716	dbrole,readwriter-read/postgres					
Rel Filenode	28429	Row level security	FALSE	Age	1775126	dbrole,offlineship/postgres					
Tablespace OID		Force RLS	FALSE	Min XID	1	Options					
Toast Table OID	28430	Is Populated	TRUE			No data					
Tuple	Value	Vacuum	Value	Analyze	Value	Blocks I/O	Value				
Seq Scan	1	Vacuums	0	Analyze	0	Heap Blocks Read	0				
Seq Scan Tuples	0	Vacuum Count	0	Analyze Count	0	Heap Blocks Hit	0				
Index Scan	0	Last Vacuum		Last Analyze		Index Blocks Read	0				
Index Scan Tuples	0	Since Last Vacuum		Last Analyzer Elapsed		Index Blocks Hit	0				
Inserted Tuples	0	Auto/Vacuum Count	0	AutoAnalyzer Count	0	Toast Blocks Read	0				
Deleted Tuples	0	Last AutoVacuum		Last AutoAnalyze Elapsed		Toast Blocks Hit	0				
Updated Tuples	0	Since Last AutoVacuum		Modify Since Analyze	0	Toast Index Blocks Read	0				
Not Updated Tuples	0	Inserts since Last Vacuum	0			Toast Index Blocks Hit	0				
Columns											
No	Drop	Name	Type	NOT NULL	Default	TypeID	Storage	Align	Len	Dir	
1	OK	cis	text	NOT NULL		25	Extended	L4	8	N	
2	OK	ckard	text	MULTIPLE		25	Extended	L4	8	N	
3	OK	index	integer	MULTIPLE		22	Plain	L4	4	N	
4	OK	status	pgsql_status	NOT NULL	'Unknown':pgsql status	20374	Plain	D8	4	N	
5	OK	cime	timestamp with time zone	NOT NULL	now()	1164	Plain	D8	8	N	
6	OK	rtime	timestamp with time zone	NOT NULL	now()	1164	Plain	D8	8	N	
Indexes											
Schema	Index	Define				Index ID	Table ID	Scan	Tap Read	Tap Fetch	
piggy	cluster_pkkey					20433	20433	0	0	0	
CREATE UNIQUE INDEX cluster_pkkey ON piggy.cluster USING btree (cis)											



PGSQL Query

The image shows a screenshot of the pgCAT interface, which is a monitoring tool for PostgreSQL databases. At the top, there's a header with the title 'PGCAT / PGCAT Query' and a status bar showing 'Last 1 hour' and various metrics. The main area has tabs for Overview, Cluster, Instance, and Database.

Overview Tab: Shows database statistics like CPU, memory, and disk usage, along with a pie chart of 'Database Time Spent (%)' for different operations.

Cluster Tab: Displays the 'Cluster' tab with a table for 'pg-msft' showing instance details (pg-ms-ft-1), database (putong-sms), and query identifier (-328823109287227845).

Instance Tab: Shows detailed performance metrics for the current query (-328823109287227845). It includes tables for 'Users' (dbuser_monitor), 'Exec Time' (plans, calls, rows, shared_blocks_read, shared_blocks_written, local_blocks_hit, local_blocks_read, local_blocks_written, temp_blocks_read, temp_blocks_written, blk_read_time, blk_write_time, wal_records, wal_fpi, wal_bytes), and 'QPS' and 'RT' charts.

Database Tab: Contains a table for 'Statements Stat for -328823109287227845' showing statistics for user dbuser_monitor across various metrics.

Monitoring Dashboard: On the right side, there's a large dashboard with numerous charts and graphs displaying real-time system health, including CPU usage, memory pressure, and network activity over time.

PGCAT

自动为新数据库注册数据源

通过Grafana直接浏览系统目录

带来无限新可能性

The screenshot shows the PGCAT configuration interface. At the top, there's a navigation bar with icons for gear (Configuration), magnifying glass (Search), plus sign (Add), user (Users), team (Teams), plugin (Plugins), preferences (Preferences), and API keys. The main title is "Configuration" under "Organization: Pigsty". Below the title, there are tabs for "Data sources", "Users", "Teams", "Plugins", "Preferences", and "API keys". A search bar at the top right says "Search by name or type" and a blue button says "Add data source". The "Data sources" section lists the following items:

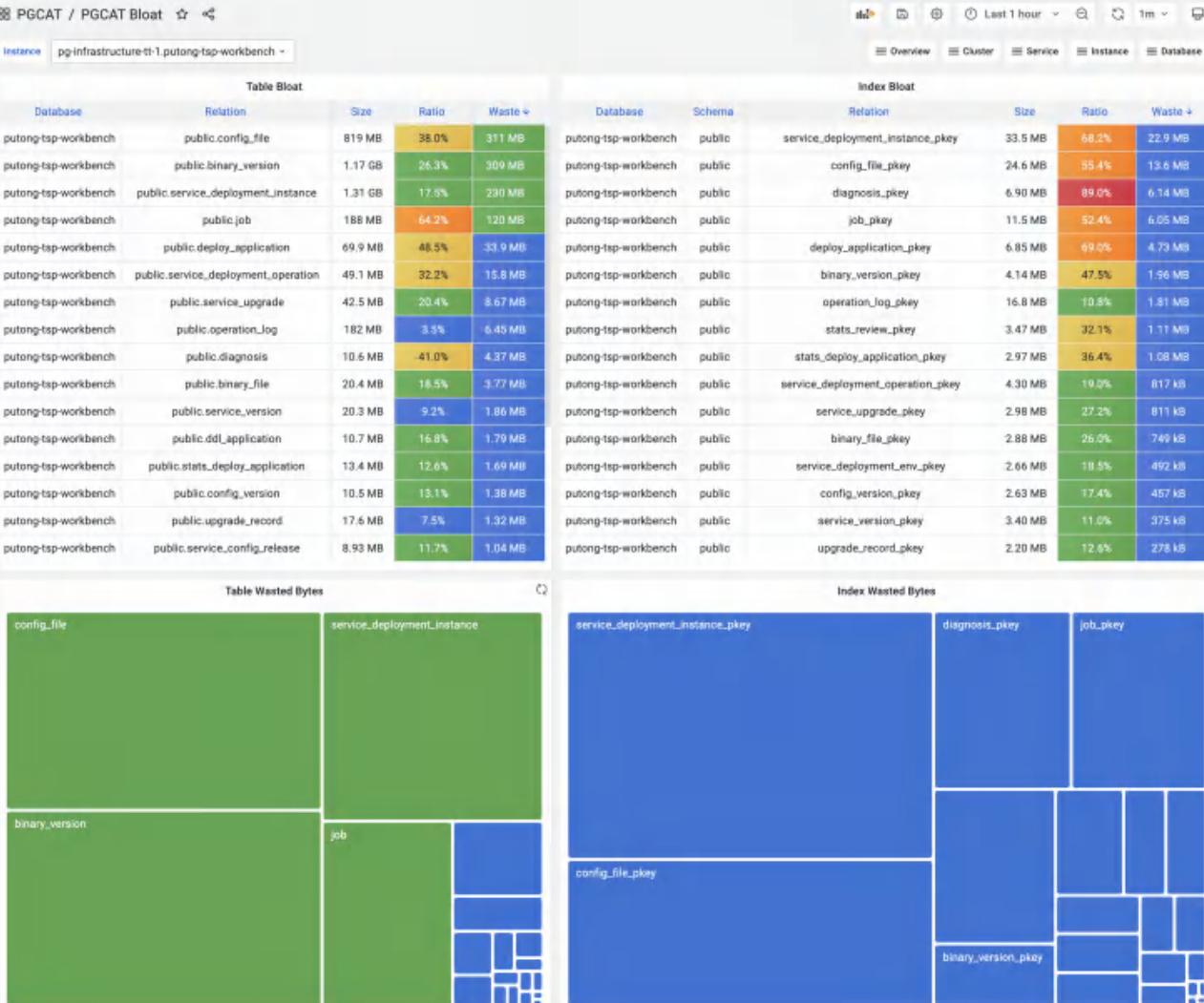
- Loki | http://127.0.0.1:3100
- Meta | PostgreSQL | 127.0.0.1:5432
- pg-meta-1.meta | PostgreSQL | 10.10.10.10:5432
- pg-test-1.test | PostgreSQL | 10.10.10.11:5432
- pg-test-2.test | PostgreSQL | 10.10.10.12:5432
- pg-test-3.test | PostgreSQL | 10.10.10.13:5432
- Prometheus | http://127.0.0.1:9090 | default
- Prometheus-PgSQL | Prometheus | http://127.0.0.1:9090

At the bottom, there are links for Documentation, Support, Community, Open Source, and version v0.0.5 (cbf12aa5001).

