

O'REILLY®

Wydanie II

Prometheus w pełnej gotowości

Jak monitorować pracę infrastruktury
i wydajność działania aplikacji



Julien Pivotto
Brian Brazil

Helion 

Tytuł oryginału: Prometheus: Up & Running:
Infrastructure and Application Performance Monitoring, 2nd Edition

Tłumaczenie: Robert Górczyński

ISBN: 978-83-289-0529-0

© 2024 Helion S.A.

Authorized Polish translation of the English edition of *Prometheus: Up & Running, 2E*
ISBN 9781098131142 © 2023 Julien Pivotto.

This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from the Publisher.

Wszelkie prawa zastrzeżone. Nieautoryzowane rozpowszechnianie całości lub fragmentu niniejszej publikacji w jakiegokolwiek postaci jest zabronione. Wykonywanie kopii metodą kserograficzną, fotograficzną, a także kopiowanie książki na nośniku filmowym, magnetycznym lub innym powoduje naruszenie praw autorskich niniejszej publikacji.

Wszystkie znaki występujące w tekście są zastrzeżonymi znakami firmowymi bądź towarowymi ich właścicieli.

Autor oraz wydawca dołożyli wszelkich starań, by zawarte w tej książce informacje były kompletne i rzetelne. Nie biorą jednak żadnej odpowiedzialności ani za ich wykorzystanie, ani za związane z tym ewentualne naruszenie praw patentowych lub autorskich. Autor oraz wydawca nie ponoszą również żadnej odpowiedzialności za ewentualne szkody wynikłe z wykorzystania informacji zawartych w książce.

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<https://helion.pl/user/opinie/prom2>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Pliki z przykładami omawianymi w książce można znaleźć pod adresem:

<https://ftp.helion.pl/przyklady/prom2.zip>

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 230 98 63

e-mail: helion@helion.pl

WWW: <https://helion.pl> (księgarnia internetowa, katalog książek)

Printed in Poland.

- [Kup książkę](#)
- [Poleć książkę](#)
- [Oceń książkę](#)

- [Księgarnia internetowa](#)
- [Lubię to! » Nasza społeczność](#)

Spis treści

Wprowadzenie	11
--------------------	----

Część I. Wprowadzenie 15

1. Czym jest Prometheus? 17

Czym jest monitorowanie?	18
Krótka i niepełna historia monitorowania	19
Kategorie monitorowania	21
Architektura systemu Prometheus	25
Biblioteki klienta	25
Komponenty eksportujące	27
Odkrywanie usług	27
Zbieranie danych	28
Pamięć masowa	29
Panele sterowania	29
Reguły rejestrowania i ostrzeżeń	29
Zarządzanie ostrzeżeniami	30
Długoterminowa pamięć masowa	31
Czym Prometheus nie jest?	31

2. Rozpoczęcie pracy z systemem Prometheus 32

Uruchamianie systemu Prometheus	32
Używanie przeglądarki wyrażenia	35
Uruchamianie komponentu Node Exporter	39
Ostrzeganie	43

3. Instrumentacja	53
Prosty program	53
Licznik	55
Zliczanie wyjątków	57
Zliczanie wielkości	58
Miernik	59
Używanie miernika	59
Wywołanie zwrotne	61
Podsumowanie	62
Histogram	63
Kubelki	64
Instrumentacja testów jednostkowych	67
Podejścia w zakresie instrumentacji	68
Dlaczego należy stosować instrumentację?	68
Jak daleko powinna sięgać instrumentacja?	70
Dlaczego należy nadawać nazwy wskaźnikom?	71
4. Ekspozycja	74
Python	75
WSGI	75
Twisted	76
Wiele procesów i Gunicorn	76
Go	80
Java	81
HTTPServer	81
Servlet	82
Pushgateway	83
Most	86
Analizator składni	87
Format ekspozycji tekstu	88
Typy wskaźników	88
Etykiety	89
Cytowanie znaków	90
Znaczniki czasu	90
Sprawdzenie wskaźników	90
OpenMetrics	91
Typy wskaźników	91
Etykiety	92
Znaczniki czasu	93

5. Etykiety	94
Czym są etykiety?	94
Etykiety instrumentacji i systemów przeznaczonych do monitorowania	95
Instrumentacja	95
Wskaźnik	97
Wiele etykiet	97
Element potomny	98
Agregacja	100
Wzorce etykiet	101
Enum	101
Info	103
Kiedy używać etykiet?	105
Licznosc	106
6. Tworzenie paneli sterowania za pomocą Grafany	109
Instalacja	109
Źródło danych	110
Panele sterowania	113
Unikanie ścian wykresów	113
Panel szeregu czasowego	114
Kontrolki czasu	116
Panel Stat	118
Panel Table	119
Panel State Timeline	120
Zmienne szablonu	121

Część III. Monitorowanie infrastruktury **127**

7. Node Exporter	129
Komponent pobierający dane dotyczące procesora	130
Komponent pobierający dane systemu plików	131
Komponent pobierający dane dysku	132
Komponent pobierający dane netdev	133
Komponent pobierający dane meminfo	133
Komponent pobierający dane hwmon	134
Komponent pobierający dane stat	135
Komponent pobierający dane uname	135
Komponent pobierający dane systemu operacyjnego	136
Komponent pobierający dane loadavg	136
Komponent pobierający dane dotyczące obciążenia	137

Komponent pobierający dane z pliku tekstowego	137
Używanie komponentu pobierającego dane z pliku tekstowego	139
Znaczniki czasu	140
8. Mechanizm wykrywania usług	141
Mechanizmy wykrywania usług	142
Podejście statyczne	143
Bazujący na pliku mechanizm wykrywania usług	144
Bazujący na HTTP mechanizm wykrywania usług	147
Consul	148
EC2	149
Zmiana etykiety	150
Wybór danych, które mają być pobierane	151
Etykiety systemów przeznaczonych do monitorowania	154
Jak pobierać dane?	162
metric_relabel_configs	164
Kolizje etykiet i opcja honor_labels	166
9. Kontenery i Kubernetes	167
cAdvisor	167
Procesor	168
Pamięć	169
Etykiety	169
Kubernetes	170
Działanie w Kubernetes	170
Mechanizm wykrywania usług	172
Komponent kube-state-metrics	181
Wdrożenia alternatywne	182
10. Najczęściej używane komponenty eksportujące	183
Consul	183
MySQLd	185
Komponent eksportujący Grok	186
Czarna skrzynka	189
Próbkowanie za pomocą ICMP	190
Próbkowanie za pomocą TCP	195
Próbkowanie za pomocą HTTP	197
Próbkowanie za pomocą DNS	199
Konfiguracja systemu Prometheus	201

11. Współpraca z innymi systemami monitorowania	205
Inne systemy monitorowania	205
InfluxDB	207
StatsD	208
12. Tworzenie komponentu eksportującego	211
Telemetria narzędzia Consul	211
Niestandardowy komponent pobierający dane	214
Etykiety	219
Wskazówki	220

Część IV. PromQL **223**

13. Wprowadzenie do PromQL	225
Podstawy agregacji	225
Miernik	225
Licznik	227
Podsumowanie	227
Histogram	229
Selektory	230
Dopasowania	231
Wektor natychmiastowy	232
Wektor zakresu	233
Podzapytanie	235
Przesunięcie	236
Modyfikator at	237
API HTTP	237
Punkt końcowy query	237
Punkt końcowy query_range	239
14. Operatory agregacji	243
Grupowanie	243
Klauzula without	244
Klauzula by	245
Operatory	246
sum	246
count	247
avg	248
group	248
stddev i stdvar	249

min i max	250
topk i bottomk	250
quantile	251
count_values	252
15. Operatory binarne	254
Praca z wartościami skalarnymi	254
Operatory arytmetyczne	254
Operator trygonometryczny	256
Operatory porównania	256
Dopasowanie wektora	258
Dopasowanie typu jeden do jednego	258
Dopasowanie typu wiele do jednego i group_left	260
Wiele do wielu i operatory logiczne	263
Kolejność wykonywania operatorów	267
16. Funkcje	268
Zmiana typu	268
vector()	268
scalar()	269
Funkcje matematyczne	270
abs()	270
ln(), log2() i log10()	270
exp()	271
sqrt()	271
ceil() i floor()	271
round()	272
clamp(), clamp_max() i clamp_min()	272
sgn()	273
Funkcje trygonometryczne	273
Data i godzina	274
time()	274
minute(), hour(), day_of_week(), day_of_month(), day_of_year(), days_in_month(), month() i year()	275
timestamp()	276
Etykiety	276
label_replace()	277
label_join()	277
Brakujące szeregi czasowe oraz funkcje absent() i absent_over_time()	278
Sortowanie za pomocą sort() i sort_desc()	279
Histogram za pomocą funkcji histogram_quantile()	279

Liczniki	280
rate()	280
increase()	282
irate()	282
resets()	283
Zmiana mierników	283
changes()	284
deriv()	284
predict_linear()	284
delta()	285
idelta()	285
holt_winters()	285
Agregacja na przestrzeni czasu	286
17. Reguły rejestrowania	288
Używanie reguł rejestrowania	288
Kiedy używać reguł rejestrowania?	291
Zmniejszenie liczności	291
Tworzenie funkcji wektora zakresu	292
Reguły dla API	293
Jak nie używać reguł?	294
Nazewnictwo reguł rejestrowania	295
<hr/>	
Część V. Ostrzeżenie	299
18. Ostrzeżenie	301
Reguły ostrzeżenia	302
for	304
Etykiety ostrzeżenia	306
Adnotacje i szablony	308
Jak wygląda dobre ostrzeżenie?	311
Konfigurowanie menedżera ostrzeżeń Alertmanager	
w oprogramowaniu Prometheus	312
Etykiety zewnętrzne	313
19. Menedżer ostrzeżeń Alertmanager	315
Potok powiadomienia	315
Plik konfiguracyjny	316
Drzewo routingu	317
Odbiorcy	324
Wstrzymywanie	333
Interfejs sieciowy menedżera ostrzeżeń Alertmanager	334

20. Zapewnienie bezpieczeństwa po stronie serwera	337
Funkcje bezpieczeństwa dostarczane przez system Prometheus	337
Włączenie obsługi TLS	338
Opcje zaawansowane TLS	339
Włączenie uwierzytelniania podstawowego	340
21. Zebranie wszystkiego w całość	342
Planowanie wdrożenia	342
Rozbudowa systemu Prometheus	343
Podejście globalne z wykorzystaniem federacji	345
Długoterminowa pamięć masowa	348
Uruchamianie systemu Prometheus	350
Sprzęt	350
Zarządzanie konfiguracją	352
Sieci i uwierzytelnianie	354
Planowanie pod kątem awarii	355
Klastrowanie Alertmanager	358
Monitorowanie meta i cross	359
Zarządzanie wydajnością działania	360
Wykrywanie problemu	360
Wyszukiwanie kosztownych wskaźników i systemów przeznaczonych do monitorowania	361
Zmniejszenie obciążenia	362
Sharding poziomy	363
Zarządzanie zmianami	364
Uzyskiwanie pomocy	365
Skorowidz	367

Czym jest Prometheus?

