

Podnoszenie wydajności platformy konwersacyjnej AI dzięki serwerom dedykowanym Advance



+120
fizycznych węzłów
Elasticsearch



+70TB
zmigrowanych
danych



6 razy
lepsza wydajność

W SKRÓCIE

SentiOne, uznana przez Deloitte za jedną z najszybciej rozwijających się firm technologicznych w Europie Środkowej, dostarcza unikalną platformę konwersacyjną AI, która umożliwia monitoring internetu, analizę zaangażowania odbiorców i automatyzację obsługi klienta we wszystkich kanałach online.

Podczas gdy większość firm, które oferują boty konwersacyjne oparte na sztucznej inteligencji, zмага się z brakiem danych do uczenia maszynowego, SentiOne wykorzystuje rzeczywiste dyskusje internautów do treningu silnika NLU (Natural Language Understanding). Dzięki obszernemu zestawowi danych, zebranemu podczas rozwijania narzędzi social listeningu, SentiOne jest w stanie wytrenować wyjątkowo dokładne algorytmy głębokiego uczenia maszynowego.

„Współpracując z OVHcloud, możemy płynnie skalować naszą działalność i przetwarzać terabajty danych, dostarczając naszym klientom wartościową wiedzę o ich grupach docelowych.”

Michał Brzezicki, Co-Founder, Chief Technical Officer w SentiOne

Zespół ekspertów IT firmy SentiOne zidentyfikował kilka potencjalnych przyczyn, które mogą prowadzić do degradacji klastra Elasticsearch, w tym niewystarczającą moc procesora, niedostateczną wydajność pamięci masowej, awarie sprzętu, konfigurację oprogramowania oraz problemy z siecią i komunikacją pomiędzy serwerami.

Sieć

Elasticsearch jest bardzo podatny na problemy z łącznością. Gdy pojawią się problemy ze stabilnością połączenia pomiędzy dwoma węzłami w klastrze, cały klastr spowalnia. Aby utrzymać niezawodną i bezpieczną komunikację pomiędzy węzłami, SentiOne zastosowała vRack: rozwiązanie sieci prywatnej, opracowane przez OVHcloud.

Sieć vRack pozwoliła firmie SentiOne na utworzenie stabilnych połączeń pomiędzy węzłami, niezależnych od sieci publicznej. Problemy z siecią zostały niemal natychmiast wyeliminowane po uruchomieniu tego rozwiązania.

Aby móc dogłębnie monitorować stan infrastruktury, zespół SentiOne wykorzystuje Grafanę z wtyczką Sensu do przechowywania wszystkich wartości statusu systemów i klastra Elasticsearch. Zarówno w przypadku interfejsów publicznych, jak i vRack, nie odnotowano żadnych anomalii, więc zespół wykluczył sieć jako źródło problemów z wydajnością.

Procesory

Podczas dalszego poszukiwania potencjalnych przyczyn problemów z wydajnością, zespół badał wykorzystanie procesora i zarządzanie wątkami. Ogólne wykorzystanie procesora w klastrze zbliżało się do 80%, co zostało uznane za właściwe, ponieważ wskazywało na prawidłowe wykorzystanie zasobów z punktu widzenia kosztów. Zespół zauważył jednak, że klastr był znacznie wolniejszy podczas rebalansowania.

W okresie spowolnienia zużycie procesora było niskie i ograniczone do jednego rdzenia. Głębsza analiza wykazała, że większość czasu procesora była poświęcona operacjom IO. Ponieważ zespół już wcześniej wykluczył sieć jako przyczynę, teraz zaczął analizować możliwe problemy z pamięcią masową.

Przestrzeń dyskowa

Utworzenie dashboardu Grafana dla metryki IO Time z wtyczki Sensu dostarczyło zespołowi SentiOne odpowiedzi. Powstały w Grafanie wykres jasno pokazał, ile czasu procesor poświęcał na operacje IO w danym okresie. Okazało się, że liczba ta sięgała prawie 100%!