PGCAT Bloat

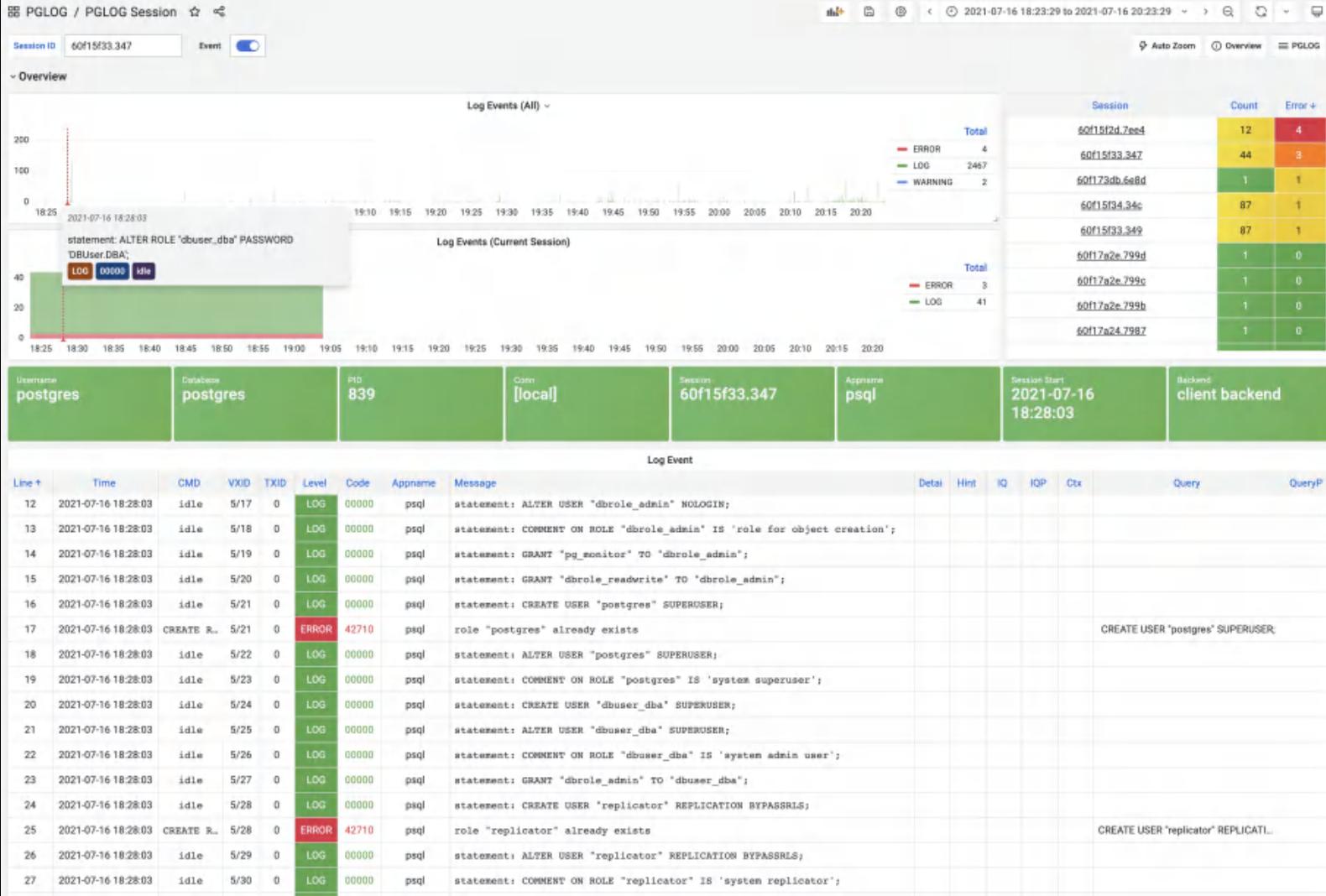
表与索引膨胀总览

从Catalog获取信息



PGLOG Session

- 单条连接详情
- Annotation



Log Events

Summary based on pgbadger



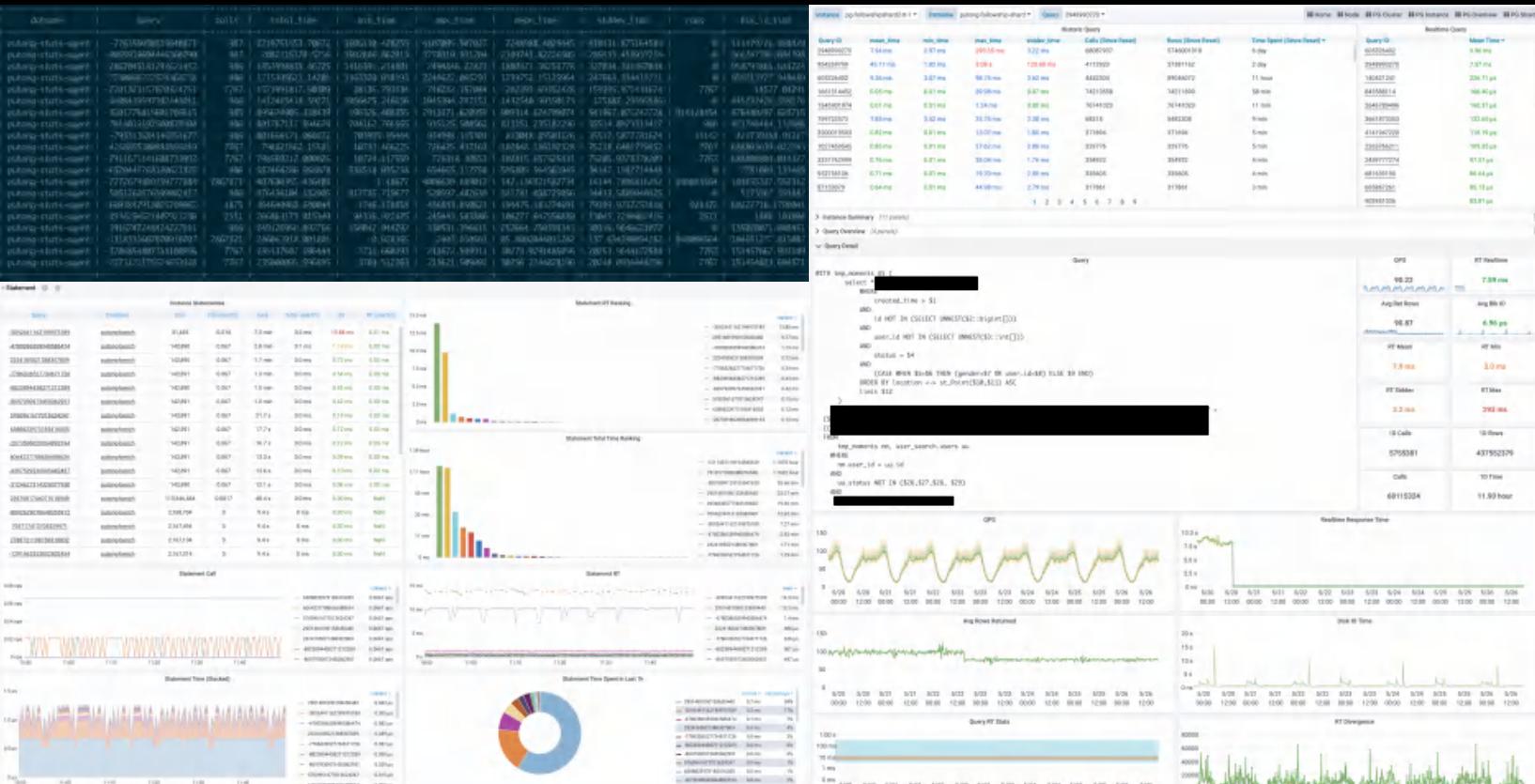
慢速Query TOP (N)

序号	耗时	Query文本信息
1	19ms	query Text: SELECT id, user_id, ab_group_id, status::text, to_char(created_time + 'MF_J6f6tYd')::text, 'YYYY-MM-DD HH24:MI:SS.MS' AS created_time, to_char(updated_time + 'MF_J6f6tYd')::interval, 'YYYY-MM-DD HH24:MI:SS.MS' AS updated_time, is_guest::text FROM public.ab_users WHERE l = 1 AND (id >= 224566432 <= 6487284432) Seq Scan ON public.ab_users COST = 0.0..19404505.35 ROWS = 364556341 width = (48) = 0.137... rows=44329 Rows = 372514945 loops = 1 Output: id, user_id, ab_group_id, (status)::text, to_char((created_time + '1muufxa44y')::interval), 'YYYY-MM-DD HH24:MI:SS.MS'::text, to_char((updated_time + '1muufxa44y')::interval), 'YYYY-MM-DD HH24:MI:SS.MS'::text, is_guest::text Filter: ((ab_users.id >= 'Gu5SDGQCVa'::bigint) AND (ab_users.id <= 'P1DvCoiml'::bigint)) ROWS Removed BY Filter: 7158213; Sets: 2020-10-14 01:16 - Database: [REDACTED] - User: [REDACTED] - Remote: 10. [REDACTED] - Application: Postgres Params: 2020-10-14 01:16:16 - Database: [REDACTED] - User: [REDACTED] - Remote: 10. [REDACTED] - Application: Postgres

Observability

Original Statements Metrics

Query Statements Dashboards



Postgres Observability



PG 14

new metrics are

READY!

v Session Time (PG14)

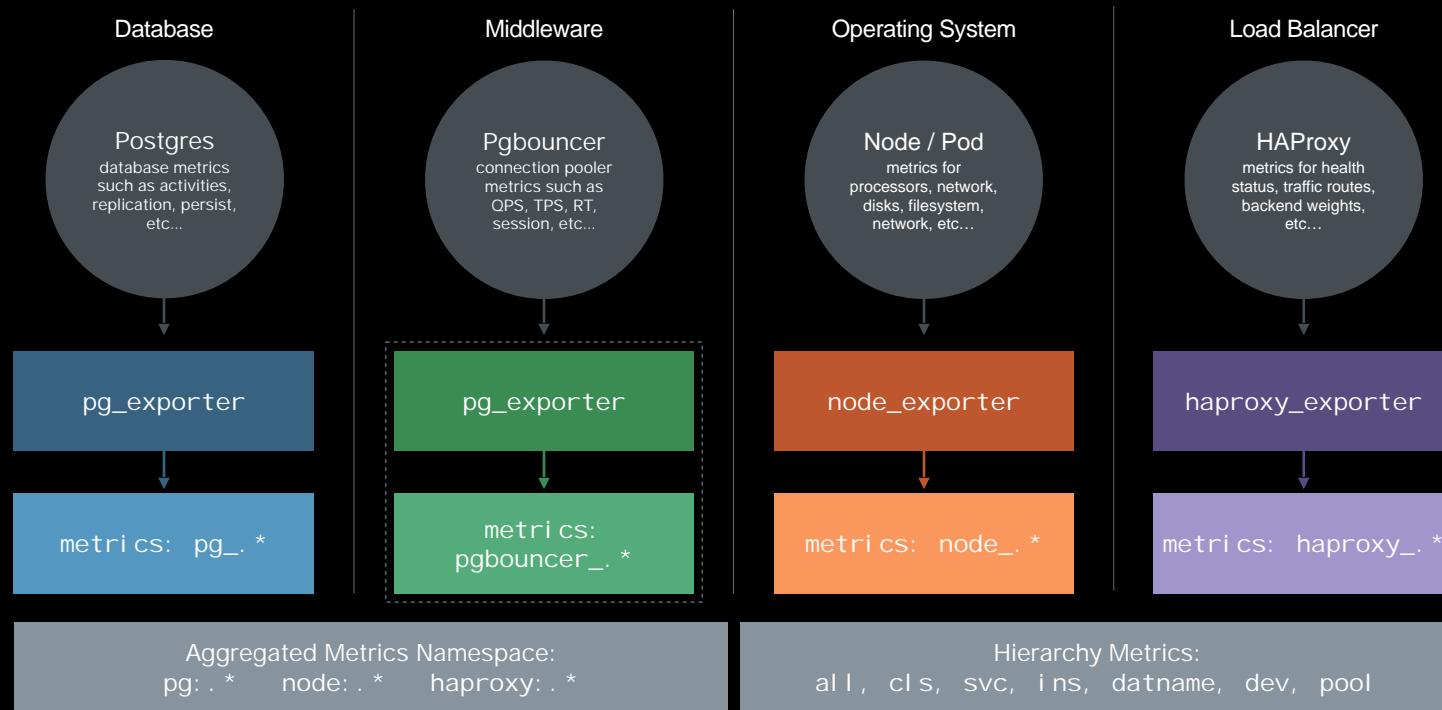
New Sessions 1m (PG14)

Sessions Failure 1m (PG14)