Prometheus to bazujący na wskaźnikach system monitorowania, który jest dostępny jako oprogramowanie typu open source. Oczywiście Prometheus nie jest jedynym rozwiązaniem tego typu, więc być może zastanawiasz się nad tym, co go wyróżnia spośród innych.

Prometheus zajmuje się wykonywaniem tylko jednego zadania i robi to świetnie. Posiada prosty i jednocześnie oferujący potężne możliwości model danych i język wykonywania zapytań, który pozwala na sprawdzanie wydajności działania aplikacji i infrastruktury. Nie wykracza poza obszar wskaźników i nie próbuje rozwiązywać problemów, pozostawiając te działania innym narzędziom, które są lepiej przystosowane do tego celu.

Spółeczność i ekosystem oprogramowania Prometheus znacznie się zwiększyły od chwili, gdy na początku w 2012 roku w społeczności SoundCloud pracowała zaledwie garstka programistów. System Prometheus został opracowany w języku Go i pozostaje dostępny na licencji Apache 2.0. Setki osób uczestniczą w pracy nad projektem, który nie jest kontrolowany przez żadną firmę. Zawsze trudno jest ustalić, ile użytkowników korzysta z projektu typu open source, ale w 2022 roku szacowano, że setki tysięcy organizacji używało systemu Prometheus w środowisku produkcyjnym. W 2016 roku projekt Prometheus stał się drugim¹ projektem należącym do fundacji Cloud Native Computing Foundation (CNCF).

Na potrzeby koordynacji własnego kodu istnieją biblioteki klienckie opracowane dla praktycznie wszystkich najpopularniejszych języków programowania i środowisk uruchomieniowych, takich jak Go, Java/JVM, C#/.Net, Python, Ruby, Node.js, Haskell, Erlang i Rust. Wiele popularnych aplikacji, np. Kubernetes, Docker, Envoy i Vault, jest już dostarczanych z bibliotekami klienckimi dla oprogramowania Prometheus. Dla innych firm tworzących oprogramowanie udostępniające wskaźniki w formacie innym niż zgodny z systemem Prometheus dostępne są setki rozwiązań integracyjnych. To są tzw. komponenty eksportujące, które obejmują produkty takie jak HAProxy, MySQL, PostgreSQL, Redis, JMX, SNMP, Consul i Kafka. Przyjaciel Briana dodał nawet obsługę monitorowania serwerów Minecrafta, ponieważ dla niego ważna jest osiągnięta liczba klatek na sekundę podczas gry.

¹ Pierwszy był produkt Kubernetes.

Format w postaci zwykłego tekstu² oznacza, że przekazywanie wskaźników do systemu Prometheus jest bardzo łatwe. Inne systemy monitorowania, zarówno komercyjne, jak i typu open source, zapewniają obsługę tego formatu. Dzięki temu w tych rozwiązaniach można skoncentrować się bardziej na ich podstawowej funkcjonalności, a nie na konieczności poświęcania czasu i wysiłku, aby zapewnić obsługę każdego elementu oprogramowania, który użytkownik będzie chciał monitorować.

Model danych identyfikuje poszczególne szeregi czasowe nie tylko za pomocą nazwy, ale również dzięki nieuporządkowanemu zbiorowi par klucz i wartość określanych mianem etykiet. Język zapytań PromQL pozwala na agregację między tymi etykietami. Dzięki temu można przeprowadzać analizy nie tylko dla poszczególnych procesów, ale także dla centrów danych i usług bądź na podstawie innych zdefiniowanych etykiet. Następnie dane można przedstawić na wykresach dzięki użyciu systemów takich jak Grafana (<https://github.com/grafana/grafana>) i Perses (<https://github.com/perses/perses>).

Ostrzeżenia mogą być zdefiniowane za pomocą dokładnie tego samego języka zapytań PromQL, który jest używany do przygotowywania wykresów. Jeżeli coś można przedstawić na wykresie, można również to wykorzystać w ostrzeżeniach. Dzięki etykietom obsługa techniczna ostrzeżeń staje się łatwiejsza, ponieważ można utworzyć pojedyncze ostrzeżenie obejmujące wszystkie możliwe wartości etykiet. W niektórych systemach monitorowania istnieje możliwość pojedynczego tworzenia ostrzeżeń dla każdego urządzenia i aplikacji. W związku z tym mechanizm wykrywania usług może automatycznie ustalić, które aplikacje i urządzenia powinny być pozyskane ze źródeł takich jak Kubernetes, Consul, Amazon Elastic Compute Cloud (EC2), Azure, Google Compute Engine (GCE) i OpenStack.

Mimo tych wszystkich funkcji i zalet Prometheus jest efektywny i prosty w użyciu. Pojedynczy serwer systemu Prometheus może pobierać miliony próbek na sekundę. Serwer ma postać pojedynczego pliku binarnego ze statycznie dołączonymi bibliotekami oraz pliku konfiguracyjnego. Wszystkie komponenty systemu Prometheus mogą być uruchamiane w kontenerach i unikają wykonywania jakichkolwiek zadań utrudniających działanie narzędziom przeznaczonym do zarządzania konfiguracją. Prometheus został opracowany z myślą o jego zintegrowaniu z istniejącą infrastrukturą i używaniu razem z nią, a nie jako platforma zarządzania.

Skoro już mniej więcej wiesz, czym jest Prometheus, warto zrobić krok wstecz i wyjaśnić, co oznacza pojęcie „monitorowanie”, aby tym samym zapewnić pewien kontekst. Następnie przedstawimy podstawowe koncepcje kryjące się za systemem Prometheus i wyjaśnimy, czym Prometheus nie jest.

Czym jest monitorowanie?

W szkole średniej jeden z nauczycieli Briana powiedział mi, że jeśli zapyta 10 ekonomistów o znaczenie określenia ekonomia, to otrzyma 11 odpowiedzi. W przypadku monitorowania mamy do czynienia z podobnym brakiem zgody co do faktycznego znaczenia tego terminu. Gdy Brian wspomina innym osobom, czym się zajmuje, ludzie są przekonani, że jego praca obejmuje praktycznie wszystko od pilnowania temperatury w fabrykach do śledzenia pracowników, aby dowiedzieć się, kto korzysta z Facebooka w pracy, oraz wykrywać intruzów w sieci firmowej.

² Poza formatem zwykłego tekstu została opracowana jeszcze bardziej standaryzowana i nieco odmienna wersja formatu tekstowego systemu Prometheus, która nosi nazwę OpenMetrics.

Prometheus nie został opracowany do wykonywania tego rodzaju zadań³. Ten produkt został zbudowany, aby pomóc twórcom oprogramowania i administratorom w nadzorowaniu produkcyjnych systemów komputerowych, np. aplikacji, narzędzi, baz danych i sieci kryjących się za popularnymi witrynami internetowymi.

Czym jest więc monitorowanie w takim kontekście? Warto zawęzić ten rodzaj monitorowania operacyjnego systemu komputerowego do czterech kwestii.

Powiadamianie

Wiedza o tym, kiedy dzieje się coś złego, to zwykle najważniejszy powód, dla którego stosuje się monitorowanie. Oczekuje się, że w przypadku problemu system monitorowania poinformuje o tym człowieka, który będzie mógł odpowiednio zareagować.

Debugowanie

Po poinformowaniu człowiek musi przyrzeć się problemowi w celu znalezienia jego podstawowej przyczyny i ostatecznie go usunąć.

Trendy

Powiadamianie i debugowanie zwykle odbywają się w czasie wyrażonym w minutach bądź godzinach. Wprawdzie mniej pilna, ale równie użyteczna jest możliwość sprawdzenia, jak systemy są używane na przestrzeni czasu i jak się wówczas zmieniają. Istniejące trendy mogą wpływać na podejmowanie decyzji projektowych, a także planowanie pojemności.

Instalowanie systemu

Gdy do dyspozycji masz jedynie młotek, wtedy wszystko zaczyna przypominać gwoździe. Ostatecznie wszystkie systemy monitorowania danych są tak naprawdę potokami przetwarzania informacji. Czasami wygodniejsze będzie przeznaczenie fragmentu systemu monitorowania do innego celu niż tworzenie oddzielnego rozwiązania. W takim przypadku nie mamy do czynienia z czystym monitorowaniem, ale to jest na tyle często spotykane podejście, że zdecydowaliśmy się na jego omówienie.