Wyjaśniło to również, dlaczego klastr zwalniał podczas rebalansowania. Gdy węzeł, z którego kopiowano shard znajdował się pod dużym obciążeniem IO, skopiowanie shardu tylko pogarszało sytuację.

Konfiguracja klastra Elasticsearch

Infrastruktura SentiOne opierała się na węzłach "gorących i zimnych". Węzły gorące zawierały popularne, regularnie sprawdzane indeksy, natomiast węzły zimne przechowywały rzadko używane dane. Choć dodanie większej liczby replik dla indeksów gorących w celu zrównoważenia obciążenia mogłoby wydawać się oczywistym rozwiązaniem, okazało się, że tak nie jest. Po pierwsze, zapytania były przypisywane losowo do shardów¹, więc jeśli ciężkie zapytanie było powtarzane wielokrotnie (np. w postaci różnych agregacji tych samych danych), to wszystkie repliki były przeciążone podobnymi lub identycznymi obliczeniami. Dodatkowo, wersja Elasticsearch zainstalowana na klastrze SentiOne w tym czasie nie zapewniała funkcji Adaptive Replica Selection (ARS). Dodanie nowych węzłów do klastra nie wpłynęłoby zatem na zwiększenie wydajności, jeśli zespół nie zaktualizował oprogramowania.

Po zidentyfikowaniu czynników, które powodowały pogorszenie wydajności klastra Elasticsearch, zespół IT mógł rozpocząć planowanie zmian w oprogramowaniu, architekturze i sprzęcie, aby uzyskać większą responsywność platformy.



¹ W Elasticsearch indeks przechowujący dokumenty, może być podzielony na wiele shardów. Każdy shard to w pełni funkcjonalny i niezależny „index”, który może być przechowywany na dowolnym węźle (ang. node) klastra. Dzięki shardom system jest skalowalny poziomo, a operacje na danych mogą być rozproszone i wykonywane równolegle, co zwiększa wydajność i skraca czas odpowiedzi.

ROZWIĄZANIE

Wymiana sprzętu i skalowanie klastrów poziomo wydawało się najbardziej oczywistym rozwiązaniem, ale ze względu na znaczny, ale ograniczony budżet, SentiOne rozpoczęła od uporządkowania wykorzystania klastra.

Zanim zespół był gotowy do migracji do nowej infrastruktury serwerów dedykowanych, należało wykonać jeszcze kilka innych działań.

Inteligentne partycjonowanie danych

W pierwszej kolejności, inżynierowie SentiOne przeanalizowali zapytania użytkowników i ich wpływ na wydajność klastra. Zależność między czasem zapytania a liczbą zapytań i ich wielkością jest prosta. Im większe shardy, tym dłużej trwa każde zapytanie, a im więcej shardów jest odpytywanych w ramach zapytania, tym dłużej trwa cały proces.

Aby skrócić czas wyszukiwania, SentiOne zdecydowała się podzielić dane na mniejsze shardy, nie większe niż 30 GB. Po wywołaniu zapytania, dane pojedynczego shardu mogą zmieścić się w pamięci cache systemu plików w RAM, przynajmniej w teorii.

Po zbadaniu, które przedziały czasowe są przeszukiwane przez użytkowników najczęściej, SentiOne zmieniła również sposób podziału danych na indeksy. Dane podzielono w zależności od zachowania klientów, minimalizując liczbę shardów w każdym zapytaniu, a zakres indeksów skrócono z miesięcznych do tygodniowych.