Metrics Source

metrics data flow



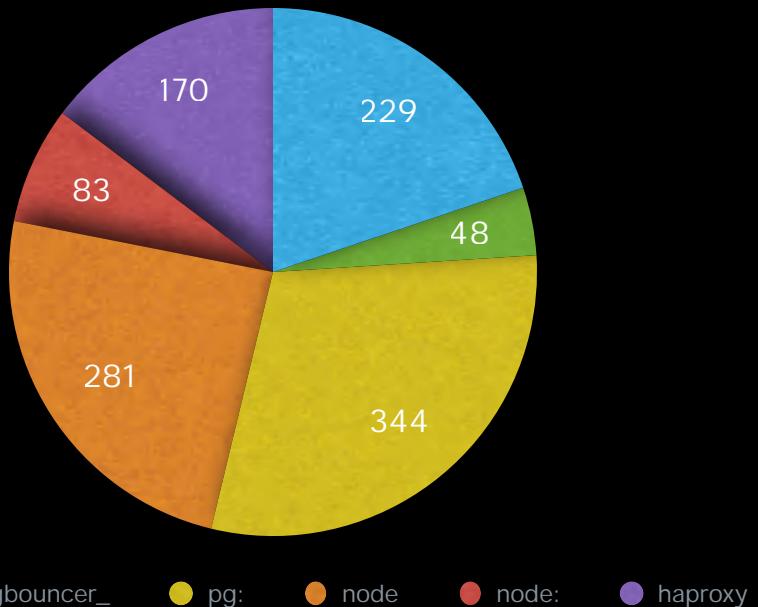
Metrics Ratio

metric type for a cluster

1155 metrics per env

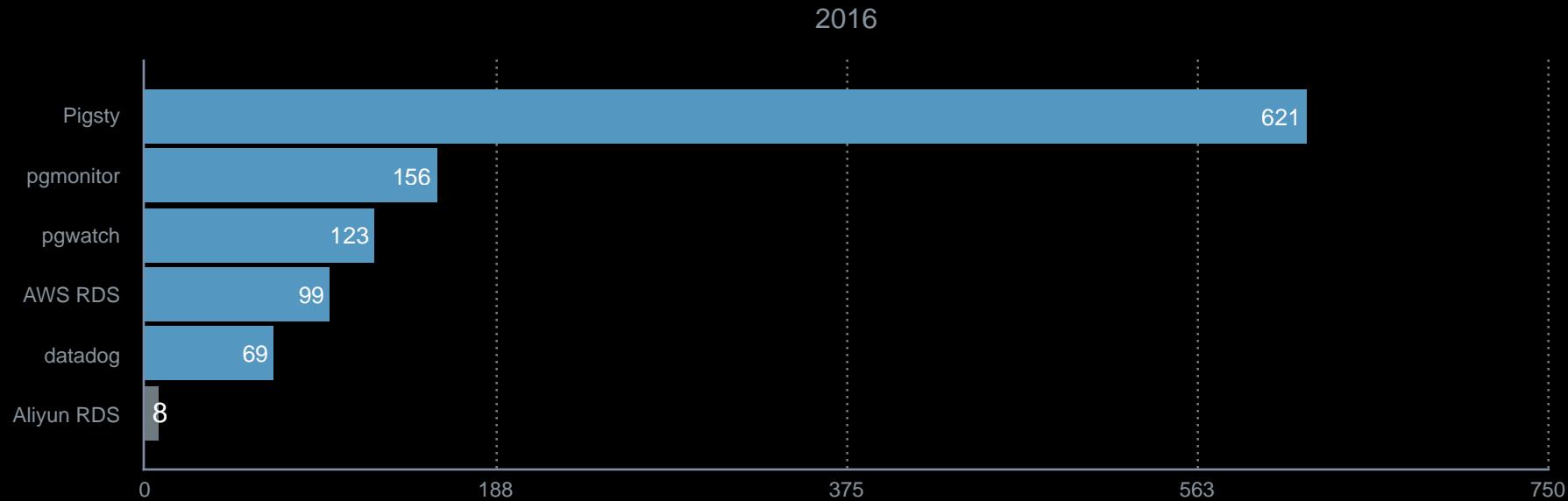
3k~10k+ time-series per instance

metrics type count



Metrics Coverage

metrics **numbers** comparing to others



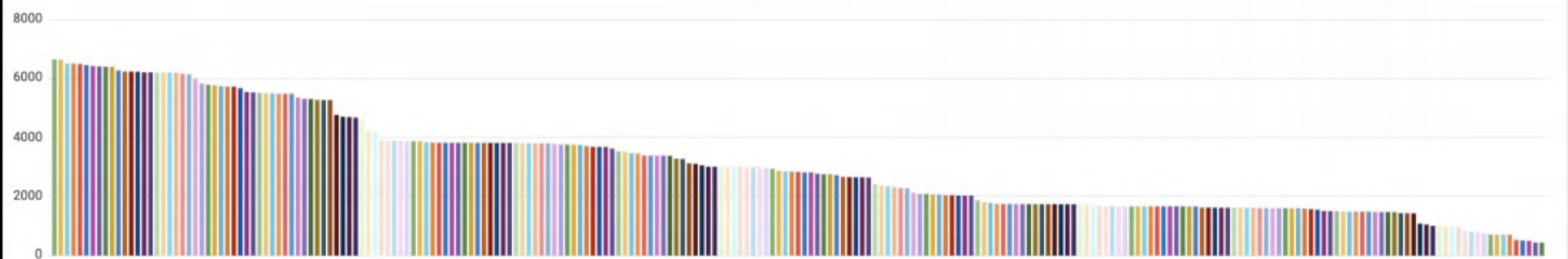
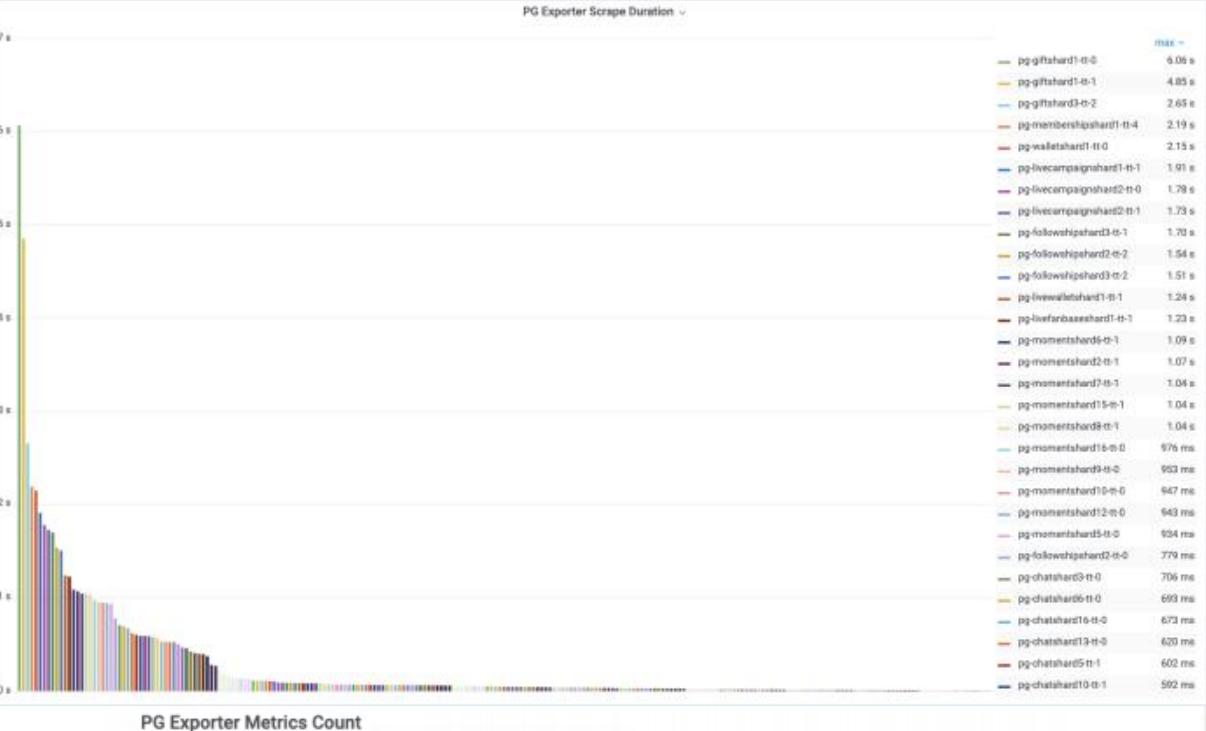
metrics related to database are counted (i.e: node metrics are excluded)

Overhead

Typical Scrape Duration: 10~100ms

Typical Scrape Count: 1500 ~ 4500

Typical Scrape Interval: 15s



Provisioning



Clearly, we must break away from the sequential and not limit the computers. We must state definitions and provide for priorities and descriptions of data. We must state relationships, not procedures.

—Grace Murray Hopper, *Management and the Computer of the Future* (1962)

ONE **Node** TO RULE
THEM ALL

Interface

Clearly, we must break away from the sequential and not limit the computers. We must state definitions and provide for priorities and descriptions of data. We must state relationships, not procedures.

—Grace Murray Hopper, *Management and the Computer of the Future* (1962)

我们必须跳出电脑指令序列的窠臼。叙述定义、描述元数据、梳理关系，而不是编写过程。

—— Grace Murray Hopper, 未来的计算机及其管理 (1962)

管控

数据库是管理数据的软件， 管控系统是管理数据库的软件

DBMS， 包括数据库内核， 和配套软件。

云厂商卖开源PG， 卖的不是内核， 而是 数据库管控能力

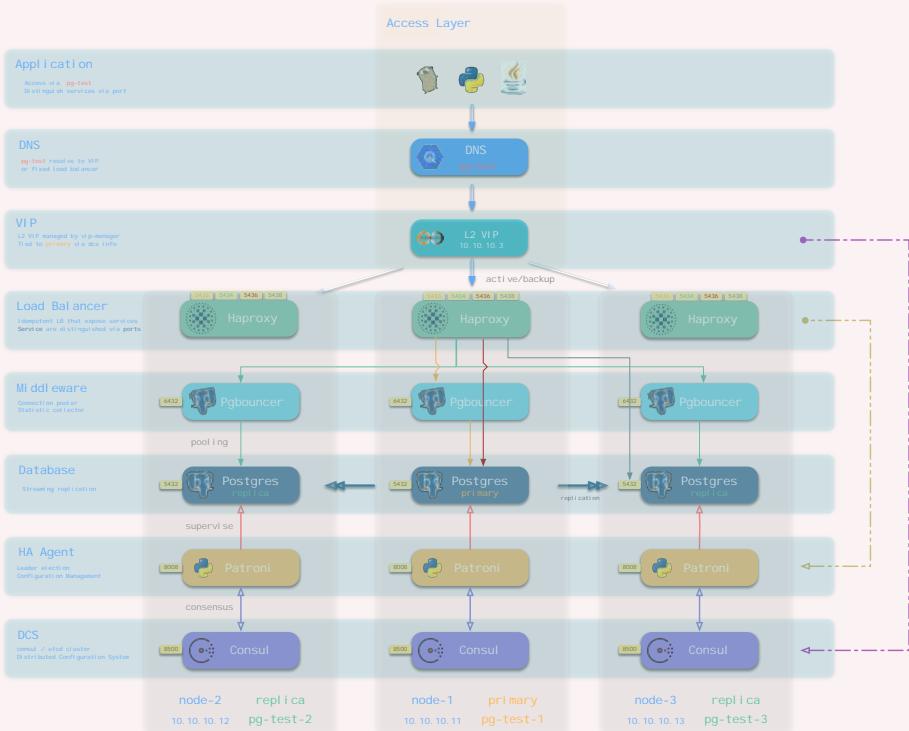
Provisioning

定义一套集群，只要几行代码

```
pg-test: # 定义一个新的 pg-test 集群
hosts: # 包含有一主两从三个实例，部署在三个节点上
  10.10.10.11: { pg_seq: 1, pg_role: primary }
  10.10.10.12: { pg_seq: 2, pg_role: replica }
  10.10.10.13: { pg_seq: 3, pg_role: replica }
vars: # 集群名称，使用的VIP地址，业务用户与数据库
  pg_cluster: pg-test
  vip_address: 10.10.10.2
  pg_users: [ { name: test } ]
  pg_databases: [ { name: test } ]
```

从配置到实现，只要一行命令

```
pgsql.yml -l pg-test # 创建出该集群
```



traffic primary replica default offline supervis e consensus inquiry check replication

primary	postgres://test@pg-test:5432/test	non-interactive read-write access via a pool
replica	postgres://test@pg-test:5434/test	non-interactive read-only access via a pool
default	postgres://dbuser_db@pg-test:5436/testdb	interactive primary direct access (DDL, DML, Admin)
offline	postgres://dbuser_stats@pg-test:5438/testdb	interactive offline direct access (ETL/SAGA/Personal)

Interface

Idempotent Ansible Playbooks

infra.yml	# 部署/修改基础设施
pgsql.yml	# 部署/修改数据库集群/实例
pgsql-remove.yml	# 下线数据库集群/实例
pgsql-createuser.yml	# 创建/修改业务角色&用户
pgsql-createdb.yml	# 创建/修改业务数据库

./infra.yml --tags=environ	# 重新在管理节点上配置环境
./infra.yml --tags=repo -e repo_rebuild=true	# 强制重新创建本地源
./infra.yml --tags=repo_upstream	# 加入上游YumRepo
./infra.yml --tags=prometheus	# 重新创建Prometheus
./infra.yml --tags=nginx_config,nginx_restart	# 重新生成Nginx配置文件
# 基础设施初始化	
./pgsql.yml --tags=infra	# 完成基础设施的初始化，包括机器节点初始化与DCS部署
./pgsql.yml --tags=node	# 完成机器节点的初始化，通常不会影响运行中数据库实例
./pgsql.yml --tags=dcs	# 完成DCS: consul/etcfd的初始化
./pgsql.yml --tags=dcs -e dcs_exists_action	# 完成consul/etcfd的初始化，强制抹除
# 数据库初始化	
./pgsql.yml --tags=pgsql	# 完成数据库部署：数据库、监控、服务
./pgsql.yml --tags=postgres	# 完成数据库部署
./pgsql.yml --tags=monitor	# 完成监控的部署
./pgsql.yml --tags=service	# 完成负载均衡的部署，(Haproxy & VIP)
./pgsql.yml --tags=register	# 将服务注册至基础设施

```
#-----  
- name: Infra Init      # init infra on common database node  
  become: yes  
  hosts: all  
  gather_facts: no  
  tags: infra  
  roles:  
    - role: node        # init common database node  
      tags: node  
  
    - role: consul      # init dcs:consul clients  
      tags: [ dcs , consul ]  
  
#-----  
- name: Psql Init       # init postgres cluster/instance  
  become: yes  
  hosts: all  
  gather_facts: no  
  tags: postgresql  
  roles:  
    - role: postgres    # init postgres pgbounce patroni  
      tags: postgres  
  
    - role: monitor     # init monitor exporters  
      tags: monitor  
  
    - role: service     # init service , lb , vip  
      tags: service  
  
    - role: register    # register cluster/instance to infra  
      tags: register
```

```
[05-23 16:53:03 CST] 24180 vonng@vonng-mac:~/pigsty master [+!?]
```

```
$ ./pigsty
```