W zależności od tego, z kim rozmawiasz i jakie ta osoba ma doświadczenie, tylko część z wymienionych zadań może ona uznać za monitorowanie. To prowadzi do wielu odbywających się w kółko dyskusji dotyczących monitorowania, których efektem jest jedynie frustracja uczestników. Aby pomóc Ci w zrozumieniu perspektywy innych, w kolejnym punkcie zamieszczono nieco historii dotyczącej monitorowania.

Krótką i niepełną historią monitorowania

W ostatnich latach można zauważyć zmianę w monitorowaniu polegającą na przejściu w kierunku narzędzi takich jak Prometheus. Przez długi czas dominującym rozwiązaniem było pewnego rodzaju połączenie Nagios i Graphite lub ich różne warianty.

³ Monitorowanie temperatury urządzeń i centrów danych jednak nie jest tak rzadko spotykane. Wiemy również o kilku użytkownikach, którzy dla rozrywki wykorzystują system Prometheus do śledzenia pogody.

Wspominając Nagios, mamy na myśli wszelkie oprogramowanie zaliczane do tej samej obszernej rodziny oprogramowania, czyli m.in. Icinga, Zmon i Sensu. Ich działanie polega przede wszystkim na regularnym wykonywaniu skryptów określanych mianem *sprawdzeń*. Jeżeli sprawdzenie kończy się niepowodzeniem i powoduje zwrot niezerowego kodu stanu, wówczas następuje wygenerowanie powiadomienia. Oprogramowanie Nagios zostało początkowo utworzone przez Ethana Galstada w 1996 roku jako aplikacja MS-DOS przeznaczona do wykonywania żądań ping. Najpierw w 1999 roku zostało wydane pod nazwą NetSaint, zaś później w 2002 roku jego nazwa została zmieniona na Nagios.

Historia Graphite wymaga cofnięcia się do 1994 roku, w którym to Tobias Oetiker utworzył skrypt Perla. Następnie w 1995 roku ten skrypt stał się programem Multi Router Traffic Grapher, MRTG 1.0. Zgodnie z nazwą jego podstawowym przeznaczeniem było monitorowanie sieci za pomocą protokołu SNMP (ang. *simple network management protocol*). Dzięki wykonywaniu skryptów ten program mógł również pobierać wskaźniki⁴. W 1997 roku nastąpiły duże zmiany związane z przepisaniem pewnych fragmentów kodu źródłowego w języku C oraz utworzeniem bazy danych Round Robin Database (RRD) przeznaczonej do przechowywania danych wskaźników. W efekcie nastąpiła znaczna poprawa wydajności działania, zaś RRD stała się punktem wyjścia dla innych narzędzi, m.in. Smokeping i Graphite.

Począwszy od 2006 roku, Graphite używa do przechowywania wskaźników produktu o nazwie Whisper, którego budowa jest podobna do bazy danych RRD. Program Graphite sam nie pobiera danych, raczej je otrzymuje przekazywane przez kolekcję narzędzi takich jak collectd i StatsD, utworzonych w odpowiednio 2005 i 2010 roku.

Jak widzisz, generowanie wykresów i powiadamianie były kiedyś oddzielnymi zadaniami, wykonywanymi przez całkowicie odmienne narzędzia. Wprawdzie można było utworzyć skrypt sprawdzenia w celu analizy zapytania w Graphite i generować powiadomienia na podstawie tego skryptu, ale większość sprawdzeń dotyczy nieoczekiwanych stanów, takich jak niedziałający proces.

Inną pozostałością po tamtej epoce jest praktycznie ręczne podejście do administrowania usługami komputerowymi. Usługi były wdrażane w poszczególnych komputerach i starannie obsługiwane przez administratorów systemów. Powiadomienia dotyczące potencjalnych problemów były przekazywane do odpowiednich inżynierów. Gdy popularność zyskała chmura i technologie związane z natywną chmurą, takie jak EC2, Docker i Kubernetes, traktowanie poszczególnych komputerów i usług jak zwierząt domowych i poświęcanie każdej z nich oddzielnej uwagi okazało się niemożliwe. Raczej są one monitorowane i administrowane jako grupa. Tak jak branża informatyczna odeszła od zarządzania ręcznego i przeszła do narzędzi typu Chef i Ansible, a obecnie zaczyna korzystać z technologii takich jak Kubernetes, tak podobne zmiany są zauważalne także w obszarze monitorowania. To oznacza przejście od sprawdzania poszczególnych procesów w pojedynczych komputerach do monitorowania na podstawie stanu usługi jako całości.

⁴ Brian ma miłe wspomnienia z początku 2000 roku związane z konfiguracją MRTG i tworzeniem skryptów przeznaczonych do informowania o temperaturze i poziomie użycia sieci moich komputerów domowych.

Powracamy do teraźniejszości. Technologia OpenTelemetry została opracowana na podstawie dwóch innych projektów typu open source, OpenCensus i OpenTracing. OTel⁵ to specyfikacja i zbiór komponentów, które mają zaoferować projektom wbudowaną telemetrię. To jest komponent wskaźników zgodny z systemem Prometheus, wzbogacony o komponent pobierający OpenTelemetry⁶, odpowiedzialny za pobieranie wskaźników i ich dostarczanie serwerowi systemu Prometheus.

Zwróć uwagę, że nie wspomnieliśmy tutaj o rejestrowaniu danych, śledzeniu i profilowaniu. W przeszłości dzienniki zdarzeń były uznawane za coś, względem czego można ręcznie używać poleceń takich jak `tail`, `grep` i `awk`. Możesz mieć doświadczenie w pracy z narzędziami analizy typu AWStats w celu generowania raportów co godzinę bądź codziennie. W ostatnich latach można zauważyć znacznie częstsze używanie dzienników zdarzeń w trakcie monitorowania np. z zastosowaniem Elasticsearch, Logstash i Kibana (ELK) oraz stosu OpenSearch. Śledzenie i profilowanie, ogólnie rzecz biorąc, odbywa się za pomocą innych stosów oprogramowania: Zipkin i Jaeger zostały opracowane na potrzeby śledzenia, podczas gdy Parca i Pyroscope są używane do nieustannego profilowania.

Skoro przedstawiliśmy nieco historii dotyczącej tworzenia wykresów i powiadamiania, warto teraz zobaczyć, jak wskaźniki i dzienniki zdarzeń wpasowują się w tym obszarze. Czy istnieje więcej kategorii monitorowania niż te dwie wymienione?

Kategorie monitorowania

Ostatecznie większość monitorowania dotyczy tego samego: zdarzeń. Te zdarzenia to może być np.:

- otrzymanie żądania HTTP,
- udzielenie odpowiedzi HTTP 400,
- rozpoczęcie wykonywania funkcji,
- osiągnięcie bloku `else` konstrukcji `if`,
- opuszczenie funkcji,
- zalogowanie użytkownika,
- zapisanie danych na dysku,
- odczytanie danych z sieci,
- żądanie od jądra większej ilości pamięci.

Wszystkie zdarzenia posiadają kontekst. Żądanie IP będzie miało adres IP źródłowy i docelowy, żądany adres URL, zdefiniowane ciasteczka oraz użytkownika wykonującego dane żądanie. Odpowiedź HTTP określi, ile czasu trwało udzielenie tej odpowiedzi, jaki jest kod stanu HTTP i jaka jest długość tej odpowiedzi. Zdarzenia obejmujące funkcje posiadają stos wywołań funkcji znajdujących się powyżej oraz wszystko, co wywołało tę część stosu, np. żądanie HTTP.

Posiadanie całego kontekstu dla wszystkich zdarzeń będzie niezwykle pomocne w przypadku debugowania i poznawania wydajności działania systemu w kategoriach zarówno technicznych, jak

⁵ OTel to nieformalna nazwa dla OpenTelemetry.

⁶ W chwili powstawania książki programiści systemu Prometheus ustalili w trakcie spotkania, że w przyszłości serwer Prometheus będzie zapewniał natywną obsługę protokołu OpenTelemetry. Jednak nie zapadła jeszcze żadna konkretna decyzja, kiedy i jak to będzie się odbywać.

i biznesowych. Jednak przetwarzanie i przechowywanie takiej ilości danych nie jest praktyczne. Dlatego też mamy z grubsza cztery podejścia pozwalające na zmniejszenie ilości danych do poziomu znacznie bardziej praktycznego: profilowanie, śledzenie, rejestrowanie danych i wskaźniki.

Profilowanie

Profilowanie polega na tym, że nie można mieć przez cały czas pełnego kontekstu dla wszystkich zdarzeń, a jedynie pewien fragment kontekstu i to przez ograniczony czas.

Przykładem narzędzia profilowania jest polecenie `tcpdump`. Pozwala ono rejestrować ruch sieciowy na podstawie określonego filtru. Wprawdzie to jest niezwykle ważne narzędzie debugowania, ale tak naprawdę nie będzie działało przez cały czas, ponieważ bardzo szybko doprowadziłoby to do zapełnienia dysku.

Innym przykładem mogą być kompilacje w trybie debugowania plików binarnych pozwalające na śledzenie danych profilowania. Dostarczają one mnóstwa użytecznych informacji, ale spadek wydajności działania podczas zbierania tych wszystkich danych, np. czas wykonywania poszczególnych funkcji, powoduje, że ogólnie rzecz biorąc, nieustanne działanie tego rodzaju narzędzia w środowisku produkcyjnym nie będzie praktyczne.

W jądrze systemu Linux technologia eBPF (ang. *enhanced berkeley packet filters*) pozwala na szczegółowe monitorowanie zdarzeń jądra, począwszy od operacji systemu plików, po działalność sieciową. Ten mechanizm zapewnia dostęp na poziomie, który wcześniej był niedostępny. eBPF oferuje także jeszcze inne zalety, np. przenośność i bezpieczeństwo. Warto zapoznać się z informacjami, które na temat eBPF opublikował Brendan Gregg na stronie <https://www.brendangregg.com/ebpf.html>.