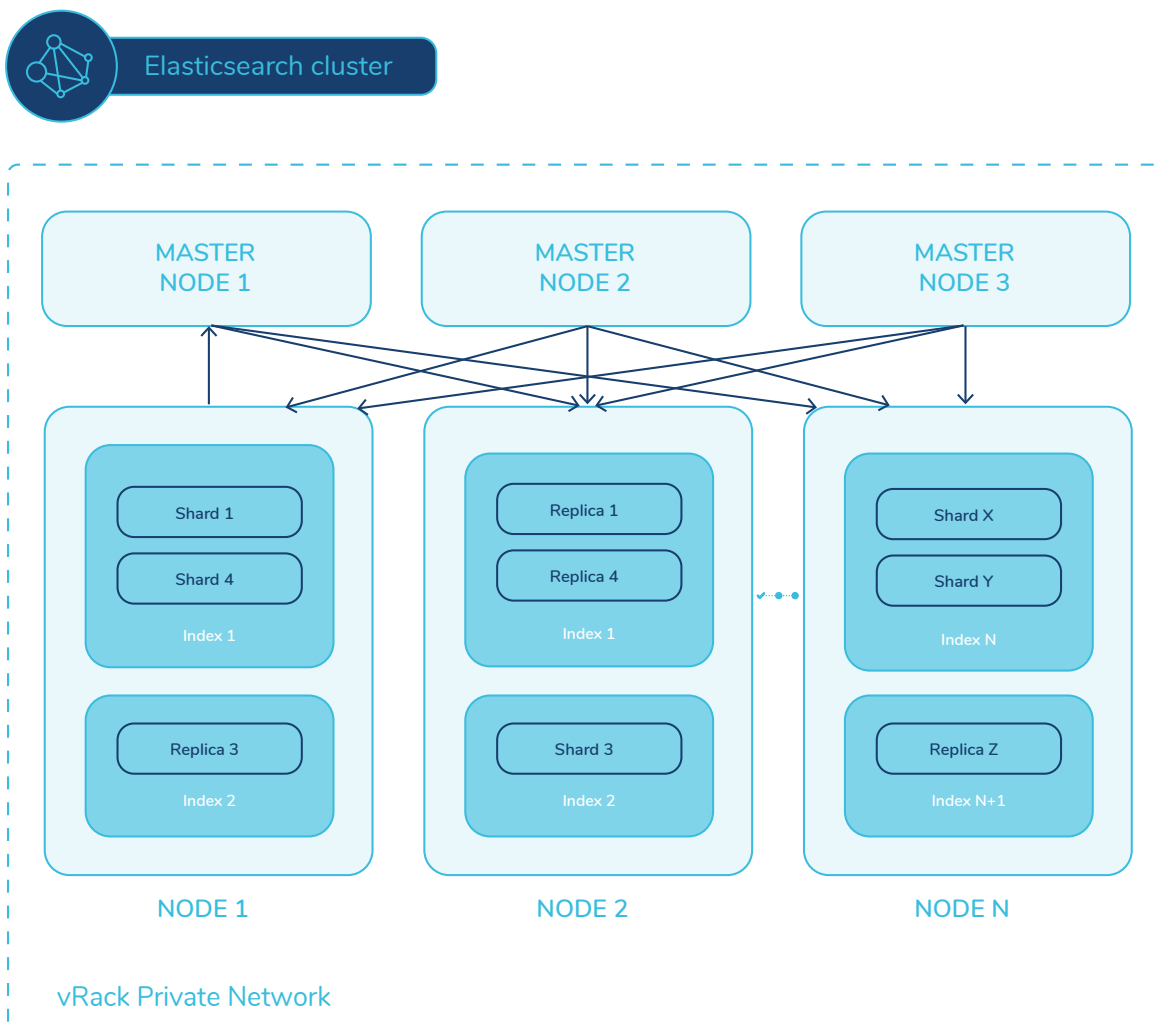
Aplikacja

Zespół SentiOne zrewidował wszystkie aspekty, związane z zapytaniami do klastra Elasticsearch, zoptymalizował wywołania i zmienił interfejs użytkownika, tak aby użytkownicy nie odczuli, że system "zamiera" podczas zapytań. Wprowadzili również kontrolę, która blokuje zbyt ciężkie zapytania, w oparciu o czas ich wykonania w ruchomych ramach czasowych. Poprzednio, pojedyncze ciężkie zapytanie mogło pogorszyć wydajność całego klastra, więc takie zabezpieczenie było koniecznością, nawet jeśli ogranicza niektórych użytkowników.