NAME

pigsty -- Pigsty Command-Line Interface v0.9

SYNOPSIS

```
meta          setup meta nodes
node          setup database nodes
pgsql         setup postgres clusters
infra         setup infrastructure
clean         clean pgsql clusters
config        mange pigsty config file
serve         run pigsty API server
demo          setup local demo
log           watch system log
pg            pg operational tasks

init          fetch|repo|cache|enable
tunel         dcs|remove|ping|bash|ssh|admin
node          dcs|postgres|template|business|monitor|service
dns           alldns|prometheus|grafana|loki|haproxy|target
service       allservice|monitor|postgres|dcs
edit          edit|info|dump|path
start         init|start|stop|restart|reload|status
up            init|up|new|clean|start|dns
query         query|postgres|patroni|pgbouncer|message
user          db|svc|hba|log|psql|deploy|backup|restore|vacuum|
```

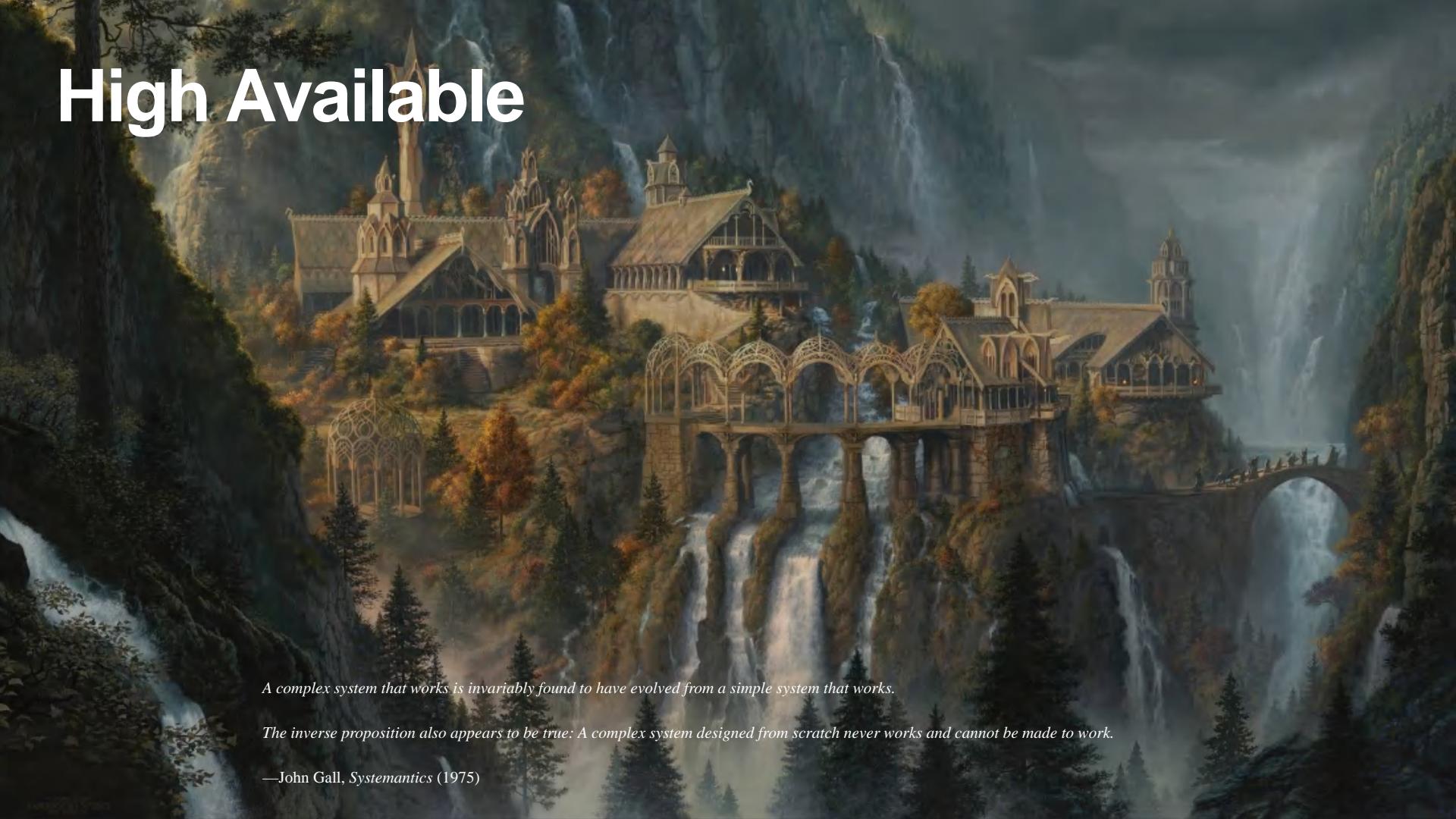
EXAMPLES

1. infra summary
pigsty infra
2. pgsql clusters summary
pigsty pgsql
3. pigsty nodes summary
pigsty node
4. create pgsql cluster 'pg-test'
pigsty pgsql init -l pg-test
5. add new instance 10.10.10.13 of cluster 'pg-test'
pigsty pgsql init -l 10.10.10.13
6. remove cluster 'pg-test'
pigsty clean -l pg-test
7. create user dbuser_vonng on cluster 'pg-test'
pigsty pg user dbuser_vonng -l pg-test
8. create database test2 on cluster 'pg-test'
pigsty pg db test -l pg-test

CLI vs GUI

The screenshot shows the Pigsty web interface with two main sections: 'pg-meta' and 'pg-test'.
pg-meta Cluster:
- Instances: One instance named 'primary' (status: up, IP: 0.0.0.10.10.10, pg_online: true).
- Users: Three users: 'dbuser_root' (default production read-write user), 'dbuser_grafana' (default readonly access for grafana datasource), and 'dbuser_pgsty' (pigsty admin user).
- Databases: One database named 'meta' (pigsty meta database).
pg-test Cluster:
- Instances: Three instances named 'primary' (IP: 10.10.10.11), 'replica' (IP: 10.10.10.12), and 'offsite' (IP: 0.0.0.10.10.10.13).
- Users: Two users: 'test' (default admin user for test database) and 'dbuser_test' (default test user for production usage).
- Databases: One database named 'test'.

High Available



A complex system that works is invariably found to have evolved from a simple system that works.

The inverse proposition also appears to be true: A complex system designed from scratch never works and cannot be made to work.

—John Gall, *Systemantics* (1975)

Architecture

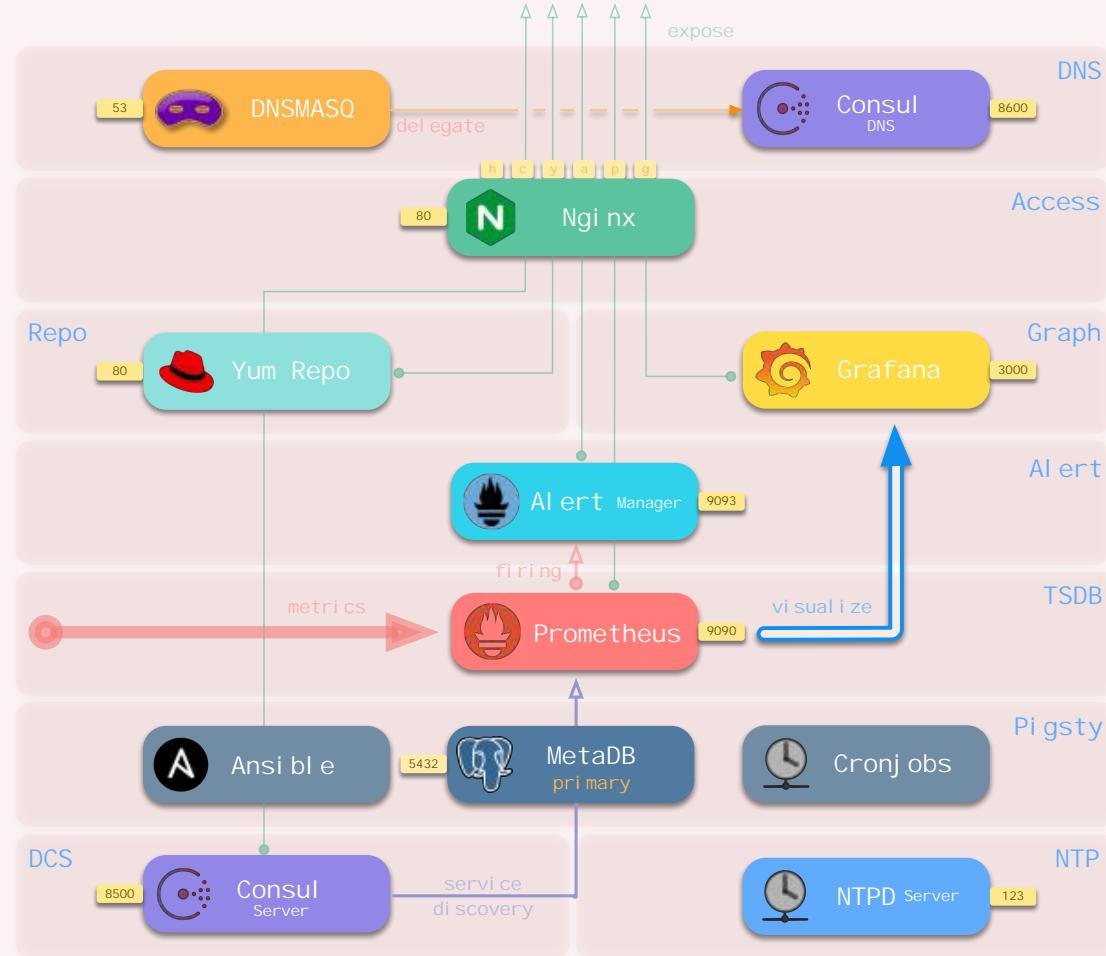
A complex system that works is invariably found to have evolved from a simple system that works.

The inverse proposition also appears to be true: A complex system designed from scratch never works and cannot be made to work.