Profilowanie służy głównie do taktycznego debugowania. Jeżeli będzie prowadzone na dłuższą metę, wówczas ilość danych będzie musiała być zmniejszona, aby mieścić się w jednej z innych kategorii monitorowania. W przeciwnym razie trzeba będzie przejść do **nieustannego profilowania**, które umożliwia zbieranie danych w dłuższym czasie.

Nowością w nieustannym profilowaniu jest to, że w celu zmniejszenia ilości danych i zapewnienia względnie niskiego obciążenia zmniejszeniu ulega częstotliwość profilowania. Jednym z zyskujących popularność narzędzi nieustannego profilowania jest bazujące na technologii eBPF oprogramowanie Parca (<https://www.parca.dev/>). Jego agent używa częstotliwości wynoszącej 19 Hz⁷. W efekcie to narzędzie próbuje pobierać statycznie ważne dane na przestrzeni minut zamiast sekund, a jednocześnie dostarcza dane niezbędne do zrozumienia, jak jest używany czas procesora w infrastrukturze, oraz pomaga w poprawieniu efektywności działania aplikacji tam, gdzie to jest konieczne.

Śledzenie

Śledzenie zwykle nie obejmuje wszystkich zdarzeń, a raczej ich pewną część, np. jedno na sto zdarzeń wywoływanych w interesującej nas funkcji. Śledzenie zanotuje funkcje na interesującym Cię stosie wywołań oraz dostarczy informacje o czasie ich wykonywania. Dzięki temu można

⁷ Porównaj to z częstotliwością 100 Hz środowiska uruchomieniowego Go lub nawet z częstotliwością 10 000 Hz w Chromium.

sprawdzić, co w programie zabiera najwięcej czasu i które ścieżki wykonywania kodu mają największy wpływ na ogólne opóźnienie działania aplikacji.

Zamiast wykonywać migawki stosu wywołań w interesujących Cię miejscach, niektóre systemy śledzą i rejestrują czas wywoływania wszystkich funkcji znajdujących się na stosie poniżej interesującej Cię funkcji. Na przykład mogą być zebrane dane jednego na sto żądań HTTP, a na ich podstawie można ustalić, ile czasu zajmuje aplikacji komunikacja z backendem takim jak bazy danych i bufory. Dzięki temu można sprawdzić, jak zmienia się czas obsługi w zależności od czynników typu użycie buforowania.

Systemy śledzenia rozproszonego stosują jeszcze inne podejście. Śledzenie odbywa się między różnymi procesami przez dołączanie unikatowych identyfikatorów do żądań przekazywanych między procesami za pomocą zdalnych wywołań procedur (ang. *remote procedure calls*). Dane zebrane z różnych procesów i urządzeń mogą być połączone ze sobą dzięki identyfikatorowi żądania. To jest ważne narzędzie przeznaczone do debugowania rozproszonego architektury mikrousług. Do technologii na tym obszarze zaliczamy OpenZipkin i Jaeger.

W przypadku śledzenia próbkowanie pozwala na zbieranie rozsądnej ilości danych i nie powoduje nadmiernego obciążenia.

Rejestrowanie danych

Rejestrowanie danych pobiera określony zbiór zdarzeń i rejestruje pewien kontekst dla każdego z tych zdarzeń. Na przykład mogą być rejestrowane wszystkie przychodzące żądania HTTP bądź wszystkie wychodzące wywołania do bazy danych. Aby uniknąć użycia zbyt wielu zasobów, zwykle stosuje się ograniczenie do mniej więcej 100 pól dla wpisu dziennika zdarzeń. Po przekroczeniu tej granicy problem może stanowić przepustowość i brak miejsca w systemie pamięci masowej.

Na przykład, jeśli serwer obsługuje 1000 żądań na sekundę, wówczas wpis dziennika zdarzeń zawierający 100 pól, z których każde ma tylko 10 bajtów, będzie wymagał 1 megabajta danych na sekundę. To wcale nie jest mało w przypadku karty sieciowej o przepustowości 100 Mb. Ponadto to oznacza zapisanie dziennie 84 GB danych na potrzeby rejestrowania zdarzeń.

Dużą zaletą rejestrowania danych jest zwykle brak próbkowania zdarzeń, więc nawet w przypadku ograniczonej liczby pól praktyczne będzie określenie, jak powoli wykonywane zapytania wpływają na określonego użytkownika wykorzystującego konkretny punkt końcowy API.

Podobnie jak monitorowanie ma odmienne znaczenia dla różnych osób, tak samo jest w przypadku rejestrowania danych. W zależności od tego, komu zadasz pytanie, otrzymasz odmienne odpowiedzi, co może prowadzić do dezorientacji. Różne rodzaje rejestrowania danych mają odmienne zastosowania oraz wymagania dotyczące trwałości i przechowywania. Jak możesz się przekonać, istnieją cztery ogólne i nieco nakładające się kategorie.

Transakcyjne dzienniki zdarzeń

To są rekordy biznesowe o znaczeniu krytycznym, którym za wszelką cenę trzeba zapewnić bezpieczeństwo, prawdopodobnie na zawsze. Do tej kategorii zaliczamy wszystko, co dotyczy pieniędzy lub jest używane dla dotyczącej użytkownika funkcjonalności o znaczeniu krytycznym.

Dzienniki zdarzeń żądań

Jeżeli śledzisz każde żądanie HTTP lub wywołanie bazy danych, wówczas masz do czynienia z dziennikiem zdarzeń żądań. Mogą być one przetwarzane w celu implementacji funkcjonalności dotyczącej użytkownika bądź też na potrzeby wewnętrznych optymalizacji. Ogólnie rzecz biorąc, nie chcesz ich utracić, ale jeśli coś się z nimi stanie, to jeszcze nie będzie końcem świata.

Dzienniki zdarzeń aplikacji

Nie wszystkie dzienniki zdarzeń dotyczą żądań, część z nich dotyczy procesów. Typowe są komunikaty początkowe, wykonywane w tle zadania konserwacyjne bądź inne wpisy na poziomie procesu. Te dzienniki zdarzeń są często przeglądane bezpośrednio przez człowieka, więc w przypadku normalnego działania aplikacji należy unikać umieszczania w nich więcej niż tylko kilku wpisów na minutę.

Dzienniki zdarzeń debugowania

Dzienniki zdarzeń debugowania zwykle są bardzo szczegółowe, a tym samym kosztowne do utworzenia i przechowywania. Bardzo często są używane w ściśle określonych sytuacjach debugowania i mają tendencję do profilowania, co wynika z ilości danych. Wymagania dotyczące niezawodności i przechowywania są małe, zaś same dzienniki zdarzeń debugowania mogą nawet nie opuszczać urządzenia, w którym zostały wygenerowane.

Traktowanie poszczególnych rodzajów dzienników zdarzeń w ten sam sposób może być najgorszym z możliwych rozwiązań — otrzymujesz ogromną ilość dzienników zdarzeń debugowania w połączeniu z wyjątkowymi wymaganiami w zakresie niezawodności dotyczącymi dzienników zdarzeń transakcji. Dlatego też wraz ze wzrostem systemu należy zaplanować oddzielenie dzienników zdarzeń debugowania, aby mogły być obsługiwane oddzielnie.

Przykładami systemów rejestrowania danych są stos ELK, OpenSearch, Grafana, Loki i Graylog.

Wskaźniki

Wskaźniki w dużej mierze ignorują kontekst i zamiast tego śledzą agregacje na przestrzeni czasu różnego rodzaju zdarzeń. W celu zapewnienia rozsądnego poziomu użycia zasobów konieczne jest ograniczenie ilości śledzonych danych: 10 000 dla procesu wydaje się rozsądną górną granicą, o której warto pamiętać.

Przykładami rodzajów wskaźników, które możesz mieć, będą liczba otrzymanych żądań HTTP, ilość czasu poświęconego na obsługę żądań, a także liczba aktualnie przetwarzanych żądań. Dzięki wyłączeniu wszelkich informacji o kontekście ilość danych koniecznych do przechowywania i przetwarzania pozostaje na rozsądnym poziomie.

To oczywiście nie oznacza, że kontekst zawsze jest ignorowany. W przypadku żądania HTTP można zdecydować o zachowaniu wskaźników dla poszczególnych ścieżek adresu URL. Pamiętaj o zalecanej granicy 10 000 wskaźników, ponieważ każda oddzielna ścieżka jest uznawana za wskaźnik. Używanie kontekstu w postaci np. adresu e-mail użytkownika nie będzie rozsądnym rozwiązaniem, ponieważ użytkownik może ich mieć nieograniczoną ilość⁸.

⁸ Adres e-mail zalicza się do danych wrażliwych, co z kolei wiąże się z aspektami zapewnienia zgodności i prywatności, których najlepiej będzie unikać w monitorowaniu.

Wskaźników można używać do śledzenia opóźnienia i ilości danych, które są obsługiwane przez poszczególne podsystemy aplikacji. Dzięki temu łatwiej jest znaleźć przyczynę wszelkich spowolnień działania. Wprawdzie dzienniki zdarzeń nie mogą rejestrować wielu pól, ale gdy już wiadomo, który podsystem odpowiada za dany problem, wówczas dzienniki zdarzeń mogą pomóc w ustaleniu, którym dokładnie żądaniom użytkownika należy się przyjrzeć.

W tym miejscu jest najbardziej widoczny kompromis między dziennikami zdarzeń i wskaźnikami. Wskaźniki pozwalają na zbieranie danych dotyczących zdarzeń we wszystkich procesach, choć ogólnie rzecz biorąc, z nie więcej niż jednym bądź dwoma polami kontekstu o ograniczonej liczności. Z kolei dzienniki zdarzeń pozwalają na zbieranie informacji o wszystkich zdarzeniach jednego typu, ale mogą śledzić sto pól kontekstu o ograniczonej liczności. Trzeba koniecznie zrozumieć kwestie dotyczące liczności i ograniczeń nakładanych na wskaźniki, do tego jeszcze powrócimy w dalszej części książki.