Zmiany w architekturze klastra

Wraz z nową wersją klastra SentiOne zrezygnowała z koncepcji "gorących i zimnych" nodów z kilku powodów. Po pierwsze, trudno jest dokonać wyraźnego rozróżnienia pomiędzy indeksami gorącymi i zimnymi. W którym momencie można już nazwać indeks "gorącym"? Po drugie, węzły zimne miały bardzo niskie zużycie IO przy skokowym wykorzystaniu procesora, podczas gdy węzły gorące były przeciążone przez IO, ale miały niewykorzystane zasoby CPU. Było to oczywiste marnotrawstwo zasobów, a tym samym strata pieniędzy. Znalezienie właściwej równowagi pomiędzy ciepłymi i zimnymi węzłami okazało się trudnym zadaniem, prowadzącym do niepotrzebnych kosztów.





Nowa platforma sprzętowa

Poprzednia infrastruktura składała się z trzech węzłów głównych, 82 węzłów gorących i 42 węzłów zimnych, wyposażonych w to samo rozwiązanie pamięci masowej: dwa dyski Intel® SSD DC S4500. Głównym celem migracji sprzętowej była poprawa wydajności dyskowej, dlatego firma SentiOne wybrała dla swoich nowych węzłów serwery serii Advance - w szczególności modele ADV-2 z dyskami NVMe.

Technologia dysków NVMe to połączenie dysków półprzewodnikowych, magistrali PCIe i protokołu NVMe zaprojektowana tak, aby zniwelować lukę w wydajności pomiędzy mocą obliczeniową a pamięcią masową.

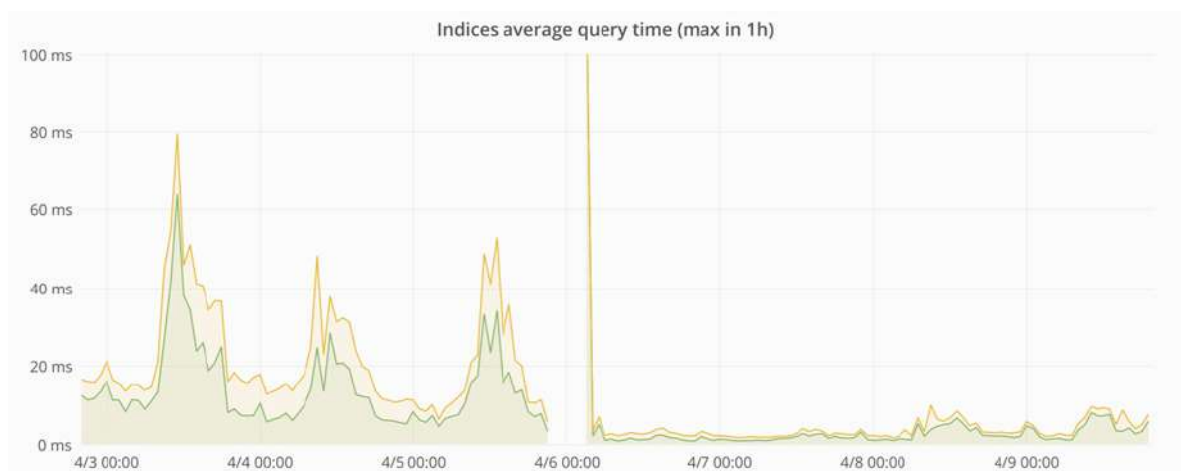
Interfejs pomiędzy dyskami SSD ma znaczący wpływ na całkowite opóźnienie, jakiego użytkownik będzie doświadczał na poziomie aplikacji, natomiast dwa elementy - magistrala (PCIe) i protokół (NVMe) - mają wpływ na ogólną wydajność. Łącze PCIe 3.0 x4 oferuje prawie pięciokrotnie większą przepustowość niż najszybszy interfejs SATA i prawie trzykrotnie większą niż najlepszy interfejs SAS. Protokół NVMe oferuje do 64 tys. kolejek, dzięki czemu interfejs jest w stanie poradzić sobie z dużą liczbą wątków I/O generowanych przez procesory wielordzeniowe.