—John Gall, *Systemantics* (1975)

Infrastructure

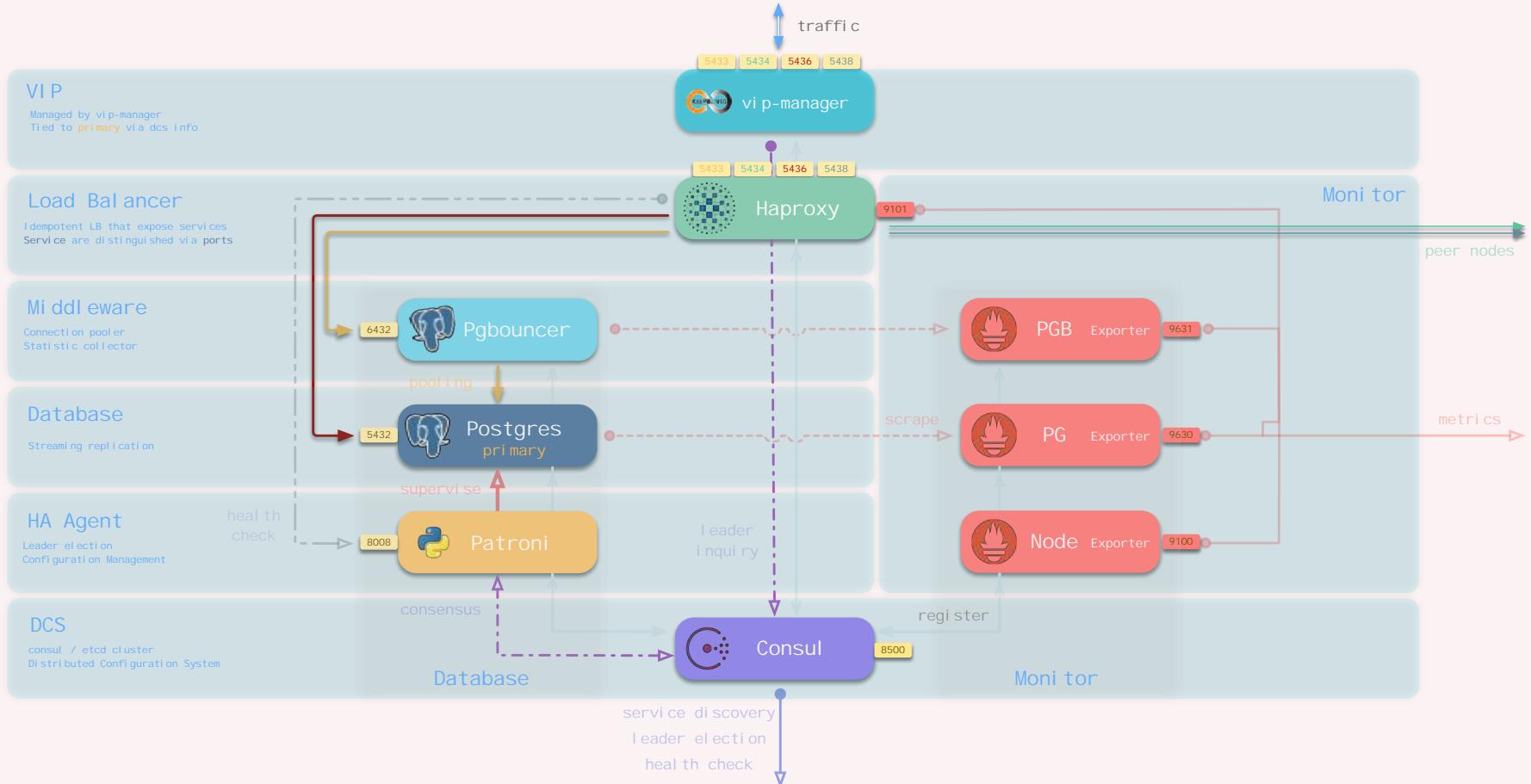
组件	端口	默认域名	说明
Grafana	3000	g.pigsty	Pigsty监控系统图形界面
Prometheus	9090	p.pigsty	监控时序数据库
AlertManager	9093	a.pigsty	报警聚合管理组件
Consul	8500	c.pigsty	分布式配置管理，服务发现
Consul DNS	8600	-	Consul提供的DNS服务
Nginx	80	pigsty	所有服务的入口代理
Yum Repo	80	yum.pigsty	本地Yum源
Haproxy Index	80	h.pigsty	所有Haproxy管理界面的访问代理
NTP	123	n.pigsty	环境统一使用的NTP时间服务器
Dnsmasq	53	-	环境统一使用的DNS域名解析服务器
Loki	3100	-	实时日志收集基础设施（选装）



Versatile Application Runtime

Database Node

以最小复杂度实现必须功能

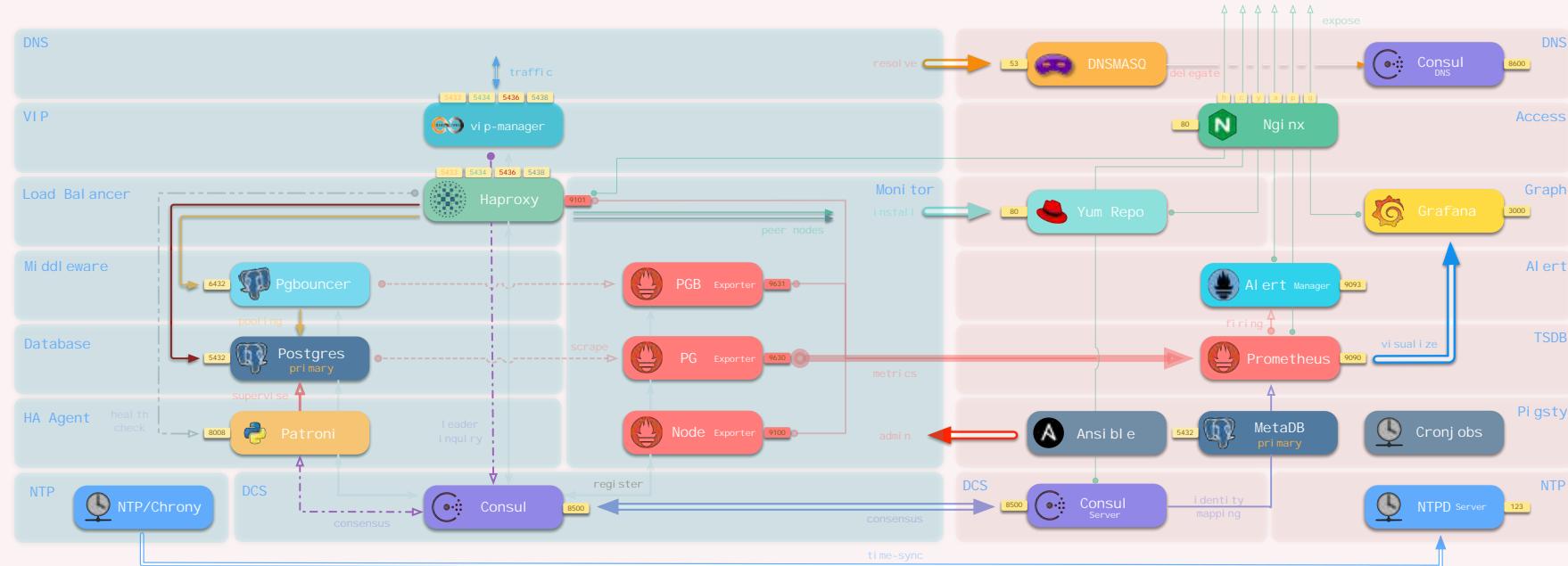


Infrastructure & Node

数据库集群与基础设施松耦合，通过注册的方式关联
基础设施升级&故障不会影响数据库集群

(除高可用使用的DCS外，维护模式下数据库集群亦免受DCS影响)

组件	端口	说明
Postgres	5432	Postgres数据库服务
PgBouncer	5432	PgBouncer连接池服务
Patroni	8008	Patroni高可用组件
Cronal	8500	分布式配置管理，服务器发现和Consul的本地Agent
Haproxy Primary	5433	集群调度服务（主库连接池）代理
Haproxy Replica	5434	集群只读服务（从库连接池）代理
Haproxy Default	5435	集群主从直连服务（用于管理、DDL/DML变更）
Haproxy Offline	5438	集群断线恢复服务（直连单线实例，用于ETL、交互式查询）
Haproxy service	5439	集群提供的外部新定义服务将被依次分配端口
Haproxy Admin	8101	Haproxy管理指标与流量管理页面
PG Exporter	9630	Postgres监控指标导出器
PGBouncer Exporter	9631	PgBouncer监控指标导出器
Node Exporter	9100	机器节点监控指标导出器
promtail	9080	实时收集Postgres, PgBouncer, Patroni日志（选装）
Consul DNS	36000	Consul提供的DNS服务
ip-manager	8	将VIP绑定至集群主库上



Database Cluster

基于DCS

Fencing & Softdog

自动故障切换

流量路由自动修复

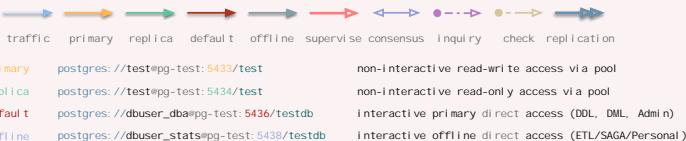
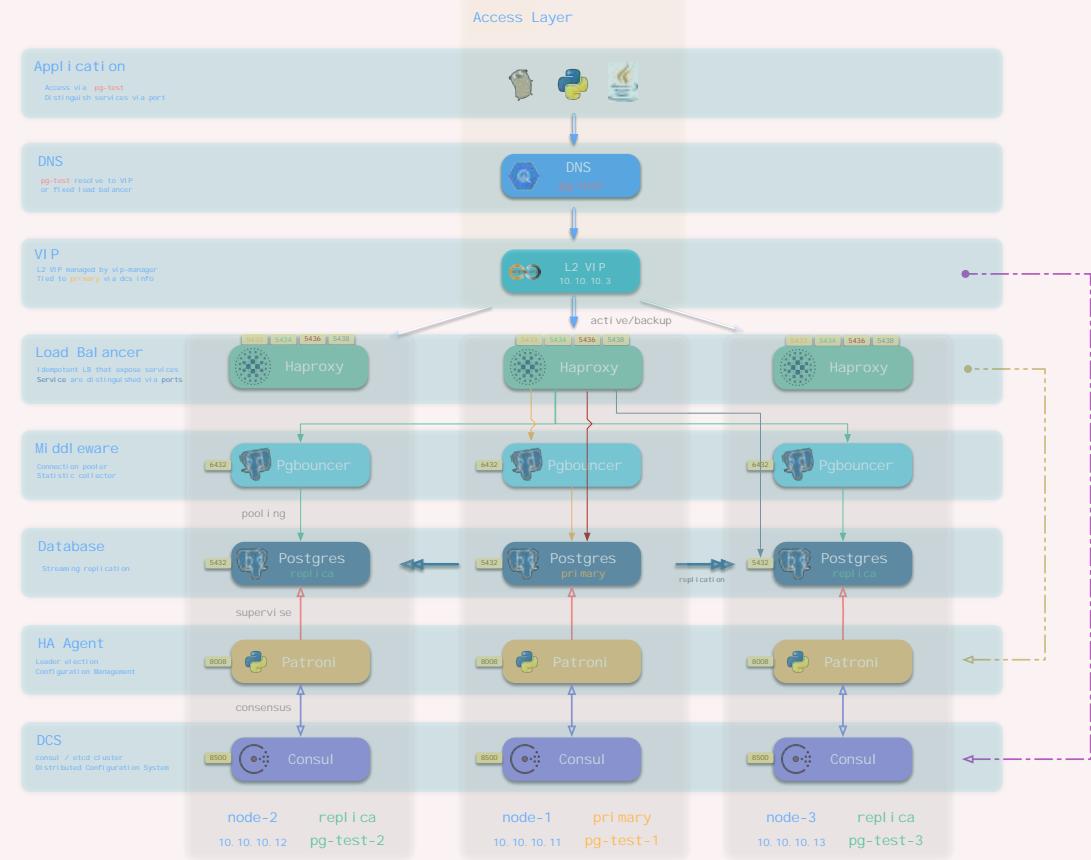
实例提供幂等的服务