Jako system monitorowania na podstawie wskaźników Prometheus został opracowany do śledzenia ogólnego stanu systemu, sposobu jego działania oraz wydajności działania, a nie poszczególnych zdarzeń. Ujmując rzecz inaczej, można powiedzieć, że dla systemu Prometheus ważne jest, że w ostatniej minucie zostało wykonanych 15 żądań, których obsługa zajęła 4 sekundy, spowodowała 40 odwołań do bazy danych, 17-krotne użycie bufora i dokonanie dwóch zakupów przez klientów. Koszt i ścieżki kodu poszczególnych wywołań to kwestie związane z profilowaniem lub rejestrowaniem danych.

Skoro już wiesz, jak Prometheus mieści się ogólnie w obszarze monitorowania, zapoznaj się z różnymi komponentami tego systemu.

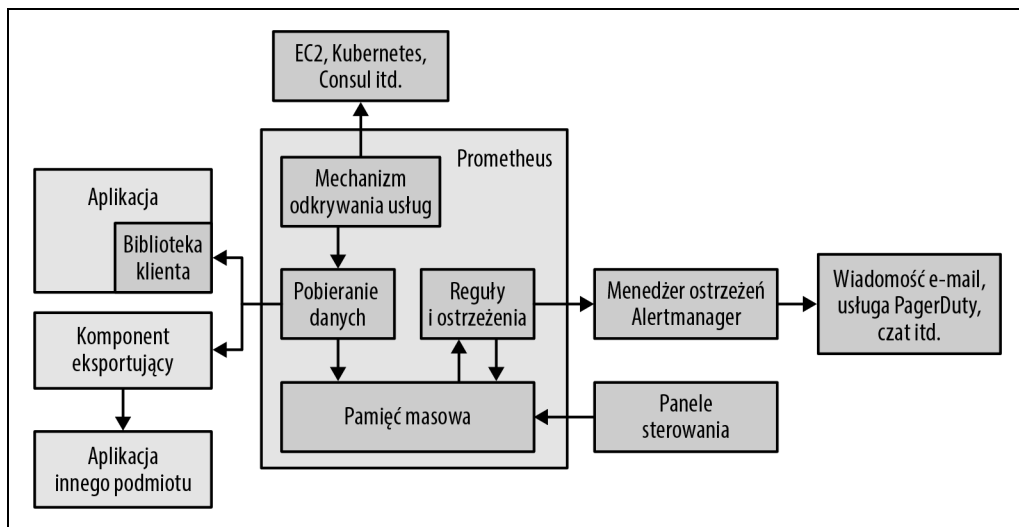
Architektura systemu Prometheus

Na rysunku 1.1 pokazaliśmy ogólną architekturę systemu Prometheus. Komponent odkrywania usług pozwala na znalezienie systemów, z których będą pobierane dane. Celem monitorowania mogą być Twoje własne aplikacje bądź też opracowane przez podmioty zewnętrzne aplikacje, z których dane będą pobierane za pomocą komponentów eksportujących. Zebrane dane są przechowywane w pamięci masowej i można ich używać w panelach sterowania, wykorzystując do tego PromQL, bądź przekazywać ostrzeżenia do menedżera ostrzeżeń Alertmanager, który skonwertuje je na wiadomości dla pagera, wiadomości e-mail bądź innego rodzaju powiadomienia.

Biblioteki klienta

Wskaźniki zwykle nie pojawiają się w magiczny sposób z aplikacji, ktoś musi zapewnić możliwość ich dostarczenia. W tym miejscu do gry wchodzi biblioteki klienta. Zwykle dzięki zaledwie dwóm lub trzem wierszom kodu można zarówno zdefiniować wskaźnik, jak i dodać niezbędną instrumentację w kodzie, nad którym masz kontrolę. To nosi nazwę instrumentacji bezpośredniej.

Biblioteki klienta są dostępne dla wszystkich najważniejszych języków programowania i środowisk uruchomieniowych. Projekt Prometheus dostarcza oficjalne biblioteki klienta dla języków Go, Python, Java/JVM, Ruby i Rust. Dostępnych jest również wiele bibliotek opracowanych przez podmioty zewnętrzne i przeznaczonych dla języków takich jak C#/.Net, Node.js, Haskell i Erlang.



Rysunek 1.1. Architektura systemu Prometheus

Oficjalne kontra nieoficjalne

Nie zniechęcaj się tym, że rozwiązania w zakresie integracji, np. biblioteki klienta, są nieoficjalne bądź zostały opracowane przez podmioty zewnętrzne. Istnieją setki aplikacji i systemów, z którymi być może chcesz przeprowadzać integrację, więc nie ma możliwości, aby zespół tworzący system Prometheus miał czas i doświadczenie wymagane do zapewnienia integracji z tymi wszystkimi aplikacjami i systemami. Dlatego też w ekosystemie oprogramowania Prometheus większość rozwiązań dotyczących integracji została przygotowana i dostarczona przez podmioty zewnętrzne. Aby zachować względną spójność i oczekiwany sposób działania, dostępne są wskazówki wyjaśniające, jak należy opracowywać tego rodzaju rozwiązania w zakresie integracji z aplikacjami i systemami.

Biblioteki klienckie dbają o wszystkie najdrobniejsze szczegóły związane z m.in. zapewnieniem bezpieczeństwa wątków, gromadzeniem danych oraz tworzeniem tekstu systemu Prometheus i/lub udostępnianiem formatu OpenMetrics w odpowiedzi na żądania HTTP. Monitorowanie na podstawie wskaźników nie śledzi poszczególnych zdarzeń, więc poziom użycia pamięci przez bibliotekę klienta nie zwiększa się wraz z liczbą zdarzeń koniecznych do obsłużenia. Zamiast tego poziom użycia pamięci jest powiązany z liczbą posiadanych wskaźników.

Jeżeli jedna z zależności biblioteki aplikacji ma infrastrukturę, zostanie ona automatycznie pobrana. Dlatego też poprzez instrumentację kluczowej biblioteki, np. klienta RPC, można zapewnić instrumentację dla wszystkich aplikacji.

Niektóre wskaźniki, np. poziom użycia procesora i dane statystyczne działania mechanizmu usuwania nieużytków, są zwykle dostarczane standardowo przez biblioteki klientów, w zależności od samej biblioteki i środowiska uruchomieniowego.

Biblioteki klientów nie są ograniczone do dostarczania wskaźników jedynie w formatach tekstowych systemu Prometheus i OpenMetrics. Prometheus to otwarty ekosystem i te same API używane w celu generowania danych w formacie tekstowym mogą być wykorzystane także do tworzenia wskaźników w innych formatach bądź do wprowadzania danych w innych systemach instrumentacji. Podobnie istnieje możliwość pobierania wskaźników z innych systemów instrumentacji i dostarczania biblioteczki klienta, jeśli jeszcze wszystko nie zostało do końca skonwertowane na postać instrumentacji systemu Prometheus.

Komponenty eksportujące

Nie wszystkie uruchamiane fragmenty kodu będą kodem, nad którym masz kontrolę lub do którego zachowujesz dostęp. Dlatego też dodawanie bezpośredniej instrumentacji nie zawsze będzie dostępne. Na przykład istnieje małe prawdopodobieństwo, że jądro systemu operacyjnego rozpocznie wkrótce generowanie przez HTTP wskaźników sformatowanych w formacie systemu Prometheus.

Tego rodzaju oprogramowanie często posiada pewne interfejsy, za pomocą których można uzyskać dostęp do wskaźników. To może być format tymczasowy, wdrażany dla aplikacji, z której mają być pobierane wskaźniki.

Komponent eksportujący pobiera żądania od systemu Prometheus, zbiera niezbędne dane z aplikacji, konwertuje je na odpowiedni format, a następnie przekazuje te dane w odpowiedzi przygotowanej dla oprogramowania Prometheus. Komponent eksportujący można traktować jako rodzaj proxy typu jeden do jednego, który konwertuje dane między interfejsem wskaźników aplikacji i formatem systemu Prometheus.

W przeciwieństwie do bezpośredniej instrumentacji stosowanej w przypadku kodu, nad którym zachowujesz kontrolę, komponenty eksportujące używają innego stylu instrumentacji, znanego jako **niestandardowe komponenty pobierające dane** lub *ConstMetrics*⁹.

Dobłą wiadomością jest to, że za względu na wielkość społeczności systemu Prometheus potrzebny Ci komponent eksportujący prawdopodobnie już istnieje, a możliwość jego wykorzystania będzie wymagała tylko niewielkiej ilości pracy z Twojej strony. Jeżeli taki komponent nie dostarcza interesującego Cię wskaźnika, zawsze możesz poprosić o jego uzupełnienie i tym samym usprawnienie komponentu dla następnej osoby, która będzie chciała go używać.

Odkrywanie usług

Gdy wszystkie aplikacje zostały przygotowane i działają już komponenty eksportujące, Prometheus musi je znaleźć. Dzięki temu Prometheus będzie wiedział, co ma monitorować, i będzie w stanie zauważyć, że coś, co ma być monitorowane, nie odpowiada. W przypadku środowisk dynamicznych nie można po prostu jednokrotnie dostarczyć listy aplikacji i komponentów eksportujących, ponieważ będzie ona nieaktualna. W tym miejscu do gry wchodzi mechanizm odkrywania usług.

⁹ Pojęcie *ConstMetrics* jest potoczne i wiąże się z funkcją `MustNewConstMetric` biblioteki klienta dla języka programowania Go, która jest używana do wygenerowania wskaźników przez komponenty eksportujące utworzone w kodzie Go.