Aby wdrożyć nowy klaster 121 węzłów danych i trzech węzłów głównych, firma SentiOne poświęciła dwa miesiące na rozwój, testowanie i planowanie. Chociaż migracja i aktualizacja Elasticsearch nie była procesem bezproblemowym i wymagała sześciogodzinnego przestoju klastra, to jednak przyniosła oczekiwane korzyści.

KORZYŚCI

Wszystostronna platforma SentiOne monitoruje miliardy dyskusji prowadzonych na tysiącach stron i portalach internetowych. Aby zapewnić przedsiębiorstwom informacje na temat odbiorców i zautomatyzować obsługę klienta z pomocą AI, firma zbiera, przetwarza i analizuje ogromne ilości danych. Dlatego też, wydajność pamięci masowej odgrywa tu kluczową rolę.

Nowy klaster, zbudowany na serwerach dedykowanych ADV-2 i wyposażonych w dyski NVMe, stał się przełomowym rozwiązaniem dla SentiOne. Migracja sprzętowa przyniosła natychmiastowy wzrost wydajności - klaster jest obecnie sześciokrotnie szybszy.



Średni czas wyszukiwania w dwóch losowo wybranych indeksach, z około 6-krotnym wzrostem wydajności.

Aktualizacja wersji Elasticsearch - która obejmowała włączenie funkcji ARS - pozwoliła na bardziej równomierną pracę klastra, a w okresach większego obciążenia dzieli je równo pomiędzy wiele węzłów. Zespół zmierzył, że podczas regularnych operacji obciążenie zbliża się do 50%, ale nie przekracza 80% podczas testów stresowych, co daje SentiOne odpowiednią rezerwę na potrzeby dalszego skalowania i wdrażania nowych funkcji.

„Pod względem budżetowym, proces migracji okazał się bardzo korzystny. Koszty utrzymania klastra zmniejszyły się o około 5%, przy znacznej poprawie wydajności.”

Michał Brzezicki, Co-Founder, Chief Technical Officer w SentiOne



OVHcloud™ jest globalnym dostawcą usług w chmurze, specjalizującym się w dostarczaniu wydajnych i przystępnych cenowo rozwiązań zapewniających sprawne zarządzanie, bezpieczeństwo i skalowalność danych. OVHcloud to rozsądna alternatywa dla rozwiązań hostingu, poczty elektronicznej, serwerów bare metal, chmury prywatnej, hybrydowej i publicznej. Grupa zarządza 30 centrami danych w 12 lokalizacjach na 4 kontynentach. Produkuje własne serwery, buduje centra danych i wdraża globalną sieć światłowodową, co pozwala jej osiągać maksymalną wydajność. Naszą filozofią jest przełamywanie schematów. Dzięki innowacjom firma OVHcloud zapewnia wolność i bezpieczeństwo, a także gotowa jest wspierać przedsiębiorstwa w wyzwaniach związanych z danymi – tych dzisiejszych i tych przyszłych. Wykorzystując 20-letnie doświadczenie i czerpiąc z europejskich wartości, OVHcloud angażuje się w rozwój odpowiedzialnych technologii. Grupa staje się siłą napędową kolejnej ewolucji chmurowej.



[ovh.com](https://www.ovh.com)  [OVH](https://twitter.com/OVH)  [ovhcom](https://www.facebook.com/ovhcom)  [OVH](https://www.linkedin.com/company/ovh)