集群可自动从常见故障中恢复，
实例成员对外表现等价，
只要仍有任意实例存活。

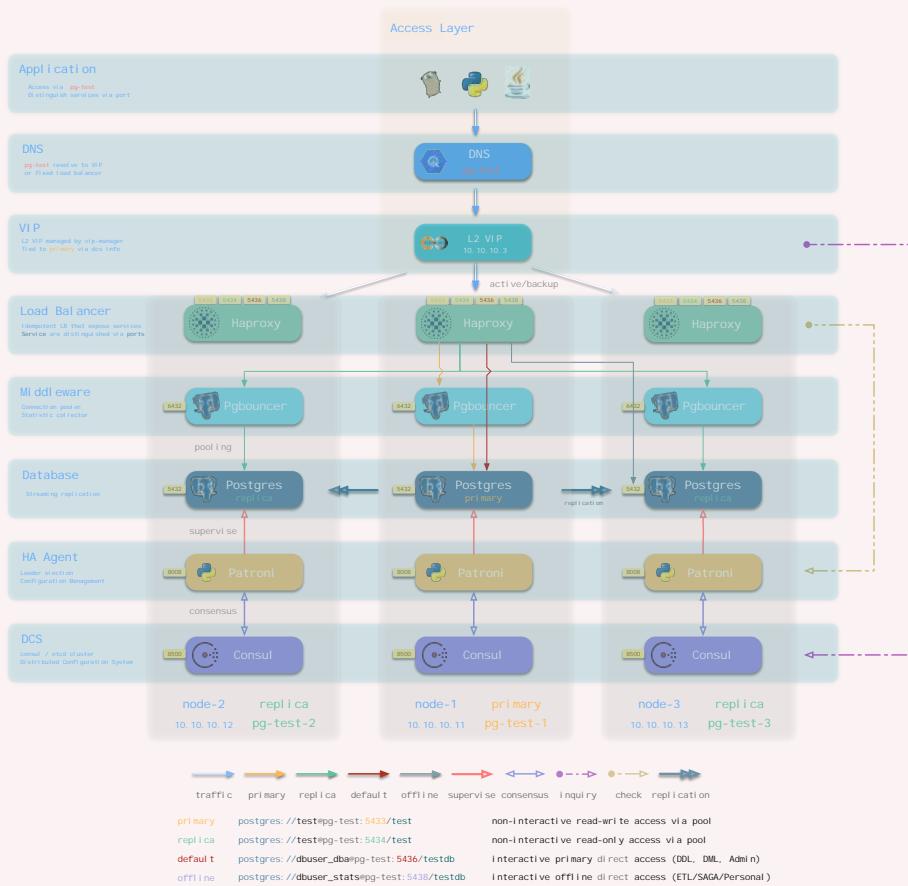
集群对外提供的各种**服务**就不会停止。

从库故障基本无影响。

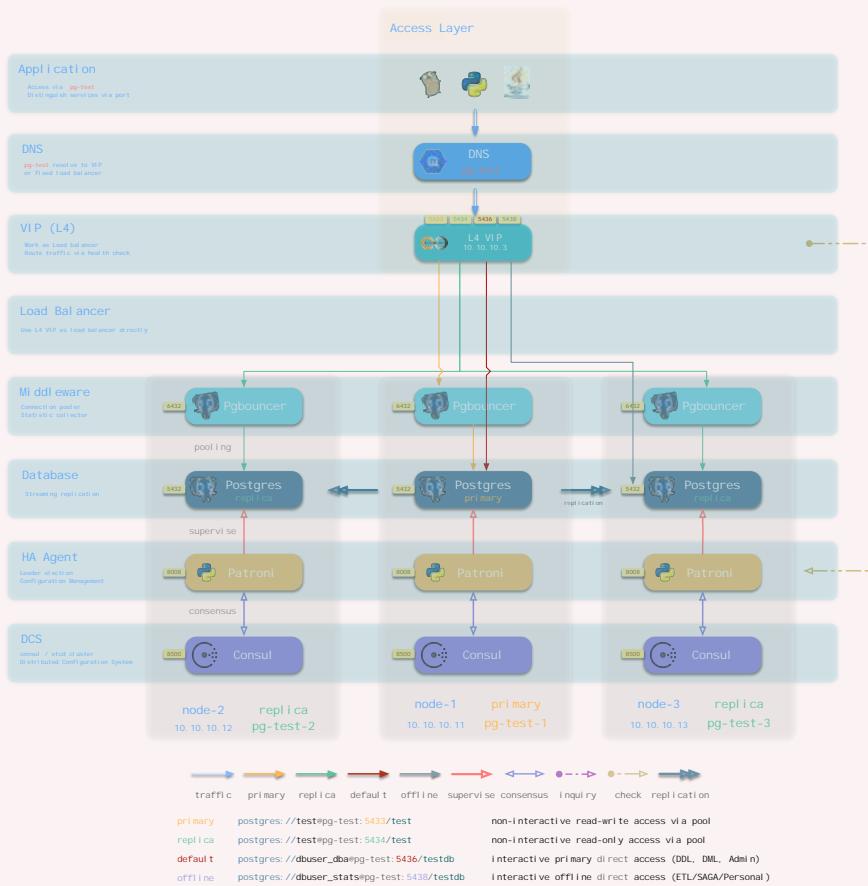
主库故障切换时，读写流量影响约15秒



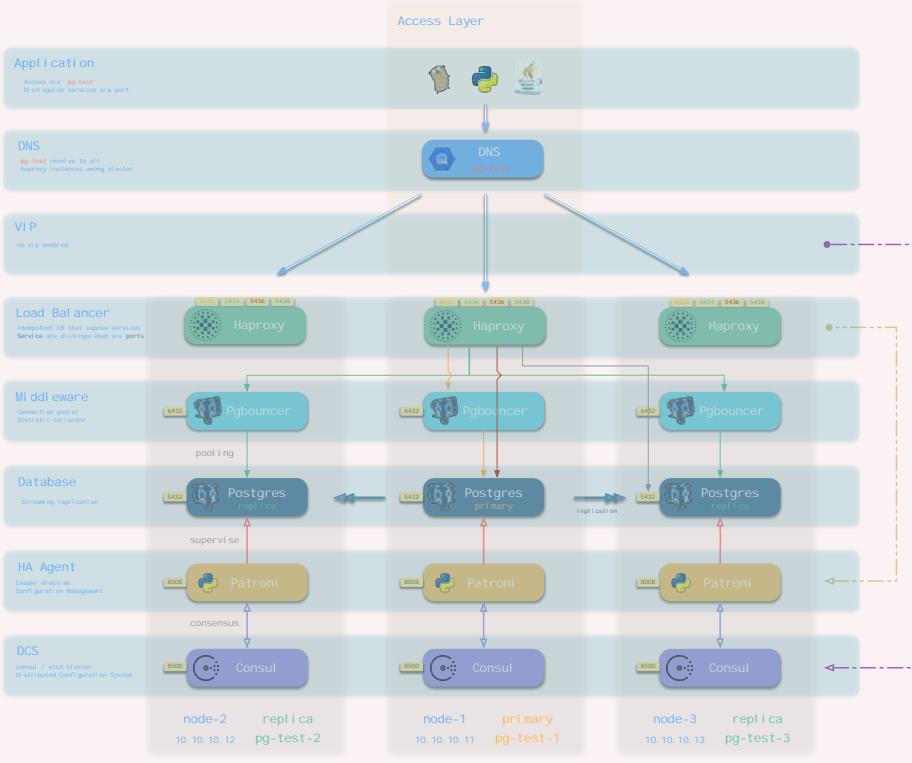
DNS + L2 VIP + HAProxy



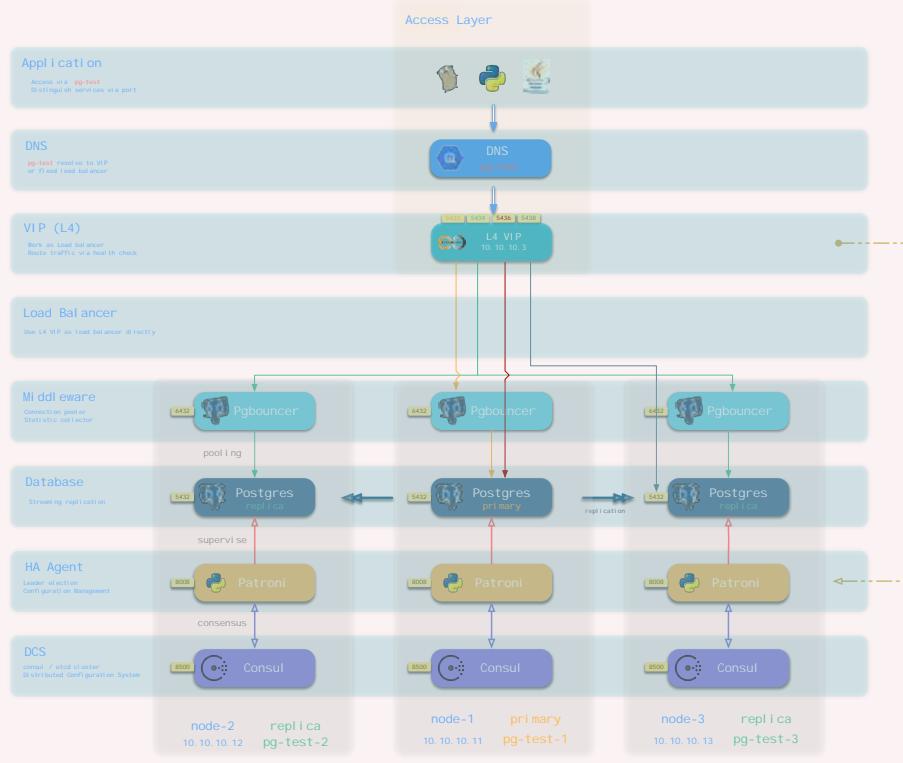
DNS + L4 VIP



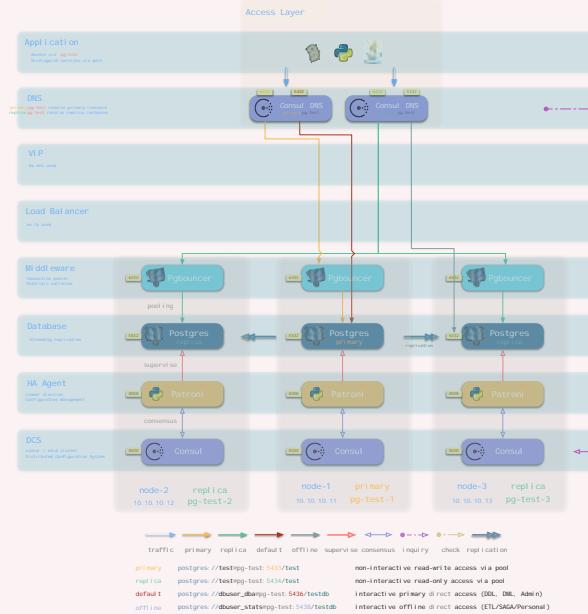
DNS + HAProxy



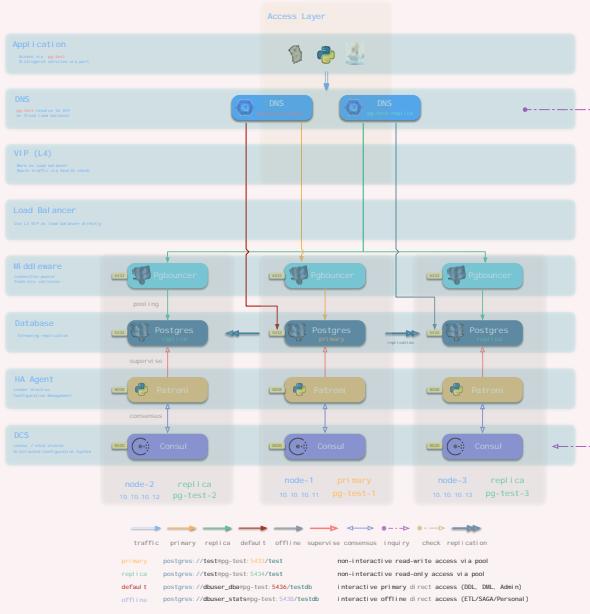
DNS + L4 VIP



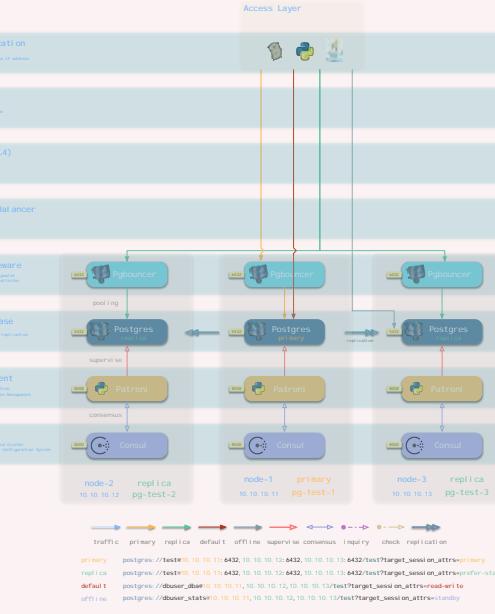
Consul DNS SD