Prawdopodobnie masz pewną bazę danych urządzeń, aplikacji i wykonywanych przez nie zadań. To może być baza danych Chef, plik zasobów dla Ansible, bazujący na znacznikach w egzemplarzu EC2, w etykietach i adnotacjach w Kubernetes bądź po prostu w dokumentacji wiki.

Prometheus zapewnia integrację z wieloma powszechnie stosowanymi mechanizmami odkrywania usług, takimi jak Kubernetes, EC2 i Consul. Istnieje również ogólne rozwiązanie w zakresie integracji (więcej informacji na ten temat znajdziesz w rozdziale 8.).

Mimo to problem wciąż pozostaje. Nawet jeśli Prometheus posiada listę urządzeń i usług, wcale nie oznacza to, że wiadomo jak one wpasowują się w architekturę. Na przykład możesz używać znacznika EC2 Name¹⁰ w celu wskazania aplikacji uruchomionych w urządzeniu, podczas gdy inni mogą używać znacznika o nazwie app.

Każda organizacja używa nieco odmiennego rozwiązania. Prometheus pozwala skonfigurować sposób, w jaki metadane z mechanizmu odkrywania usług będą mapowane na systemy przeznaczone do monitorowania i ich etykiety, używając do tego procesu określanego mianem *relabeling*.

Zbieranie danych

Odkrywanie usług i proces relabeling dostarczają listę systemów przeznaczonych do monitorowania. Teraz Prometheus musi pobrać z nich wskaźniki. W tym celu Prometheus wykonuje żądanie HTTP określane mianem *scrape*. Odpowiedź udzielona na to żądanie zostaje przetworzona i zapisana w pamięci masowej. Dodanych zostaje również wiele użytecznych wskaźników, m.in. informujących o czasie trwania tej operacji oraz o tym, czy zakończyła się sukcesem. Zbieranie danych odbywa się regularnie, zwykle w przedziale czasu od 10 do 60 sekund dla każdego monitorowanego systemu.

Pobieranie kontra przekazywanie

Prometheus to system bazujący na pobieraniu danych. Decyduje o tym, kiedy i co pobrać, bazując na konfiguracji. Istnieją również systemy bazujące na przekazywaniu, w których to monitorowany system decyduje, czy może być monitorowany i jak często.

W internecie toczony są żarliwe debaty dotyczące tych dwóch podejść projektowych, które często są podobne do debat toczonych między zwolennikami edytorów Vim i Emacs. Warto w tym miejsc wyjaśnić, że każde z tych podejść ma swoje wady i zalety. Ostatecznie nie ma znaczenia, które z nich będzie wykorzystane.

Jako użytkownik systemu Prometheus musisz zrozumieć, że mechanizm pobierania danych został wbudowany w program i podejmowanie prób mających na celu użycie podejścia bazującego na przekazywaniu danych będzie w najlepszym razie nierozsądne.

¹⁰ Znacznik EC2 Name to wyświetlana nazwa egzemplarza EC2 w konsoli internetowej EC2.

Pamięć masowa

Prometheus przechowuje dane lokalnie w niestandardowej bazie danych. Zapewnienie niezawodności działania systemowi rozproszonemu jest wyzwaniem, więc Prometheus nie próbuje stosować żadnych form klastrowania. Pomijając kwestie niezawodności, to również ułatwia działanie systemu Prometheus.

Na przestrzeni lat pamięć masowa była wielokrotnie zmieniana, zaś wydanie Prometheus 2.0 posiada pamięć masową w trzeciej postaci. Ten system pamięci masowej jest w stanie obsłużyć pobieranie milionów próbek na sekundę, a tym samym pozwala, aby pojedynczy serwer systemu Prometheus był w stanie monitorować tysiące urządzeń. Zastosowany algorytm kompresji pozwala w rzeczywistych danych zmniejszyć do 1,3 bajta ilość danych w próbce. Zaleca się używanie pamięci masowej bazującej na napędach SSD, choć to nie jest ściśle wymagane.

Panele sterowania

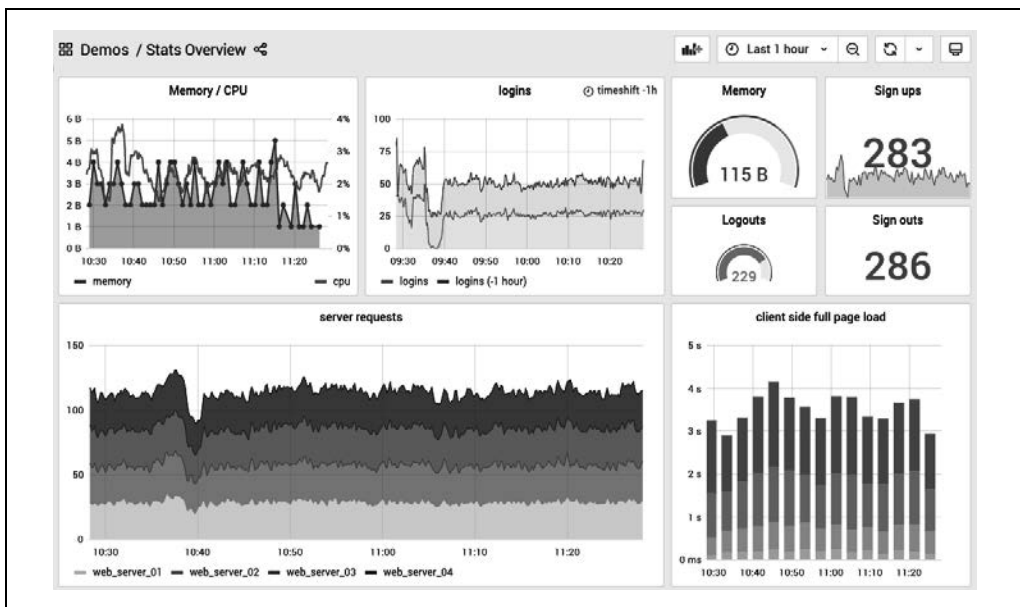
Prometheus posiada pewną liczbę API HTTP pozwalających na zarówno żądanie niezmodyfikowanych danych, jak i analizowanie zapytań PromQL. Te API można wykorzystać do generowania wykresów i paneli sterowania. Standardowo Prometheus oferuje tzw. **przeglądarkę wyrażeń** korzystającą z tych API oraz odpowiednią do wykonywania zapytań tymczasowych i eksploracji danych, choć to nie jest ogólny system paneli sterowania.

Na potrzeby zadań związanych z panelami sterowania zaleca się użycie produktu o nazwie Grafana. Oferuje on wiele funkcjonalności, np. oficjalną obsługę systemu Prometheus jako źródła danych. Grafana potrafi wygenerować wiele różnych paneli sterowania, takich jak pokazane na rysunku 1.2. Grafana obsługuje komunikację z wieloma serwerami Pythona, nawet z poziomu pojedynczego panelu sterowania.

Reguły rejestrowania i ostrzeżeń

Wprawdzie PromQL i silnik pamięci masowej zapewniają potężne możliwości i są efektywne, to jednak agregacja „w locie” wskaźników pochodzących z tysięcy urządzeń w trakcie każdego zadania generowania wykresu może spowodować, że ta operacja stanie się nieco powolna. Reguły rejestrowania pozwalają na regularne wykonywanie wyrażeń PromQL oraz umieszczanie ich wyników w silniku pamięci masowej.

Reguły ostrzeżeń to inna postać reguł rejestrowania. One również pozwalają na regularne wykonywanie wyrażeń PromQL, a ich wyniki stają się ostrzeżeniami, które następnie są przekazywane to tzw. **menedżera ostrzeżeń Alertmanager**.



Rysunek 1.2. Panel sterowania wygenerowany przez produkt Grafana (<https://play.grafana.org/d/cL5pLH7Wz/stats-overview?theme=light&orgId=1>)

Zarządzanie ostrzeżeniami

Menedżer ostrzeżeń Alertmanager otrzymuje ostrzeżenia z serwera systemu Prometheus i zamienia je na powiadomienia. Powiadomienia mogą obejmować wiadomości e-mail, aplikacje czatu, takie jak Slack, a także wiadomości dla usługi, takiej jak PagerDuty.

Działanie menedżera ostrzeżeń Alertmanager jest nieco bardziej złożone niż jedynie bezwarunkowa zamiana ostrzeżenia na powiadomienie w ramach systemu typu „jeden do jednego”. Powiązane ze sobą ostrzeżenia mogą być agregowane na postać pojedynczego powiadomienia, co tym samym ograniczy je w celu zmniejszenia burzy pagerów¹¹, zaś dla poszczególnych zespołów mogą być skonfigurowane różne trasy i dane wyjściowe powiadomień. Powiadomienia można również wyciszać, np. w celu odłożenia na bok informacji o problemie, o istnieniu którego wiadomo i do usunięcia którego została już zaplanowana obsługa techniczna.

Rola menedżera ostrzeżeń Alertmanager kończy się na przekazaniu powiadomień. W celu zarządzania odpowiedzią człowieka na incydent należy używać usług takich jak PagerDuty bądź systemu bazującego na zgłoszeniach.

¹¹ Wspomniany **pager** to urządzenie pozwalające na przekazywanie powiadomień inżynierowi, od którego oczekuje się, że sprawdzi dane zgłoszenie. Wprawdzie istnieje możliwość przekazania zgłoszenia przez tradycyjny pager, ale obecnie prawdopodobnie zostanie użyta wiadomość SMS bądź połączenie głosowe. Burza pagerów to sytuacja, w której otrzymujesz dużą liczbę powiadomień w krótkich odstępach czasu.



Konfiguracja ostrzeżeń i ich poziomów progowych odbywa się w oprogramowaniu Prometheus, a nie w menedżerze ostrzeżeń Alertmanager.

Długoterminowa pamięć masowa

Skoro Prometheus przechowuje dane jedynie w urządzeniu lokalnym, istnieje ograniczenie dotyczące ilości miejsca na dysku dostępnego w tym urządzeniu¹². Wprowadzie najczęściej interesują Cię dane mniej więcej z ostatniego dnia, ale na potrzeby długoterminowego planowania pojemności pożądanym będzie dłuższy okres przechowywania danych.

Prometheus nie oferuje rozwiązania klastrowanego w zakresie przechowywania danych w wielu urządzeniach, choć istnieją API zdalnego odczytu i zapisu pozwalające na wykorzystanie innych systemów w podanej roli. Umożliwiają one wykonywanie zapytań PromQL do danych zarówno lokalnych, jak i zdalnych.

Czym Prometheus nie jest?

Skoro już wiesz, jak Prometheus wpasowuje się w szerszy obszar monitorowania, i znasz jego podstawowe komponenty, warto zapoznać się z przykładowymi przypadkami użycia, w których Prometheus nie okazuje się najlepszym wyborem.

Prometheus to system bazujący na wskaźnikach, a tym samym nie nadaje się do przechowywania dzienników zdarzeń lub poszczególnych zdarzeń. Nie będzie również najlepszym wyborem w przypadku danych o dużej liczności, takich jak adresy e-mail i nazwy użytkowników.

Prometheus został opracowany na potrzeby monitorowania operacyjnego, gdzie małe niedokładności i stan wyścigu ze względu na czynniki takie jak algorytmy szeregowania procesora i nieudane zebranie danych są codziennością. Używanie systemu Prometheus wiąże się z pewnymi kompromisami. Oprogramowanie woli dostarczać dane poprawne w 99,9% niż przerwać monitorowanie w oczekiwaniu na otrzymanie doskonałych danych. Dlatego też należy zachować ostrożność podczas używania systemu Prometheus w przypadkach obejmujących pracę z pieniędzmi lub systemami rozliczeniowymi.

W następnym rozdziale wyjaśnimy, jak można uruchamiać system Prometheus i przeprowadzać za jego pomocą proste monitorowanie.

¹² Jednakże nowoczesne urządzenia mogą lokalnie przechowywać ogromne ilości danych, a tym samym może nie zachodzić potrzeba używania oddzielnych klastrowanych systemów pamięci masowej.

A

adnotacja ostrzeżenia, 309
agregacja, 100, 225, 243, 286
akcja
 drop, 153
 hashmod, 361
 keep, 152
 labeldrop, 165
 labelkeep, 165
 labelmap, 159
 lowercase, 160
 replace, 156
 uppercase, 160
aliasing, 116
analizator składni, 87
API
 EC2, 150
 HTTP, 237, 288
architektura
 federacji globalnej, 345
 komponentu Pushgateway, 84
 klastrowania menedżera ostrzeżeń
 Alertmanager, 358
 systemu Prometheus, 26
asercje przewidywania, lookahead assertions, 154
awarie, 355

B

bezpieczeństwo, 337
biblioteka klienta, 25, 53
Blackbox Exporter, 190
 konfiguracja systemu Prometheus, 201
 próbkiwanie
 DNS, 199
 HTTP, 197

ICMP, 190
 TCP, 195
 przekroczenie czasu oczekiwania, 203
brakujące szeregi czasowe, 278

C

cAdvisor, 167
 etykiety, 169
 pamięć, 169
 procesor, 168
Consul, 148, 183
 telemetria narzędzia, 211
crossmonitorowanie, 359
cytowanie znaków, 90
czarna skrzynka, 189
czas trwania, 235

D

data i godzina, 274
debugowanie, 19
DNS, 199
dopasowanie
 etykiety, 41
 group_left, 260
 negatywne równości, 231
 negatywne wyrażenia regularnego, 231
 równości, 231
 typu jeden do jednego, 258
 typu wiele do jednego, 260
 typu wiele do wielu, 263
 wektora, 258
 wyrażenia regularnego, 231
dzienniki zdarzeń
 aplikacji, 24
 debugowania, 24

dzienniki zdarzeń
transakcyjne, 23
zadań, 24

E

EC2, Elastic Compute Cloud, 149
ekspozycja, 74
element potomny, 98
enum, 101
etykiety, 89, 92, 94, 219
 __address__, 159
 __name__, 97
funkcje, 276
instance, 159
instrumentacji, 95
job, 159
kolizje, 166
liczność, 106
o pustej wartości, 157
opcja honor_labels, 166
ostrzeżenia, 306, 308
systemów przeznaczonych
 do monitorowania, 95, 154
tymczasowe, 162
typu klucz-wartość, 161
używanie, 105
wielkość znaków, 160
wzorce, 101
zarezerwowane, 97
zewnętrzne, 313
złączanie elementów listy, 161
zmienianie, 150

F

federacja, 345
 architektura, 345
format
 ekspozycji tekstu, 88
 OpenMetrics, 79, 91
 YAML, 34, 288, 316
funkcja
 abs(), 270
 absent(), 278
 absent_over_time(), 278
 ceil(), 271
 changes(), 284
 clamp(), 272

clamp_max(), 272
clamp_min(), 272
day_of_month(), 275
day_of_week(), 275
day_of_year(), 275
days_in_month(), 275
delta(), 285
deriv(), 284
exp(), 271
floor(), 271
histogram_quantile(), 279
holt_winters(), 285
hour(), 275
idelta(), 285
increase(), 282
irate(), 282
label_join(), 277
label_replace(), 277
ln(), 270
minute(), 275
month(), 275
predict_linear(), 284
rate(), 280
resets(), 283
round(), 272
scalar(), 269
sgn(), 273
sort(), 279
sort_desc(), 279
sqrt(), 271
time(), 274
timestamp(), 276
vector(), 268
year(), 275

funkcje
 agregacji, 286
 bezpieczeństwa, 337
 dotyczące daty i godziny, 274
 dotyczące zmiany typu, 268
 matematyczne, 270
 trygonometryczne, 273

G

GIL, Global Interpreter Lock, 76
Go, 80
Grafana, 37, 109
 instalacja, 109
 kontrolki czasu, 116

- ograniczanie liczby wykresów, 114
- panel
 - Stat, 118
 - State Timeline, 120
 - szeregu czasowego, 114
 - Table, 119
- panele sterowania, 113
- zmienne szablonu, 121
- źródło danych, 110
- Graphite Exporter, 205, 210
- Grok, 186
- grupowanie, 243
- Gunicorn, 76

H

- histogram, 63, 229
 - częstotliwości, 252
 - funkcja `histogram_quantile()`, 279
 - kumulacyjny, 65
- HTTP, HyperText Transfer Protocol, 197
- HTTPServer, 81

I

- ICMP, Internet Control Message Protocol, 190
- import zależności, 55
- InfluxDB, 207
 - Exporter, 205
- info, 103
- instalowanie systemu, 19
- instrumentacja, 53, 56, 70, 95
 - bezpośrednia, 53
 - biblioteki, 69
 - testów jednostkowych, 67
 - usługi, 68
- interfejs użytkownika, 35

J

- Java, 81
 - HTTPServer, 81
 - Servlet, 82
- język
 - Go, 80
 - Java, 81
 - PromQL, 223
 - Python, 75
- JMX Exporter, 206

K

- klastrowanie Alertmanager, 358
- klauzula
 - by, 245
 - on, 260
 - without, 244
- komponent
 - cAdvisor, 167
 - eksportujący, 27
 - Blackbox Exporter, 190
 - Consul Exporter, 183
 - Graphite Exporter, 205, 210
 - Grok, 186
 - InfluxDB Exporter, 205
 - JMX Exporter, 206
 - MySQLd, 185
 - Node Exporter, 39, 130
 - NRPE Exporter, 207
 - numerów portów, 186
 - StatsD Exporter, 208
 - tworzenie, 211
 - Windows Exporter, 39, 129
 - InfluxDB, 207
 - kube-state-metrics, 181
 - pobierający dane
 - dotyczące obciążenia, 137
 - dotyczące procesora, 130
 - dysku, 132
 - hwmon, 134
 - loadavg, 136
 - meminfo, 133
 - netdev, 133
 - niestandardowy, 102, 214
 - stat, 135
 - systemu operacyjnego, 136
 - systemu plików, 131
 - uname, 135
 - z pliku tekstowego, 137
 - powiadamiający, 324
 - Pushgateway, 83
 - komunikator Slack, 327, 331
 - kontenery, 167
 - kopia zapasowa, 348
 - kubelki, 64
 - Kubernetes, 170
 - działanie systemu Prometheus, 170
 - mechanizmy wykrywania usług, 172
 - kwantyle, 66

L

liczniki, 38, 55, 227
funkcje, 280
listy, 161

M

mapowanie metadanych, 150
mechanizm wykrywania usług, 141
bazujący na HTTP, 147
bazujący na pliku, 144
Consul, 148
EC2, 149
Kubernetes
endpointslice, 175
ingress, 180
komponent kube-state-metrics, 181
node, 172
pod, 180
service, 175
mapowanie metadanych, 150
podejście statyczne, 143
menedżer ostrzeżeń Alertmanager, 29, 47, 306, 315
architektura, 301
architektura klastrowania, 358
drzewo routingu, 317
interfejs sieciowy, 334
konfiguracja, 312
odbiorcy, 324
plik konfiguracyjny, 316
potok powiadomień, 315
szablony powiadomień, 327
metamonitorowanie, 359
mierniki, 38, 59, 225
funkcje, 283
modyfikator
@, 237
bool, 257
offset, 236
monitorowanie, 18, 19
aplikacji, 51
cross, 359
infrastruktury, 127
kategorie, 21
meta, 359
typu „czarna skrzynka”, 189

most, 86
Graphite, 86
MySQLd, 185

N

narzędzie
Ansible, 20, 143
AWStats, 21
Chef, 20, 138
collectd, 20
Consul, 148, 183, 211
envsubst, 353
Grafana, 109
Graphite, 20, 94
Gunicorn, 76
Jaeger, 23
kubectl, 171
Maven, 81, 83
Minikube, 170
Netbox, 141
OpenZipkin, 23
Parca, 22
Promdash, 110
Promtool, 91, 338
Pulumi, 352
Pyrra, 66
sed, 353
Smokeping, 20
StatsD, 20
tcpdump, 22
nieustanne profilowanie, 22
Node Exporter, 129
uruchamianie komponentu, 39
NRPE Exporter, 207