Pure DNS



IP + Smart Client



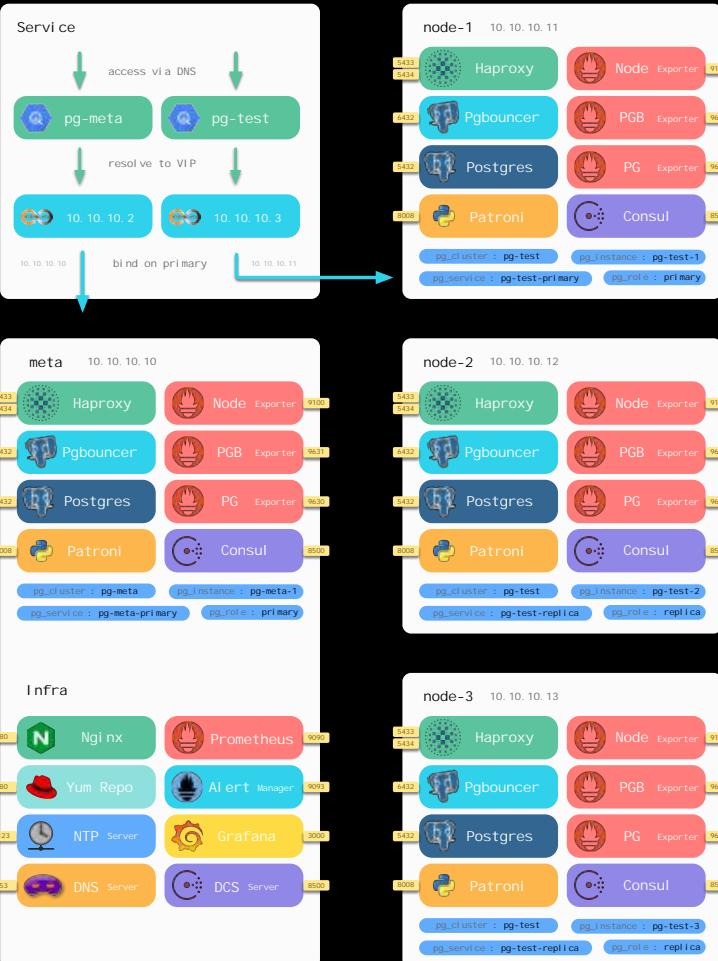
Sandbox

本地运行，一键拉起，低配置无压力

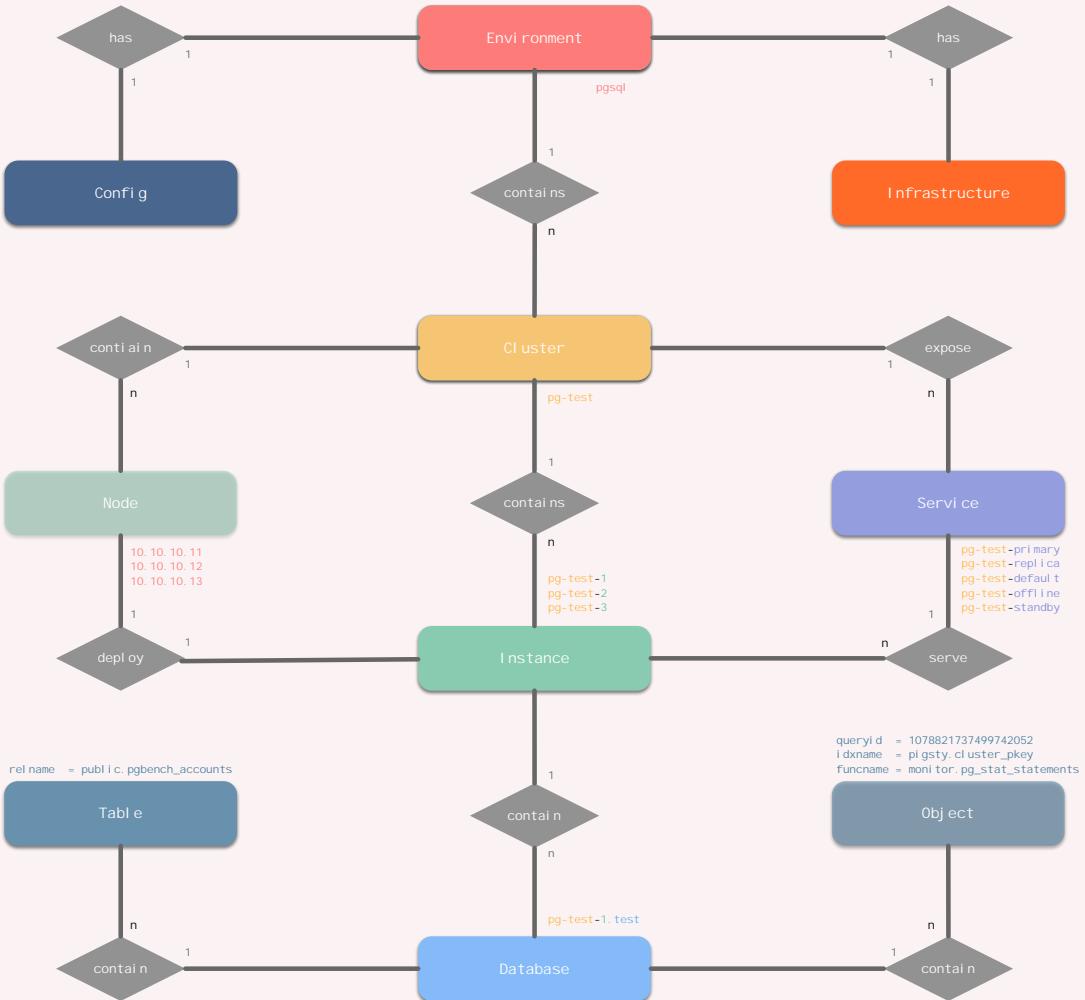
Jupyterlab, Grafana, Echarts
方便易用，助力数据分析

确定性的环境，便于交付演示案例

Just like Minikube for Kubernetes



Domain Extended



Analytic

数据分析
与可视化



*The best vision is **insight**.*

— Malcolm Forbes

Pigsty for Analysis

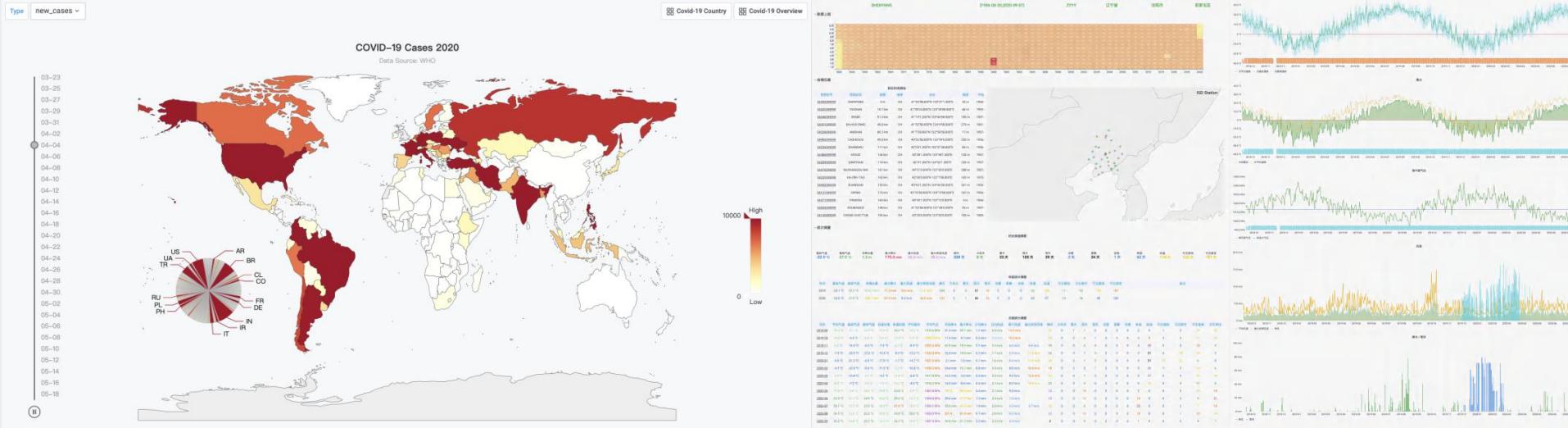
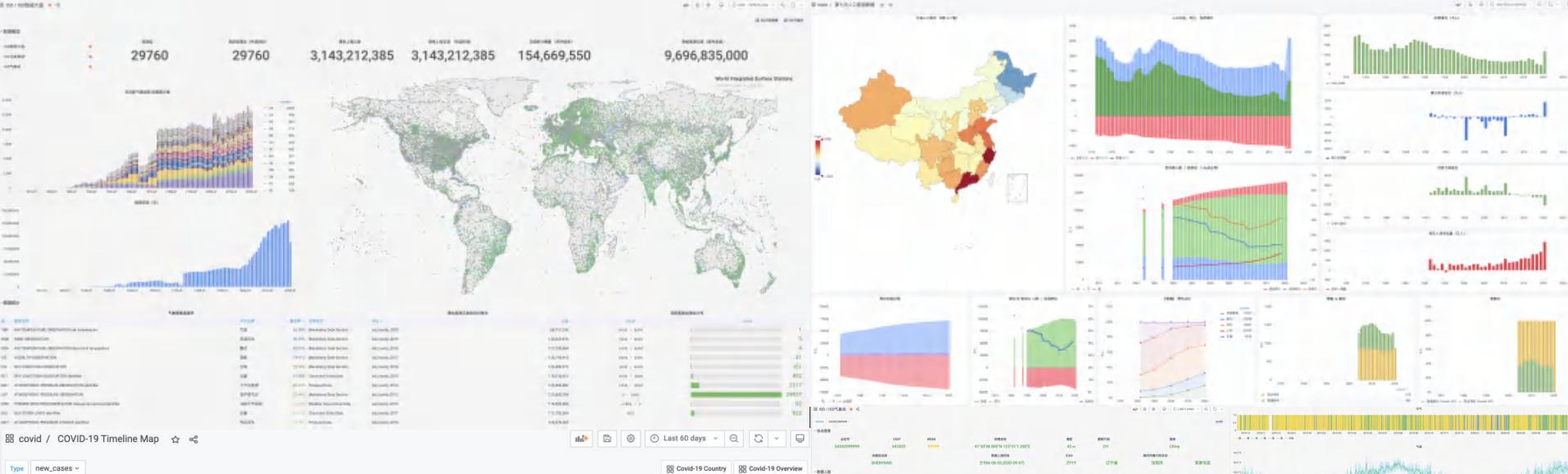
专注于数据处理、**分析**与**可视化**