O

obciążenia
zmniejszanie, 362
odchylenie standardowe, 249
odkrywanie usług, 27
odwołania wsteczne, backreference, 154
OpenMetrics, 79, 91
operacje, 296
operator
agregacji, 243
avg, 248

- bottomk, 250
- count, 247
- count_values, 252
- group, 248
- max, 250
- min, 250
- quantile, 251
- stddev, 249
- stdvar, 249
- sum, 246
- topk, 250
- and, 266
- atan2, 256
- or, 264
- unless, 265
- operatory
 - arytmetyczne, 254
 - binarne, 254
 - kolejność wykonywania, 267
 - logiczne, 263
 - porównania, 256
- ostrzeganie, 29, 43, 299, 301, 311
 - adnotacja, 309
 - na stronie Alerts, 47

P

- pager, 30
- pamięć masowa, 29
 - długoterminowa, 31, 348
- panele sterowania, 29, 109
- platforma
 - Kubernetes, 170
- plik konfiguracyjny operacji pobierania danych, 84
- pobieranie danych, 28, 162
- podsumowanie, 62, 227
- podzapytanie, 235
- pomoc, 365
- potok powiadomienia, 315
- powiadamianie, 19, 49, 315
- powiadomienia rozwiązane, 332
- powtarzalność, 69
- poziom, 296
- profilowanie, 22
- Promdash, 110
- Prometheus, 17, 31
 - architektura systemu, 25
 - awarie, 355

- długoterminowa pamięć masowa, 348
- planowanie wdrożenia, 342
- rozbudowa systemu, 343
- uruchamianie systemu, 32, 350
- uwierzytelnianie, 354
- uzyskiwanie pomocy, 365
- wykorzystanie federacji, 345
- wymagania sprzętowe, 350
- zarządzanie konfiguracją, 352
- zarządzanie wydajnością działania, 360
- PromQL, Prometheus Query Language, 223, 225
 - agregacja, 225
 - API HTTP, 237
 - czas trwania, 235
 - dopasowanie wektora, 258
 - funkcje, 268
 - operatory agregacji, 243
 - operatory binarne, 254
 - priorytet operatorów, 267
 - reguły rejestrowania, 288
 - selektory, 230

Promtool, 91

protokół

- HTTP, 197
- ICMP, 191
- InfluxDB, 205, 207
- SNMP, 206
- StatsD, 209
- TCP, 195

przełęczarka wyrażenia, 35

przekazywanie danych, 28

przesunięcie, 236

punkt końcowy

- query, 237
- query_range, 239
- zdalnego zapisu, 349

Pushgateway, 83

- architektura komponentu, 84

Python, 75

- CPython, 76

- Twisted, 76

- WSGI, 75

R

regresja najmniejszych kwadratów, 284

reguły

- ostrzegania, 29, 302
- etykiety ostrzeżenia, 306

reguły

- opcja annotations, 309
- opcja for, 304
- rejestrowania, 29, 288
 - dla API, 293
 - funkcja wektora zakresu, 292
 - informacje o stanie Rules, 290
 - nazewnictwo, 295
 - używanie, 288, 291, 294
 - zmniejszenie liczności, 291

rejestr

- domyślny, 56, 74
- niestandardowy, 85
- rejestrwanie danych, 23
- role komputera, 103
- routing, 317

S

selektory, 230

Servlet, 82

serwer

- zapewnienie bezpieczeństwa, 337

sharding

- pionowy, 344
- poziomy, 363

skalar, 254

SLA, service level agreement, 64

SNMP, 206

sortowanie, 279

sprawdzenie wskaźników, 90

StatsD, 208

- Exporter, 208

sterta, 73

systemy monitorowania, 205

- InfluxDB, 207
- StatsD, 208

szablony

- konsoli, 110
- powiadomień, 310, 327

szyfrowanie TLS, 338

Ś

ścieżka dostępu, 76

śledzenie, 22

średnia

- arytmetyczna, 248
- geometryczna, 248

T

tabela wyjątku, 106

TCP, Transmission Control Protocol, 195

telemetria narzędzia Consul, 211

testy jednostkowe, 67

TLS, 338

- opcje zaawansowane, 339

trendy, 19

tryb wieloprotocowy, 79

Twisted, 76

tworzenie

- funkcji wektora zakresu, 292
- komponentu eksportującego, 211
- kopii zapasowej, 348
- paneli sterowania, 109
- typy wskaźników, 88, 91

U

uwierzytelnianie, 354

- podstawowe, 340

W

wariancja standardowa, 249

wdrożenie systemu Prometheus, 342

wektor

- natychmiastowy, 232
- zakresu, 233

Windows Exporter, 39, 129

WSGI, Web Server Gateway Interface, 75

wskaźniki, 24, 55, 97, 296

- info, 103
- nadawanie nazwy, 71
- sprawdzanie, 90
- typy, 88, 91

współczynnik

- trendu, 286
- wygładzania, 286

wstrzymywanie, 333

wydajność systemu, 360

wykres

- częstotliwości żądań, 57
 - poziomu użycia pamięci, 43
 - ruchu sieciowego, 44
 - wskaźnika, 39
 - wyrażenia, 38
- wykrywanie usług, 142

- wyliczenia, 101
- wrażenia regularne, 154
- wyszukiwanie
 - kosztownych wskaźników, 361
 - systemów do monitorowania, 361
- wywołania zwrotne, 61

Z

- zaczep sieciowy, 326
- zarządzanie
 - ostrzeżeniami, 30
 - zmianami, 364
- zbieranie danych, 28
- zdalny
 - odczyt, remote read, 349
 - zapis, remote write, 349
- zliczanie
 - wielkości, 58
 - wyjątków, 57

- zmiana
 - etykiet, 150
 - akcja hashmod, 361
 - akcja labeldrop, 165
 - akcja labelkeep, 165
 - akcja labelmap, 159
 - akcja replace, 156
 - opcja metric_relabel_configs, 164, 166
 - wielkość znaków, 160
 - wskaźników, 164
 - wrażenia regularne, 154
 - mierników, 283
 - typu, 268
- znaczniki czasu, 90, 93, 140

Ż

- źródło danych, 110

PROGRAM PARTNERSKI

— GRUPY HELION —

1. ZAREJESTRUJ SIĘ
2. PREZENTUJ KSIĄŻKI
3. ZBIERAJ PROWIZJĘ

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

<http://program-partnerski.helion.pl>

GRUPA
Helion 

W tej książce znajdziesz bezcenne wskazówki dotyczące wdrażania serwera Prometheus i jego używania w rzeczywistych rozwiązaniach.

Julius Volz, współtwórca oprogramowania Prometheus

Monitorowanie wydajności systemu jest ważnym zadaniem i nie polega tylko na obserwacji pracy procesora. Należy orientować się pod jakim obciążeniem pracuje baza danych, czy przepustowość urządzeń sieciowych jest wystarczająca i jaki jest koszt niepełnego użycia bufora. Warto też wiedzieć, czy stopień wykorzystania droższych elementów uzasadnia ich zastosowanie i utrzymywanie.

To drugie, zaktualizowane wydanie przewodnika po systemie Prometheus, pozwoli zrozumieć, dlaczego ten system open source w ostatnich latach wciąż zyskuje na popularności. Znajdziesz w nim wyczerpujące wprowadzenie do tego oprogramowania, a także wskazówki dotyczące monitorowania aplikacji i infrastruktury, tworzenia wykresów, przekazywania ostrzeżeń, bezpośredniej instrumentacji kodu i pobierania wskaźników pochodzących z systemów zewnętrznych. Zrozumiesz zasady konfiguracji systemu Prometheus, komponentu Node Exporter i menedżera ostrzeżeń Alertmanager. Zapoznasz się też z nowymi funkcjonalnościami języka PromQL. Dokładnie zaprezentowano tu również zagadnienia bezpieczeństwa po stronie serwera, w tym mechanizm TLS i uwierzytelniania podstawowego.

Dzięki najlepszym praktykom i wskazówkom dotyczącym instrumentacji w kodzie ta książka pomoże Ci w niezawodnym monitorowaniu usług!

TJ Hoplock, starszy inżynier monitorowania, NS1

W książce między innymi:

- czym jest Prometheus i jak wygląda jego architektura
- monitorowanie infrastruktury za pomocą komponentów takich jak Node Exporter
- instrumentacja w kodzie aplikacji
- tworzenie paneli sterowania za pomocą Grafany
- współpraca oprogramowania Prometheus i Kubernetes

Julien Pivotto jest współtwórcą serwera Prometheus i ekosystemu CNCF. Zajmuje się tym projektem od 2017 roku.

Brian Brazil jest programistą systemu Prometheus. Słynie z głębokiej znajomości tego oprogramowania.

Helion
helion.pl
HELION SA
ul. Kościuszki 1c
44-100 Gliwice
tel. 32 250 99 63
helion@helion.pl

KOD KORZYŚCI
Sięgnij po więcej! ▶



ISBN 978-83-289-0529-0



Cena: 89,00 zł