使用**Grafana**承载”前端”与”后端”的主体功能。

使用**PostgreSQL**作为主要数据库，存储数据并实现业务逻辑

使用**Echarts**与**Grafana**实现数据可视化与交互功能

采用统一的格式，可以在**Pigsty**环境中一键安装运行

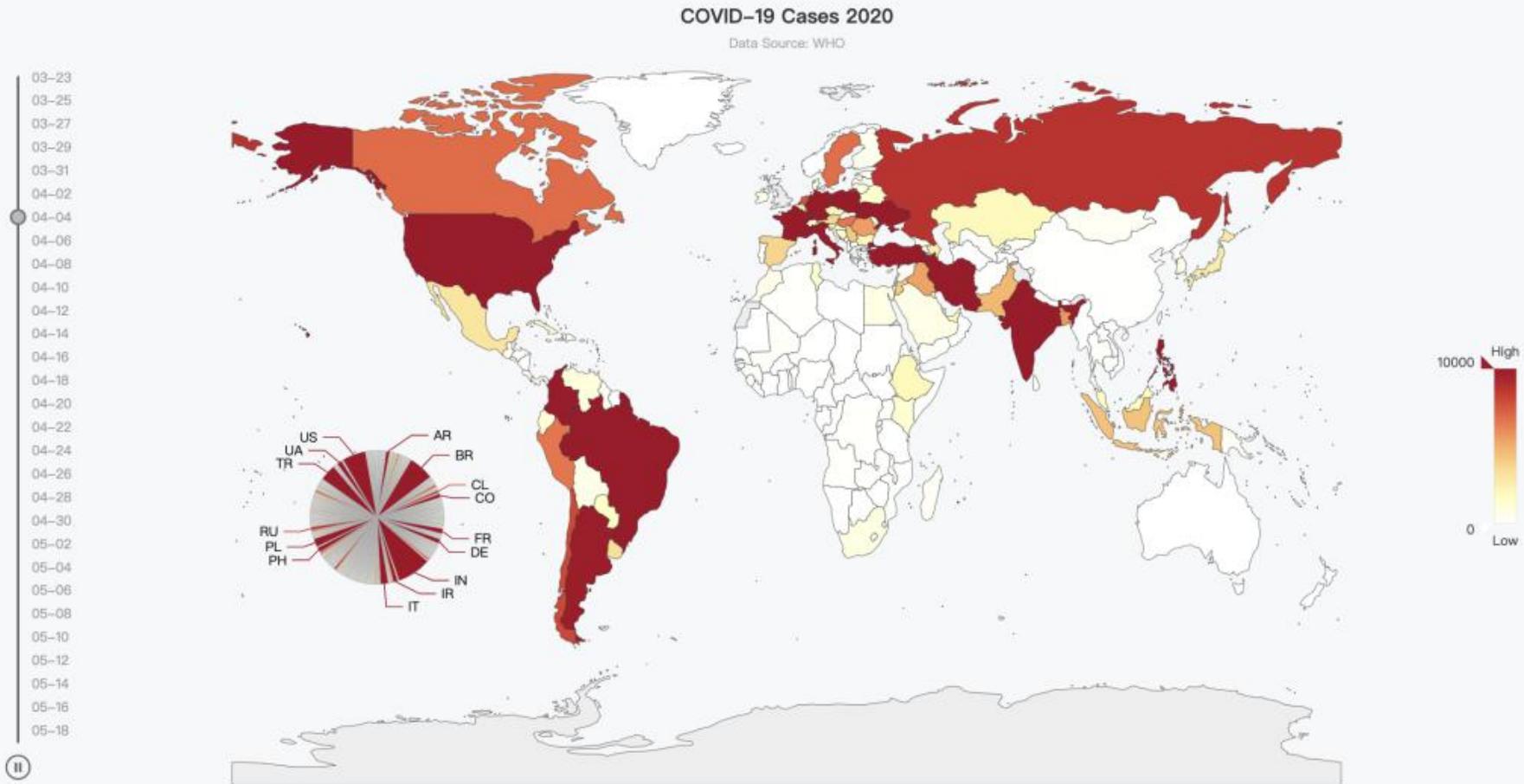


Type new_cases

Grafana太丑，Echarts来凑，100行代码解决所有

Covid-19 Country

Covid-19 Overview



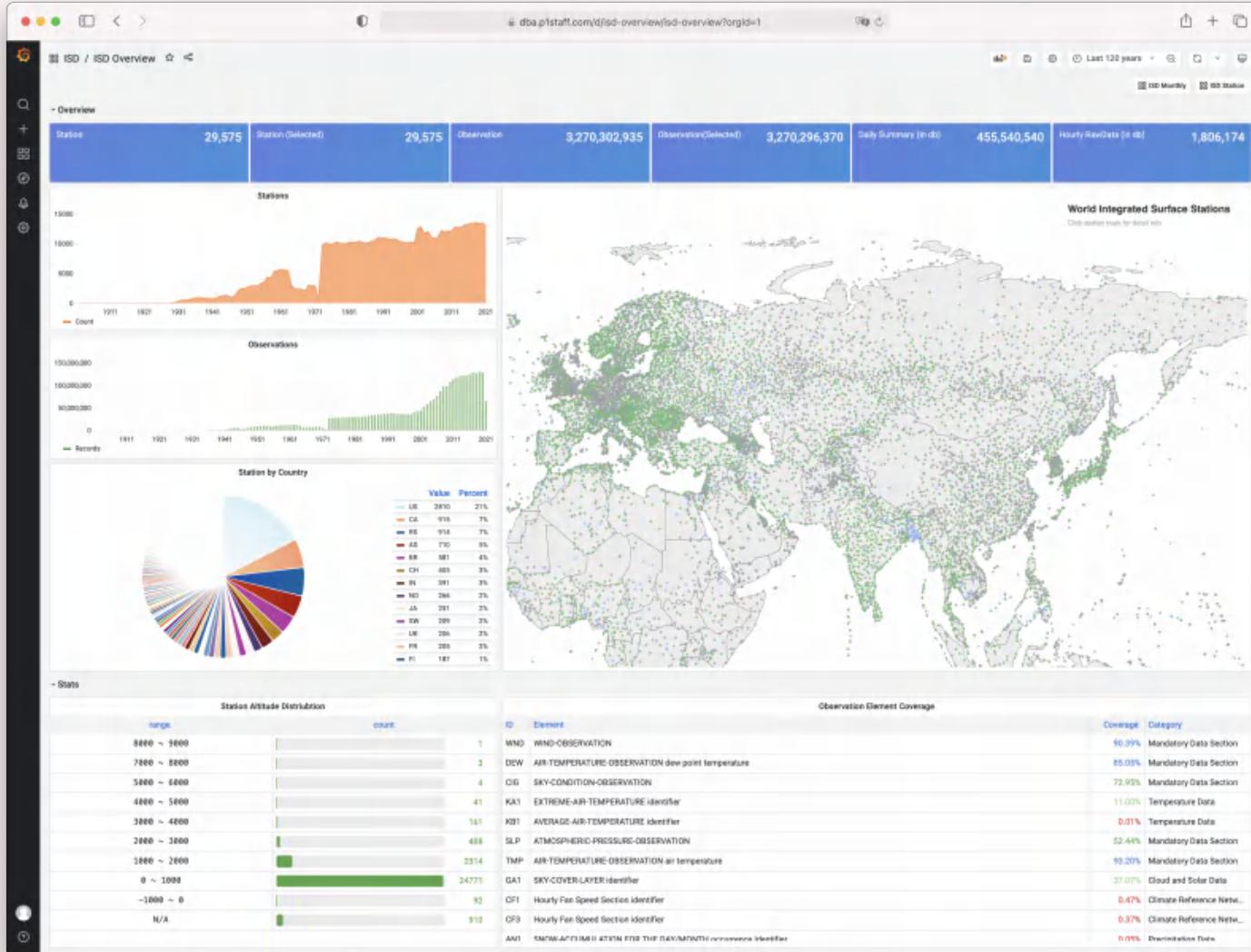
样例：ISD

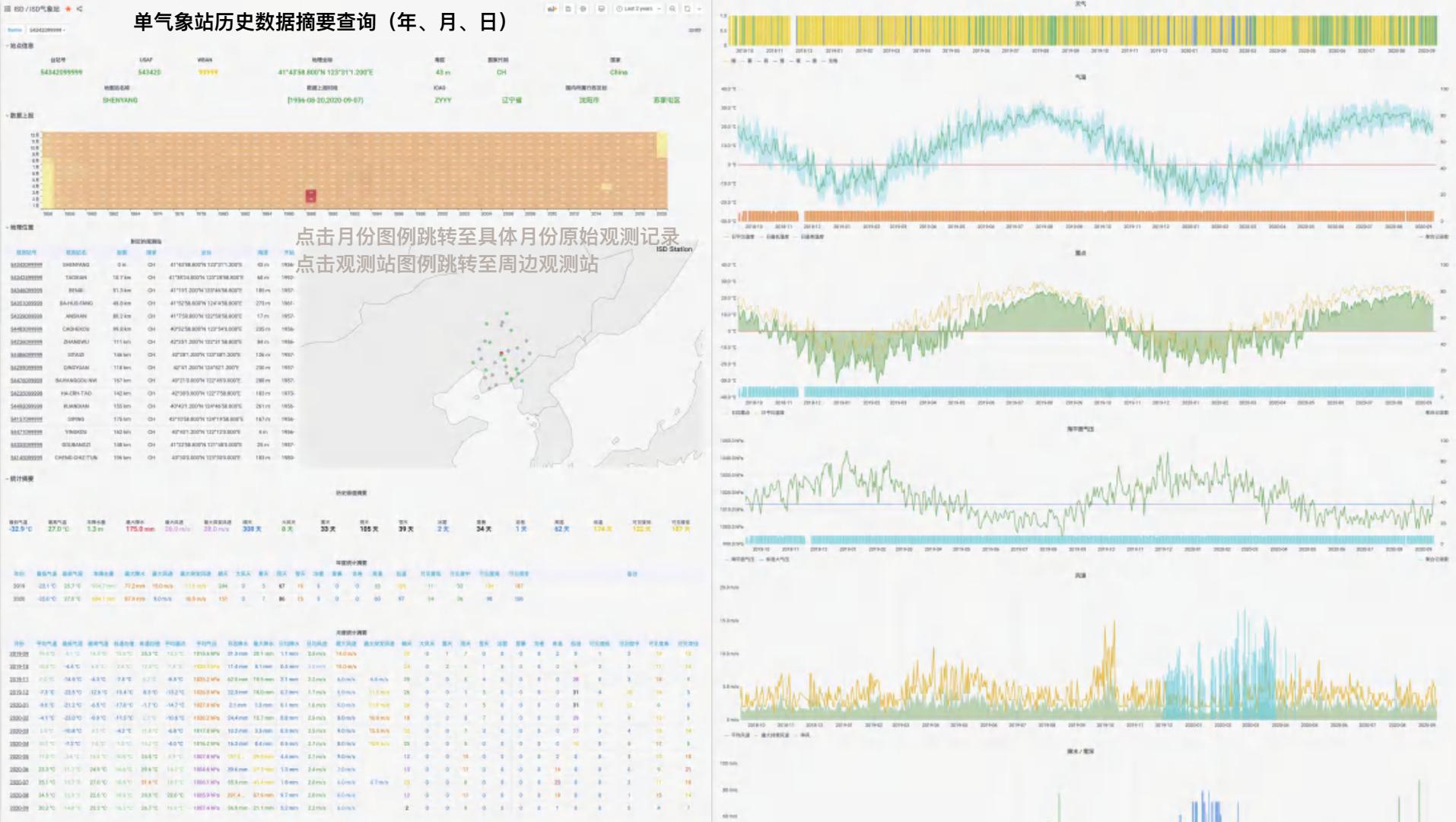
全球地表气象站 历史数据查询

3万个地表气象站

120年来小时级

原始气象观测数据





More

自带应用：Pigsty CMDB

使用PG管理PG

将静态配置文件

替换为动态CMDB查询

从而与外部系统与平台集成



More

Beta应用：Redis

使用同样的基础设施
部署、监控Redis



Open Source



*A system cannot be **successful** if it is too strongly influenced by a single person. Once the initial design is complete and fairly robust, the real test begins as people with many different viewpoints undertake their own experiments.*

— Donald Knuth

Pigsty Timeline

2021-07-31 v1.0.0 GA



Next Episode of Pigsty

更多功能
平台化

原生支持 部署Citus高可用水平分片集群组
支持使用 etcd 替代完整的 consul 功能
添加离线从库，延时从库部署支持

集中式备份/恢复方案
安全性强化：严格模式
巡检、日报、报警分析

管理平台 & CMDB
集成模式变更方案
集成数据迁移方案

监控系统
应用案例

更多、更全面、更专业的监控面板
更多、更实用的专题应用
集成pgbager日志分析报告
集成PEV慢查询可视化分析工具

集成黑盒监控 blackbox_exporter
集成更多日志衍生指标：mtail
完善的PostgreSQL系统目录浏览器：PGCAT

系统支持
云原生改造

支持ARM架构，支持更多Linux发行版：SUSE (✓) , Debian, Ubuntu
容器化方案，Kubernetes 支持，Kubernetes Operator

管理更多
开源数据库

支持部署、监控更多种类的开源数据库：Redis (✓) , MongoDB, MySQL,

Open Source

第一时间获取反馈，了解用户的痛点、痒点与爽点，解决**真正的问题**

让用户能用**好数据库，用好**数据库

降低门槛，扩大用户群，促进PostgreSQL生态繁荣发展。

提供一种“自主可控”的选择，让中小企业，个人用户更有**自由选择**的底气

THANKS